

Le Frido 2019,
volume 4
Laurent Claessens

Plusieurs extensions et versions de ce livre.

1. La version courante, régulièrement mise à jour et qui deviendra petit à petit le Frido 2020. Téléchargeable sur

<https://laurent.claessens-donadello.eu/pdf/lefrido.pdf>

2. La version la plus complète, contenant des exercices ainsi que de la mathématique de niveau recherche sur

<https://laurent.claessens-donadello.eu/pdf/giulietta.pdf>

3. Et bien entendu les sources $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$

<https://github.com/LaurentClaessens/mazhe>

Copyright 2011-2019 Laurent Claessens, Carlotta Donadello, Lilian Besson, Bertrand Desmons, and many other contributors. A complete list could be retrieved from the git log. Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.3 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license is included in the chapter entitled “GNU Free Documentation License”.

(c) 2015 David Revoy pour les illustrations de couverture CC-BY,
<https://www.peppercarrot.com/>

ISBN : 979-10-97085-21-6

Chapitre 29

Séries de Fourier

29.1 Densité des polynômes trigonométriques

29.1.1 Convergence pour les fonctions continues (via Weierstrass)

Le résultat fondamental qui nous permet d'utiliser les polynômes trigonométriques comme base pour les fonctions *continues* périodiques est le suivant. Notons que pour les fonctions non continues, il y a encore du travail.

Lemme 29.1.

Si $f: \mathbb{R} \rightarrow \mathbb{C}$ est une fonction continue 2π -périodique et si $\epsilon > 0$, alors il existe un polynôme trigonométrique P tel que $\|f - P\|_\infty \leq \epsilon$.

Démonstration. Nous allons utiliser le théorème de Stone-Weierstrass 13.303. Soit le compact Hausdorff

$$S^1 = \{z \in \mathbb{C} \text{ tel que } |z| = 1\}, \quad (29.1)$$

et $C(S^1, \mathbb{C})$ l'algèbre des fonctions continues de S^1 vers \mathbb{C} . Il suffit de vérifier que les polynômes trigonométriques vérifient les hypothèses du théorème de Stone-Weierstrass. Un polynôme trigonométrique est un polynôme en z et \bar{z} défini sur S^1 .

- (1) Le polynôme constant est dans l'algèbre, ok.
- (2) Pour la séparation des points, le polynôme trigonométrique $x \mapsto e^{ix}$.
- (3) Si P est un polynôme en z et \bar{z} , alors \bar{P} l'est encore.

Donc si $\epsilon > 0$ et $\tilde{f} \in C(S^1, \mathbb{C})$ sont donnés, il existe un polynôme trigonométrique P tel que

$$\sum_t |\tilde{f}(e^{it}) - P(t)| < \epsilon. \quad (29.2)$$

Soit $f: \mathbb{R} \rightarrow \mathbb{C}$ une fonction continue 2π -périodique. Nous considérons $\tilde{f} \in C(S^1, \mathbb{C})$ donnée par $\tilde{f}(e^{it}) = f(t)$. Alors $\sup_t |f(t) - P(t)| \leq \epsilon$. \square

29.1.2 Convergence pour les fonctions continues (via Fejér)

Si nous ne voulons pas passer par le gros théorème de Stone-Weierstrass pour prouver la densité des polynômes trigonométriques dans $(C_{2\pi}^0, \|\cdot\|_\infty)$, nous pouvons passer par le gros théorème de Fejér. C'est ce que nous faisons maintenant.

Le **noyau de Dirichlet** est la fonction

$$D_n(t) = \sum_{k=-n}^n e^{ikt}. \quad (29.3)$$

Le **noyau de Fejér** est la moyenne de Cesaro des noyaux de Dirichlet :

$$F_n(t) = \frac{1}{n} \sum_{k=0}^{n-1} D_k(t). \quad (29.4)$$

Lemme 29.2.

Le noyau de Dirichlet s'exprime sous la forme

$$D_n(t) = \sum_{k=-n}^n e^{-ikt} = \frac{\sin\left(\frac{2n+1}{2}t\right)}{\sin(t/2)} \quad (29.5)$$

Note : ce noyau n'est pas positif.

Démonstration. Nous commençons par mettre en facteur le premier terme :

$$D_n(t) = \sum_{k=-n}^n e^{int} = e^{-int} \sum_{k=0}^{2n} e^{ikt}. \quad (29.6)$$

En utilisant la formule de la somme géométrique,

$$D_n(t) = e^{-int} \frac{1 - (e^{it})^{2n+1}}{1 - e^{it}} \quad (29.7a)$$

$$= e^{-int} \frac{1 - e^{(2n+1)it}}{1 - e^{it}} \quad (29.7b)$$

$$= e^{-int} \frac{e^{(2n+1)it/2} e^{-(2n+1)it/2} - e^{(2n+1)it/2}}{e^{i\frac{t}{2}} e^{-it/2} - e^{it/2}} \quad (29.7c)$$

$$= \frac{(-2i) \sin\left(\frac{2n+1}{2}t\right)}{(-2i) \sin\left(\frac{t}{2}\right)}. \quad (29.7d)$$

□

Théorème 29.3 (Théorème de Dirichlet).

Soit f une fonction 2π -périodique et C^1 par morceaux. Pour tout $x \in \mathbb{R}$ nous posons

$$s_n(x) = \sum_{k=-n}^n c_k(f) e^{ikx}. \quad (29.8)$$

Alors nous avons

$$\lim_{n \rightarrow \infty} s_n(x) = \frac{f(x^+) + f(x^-)}{2}. \quad (29.9)$$

Lemme 29.4.

Le noyau de Fejér s'exprime sous la forme

$$F_n(t) = \frac{1}{n} \left(\frac{\sin \frac{nt}{2}}{\sin \frac{t}{2}} \right)^2. \quad (29.10)$$

Note : ce noyau est positif. C'est important parce qu'on s'en sert dans la preuve du théorème de Fejér.

Démonstration. L'astuce est de noter $\sin(x) = \Im(e^{ix})$ et de repartir du résultat à propos du noyau

de Dirichlet. En utilisant encore la formule de la série géométrique partielle¹,

$$F_n(t) = \frac{1}{n \sin(t/2)} \Im \sum_{k=0}^{n-1} e^{(2k+1)it/2} \tag{29.11a}$$

$$= \frac{1}{n \sin(t/2)} \Im e^{\frac{it}{2}} \sum_{k=0}^{n-1} e^{ikt} \tag{29.11b}$$

$$= \frac{1}{n \sin(t/2)} \Im e^{\frac{it}{2}} \left(\frac{1 - e^{nit}}{1 - e^{it}} \right) \tag{29.11c}$$

$$= \frac{1}{n \sin(t/2)} \Im e^{it/2} \frac{e^{\frac{nit}{2}} \left(e^{-\frac{int}{2}} - e^{\frac{nit}{2}} \right)}{e^{\frac{it}{2}} \left(e^{-it/2} - e^{it/2} \right)} \tag{29.11d}$$

$$= \frac{1}{n \sin(t/2)} \underbrace{\Im e^{\frac{nit}{2}}}_{\sin(nt/2)} \frac{\sin\left(\frac{nt}{2}\right)}{\sin\left(\frac{t}{2}\right)} \tag{29.11e}$$

$$= \frac{1}{n} \left(\frac{\sin \frac{nt}{2}}{\sin \frac{t}{2}} \right)^2. \tag{29.11f}$$

□

Théorème 29.5 (Fejèr).

Soit $f : \mathbb{R} \rightarrow \mathbb{C}$ une fonction continue et 2π -périodique. Pour tout $k \in \mathbb{Z}$ nous notons

$$\begin{aligned} e_k : \mathbb{R} &\rightarrow \mathbb{C} \\ x &\mapsto e^{ikx}. \end{aligned} \tag{29.12}$$

Pour chaque $n \in \mathbb{N}$ nous posons

$$D_n = \sum_{k=-n}^n e_k \qquad \tilde{S}_n(f) = \sum_{k=-n}^n c_k(f) e_k \tag{29.13a}$$

$$F_n = \frac{D_0 + \dots + D_{n-1}}{n} \qquad \tilde{F}_n = \sigma_n(f) = \frac{1}{n} \sum_{k=0}^{n-1} S_k(f). \tag{29.13b}$$

Alors

- (1) $\frac{1}{2\pi} \int_{-\pi}^{\pi} F_n(t) dt = 1$.
- (2) Pour tout $\alpha \in]0, \pi[$, F_n converge uniformément vers 0 sur $[-\pi, \pi] \setminus [-\alpha, \alpha]$.
- (3) La suite \tilde{F}_n converge uniformément sur \mathbb{R} vers f .
- (4) Le système trigonométrique $\{e_k\}_{k \in \mathbb{Z}}$ est total pour l'espace $(C^0(S^1), \|\cdot\|_{\infty})$ des fonctions continues 2π -périodiques.

Démonstration. Un calcul usuel montre que

$$\int_{-\pi}^{\pi} e_l(t) dt = \begin{cases} 0 & \text{si } l \neq 0 \\ 2\pi & \text{si } l = 0 \end{cases} \tag{29.14}$$

Nous avons alors

$$\frac{1}{2\pi} \int_{-\pi}^{\pi} F_n(t) dt = \frac{1}{2\pi} \frac{1}{n} \sum_{k=0}^{n-1} \sum_{l=-k}^k \underbrace{\int_{-\pi}^{\pi} e_l(t) dt}_{2\pi \delta_l} = \frac{1}{n} \sum_{k=0}^{n-1} 1 = 1. \tag{29.15}$$

Cela prouve déjà le premier point.

1. Voir l'exemple 12.75.

Pour le second point, en partant de l'expression (29.10) et en considérant $x \in [-\pi, \pi] \setminus [-\alpha, \alpha]$ (ce qui nous évite l'annulation du dénominateur),

$$|F_n(x)| \leq \frac{1}{(n+1)\sin^2(\alpha/2)}, \quad (29.16)$$

et donc $F_n \rightarrow 0$ uniformément sur l'ensemble considéré.

Nous passons maintenant à cette histoire de convergence uniforme de la moyenne de Cesaro vers f . Pour tout $n \in \mathbb{N}$ nous avons

$$\tilde{D}_n(x) = \frac{1}{2\pi} \sum_{k=-n}^n \left(\int_{-\pi}^{\pi} f(t)e^{-ikt} dt \right) e^{ikx} \quad (29.17a)$$

$$= \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) \sum_{k=-n}^n e_k(x-t) \quad (29.17b)$$

$$= \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) D_k(x-t). \quad (29.17c)$$

Par conséquent, en effectuant le changement de variable $u = x - t$ et la périodicité,

$$\tilde{F}_n(x) = \int_{-\pi}^{\pi} f(t) F_n(x-t) dt \quad (29.18a)$$

$$= - \int_{x+\pi}^{x-\pi} f(x-u) F_n(u) du \quad (29.18b)$$

$$= \int_{-\pi}^{\pi} f(x-u) F_n(u) du. \quad (29.18c)$$

Nous prouvons à présent l'uniforme continuité. Soit $\epsilon > 0$; étant donné que f est continue et 2π -périodique, elle est uniformément continue et nous considérons $\delta > 0$ tel que $|x - y| < \delta$ implique $|f(x) - f(y)| < \epsilon$. Soit M un majorant de $|f|$ sur \mathbb{R} . L'équation (29.18) nous donne

$$|f(x) - \tilde{F}_n(x)| = \left| \frac{1}{2\pi} \int_{-\pi}^{\pi} (f(x-t) - f(x)) F_n(t) dt \right| \quad (29.19a)$$

$$\leq \frac{1}{2\pi} \int_{\delta \leq |t| \leq \pi} |2M F_n(t)| dt + \frac{1}{2\pi} \int_{-\delta}^{\delta} \epsilon |F_n(t)| dt \quad (29.19b)$$

$$\leq \frac{2M}{2\pi} \int_{\delta \leq |t| \leq \pi} F_n(t) dt + \epsilon' \quad (29.19c)$$

Pour obtenir (29.19a) nous avons pu rentrer $f(x)$ dans l'intégrale en utilisant le premier point. Pour obtenir (29.19c) nous avons d'abord utilisé la positivité de F_n (lemme 29.4) pour enlever les valeurs absolues, et nous avons ensuite utilisé le fait que son intégrale valait 2π .

Étant donné que $F_n \rightarrow 0$ uniformément sur $[-\pi, \pi] \setminus [-\alpha, \alpha]$, il existe un N tel que

$$\int_{\delta \leq |t| \leq \pi} F_n(t) dt \leq \epsilon \quad (29.20)$$

dès que $n > N$. Le résultat découle.

Pour le point (4), il suffit de remarquer que chacun des \tilde{F}_n est une combinaison finie d'éléments du système trigonométrique. \square

29.1.3 Densité dans L^p

Nous venons de voir (de deux façons différentes) que les polynômes trigonométriques étaient dense dans $(C_{2\pi}^0(\mathbb{R}), \|\cdot\|_{\infty})$. Nous avons aussi déjà vu par le théorème 28.70 que ces polynômes trigonométriques étaient denses dans $L^p(S^1)$. Nous présentons à présent une autre façon de prouver cette dernière densité.

Théorème 29.6.

Les polynômes trigonométriques sont denses dans $L^p(S^1)$ pour $1 \leq p < \infty$.

Démonstration. Par les théorèmes 29.1 ou 29.5 (au choix), nous savons que les polynômes trigonométriques sont denses dans $(C_{2\pi}^0(S^1), \|\cdot\|_\infty)$. Vu que S^1 est compact, la densité est également au sens L^p . En effet si $\|f_n - f\|_\infty \leq \epsilon$, alors

$$\|f_n - f\|_\infty = \int_0^{2\pi} |f_n - f|^p \leq \int_0^{2\pi} \epsilon^p = 2\pi\epsilon^p. \quad (29.21)$$

Donc les polynômes trigonométriques sont denses dans $(C_{2\pi}^0(S^1), \|\cdot\|_p)$. Mais nous savons par (un a fortiori sur) le théorème 28.47 que les fonctions continues sont denses dans $L^p(S^1)$.

Par densité de la densité, les polynômes trigonométriques sont denses dans $L^p(S^1)$. \square

29.1.4 Suite équirépartie, critère de Weyl**Définition 29.7.**

Soit u une suite dans $[0, 1]$. Pour $0 \leq a \leq b \leq 1$ nous posons

$$X_n(a, b) = \text{Card} \{k \in \{1, \dots, n\} \text{ tel que } u_k \in [a, b]\}. \quad (29.22)$$

Nous disons que la suite u est **équirépartie** si pour tout $0 \leq a < b < 1$, on a

$$\lim_{n \rightarrow \infty} \frac{X_n(a, b)}{n} = b - a. \quad (29.23)$$

Voir aussi la remarque 37.132 sur les nombres normaux.

Proposition 29.8 (Critère de Weyl[276, 43]).

Soit (x_n) une suite dans $[0, 1[$. Les conditions suivantes sont équivalentes.

- (1) La suite (x_n) est équirépartie.
- (2) Pour toute fonction continue à valeurs réelles sur $[0, 1]$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n f(x_k) = \int_0^1 f(x) dx. \quad (29.24)$$

- (3) Pour tout $p \in \mathbb{N}^*$ nous avons

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n e^{2i\pi p x_k} = 0. \quad (29.25)$$

Démonstration. On pose

$$S_n(f) = \frac{1}{n} \sum_{k=1}^n f(x_k). \quad (29.26)$$

Une espèce de lemme Supposons connaître un ensemble de fonctions A dense dans $C^0([0, 1])$ pour toutes les fonctions duquel nous avons la limite (29.24). Alors la limite a lieu pour toute fonction de $C^0([0, 1])$. En effet, soit $f \in C^0([0, 1])$ et $g \in A$ tel que $\|f - g\|_\infty < \epsilon$. Alors

$$\left\| \frac{1}{n} \sum_{k=1}^n f(x_k) - \int_0^1 f(t) dt \right\| \leq \left\| \frac{1}{n} \sum_{k=1}^n (f(x_k) - g(x_k)) \right\| \quad (29.27a)$$

$$+ \left\| \frac{1}{n} \sum_{k=1}^n g(x_k) - \int_0^1 g(t) dt \right\| \quad (29.27b)$$

$$+ \left\| \int_0^1 g(t) dt - \int_0^1 f(t) dt \right\|. \quad (29.27c)$$

Le premier terme se majore par ϵ . Le troisième est la même majoration : $\int_0^1 (f(t) - g(t)) dt \leq \|f - g\|_\infty = \epsilon$. Par hypothèse sur l'espace A , le second terme se majore par ϵ lorsque n est grand.

(1)⇒(2) Nous supposons que la suite est équirépartie et nous commençons par montrer le résultat pour les fonctions en escalier. Soit donc la fonction en escalier $\eta(x) = c_j$ sur $a_{j-1} < x < a_j$. Sur le point a_j lui-même, la fonction η vaut soit c_j soit c_{j+1} . Nous avons

$$\frac{1}{n} \sum_{k=1}^n \eta(x_k) = \frac{1}{n} \left[\sum_{j=1}^m c_j X_n(a_j, a_{j+1}) - \sum_{j=1}^m c_j X_n(a_j, a_j) + \sum_{j=1}^m \eta(a_j) X_n(a_j, a_j) \right]. \quad (29.28)$$

À la limite $n \rightarrow \infty$, les deux derniers termes tombent² et il reste

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \eta(x_k) = \sum_{j=1}^m c_j (a_{j-1} - a_j). \quad (29.29)$$

Or par construction, pour une fonction en escalier,

$$\sum_{j=1}^m c_j (a_{j-1} - a_j) = \int_0^1 \eta. \quad (29.30)$$

Étant donné que les fonctions en escalier sont denses dans les fonctions continues, l'espèce de lemme plus haut conclut.

(2)⇒(1) Nous prouvons maintenant le sens inverse. C'est-à-dire que pour toute fonction continue sur $[0, 1]$, nous avons

$$\int_0^1 f(x) dx = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n f(x_k). \quad (29.31)$$

Nous devons en déduire que (x_n) est équirépartie. Pour ce faire, soit $x \in [0, 1[$ et $\epsilon > 0$ tel que $x + \epsilon < 1$. Nous considérons $\varphi = \mathbb{1}_{[x, 1[}$ et

$$\varphi_\epsilon(t) = \begin{cases} 0 & \text{si } t \in [0, x[\\ \frac{t-x}{\epsilon} & \text{si } t \in [x, x + \epsilon[\\ 1 & \text{si } t \geq x + \epsilon. \end{cases} \quad (29.32)$$

Cela est une fonction continue, donc

$$\lim_{n \rightarrow \infty} S_n(\varphi_\epsilon(t)) = \int_0^1 \varphi_\epsilon(t) dt = \int_x^{x+\epsilon} \frac{t-x}{\epsilon} dt + \int_{x+\epsilon}^1 1 dt = 1 - x - \frac{\epsilon}{2}. \quad (29.33)$$

Mais $\varphi_\epsilon \leq \varphi$, donc $S_n(\varphi_\epsilon) \leq S_n(\varphi)$ et donc

$$\liminf_{n \rightarrow \infty} S_n(\varphi) \geq 1 - x. \quad (29.34)$$

Notons que nous ne savons pas si la vraie limite de gauche existe ; c'est pourquoi nous prenons la limite inférieure, qui existe toujours.

Nous définissons aussi

$$\psi_\epsilon(t) = \begin{cases} 0 & \text{si } t \in [0, x - \epsilon[\\ \frac{t-x+\epsilon}{\epsilon} & \text{si } t \in [x - \epsilon, x[\\ 1 & \text{si } t > x. \end{cases} \quad (29.35)$$

C'est encore une fonction continue et nous trouvons³

$$\int_0^1 \psi_\epsilon(t) dt = 1 - x + \frac{\epsilon}{2}. \quad (29.36)$$

2. J'en profite pour mentionner que mon équation (29.28) n'est pas la même que celle de [276] dans laquelle il me semble voir une faute ; quoi qu'il en soit, les termes litigieux tombent.

3. Je recommande chaudement de dessiner les fonctions φ_ϵ et ψ_ϵ pour avoir une idée de la situation.

Vu que $\psi_\epsilon \geq \varphi$, nous avons $S_n(\psi_\epsilon) \geq S_n(\varphi)$ et donc

$$\limsup_n S_n(\varphi) \leq 1 - x. \tag{29.37}$$

Nous avons déjà obtenu que

$$1 - x \leq \liminf S_n(\varphi) \leq \limsup S_n(\varphi) \leq 1 - x, \tag{29.38}$$

donc la limite existe et vaut

$$\lim_{n \rightarrow \infty} S_n(\varphi) = 1 - x. \tag{29.39}$$

Cela est pour la fonction caractéristique $\varphi = \mathbb{1}_{[x,1[}$. Si nous prenons une fonction caractéristique $\mathbb{1}_{[a,b]}$, nous avons la même chose parce que $\mathbb{1}_{[a,b]}$ est une combinaison linéaire de fonctions du type $\mathbb{1}_{[x,1[}$.

Nous avons donc

$$\lim_{n \rightarrow \infty} S_n(\mathbb{1}_{[a,b]}) = b - a, \tag{29.40}$$

alors que le membre de gauche n'est autre que

$$S_n(\mathbb{1}_{[a,b]}) = \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{[a,b]}(x_k) = \frac{1}{n} N(n, a, b). \tag{29.41}$$

(2) ⇒ (3) Vu que $e^{2i\pi p x_k} = \cos(2\pi p x_k) + i \sin(2\pi p x_k)$ est une fonction périodique, c'est immédiat.

(3) ⇒ (2) Par linéarité, le point (2) montre que si f est un polynôme trigonométrique, alors

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n f(x_k) = \int_0^1 f(t) dt. \tag{29.42}$$

Densité des polynômes trigonométriques Il nous reste à prouver que les polynômes trigonométriques sont denses dans les fonctions continues sur $[0, 1]$. Soit une fonction continue sur $[0, 1]$ avec $f(0) = f(1)$. Alors le théorème de Stone-Weierstrass dans sa version trigonométrique (lemme 29.1) nous donne la densité.

Si $f(1) \neq f(0)$ c'est pas très grave : on peut trouver une fonction g vérifiant $g(0) = g(1)$ et $\|f - g\|_\infty \leq \epsilon$. Ensuite un polynôme trigonométrique approxime très bien g . .

□

29.2 Fonctions de Dirichlet

Définition 29.9.

Une fonction $f : \mathbb{R} \rightarrow \mathbb{C}$ est une **fonction de Dirichlet** si

- (1) elle est 2π -périodique,
- (2) elle est continue par morceaux,
- (3) pour tout $x \in \mathbb{R}$ nous avons

$$f(x) = \frac{f(x^+) + f(x^-)}{2}. \tag{29.43}$$

Nous notons \mathcal{D} l'ensemble des fonctions de Dirichlet.

Lemme 29.10 ([431]).

L'ensemble $C^0(S^1)$ est dense dans $(\mathcal{D}, \|\cdot\|_2)$.

Démonstration. Nous commençons par supposer que $f \in \mathcal{D}$ n'ait qu'un seul point de discontinuité, x_0 . Alors nous considérons la fonction f_n qui est égale à f sur $S^1 \setminus B(x_0, \frac{1}{n})$ et qui sur $B(x_0, \frac{1}{n})$ est le segment de droite joignant $f(x_0 - \frac{1}{n})$ et $f(x_0 + \frac{1}{n})$. Cela est une fonction continue, et de plus nous avons

$$|f_n(x)| \leq \|f\|_\infty \quad (29.44)$$

pour tout x . En effet si x est en dehors de $B(x_0, \frac{1}{n})$ c'est évident, et si $x \in B(x_0, \frac{1}{n})$, alors $|f_n(x)|$ est majoré soit par $f(x_0 - \frac{1}{n})$ soit par $f(x_0 + \frac{1}{n})$ suivant que le raccord affine soit croissant ou décroissant. Avec ça nous avons

$$\|f_n - f\|_2^2 = \int_{x_0 - 1/n}^{x_0 + 1/n} |f(x) - f_n(x)|^2 dx \leq \int_{x_0 - 1/n}^{x_0 + 1/n} 4\|f\|_\infty^2 dx = \frac{8\|f\|_\infty^2}{n}. \quad (29.45)$$

Et nous voyons que $\|f_n - f\|_2 \rightarrow 0$.

Si f contient plusieurs points de continuité, on fait le même coup autour de chaque point, en prenant n assez grand pour que si x_0 est un point de discontinuité, $B(x_0, \frac{1}{n})$ n'en contienne pas d'autres. \square

Notons que la densité de $C^0(S^1)$ dans $(\mathcal{D}, \|\cdot\|_\infty)$ est impossible parce qu'une limite uniforme de fonctions continues est continue.

Théorème 29.11.

Le système trigonométrique $\{e_n\}_{n \in \mathbb{Z}}$ est total dans $(\mathcal{D}, \|\cdot\|_2)$.

Démonstration. Soit $f \in \mathcal{D}$. Si elle est continue, le théorème de Fejèr 29.5 nous donne convergence uniforme sur S^1 d'une suite de polynômes trigonométriques vers f . Cette convergence est également une convergence L^2 parce que S^1 est compact.

Prenons donc $f \in \mathcal{D}$ non continue et $\epsilon > 0$ ⁴. Par le lemme 29.10, il existe une fonction $g \in C^0(S^1)$ telle que

$$\|g - f\|_2 \leq \epsilon. \quad (29.46)$$

Le théorème de Fejèr donne aussi un polynôme trigonométrique P tel que $\|P - g\|_2 < \epsilon$; nous avons alors

$$\|P - f\|_2 \leq \|P - g\|_2 + \|g - f\|_2 \leq 2\epsilon. \quad (29.47)$$

\square

Notons que cette histoire de fonctions de Dirichlet n'a pas attaqué le vrai fond du problème de la densité des polynômes trigonométriques dans $L^2(S^1)$ parce que nous restons avec une hypothèse de continuité, alors que les représentants des éléments de $L^2(S^1)$ n'ont strictement aucune régularité a priori.

29.3 Coefficients et série de Fourier

Définition 29.12.

La série de Fourier associée à f est

$$f(x) \sim \sum_{n=-\infty}^{\infty} c_n(f) e^{2\pi i \frac{n}{T} x}. \quad (29.48)$$

Cette expression est pour l'instant purement formelle. Cela ne présume ni de la convergence de la série, ni, au cas où elle serait convergente, que la limite soit f .

Pour la suite nous allons considérer des fonctions périodiques de période 2π , et les coefficients de Fourier de f (quand ils existent) sont alors

$$c_n(f) = \frac{1}{2\pi} \int_0^{2\pi} f(t) e^{-int} dt \quad (29.49)$$

4. Par exemple $\epsilon = 0.4$, mais ce n'est qu'un exemple hein. Si vous en voulez un autre, prenez p , un nombre premier puis calculez $\epsilon = 1/p$.

Proposition 29.13 ([432]).

Soit f une fonction continue et périodique telle que sa série de Fourier converge uniformément. Alors la convergence est vers f .

Démonstration. Notons d'abord que f étant continue sur $[0, 2\pi]$, elle y est bornée et L^2 . Par conséquent Parseval nous enseigne que

$$\|S_N(f) - f\|_{L^2} \rightarrow 0. \quad (29.50)$$

Cela signifie que

$$\lim_{N \rightarrow \infty} \frac{1}{2\pi} \int_0^{2\pi} |f(t) - S_N(t)|^2 dt = 0. \quad (29.51)$$

L'hypothèse de convergence uniforme nous dit que la fonction $|f(t) - S_N(t)|^2$ converge uniformément vers la fonction $|f(t) - S(t)|^2$ où nous avons écrit S la limite de S_N . En permutant la limite et l'intégrale,

$$\frac{1}{2\pi} \int_0^{2\pi} |f(t) - S(t)|^2 dt = 0, \quad (29.52)$$

ce qui signifie que la fonction $t \mapsto |f(t) - S(t)|^2$ est la fonction nulle. Nous en déduisons que $f = S$. \square

Proposition 29.14.

Soit f une fonction 2π -périodique. Si $\sum_{n \in \mathbb{Z}} |c_n(f)| < \infty$, alors pour tout $x \in \mathbb{R}$ nous avons

$$f(x) = \sum_{n \in \mathbb{Z}} c_n(f) e^{inx}. \quad (29.53)$$

De plus, la suite $(S_n f)$ converge uniformément vers f .

Démonstration. Nous posons

$$g(x) = \sum_{n \in \mathbb{Z}} c_n(f) e^{inx}. \quad (29.54)$$

Étant donné les hypothèses, la série de droite converge absolument, la fonction g est continue sur \mathbb{R} . Nous avons

$$|g(x) - (S_n f)(x)| \leq \sum_{|k| > n} |c_k(f)|, \quad (29.55)$$

mais le terme de droite tend vers zéro lorsque $n \rightarrow \infty$ parce que c'est le reste d'une série convergente. Cela signifie que $S_n f$ converge uniformément vers g .

Par ailleurs nous savons que dans L^2 nous avons la convergence $S_n f \rightarrow f$ (parce que f est continue sur le compact $[0, 2\pi]$ et donc y est bornée et L^2), ce qui signifie que $g = f$ presque partout au sens L^2 . Ces deux fonctions étant continues, elles sont égales partout. \square

Théorème 29.15.

Soit f , une fonction C^1 et 2π -périodique. Nous notons $(c_n)_{n \in \mathbb{Z}}$ la suite de ses coefficients de Fourier. Alors $(c_n) \in \ell^1(\mathbb{Z})$ et pour tout $x \in \mathbb{R}$ nous avons

$$f(x) = \sum_{n \in \mathbb{Z}} c_n(f) e^{inx}. \quad (29.56)$$

Démonstration. Soit $n \in \mathbb{Z}$. Nous posons $g(t) = f(t) e^{-int}$. Nous avons

$$0 = g(2\pi) - g(0) = \int_0^{2\pi} g'(t) dt = \int_0^{2\pi} [f'(t) e^{-int} - in f(t) e^{-int}]. \quad (29.57)$$

Du coup, $c_n(f') = in c_n(f)$. La fonction f' étant bornée (parce que continue sur $[0, 2\pi]$), elle est de carré intégrable sur $[0, 2\pi]$ et par les inégalités de Parseval (théorème 26.44) nous avons

$$\sum_{n \in \mathbb{Z}} |c_n(f')|^2 < \infty. \quad (29.58)$$

Par conséquent $(c_n(f')) \in \ell^2(\mathbb{Z})$ et a fortiori $(c_n(f'))_{n \in \mathbb{N}} \in \ell^2(\mathbb{N})$. L'inégalité de Cauchy-Schwarz nous indique alors

$$\sum_{n \in \mathbb{N}} |c_n(f)| = \sum_{n \in \mathbb{N}} \frac{1}{n} |c_n(f')| \leq \left(\sum_n \frac{1}{n^2} \right)^{1/2} \left(\sum_n |c_n(f')|^2 \right)^{1/2} < \infty. \quad (29.59)$$

Nous procédons de même pour $n < 0$. Cela prouve que

$$\sum_{n \in \mathbb{Z}} |c_n(f)| < \infty. \quad (29.60)$$

□

Corollaire 29.16.

Soient f, g deux fonctions continues et 2π -périodiques. Si $c_n(f) = c_n(g)$ alors $f = g$.

Démonstration. Dans le cas de fonctions continues, le théorème de Fejér nous enseigne que si nous posons

$$S_n(x) = \sum_{k=-n}^n c_k(f) e^{ikx} \quad (29.61)$$

alors nous avons la convergence

$$\frac{1}{N+1} \sum_{n=0}^N S_n(f)(x) \rightarrow f(x). \quad (29.62)$$

C'est-à-dire qu'une fonction continue est déterminée par ses coefficients de Fourier. □

Exemple 29.17

Considérons la fonction

$$f(x) = 1 - \frac{x^2}{\pi^2} \quad (29.63)$$

sur $[-\pi, \pi]$. Nous la développons en série trigonométrique, et étant paire il n'y a pas de sinus. Un calcul montre que

$$a_0 = \frac{4}{3} \quad (29.64)$$

et

$$a_n = (-1)^{n+1} \frac{4}{n^2 \pi^2}, \quad (29.65)$$

de telle sorte que

$$f(x) = \frac{2}{3} - \frac{4}{\pi^2} \sum_{n=1}^{\infty} (-1)^n \frac{\cos(nx)}{n^2}. \quad (29.66)$$

Nous avons $f(\pi) = 0$, mais vu le développement,

$$f(\pi) = \frac{2}{3} - \frac{4}{\pi^2} \sum_{n=1}^{\infty} \frac{1}{n^2}, \quad (29.67)$$

donc

$$\sum_{n=1}^{\infty} \frac{1}{n^2} = \frac{\pi^2}{6}. \quad (29.68)$$

△

29.3.1 Le contre-exemple que nous attendions tous

Nous montrons maintenant que la continuité et la périodicité ne sont pas suffisantes pour avoir convergence de la série de Fourier.

Proposition 29.18 ([43]).

Soit $C_{2\pi}^0$ l'ensemble des fonctions continues muni de la norme uniforme. Nous définissons

$$S_n(f)(x) = \sum_{k=-n}^n c_k(f)e^{ikx}. \quad (29.69)$$

Alors il existe $f \in C_{2\pi}^0$ tel que la suite $n \mapsto S_n(f)(0)$ soit divergente. En particulier f n'est pas la somme de sa série de Fourier.

Démonstration. Nous considérons la forme linéaire

$$\begin{aligned} l_n : C_{2\pi}^0 &\rightarrow \mathbb{C} \\ f &\mapsto S_n(f)(0) = \sum_{k=-n}^n c_k(f). \end{aligned} \quad (29.70)$$

La forme est continue Nous montrons d'abord que $\|l_n\|$ est continue en montrant que $\|l_n\| < \infty$ et en utilisant la proposition 12.25. Pour cela nous calculons un peu :

$$l_n(f) = \sum_{k=-n}^n \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t)e^{-ikt} dt = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) \sum_{k=-n}^n e^{-ikt} dt = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t)D_n(t) dt \quad (29.71)$$

où $D_n(t)$ est le noyau de Dirichlet dont nous savons une formule par le lemme 29.2. Nous avons donc

$$|l_n(f)| \leq \frac{1}{2\pi} \int_{-\pi}^{\pi} |D_n(t)| \|f\|_{\infty} dt. \quad (29.72)$$

En prenant $\|f\|_{\infty} = 1$ nous avons la borne suivante pour la norme de l_n :

$$\|l_n\| \leq \frac{1}{2n} \int_{-\pi}^{\pi} |D_n(t)| dt < \infty. \quad (29.73)$$

Notons que la convergence de l'intégrale vient de la continuité de la fonction

$$t \mapsto \frac{\sin\left(\frac{2n+1}{2}t\right)}{\sin\left(\frac{t}{2}\right)} \quad (29.74)$$

qui, elle même, se prouve avec une règle de l'Hospital :

$$\lim_{t \rightarrow 0} \frac{\sin(at)}{\sin(t)} = \lim_{t \rightarrow 0} \frac{a \cos(at)}{\cos(t)} = a. \quad (29.75)$$

Donc $D_n(t)$ a une limite bien définie pour $t \rightarrow 0$ et est alors une fonction continue sur le compact $[-\pi, \pi]$.

La norme de l_n (début) Nous avons prouvé que $\|l_n\| \leq \frac{1}{2n} \int_{-\pi}^{\pi} |D_n(t)| dt$. Nous allons à présent prouver que cela est effectivement la norme de l_n . Pour $\epsilon > 0$ nous considérons la fonction

$$\begin{aligned} f_{\epsilon} : \mathbb{R} &\rightarrow \mathbb{C} \\ x &\mapsto \frac{D_n(x)}{|D_n(x)| + \epsilon}. \end{aligned} \quad (29.76)$$

C'est une fonction continue et 2π -périodique satisfaisant $\|f_{\epsilon}\| \leq 1$ parce que le dénominateur est toujours plus grand que le numérateur. Nous nous proposons de calculer

$$l_n(f_{\epsilon}) = \sum_{k=-n}^n \frac{1}{2\pi} \int_{-\pi}^{\pi} f_{\epsilon}(t)e^{-ikt} dt. \quad (29.77)$$

Vu que $f_\epsilon(t)e^{-ikt}$ vaut en norme $|f_\epsilon(t)|$ qui est une fonction intégrable (ne dépendant pas de k) sur $[-\pi, \pi]$, le théorème de la convergence dominée 15.184 nous permet de permuter la somme et l'intégrale :

$$l_n(f_\epsilon) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{D_n(t)}{|D_n(t)| + \epsilon} \underbrace{\sum_{k=-n}^n e^{-ikt}}_{=D_n(t)} dt = \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{|D_n(t)|^2}{|D_n(t)| + \epsilon} dt. \quad (29.78)$$

Nous avons donc

$$\lim_{\epsilon \rightarrow 0} l_n(f_\epsilon) = \frac{1}{2\pi} \int_{-\pi}^{\pi} |D_n(t)| dt. \quad (29.79)$$

Mais vue l'inégalité (29.73) nous avons

$$\|l_n\| = \frac{1}{2\pi} \int_{-\pi}^{\pi} |D_n(t)| dt. \quad (29.80)$$

Notre tâche est maintenant de donner une valeur à cette intégrale.

Norme de l_n tend vers ∞ D'abord nous écrivons

$$\|l_n\| = \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{|\sin(\frac{2n+1}{2}t)|}{|\sin(t/2)|} dt, \quad (29.81)$$

ensuite nous nous souvenons que $|\sin(x)| \leq |x|$ pour tout x , ce qui nous permet de changer le dénominateur :

$$\|l_n\| \geq \frac{2}{\pi} \int_0^{\pi} \frac{|\sin(\frac{2n+1}{2}t)|}{|t|} dt \quad (29.82)$$

Nous y effectuons le changement de variable $u = \frac{2n+1}{2}t$ qui donne

$$\|l_n\| \geq \frac{2}{\pi} \int_0^{(n+\frac{1}{2})\pi} \frac{|\sin(u)|}{|u|} du. \quad (29.83)$$

Nous y reconnaissons l'intégrale (18.484) du sinus cardinal que nous savons diverger. Cela donne

$$\lim_{n \rightarrow \infty} \|l_n\| = \infty. \quad (29.84)$$

La conclusion L'espace $(C_{2\pi}^0, \|\cdot\|_\infty)$ est complet⁵, donc le théorème de Banach-Steinhaus 12.92 s'applique. Par rapport aux notations de l'énoncé de Banach-Steinhaus, nous posons

$$E = (C_{2\pi}^0, \|\cdot\|_\infty) \quad (29.85a)$$

$$F = \mathbb{R} \quad (29.85b)$$

$$H = \{l_n\}_{n \in \mathbb{N}}. \quad (29.85c)$$

Vu que la suite $(\|l_n\|)$ n'est pas bornée, il existe $f \in C_{2\pi}^0$ tel que

$$\sup_n \|l_n(f)\| = \infty. \quad (29.86)$$

Pour cette fonction nous avons

$$\sup_{n \geq 0} S_n(f)(0) = \infty, \quad (29.87)$$

et donc la série de Fourier de f ne converge pas en zéro.

□

5. Parce qu'une limite uniforme de fonctions continues est continue, théorème 13.280.

29.3.2 Inégalité isopérimétrique

Définition 29.19.

Une **courbe de Jordan** dans le plan est une application $\gamma: S^1 \rightarrow \mathbb{R}^2$ qui est continue et injective.

Une telle courbe peut évidemment être vue comme une application $\gamma: [0, 2\pi] \rightarrow \mathbb{R}^2$ telle que $\gamma(0) = \gamma(2\pi)$. En particulier il n'est jamais mauvais de se rappeler qu'on peut choisir une paramétrisation normale par la proposition 22.46.

Théorème 29.20 (Théorème de Jordan[433]).

Si γ est une courbe de Jordan, alors l'ensemble $\mathbb{R}^2 \setminus \gamma$ a exactement deux composantes connexes. L'une est bornée, l'autre non. Les deux ont γ comme frontière.

Le théorème suivant dit que parmi les courbes C^1 , le cercle a la plus grande surface possible à périmètre donné.

Théorème 29.21 (Inégalité isopérimétrique[43]).

Soit $f: S^1 \rightarrow \mathbb{C}$ une courbe de Jordan de classe C^1 . Nous notons L sa longueur et S l'aire contenue de la surface délimitée⁶ par f . Alors

- (1) Nous avons l'**inégalité isopérimétrique** : $L^2 \geq 4\pi S$.
- (2) Nous avons l'égalité $L^2 = 4\pi S$ si et seulement si la courbe donnée par f est un cercle.

Démonstration. Nous commençons par considérer un chemin dont la longueur est 2π et nous en considérons sa paramétrisation normale. Nous allons exprimer l'aire S en utilisant le théorème de Green, et plus particulièrement la formule de surface (21.256).

Si $f(s) = x(s) + iy(s)$, nous devons intégrer $y'x - x'y$, qui n'est rien d'autre que la partie imaginaire de $f'(s)\overline{f(s)}$. Donc

$$S = \frac{1}{2} \operatorname{Im} \int_0^{2\pi} f'(s)\overline{f(s)} ds \quad (29.88)$$

Nous considérons les coefficients de Fourier de f donnés par la formule (29.49) :

$$c_n(f) = \frac{1}{2\pi} \int_0^{2\pi} f(s) e^{-ins} ds. \quad (29.89)$$

Ceux de f' (qui est aussi continue sur le compact S^1 et donc tout autant L^2) sont donnés par

$$c_n(f') = in c_n(f). \quad (29.90)$$

D'autre part en vertu du théorème 22.10, la longueur de γ s'exprime en terme de l'intégrale de la norme de sa dérivée :

$$2\pi = L = \int_0^{2\pi} |f'(s)| ds = \int_0^{2\pi} |f'(s)|^2 ds \quad (29.91)$$

parce que nous avons choisi une paramétrisation normale qui vérifie automatiquement $|f'(s)| = 1$ pour tout s . L'identité de Parseval sous sa forme (26.109) appliquée à f' nous enseigne que

$$L = 2\pi = \int_0^{2\pi} |f'(s)|^2 ds = 2\pi \langle f', f' \rangle = 2\pi \sum_{n=-\infty}^{\infty} |c_n(f')|^2 = 2\pi \sum_n n^2 |c_n(f)|^2. \quad (29.92)$$

Par ailleurs le système trigonométrique étant une base hilbertienne, et les fonctions f et f' étant dans $L^2([0, 2\pi])$ (parce que continues sur un compact), elles sont égales à leurs séries de Fourier

6. C'est la partie connexe bornée de $\mathbb{C} \setminus \gamma$ dont l'existence est donnée par le théorème de Jordan 29.20.

(au sens L^2), c'est-à-dire que nous avons l'égalité (28.336). Nous avons alors

$$\langle f', f \rangle_{L^2} = \left\langle \sum_{n \in \mathbb{Z}} c_n(f') e_n, \sum_{m \in \mathbb{Z}} c_m(f) e_m \right\rangle \quad (29.93a)$$

$$= \sum_m \sum_n c_n(f') \overline{c_m(f)} \underbrace{\langle e_n, e_m \rangle}_{\delta_{m,n}} \quad (29.93b)$$

$$= \sum_{n \in \mathbb{Z}} c_n(f') \overline{c_n(f)} \quad (29.93c)$$

$$= \sum_n i n |c_n(f)|^2 \quad (29.93d)$$

où nous avons utilisé la continuité du produit scalaire pour sortir les sommes. Avec cela nous pouvons exprimer l'aire (29.88) en termes de coefficients de Fourier :

$$S = \frac{1}{2} \operatorname{Im} 2\pi \langle f', f \rangle = \pi \sum_{n \in \mathbb{Z}} n |c_n(f)|^2. \quad (29.94)$$

En utilisant les expressions (29.92) et (29.94) pour L et S , et en écrivant $L = 2\pi L$, nous avons

$$L^2 - 4\pi S = 4\pi^2 \sum_{n \in \mathbb{Z}} (n^2 - n) |c_n(f)|^2 \geq 0. \quad (29.95)$$

Cela prouve l'inégalité demandée dans le cas où $L = 2\pi$.

Si γ n'est pas de longueur 2π mais L , alors nous considérons le chemin $\sigma(t) = \frac{2\pi\gamma(t)}{L}$. Sa longueur est 2π et son aire, au vu de la formule de Green (29.88), son aire est $4\pi^2 \frac{S}{L^2}$. L'inégalité isopérimétrique appliquée au chemin σ donne alors $L^2 \geq 4\pi S$.

Le cas d'égalité s'obtient uniquement si $c_n = 0$ pour tout n différent de 0 ou 1. Dans ce cas nous avons

$$f(s) = c_0(f) + c_1(f) e^{is}, \quad (29.96)$$

qui est un cercle de centre $c_0(f)$ et de rayon $|c_1(f)|$. \square

29.3.3 À propos des coefficients

Nous considérons l'application

$$\begin{aligned} c: (L^1_{2\pi}, \|\cdot\|_1) &\rightarrow (C_0, \|\cdot\|_\infty) \\ f &\mapsto (c_n(f))_{n \in \mathbb{Z}} \end{aligned} \quad (29.97)$$

qui à une fonction 2π -périodique fait correspondre la suite (bornée) de ses coefficients de Fourier. Nous rappelons la définition

$$c_n(f) = \frac{1}{2\pi} \int_0^{2\pi} f(t) e^{-int} dt. \quad (29.98)$$

Nous allons montrer que cette application est linéaire, continue, injective et non surjective. Pour la continuité, par la linéarité il suffit de la montrer en 0. Nous devons donc montrer que si nous avons une suite de fonctions f_k qui tend vers 0 au sens L^1 , alors $c(f_k) \rightarrow 0$ au sens de la norme $\|\cdot\|_\infty$ sur l'ensemble des suites.

Si nous posons $r_k = \int_0^{2\pi} |f_k(t)| dt$, alors $r_k = \|f_k\|_1$ et nous avons $r_k \rightarrow 0$. Mais par définition

$$|c_n(f_k)| \leq r_k, \quad (29.99)$$

et donc $\|c(f_k)\|_\infty \leq r_k$. L'application c est donc continue. L'injectivité est donnée par le corollaire 29.16.

Si nous supposons que l'application c est continue, alors le théorème d'isomorphisme de Banach (28.1) nous dit que cela devrait être un homéomorphisme, c'est-à-dire que c^{-1} serait également continue. Nous allons montrer qu'il n'en est rien.

Nous considérons la suite de suite

$$(c_n)_k = \begin{cases} 1 & \text{si } k < n \\ 0 & \text{sinon.} \end{cases} \quad (29.100)$$

Ici $(c_n)_k$ est le terme numéro k de la suite n . Par injectivité de l'application qui à une fonction fait correspondre la suite de ses coefficients de Fourier, la seule fonction qui possède ces coefficients est

$$f_n(t) = \sum_{k \in \mathbb{N}} c_{n,k} e^{ikt}. \quad (29.101)$$

Étant donné que $\|f_n\|_1 = n$, la suite $(\|f_n\|_1)$ n'est pas bornée alors que la suite de suites (29.100) est bornée dans l'ensemble des suites parce que $\|c_n\|_\infty = 1$.

Chapitre 30

Transformation de Fourier

Définition 30.1.

Soit une fonction f sur \mathbb{R}^d , dont nous ne précisons pas la régularité. Sa **transformée de Fourier** est la fonction \hat{f} définie par

$$\hat{f}(\xi) = \int_{\mathbb{R}^d} f(x) e^{-i\xi \cdot x} dx \quad (30.1)$$

si elle existe.

Ce qui est bien avec cette définition est que si la formule (30.1) ne définit pas \hat{f} (parce que l'intégrale n'existe pas, par exemple), nous nous réservons le droit de définir tout de même \hat{f} par d'autres biais. Ce sera d'ailleurs l'objet du théorème 30.28 qui définira \hat{f} pour tout $f \in L^2(\mathbb{R}^d)$ alors que la formule (30.1) ne fonctionne pas sur toutes ces fonctions.

Une bonne partie de ce qui va suivre aura pour objet de déterminer des espaces de fonctions sur lesquels la transformée est bien définie, et sur lesquels elle a de bonnes propriétés.

Nous allons par ailleurs utiliser indifféremment les notations $\mathcal{F}(f)$ ou \hat{f} pour la transformée de Fourier de f . La notation \mathcal{F} est pratique pour les transformées de loooooongues expressions ainsi que pour parler de l'application « transformée de Fourier » d'un espace de fonction vers un autre.

30.2.

Nous verrons dans le théorème 30.28 que la transformée de Fourier n'est pas une isométrie de L^2 . Pour avoir une isométrie, il aurait fallu choisir des coefficients moins simples.

30.1 Transformée de Fourier sur $L^1(\mathbb{R}^d)$

Nous rappelons que les espaces L^p sont des ensembles de classes de fonctions, définition 28.8. La transformée de Fourier, comme presque tout ce qui a trait aux intégrales, passe aux classes.

Lemme 30.3.

Soit une fonction $f: \mathbb{R}^d \rightarrow \mathbb{C}$ telle que $\mathcal{F}(f)$ existe. Alors pour toute fonction $g \in [f]$ la transformée $\mathcal{F}(g)$ existe et $\hat{g} = \hat{f}$.

Démonstration. Par définition des classes, il existe une fonction $s: \mathbb{R}^d \rightarrow \mathbb{C}$ presque partout nulle telle que $g = f + s$. Soit $\xi \in \mathbb{R}^d$ fixé. La fonction $x \mapsto s(x)e^{-i\xi x}$ est presque partout nulle et donc intégrable d'intégrale nulle. La proposition 15.172 nous permet alors d'affirmer que $f + s$ est intégrable et que

$$\mathcal{F}(f + s)(\xi) = \int_{\mathbb{R}^d} (f + s)(x) e^{-i\xi x} dx = \int_{\mathbb{R}^d} f(x) e^{-i\xi x} dx + \int_{\mathbb{R}^d} s(x) e^{-i\xi x} dx = \mathcal{F}(f)(\xi). \quad (30.2)$$

□

À partir de maintenant, lorsque nous parlons de transformée de Fourier d'une fonction dans L^p , nous parlons indifféremment d'une vraie fonction ou d'une classe.

Lemme 30.4.

Si $f \in L^1(\mathbb{R}^d)$, alors \hat{f} existe.

Démonstration. Par définition, si $f \in L^1(\mathbb{R}^d)$, alors l'intégrale $\int_{\mathbb{R}^d} |f|$ existe et est finie. Alors la fonction qui arrive dans la transformée de Fourier en ξ , la fonction $s: x \rightarrow f(x)e^{-i\xi x}$ est également dans $L^1(\mathbb{R}^d)$ parce que $|s| = |f|$. \square

<+++>

Lemme 30.5.

Si $f \in L^1(\mathbb{R}^d)$ et si $g(x) = f(\lambda x)$ alors

$$\hat{g}(\xi) = \lambda^{-d} \hat{f}(\xi/\lambda). \quad (30.3)$$

Démonstration. Il s'agit de faire le changement de variable $y = \lambda x$ dans l'intégrale

$$\hat{g}(\xi) = \int_{\mathbb{R}^d} f(\lambda x) e^{-i\xi x} dx. \quad (30.4)$$

Dans le changement de variables, vient le coefficient $dx = \lambda^{-d} dy$. \square

Proposition 30.6.

La transformée de Fourier est un morphisme vis-à-vis de la convolution sur $L^1(\mathbb{R}^n)$:

$$\widehat{f * g} = \hat{f} \hat{g}. \quad (30.5)$$

Démonstration. Nous devons étudier l'intégrale

$$\widehat{f * g}(\xi) = \int_{\mathbb{R}} \left[\int_{\mathbb{R}} f(y)g(t-y) \right] e^{-it\xi} dt. \quad (30.6)$$

Ici nous avons choisi des représentants f et g dans les classes de L^1 . Montrons que f est borélienne. D'abord $f(x) = f_+(x) - f_-(x)$ où f_+ et f_- sont des fonctions positives. Afin d'alléger les notations nous supposons un instant que f est positive et nous posons

$$f_n(x) = \sum_{k=1}^{2^n} \frac{k}{n} \mathbb{1}_{f(x) \in [\frac{k}{n}, \frac{k+1}{n}[}. \quad (30.7)$$

Le fait que f soit dans L^1 implique que chacune des fonctions f_n est borélienne¹ et donc que f l'est aussi en tant que limite ponctuelle de fonctions boréliennes².

Nous allons appliquer le théorème de Fubini 15.258 à la fonction

$$\phi(x, y) = f(x)g(y)e^{-i\xi(x+y)} \quad (30.8)$$

qui est borélienne en tant que produit et composé de fonctions boréliennes. Nous avons

$$\int_{\mathbb{R}} \left(\int_{\mathbb{R}} |f(x)e^{-i\xi x}| |g(y)e^{-i\xi y}| dy \right) dx = \int_{\mathbb{R}} \left(|f(x)| \int_{\mathbb{R}} |g(y)| dy \right) dx \quad (30.9a)$$

$$= \int_{\mathbb{R}} |f(x)| \|g\|_1 \quad (30.9b)$$

$$= \|f\|_1 \|g\|_1 < \infty. \quad (30.9c)$$

Le théorème est donc applicable. D'abord nous avons :

$$\hat{f}(\xi)\hat{g}(\xi) = \left(\int_{\mathbb{R}} f(x)e^{-i\xi x} dx \right) \left(\int_{\mathbb{R}} g(y)e^{-i\xi y} dy \right) \quad (30.10a)$$

$$= \int_{\mathbb{R}} \left(\int_{\mathbb{R}} f(x)g(y)e^{-i\xi(x+y)} dy \right) dx \quad (30.10b)$$

$$= \int_{\mathbb{R}} \left(\int_{\mathbb{R}} f(x)g(t-x)e^{-i\xi t} \right) dx. \quad (30.10c)$$

1. Ceci demanderait plus de justification. Dites moi si vous savez comment justifier que les f_n soient boréliennes.

2. Le fait que f soit borélienne est une conséquence du théorème 28.152.

Jusqu'ici nous n'avons pas utilisé Fubini. Nous avons seulement introduit le nombre $\int_{\mathbb{R}} g(y)e^{-i\xi y} dy$ dans l'intégrale par rapport à x et effectué le changement de variables $y \mapsto t = x + y$. Maintenant nous appliquons le théorème de Fubini pour inverser l'ordre des intégrales :

$$\hat{f}(\xi)\hat{g}(\xi) = \int_{\mathbb{R}} \left(\int_{\mathbb{R}} f(x)g(t-x)e^{-it\xi} dx \right) dy \quad (30.11a)$$

$$= \int_{\mathbb{R}} e^{-it\xi} \left(\int_{\mathbb{R}} f(x)g(t-x) dx \right) dt \quad (30.11b)$$

$$= \int_{\mathbb{R}} e^{-it\xi} (f * g)(t) dt \quad (30.11c)$$

$$= \widehat{f * g}(\xi). \quad (30.11d)$$

□

Proposition 30.7.

Soit une fonction $f \in L^1(\mathbb{R}^d)$. Alors sa transformée de Fourier est continue.

Démonstration. Nous considérons une fonction f définie sur \mathbb{R}^d et à valeurs dans \mathbb{R} ou \mathbb{C} . Sa transformée de Fourier est donnée par

$$\hat{f}(\xi) = \int_{\mathbb{R}^d} e^{-i\xi x} f(x) dx. \quad (30.12)$$

Pour montrer que cette fonction \hat{f} est continue en ξ_0 nous considérons une suite $(\xi_n) \rightarrow \xi_0$ et nous voulons montrer que $\hat{f}(\xi_n) \rightarrow \hat{f}(\xi_0)$. Pour cela nous considérons les fonctions

$$g_n(x) = e^{-i\xi_n x} f(x) \quad (30.13)$$

qui convergent simplement vers $g(x) = e^{-i\xi x} f(x)$. Étant donné que

$$|g_n(x)| < |f(x)|, \quad (30.14)$$

le théorème de la convergence dominée donne alors

$$\lim_{n \rightarrow \infty} \int g_n(x) = \int \lim_{n \rightarrow \infty} g_n(x), \quad (30.15)$$

c'est-à-dire $\lim_{n \rightarrow \infty} \hat{f}(\xi_n) = \hat{f}(\xi)$. La fonction \hat{f} est donc continue. □

Lemme 30.8.

Pour tout $f \in L^1(\mathbb{R}^n)$ nous avons $\|\hat{f}\|_{\infty} \leq \|f\|_1$.

Démonstration. Cela est une simple vérification :

$$\hat{f}(\xi) = \int_{\mathbb{R}^n} f(x)e^{-ix\xi} dx, \quad (30.16)$$

nous avons, pour tout ξ ,

$$|\hat{f}(\xi)| \leq \int_{\mathbb{R}^n} |f(x)| dx, \quad (30.17)$$

ce qui signifie exactement $\|\hat{f}\|_{\infty} \leq \|f\|_1$. □

Lemme 30.9 (Lemme de Riemann-Lebesgue[434]).

Si f est une fonction $L^1(\mathbb{R})$ alors $\lim_{\xi \rightarrow \pm\infty} \hat{f}(\xi) = 0$.

Démonstration. Nous commençons par prouver le résultat dans le cas d'une fonction g en escalier, et plus précisément par une fonction caractéristique d'un compact $K = [a, b]$. Au niveau de la transformée de Fourier nous avons

$$\hat{\mathbb{1}}_K(\xi) = \int_a^b e^{-i\xi x} dx = -\frac{1}{i\xi}(e^{-ib\xi} - e^{-ia\xi}). \quad (30.18)$$

Par conséquent

$$|\hat{\mathbb{1}}_K(\xi)| \leq \frac{2}{|\xi|}. \quad (30.19)$$

Plus généralement si $g = \sum_{i=1}^N c_i \mathbb{1}_{K_i}$, alors

$$|\hat{g}(\xi)| \leq \frac{2}{|\xi|} \sum_{i=1}^N |c_i|, \quad (30.20)$$

et donc nous avons effectivement $\lim_{\xi \rightarrow \pm\infty} |\hat{g}(\xi)| = 0$.

Nous passons maintenant au cas général $f \in L^1(\mathbb{R})$. Étant donné que les fonctions L^1 en escalier sont denses dans L^1 , nous considérons une fonction $g \in L^1(\mathbb{R})$ en escalier telle que $\|f - g\|_1 < \epsilon$. Nous avons donc

$$\|\hat{f} - \hat{g}\|_\infty \leq \|f - g\|_1 < \epsilon. \quad (30.21)$$

Donc

$$\|\hat{f}(\xi)\| \leq \|\hat{f}(\xi) - \hat{g}(\xi)\| + |\hat{g}(\xi)|. \quad (30.22)$$

Le premier terme est plus petit que ϵ . Il nous reste à voir que

$$\lim_{\xi \rightarrow \infty} |\hat{g}(\xi)| = 0, \quad (30.23)$$

mais cela est le résultat de la première partie de la preuve. \square

Corollaire 30.10.

La transformée de Fourier d'une fonction $L^1(\mathbb{R})$ est bornée.

Démonstration. Par le corollaire 30.7, la transformée de Fourier d'une fonction L^1 est continue. Le lemme de Riemann-Lebesgue 30.9 impliquant qu'elle tend vers zéro en $\pm\infty$, elle doit être bornée. \square

30.1.1 Formule sommatoire de Poisson

Proposition 30.11 (Formule sommatoire de Poisson).

Soit $f: \mathbb{R} \rightarrow \mathbb{C}$ une fonction continue et $L^1(\mathbb{R})$. Nous supposons que

(1) *il existe $M > 0$ et $\alpha > 1$ tels que*

$$|f(x)| \leq \frac{M}{(1 + |x|)^\alpha}, \quad (30.24)$$

(2) $\sum_{n=-\infty}^{\infty} |\hat{f}(2\pi n)| < \infty$.

Alors nous avons

$$\sum_{n=-\infty}^{\infty} f(n) = \sum_{n=-\infty}^{\infty} \hat{f}(2\pi n). \quad (30.25)$$

Démonstration. **Convergence normale** Nous commençons par montrer qu'il y a convergence normale sur tout compact séparément des séries sur les $n \geq 0$ et sur les $n < 0$.

Soit K un compact de \mathbb{R} contenu dans $[-A, A]$ et $n \in \mathbb{Z}$ tel que $|n| \geq 2A$. Pour $x \in K$ nous avons

$$|x + n| \geq |n| - |x| \geq |n| - A \geq \frac{|n|}{2}. \quad (30.26)$$

Du coup nous avons un $\alpha > 1$ tel que

$$|f(x+n)| \leq \frac{M}{(1+|x+n|)^\alpha} \leq \frac{M}{\left(1+\frac{|n|}{2}\right)^\alpha}. \quad (30.27)$$

Lorsque n est grand, cela a le comportement de $M/|n|^\alpha$ et donc la série

$$\sum_{n=0}^{\infty} f(x+n) \quad (30.28)$$

est une série convergent normalement. Les deux séries (usuelles)

$$a_- = \sum_{n \leq 0} f(x+n) \quad (30.29a)$$

$$a_+ = \sum_{n > 0} f(x+n) \quad (30.29b)$$

convergent normalement.

Convergence commutative Au sens de la définition 12.100 nous avons

$$\sum_{n \in \mathbb{Z}} f(x+n) = a_+ + a_-. \quad (30.30)$$

En effet si nous prenons $J'_0 \subset \mathbb{N}$ fini tel que $|\sum_{\mathbb{N} \setminus J'_0} f(x+n) - a_+| \leq \epsilon$ et $J'_1 \in -\mathbb{N}$ tel que $|\sum_{\mathbb{N} \setminus J'_1} f(x+n) - a_-| < \epsilon$, et si nous posons $J_0 = J'_0 \cup J'_1$ alors si K est un ensemble fini de \mathbb{Z} contenant J_0 nous avons

$$\left| \sum_{n \in K} f(n+x) - (a_+ + a_-) \right| \leq \left| \sum_{n \in K^+} f(n+x) - a_+ \right| + \left| \sum_{n \in K^-} f(n+x) - a_- \right| \leq 2\epsilon \quad (30.31)$$

où K^+ sont les éléments positifs de K et K^- sont les *strictement* négatifs. Maintenant que la famille $\{f(n+x)\}_{n \in \mathbb{Z}}$ est une famille sommable, nous savons qu'elle est commutativement sommable et que la proposition 12.105 nous permet de sommer dans l'ordre que l'on veut. Nous pouvons donc écrire sans ambiguïté l'expression $\sum_{n \in \mathbb{Z}} f(x+n)$ ou $\sum_{n=-\infty}^{\infty} f(x+n)$.

re-convergence normale Nous posons donc sans complexes la série

$$F(x) = \sum_{n \in \mathbb{Z}} f(x+n) \quad (30.32)$$

qui converge tant commutativement que normalement. Notons que nous pouvons maintenant dire que la série sur \mathbb{Z} converge normalement ; pas seulement les deux séries séparément.

Continuité, périodicité Étant donné que chacune des fonctions $f(x+n)$ est continue, la convergence normale nous assure que F est continue.

De plus F est périodique de période 1 parce que

$$F(x+1) = \sum_{n=-\infty}^{\infty} f(x+1+n) = \sum_{p=-\infty}^{\infty} f(x+p) = F(x) \quad (30.33)$$

où nous avons posé $p = 1+n$.

Notons que nous n'avons pas spécialement prouvé que F n'était pas périodique avec des périodes plus petites que 1. Mais cela n'a pas d'importance ici.

Coefficients de Fourier En vertu de la définition (28.192) et de la périodicité de F ,

$$c_n(F) = \int_{-1/2}^{1/2} F(t)e^{-2\pi int} dt \quad (30.34a)$$

$$= \int_0^1 F(t)e^{-2\pi int} dt \quad (30.34b)$$

$$= \int_0^1 \sum_{n \in \mathbb{Z}} f(t+n)e^{-2i\pi nt} dt \quad (30.34c)$$

$$= \sum_{n \in \mathbb{Z}} \int_n^{n+1} f(u)e^{-2\pi i(u-n)t} du \quad (30.34d)$$

$$= \int_{-\infty}^{\infty} f(u)e^{-2\pi i nu} du \quad (30.34e)$$

$$= \hat{f}(2\pi n). \quad (30.34f)$$

où nous avons effectué le changement de variables $u = t + n$, et permuté l'intégrale et la somme en vertu du fait que la somme converge normalement.

Conclusion Étant donné l'hypothèse $\sum_{n \in \mathbb{Z}} |\hat{f}(n)| < \infty$ la proposition 29.14 nous dit que

$$F(x) = \sum_{n \in \mathbb{Z}} c_n(F)e^{2\pi inx}, \quad (30.35)$$

c'est-à-dire que

$$\sum_{n=-\infty}^{\infty} f(x+n) = \sum_{n=-\infty}^{\infty} \hat{f}(2\pi n)e^{2\pi inx}. \quad (30.36)$$

En écrivant cette égalité en $x = 0$ nous trouvons le résultat :

$$\sum_{n \in \mathbb{Z}} f(n) = \sum_{n \in \mathbb{Z}} \hat{f}(2\pi n). \quad (30.37)$$

□

Exemple 30.12

La formule sommatoire de Poisson peut être utilisée pour calculer des sommes dans l'espace de Fourier plutôt que dans l'espace direct. Nous allons montrer dans cet exemple l'égalité

$$\sum_{n=-\infty}^{\infty} e^{-\alpha n^2} = \sum_{n=-\infty}^{\infty} \sqrt{\frac{\pi}{\alpha}} e^{-\pi^2 n^2 / \alpha}. \quad (30.38)$$

Si α est grand, alors la somme de gauche est plus rapide, tandis que si α est petit, c'est le contraire.

Nous appliquons la formule sommatoire de Poisson à la fonction

$$f(x) = e^{-\alpha x^2}. \quad (30.39)$$

Nous avons

$$\hat{f}(x) = \int_{\mathbb{R}} e^{-\alpha t^2 - ixt} dt \quad (30.40a)$$

$$= e^{-x^2/4\alpha} \int_{\mathbb{R}} e^{-(\sqrt{\alpha}t + \frac{ix}{2\sqrt{\alpha}})^2} dt \quad (30.40b)$$

$$= e^{-x^2/4\alpha} \frac{1}{\sqrt{\alpha}} \int_{\mathbb{R} + \frac{ix}{2\sqrt{\alpha}}} e^{-u^2} du. \quad (30.40c)$$

Pour traiter cette intégrale nous utilisons la proposition 27.6 en considérant le chemin rectangulaire fermé qui joint les points $-R$, R , $R + ai$, $-R + ai$ et $f(z) = e^{-z^2}$. Calculons l'intégrale sur les deux côtés verticaux. Nous posons

$$\gamma_R(t) = R + tia \quad (30.41)$$

avec $t: 0 \rightarrow 1$. Nous avons

$$\int_{\gamma_R} f = \int_0^1 f(\gamma_R(t)) \|\gamma_R'(t)\| dt \quad (30.42a)$$

$$= ae^{-R^2} \int_0^1 e^{-2tRia+at^2} dt, \quad (30.42b)$$

donc en module nous avons

$$\left| \int_{\gamma_R} f \right| \leq ae^{-R^2} \int_0^1 e^{at^2} dt \leq Me^{-R^2}, \quad (30.43)$$

où M est une constante ne dépendant pas de R . Lorsque $R \rightarrow \infty$, la contribution des chemins verticaux s'annule et nous trouvons que

$$\int_{\mathbb{R}+ai} e^{-u^2} du = \int_{\mathbb{R}} e^{-u^2} du, \quad (30.44)$$

que nous pouvons utiliser pour continuer le calcul (30.40). Nous avons

$$\hat{f}(x) = \frac{e^{-x^2/4\alpha}}{\sqrt{\alpha}} \int_{\mathbb{R}} e^{-u^2} du = \sqrt{\frac{\pi}{\alpha}} e^{-x^2/4\alpha} \quad (30.45)$$

où nous avons utilisé la formule (15.831). Par conséquent ce qui rentre dans la formule sommatoire de Poisson est

$$\hat{f}(2\pi n) = \sqrt{\frac{\pi}{\alpha}} e^{-\pi^2 n^2/\alpha}. \quad (30.46)$$

△

30.2 Transformée de Fourier dans l'espace de Schwartz

La définition de la transformée de Fourier de $\varphi \in \mathcal{S}(\mathbb{R}^d)$ est

$$\hat{\varphi}(\xi) = \int_{\mathbb{R}^n} \varphi(x) e^{-ix \cdot \xi} dx. \quad (30.47)$$

Si α est un multiindice de taille m , nous notons

$$(M_\alpha f)(x) = x_{\alpha_1} \dots x_{\alpha_m} f(x). \quad (30.48)$$

Lemme 30.13 (Lemme de transfert).

Si $\varphi \in \mathcal{S}(\mathbb{R}^d)$ et si α est un multiindice, alors

$$\partial^\alpha \hat{\varphi} = (-i)^{|\alpha|} \widehat{M_\alpha \varphi}. \quad (30.49)$$

et

$$\widehat{\partial^\alpha \varphi}(\xi) = (-i)^{|\alpha|} \xi^\alpha \hat{\varphi}(\xi). \quad (30.50)$$

Démonstration. Nous considérons la fonction $h(x, \xi) = \varphi(x) e^{-ix \cdot \xi}$ dont la dérivée par rapport à ξ_i est donnée par $-i(M_i \varphi)(x) e^{-ix \cdot \xi}$. Cette fonction est majorée en norme par

$$G(x) = M_i \varphi(x), \quad (30.51)$$

qui est encore une fonction à décroissance rapide et donc parfaitement intégrable sur \mathbb{R}^d . Le théorème 18.18 nous dit donc que la dérivée de $\widehat{\varphi}$ par rapport à ξ_i existe et vaut

$$\frac{\partial \widehat{\varphi}}{\partial \xi_i}(\xi) = -i \int_{\mathbb{R}^n} x_i \varphi(x) e^{-i\xi \cdot x} = -i \widehat{M_i \varphi}(\xi). \quad (30.52)$$

En appliquant ce résultat en chaîne, nous trouvons la première formule annoncée.

Nous passons à la seconde formule annoncée. Étant donné que $\varphi \in \mathcal{S}$, ses dérivées le sont aussi et par conséquent, il n'y a pas de problèmes pour écrire

$$\widehat{\partial_{x_k} \varphi}(\xi) = \int_{\mathbb{R}^d} \frac{\partial \varphi}{\partial x_k}(x) e^{-ix \cdot \xi} dx. \quad (30.53)$$

Étant donné que

$$\frac{\partial}{\partial x_k} \left(\varphi(x) e^{-ix \cdot \xi} \right) = \frac{\partial \varphi}{\partial x_k}(x) e^{-ix \cdot \xi} - i \xi_k \varphi(x) e^{-ix \cdot \xi}, \quad (30.54)$$

notre tâche sera de prouver que

$$\int_{\mathbb{R}^d} \frac{\partial}{\partial x_k} \left(\varphi(x) e^{-ix \cdot \xi} \right) dx = 0. \quad (30.55)$$

Autrement dit, nous voulons montrer que le terme au bord d'une intégration par partie s'annule. D'abord le fait que φ soit à décroissance rapide nous assure que l'intégrale (30.55) converge. Pour chaque ξ , la fonction

$$f(x, \xi) = \frac{\partial}{\partial x_k} \left(\varphi(x) e^{-ix \cdot \xi} \right) \quad (30.56)$$

est intégrable par rapport à x . De plus, f est dans $\mathcal{S}(\mathbb{R})$ pour chacune de ses variables (les autres étant fixées). Le théorème de Fubini 15.259 nous permet alors de décomposer l'intégrale en

$$\int_{\mathbb{R}^d} f(x, \xi) dx = \int_{\mathbb{R}} \dots \int_{\mathbb{R}} f(x_1, \dots, x_d) dx_1 \dots dx_d. \quad (30.57)$$

De plus nous pouvons intégrer dans l'ordre de notre choix et nous choisissons évidemment d'intégrer d'abord par rapport à x_k . Étudions donc l'intégrale

$$\int_{\mathbb{R}} \frac{\partial}{\partial x} \left(\varphi(x) e^{-ix\xi} \right) dx = \lim_{A \rightarrow \infty} \int_{-A}^A \frac{\partial}{\partial x} \left(\varphi(x) e^{-ix\xi} \right) dx \quad (30.58)$$

dans laquelle nous avons un peu allégé les notations. Une primitive de ce qui est intégré est toute trouvée : c'est $\varphi(x) e^{-ix\xi}$, et nous pouvons utiliser le théorème fondamental du calcul intégral pour écrire que

$$\int_{-A}^A \left(\varphi(x) e^{-ix\xi} \right)' dx = \left[\varphi(x) e^{-ix\xi} \right]_{x=-A}^{x=A}. \quad (30.59)$$

Vu que φ est dans \mathcal{S} , la limite $A \rightarrow \infty$ donne zéro.

En substituant maintenant (30.54) dans (30.53) et en tenant compte du terme que nous venons de montrer s'annuler, nous avons

$$\widehat{\partial_k \varphi}(\xi) = -i \xi_k \int_{\mathbb{R}^d} \varphi(x) e^{-ix \cdot \xi} = -i \xi_k \widehat{\varphi}(\xi). \quad (30.60)$$

En recommençant la procédure $|\alpha|$ fois nous trouvons la seconde formule annoncée. \square

Proposition 30.14 ([263]).

L'espace de Schwartz est stable par transformée de Fourier. De plus l'application

$$\mathcal{F}: \mathcal{S}(\mathbb{R}^d) \rightarrow \mathcal{S}(\mathbb{R}^d) \quad (30.61)$$

est une bijection linéaire et continue.

Démonstration. La linéarité découle de celle de l'intégrale. La difficulté est de prouver que pour $\varphi \in \mathcal{S}(\mathbb{R}^d)$ nous avons bien que $\hat{\varphi} \in \mathcal{S}(\mathbb{R}^d)$ et que cette association est continue³.

Stabilité Nous devons prouver que pour tout multiindices α et β , nous avons $p_{\alpha,\beta}(\hat{\varphi}) < \infty$. Nous avons

$$\xi^\beta \partial^\alpha \hat{\varphi}(\xi) = \xi^\beta (-i)^{|\alpha|} \widehat{M_\alpha \varphi}(\xi) = (-i)^{|\alpha|+|\beta|} \widehat{\partial^\beta M_\alpha \varphi}(\xi). \quad (30.62)$$

Ensuite nous nous souvenons que $\|\hat{f}\|_\infty \leq \|f\|_1$ parce que

$$|\hat{f}(\xi)| \leq \int_{\mathbb{R}^d} |f(x) e^{-ix \cdot \xi}| = \int_{\mathbb{R}^d} |f(x)| dx = \|f\|_1. \quad (30.63)$$

Donc

$$p_{\alpha,\beta}(\hat{\varphi}) = \|\widehat{\partial^\beta M_\alpha \varphi}\|_\infty \leq \|\partial^\beta M_\alpha \varphi\|_1. \quad (30.64)$$

Du fait que φ soit dans \mathcal{S} , la dernière expression est finie. Cela prouve déjà que

$$\mathcal{F}(\mathcal{S}(\mathbb{R}^d)) \subset \mathcal{S}(\mathbb{R}^d). \quad (30.65)$$

Continuité Nous supposons avoir une suite $\varphi_n \xrightarrow{\mathcal{S}} \varphi$, et nous devons prouver que $\hat{\varphi}_n \xrightarrow{\mathcal{S}} \hat{\varphi}$. Pour alléger les notations, nous posons $f_n = \varphi_n - \varphi$. Nous avons

$$\|\hat{f}_n\|_{\alpha,\beta} = \|\xi^\beta \partial^\alpha \hat{f}_n\|_\infty \quad (30.66a)$$

$$= \|\widehat{\partial^\beta M_\alpha f_n}\|_\infty \text{ lemme 30.13.} \quad (30.66b)$$

$$\leq \|\partial^\beta M_\alpha f_n\|_1 \quad (30.66c)$$

La convergence $f_n \xrightarrow{\mathcal{S}} 0$ nous dit entre autres que $\partial^\beta M_\alpha f_n \xrightarrow{\mathcal{S}} 0$; en particulier la proposition 28.151 nous dit que $\partial^\beta M_\alpha f_n \xrightarrow{L^1} 0$, ce qui signifie, par les majorations (30.66) que

$$\|\hat{f}_n\|_{\alpha,\beta} \leq \|\partial^\beta M_\alpha f_n\|_1 \rightarrow 0, \quad (30.67)$$

ce qui prouve la continuité de transformée de Fourier dans $\mathcal{S}(\mathbb{R}^d)$.

Bijection Une preuve peut être trouvée dans [435].

□

Proposition 30.15 ([1]).

Soit $\varphi \in \mathcal{S}(\mathbb{R}^n \times \mathbb{R}^m)$ et la transformée de Fourier partielle

$$\tilde{\varphi}(x, k) = \int_{\mathbb{R}^m} e^{-iky} \varphi(x, y) dy. \quad (30.68)$$

Alors $\tilde{\varphi} \in \mathcal{S}(\mathbb{R}^n \times \mathbb{R}^m)$.

Démonstration. Il s'agit de reprendre les étapes de la partie correspondante de la preuve de la proposition 30.14. Soient des multiindices α, α', β et β' où α et β se réfèrent à la variable x tandis que α' et β' se réfèrent à la variable k .

Vu que la multiplication par $k^{\beta'}$ commute avec ∂^α nous avons

$$x^\beta k^{\beta'} \partial^\alpha \partial^{\alpha'} \tilde{\varphi}(x, k) = x^\beta k^{\beta'} \partial^\alpha (-i)^{|\alpha'|} \widehat{M_{\alpha'} \varphi}(x, k) = (-i)^{|\alpha'|+|\beta'|} x^\beta \partial^\alpha \widehat{\partial^{\beta'} M_{\alpha'} \varphi}(x, k). \quad (30.69)$$

D'autre part nous avons $\partial^\alpha \tilde{\varphi} = \widetilde{\partial^\alpha \varphi}$ parce que la fonction $\partial_x \varphi$ étant Schwartz, la fonction

$$G(y) = \sup_{x \in \mathbb{R}^n} |(\partial_x \varphi)(x, y)| \quad (30.70)$$

3. Pour rappel, en dimension infinie, il n'est pas garanti qu'une application linéaire soit continue.

est dans $L^1(\mathbb{R}^m)$ par le corollaire 28.147. Par conséquent le théorème 18.18 permet de permuter la dérivée et l'intégrale dans

$$\frac{\partial}{\partial x} \tilde{\varphi}(x, k) = \frac{\partial}{\partial x} \int_{\mathbb{R}^m} e^{-iky} \varphi(x, y) dy. \quad (30.71)$$

Dans le même ordre d'esprit mais dans difficultés de permutation de limites nous avons $M_\beta \tilde{\varphi} = \widetilde{M_\beta \varphi}$.

D'autre part nous avons encore $\|\tilde{\varphi}\|_\alpha < \infty$ parce que

$$|\tilde{\varphi}(x, k)| \leq \int_{\mathbb{R}^m} |\varphi(x, y)| dy \leq \sup_x \int_{\mathbb{R}^m} |\varphi(x, y)| dy \leq \int_{\mathbb{R}^m} |\sup_x \varphi(x, y)| dy < \infty \quad (30.72)$$

parce que φ est Schwartz et le corollaire 28.147 donne l'intégrabilité.

Donc nous avons

$$p_{(\alpha\alpha'), (\beta\beta')}(\tilde{\varphi}) = \|\partial^{\beta'} \widetilde{M_{\alpha'} M_\beta} \partial^\alpha \varphi\|_\infty < \infty. \quad (30.73)$$

Cela prouve que $\tilde{\varphi}$ est Schwartz. \square

30.2.1 Quelques transformées de Fourier

Exemple 30.16 ([43])

Soit la fonction $g_\epsilon(x) = e^{-\epsilon x^2}$. Sa transformée de Fourier sera déduite dans le lemme 30.17 en utilisant le lemme de transfert 30.13. Nous nous proposons ici de déduire de façon directe l'équation différentielle vérifiée par la transformée de Fourier de g_ϵ .

Nous posons

$$I(k) = \int_{\mathbb{R}} e^{-ikx} e^{-\epsilon x^2} dx. \quad (30.74)$$

et nous considérons la fonction

$$f(k, x) = e^{-ikx} e^{-\epsilon x^2}. \quad (30.75)$$

Elle est de classe C^1 par rapport à k , et intégrable en x pour chaque k . De plus sa dérivée

$$(\partial_k f)(k, x) = -ix e^{-ikx} e^{-\epsilon x^2} \quad (30.76)$$

vérifie $|\partial_k f| \leq x e^{-\epsilon x^2}$. La dérivée est donc majorée (uniformément en k) par une fonction intégrable. Le théorème 18.18 permet de permuter la dérivée et l'intégrale :

$$I'(k) = \int_{\mathbb{R}} -ix e^{-ikx} e^{-\epsilon x^2} dx \quad (30.77a)$$

$$= i \int_{\mathbb{R}} e^{-ikx} \frac{1}{2\epsilon} \frac{d}{dx} (e^{-\epsilon x^2}) dx \quad (30.77b)$$

$$= \frac{-i}{2\epsilon} \int_{\mathbb{R}} \frac{d}{dx} (e^{-ikx}) e^{-\epsilon x^2} dx \quad \text{par partie} \quad (30.77c)$$

$$= \frac{-k}{2\epsilon} \int_{\mathbb{R}} e^{-ikx} e^{-\epsilon x^2} dx \quad (30.77d)$$

$$= \frac{-k}{2\epsilon} I(k). \quad (30.77e)$$

D'où l'équation différentielle $I'(k) = -\frac{k}{2\epsilon} I(k)$. \triangle

Lemme 30.17 (Transformée de Fourier de la Gaussienne [436]).

La transformée de Fourier de

$$g_\epsilon : \mathbb{R}^d \rightarrow \mathbb{R} \\ x \mapsto e^{-\epsilon \|x\|^2} \quad (30.78)$$

est donnée par

$$\hat{g}_\epsilon(\xi) = \left(\frac{\pi}{\epsilon}\right)^{d/2} e^{-\|\xi\|^2/4\epsilon} \quad (30.79)$$

Démonstration. Nous commençons par la fonction $g(x) = e^{-\|x\|^2/2}$ et nous prouvons que sa transformée de Fourier est $\hat{g}(\xi) = (2\pi)^{d/2}g(\xi)$.

Réduction à la dimension 1 La fonction g est dans l'espace de Schwartz. Par le théorème de Fubini,

$$\hat{g}(\xi) = \int_{\mathbb{R}^d} \prod_{k=1}^d e^{-x_k^2} e^{-i\xi_k x_k} dx = \prod_{k=1}^d \int_{\mathbb{R}} e^{-t^2/2} e^{-\xi_k t} dt = \prod_{k=1}^d \hat{f}(\xi_k) \tag{30.80}$$

où f est la fonction d'une variable

$$f(x) = e^{-x^2/2}. \tag{30.81}$$

Notons que $f \in \mathcal{D}(\mathbb{R})$.

Une équation différentielle Voyons l'équation différentielle satisfaite par la transformée de Fourier \hat{f} de la fonction (30.81). Grâce au lemme 30.13 nous trouvons l'équation différentielle⁴

$$\xi \hat{f}(\xi) + (\hat{f})'(\xi) = 0. \tag{30.82}$$

C'est le moment d'utiliser le théorème de Cauchy-Lipschitz (18.40), appliqué à la fonction $f(t, y) = -ty$ qui est Lipschitz et continue au problème

$$\begin{cases} y' + ty = 0 \\ y(0) = y_0 \end{cases} \tag{30.83a}$$

$$\tag{30.83b}$$

possède une unique solution maximale, en l'occurrence $y(x) = y_0 e^{-x^2/2}$. En ce qui concerne la condition initiale nous avons

$$\hat{f}(0) = \int_{\mathbb{R}} e^{-x^2/2} dx = \sqrt{2\pi}. \tag{30.84}$$

par l'exemple 15.262. Donc

$$\hat{f}(\xi) = \sqrt{2\pi} e^{-\xi^2/2}. \tag{30.85}$$

En reformant le produit (30.80) nous concluons.

Nous passons maintenant à la fonction g_ϵ . Nous pouvons écrire g_ϵ sous la forme

$$g_\epsilon(x) = g(\sqrt{2\epsilon}x). \tag{30.86}$$

Utilisant successivement la transformée de Fourier de g que nous venons de calculer et 30.5 (facteur d'échelle) nous trouvons

$$\hat{g}_\epsilon(\xi) = (2\pi)^{d/2} g_\epsilon(\xi) \tag{30.87a}$$

$$\hat{g}_\epsilon(\xi) = (2\epsilon)^{-d/2} \hat{g}(\xi/\sqrt{2\epsilon}) \tag{30.87b}$$

$$= \left(\frac{\pi}{\epsilon}\right)^{d/2} e^{-|\xi|^2/4\epsilon}. \tag{30.87c}$$

Nous voyons que $\hat{g}_\epsilon \in \mathcal{S}(\mathbb{R}^d)$ (c'était gagné d'avance par la proposition 30.14). □

30.3 Suite régularisante

Définition 30.18.

Une **suite régularisante** est une suite (ρ_n) dans $L^1(\mathbb{R}^d)$ telle que

(1) pour tout n , $\rho_n \geq 0$ et $\int_{\mathbb{R}^d} \rho_n = 1$;

(2) pour tout $\alpha > 0$,

$$\lim_{n \rightarrow \infty} \int_{|t| > \alpha} \rho_n = 0. \tag{30.88}$$

4. Une façon directe de déduire cette équation différentielle est donnée dans l'exemple 30.16.

Une telle suite est régularisante parce que souvent $\rho_n \in \mathcal{D}(\mathbb{R}^d)$, ce qui donne $f * \rho_n \in C^\infty$ par le corollaire 28.58.

Proposition 30.19 ([437, 438]).

Soit une suite régularisante $\rho_n \in L^1(\mathbb{R}^d)$. Alors :

(1) Si f est continue à support compact, nous avons la convergence uniforme sur \mathbb{R}^d :

$$f * \rho_n \xrightarrow{\text{unif}} f. \quad (30.89)$$

(2) Si $g \in L^p$ ($1 \leq p < \infty$) alors

$$g * \rho_n \xrightarrow{L^p} g. \quad (30.90)$$

Démonstration. Si f est continue à support compact, elle est uniformément continue⁵, et elle est bornée. Soit $\epsilon > 0$ et $\alpha > 0$ tel que pour tout x, y tels que $\|x - y\| < \alpha$ nous ayons $|f(x) - f(y)| < \epsilon$. Nous prenons de plus n suffisamment grand pour avoir $\int_{B(0, \alpha)^c} \rho_n < \epsilon$. Nous avons alors

$$|f(x) - (f * \rho_n)(x)| = \left| \int_{\mathbb{R}^d} (f(x) - f(y)) \rho_n(x - y) dy \right| \quad (30.91a)$$

$$\leq \int_{B(x, \alpha)} \underbrace{|f(x) - f(y)|}_{\leq \epsilon} \rho_n(x - y) dy + \int_{B(x, \alpha)^c} \underbrace{|f(x) - f(y)|}_{\leq 2\|f\|_\infty} \rho_n(x - y) dy$$

$$(30.91b)$$

$$\leq \epsilon(1 + 2\|f\|_\infty). \quad (30.91c)$$

Nous avons prouvé que pour tout $\epsilon > 0$, il existe N tel que $n > N$ implique $|f(x) - (f * \rho_n)(x)| \leq \epsilon$. Cela prouve l'uniforme convergence sur \mathbb{R}^d de $f * \rho_n$ vers f .

Pour le point (2) nous considérons $g \in L^1(\mathbb{R}^d)$ et $\phi \in \mathcal{D}(\mathbb{R}^d)$. Nous avons la majoration

$$\|g * \rho_n - g\|_p \leq \|g * \rho_n - \phi * \rho_n\|_p + \|\phi * \rho_n - \phi\|_p + \|\phi - g\|_p \quad (30.92)$$

En ce qui concerne le premier terme ;

$$\|(g - \phi) * \rho_n\|_p \leq \|g - \phi\|_p \quad (30.93)$$

par la proposition 28.56. Donc

$$\|g * \rho_n - g\|_p \leq 2\|g - \phi\|_p + \|\phi * \rho_n - \phi\|_p. \quad (30.94)$$

Par la densité de \mathcal{D} dans L^p (théorème 28.47(5)) nous pouvons considérer une suite $\phi_i \xrightarrow{L^p} g$ dans $\mathcal{D}(\mathbb{R}^d)$. Pour tout i nous avons

$$\|g * \rho_n - g\|_p \leq 2\|g - \phi_i\|_p + \|\phi_i * \rho_n - \phi_i\|_p. \quad (30.95)$$

Nous effectuons la limite sur $n \rightarrow \infty$:

$$\lim_{n \rightarrow \infty} \|g * \rho_n - g\|_p \leq 2\|g - \phi_i\|_p + \underbrace{\lim_{n \rightarrow \infty} \|\phi_i * \rho_n - \phi_i\|_p}_{=0} \quad (30.96)$$

parce que le point (1) s'applique à ϕ_i . Nous effectuons ensuite la limite sur $i \rightarrow \infty$ dans

$$\lim_{n \rightarrow \infty} \|g * \rho_n - g\|_p \leq 2\|g - \phi_i\|_p \rightarrow 0. \quad (30.97)$$

□

5. Théorème de Heine 13.87.

Lemme 30.20.

Si $g_\epsilon(x) = e^{-\epsilon\|x\|^2}$ alors la suite

$$\rho_n = \frac{1}{(2\pi)^d} \hat{g}_{1/n} \quad (30.98)$$

est une suite régularisante (définition 30.18).

Démonstration. Nous savons déjà la transformée de Fourier de g_ϵ par le lemme 30.17. Nous montrons que la suite ρ_n est régularisante. Nous avons $\hat{g}_\epsilon \in L^1(\mathbb{R}^d)$ et $\hat{g}_\epsilon \geq 0$ ainsi que $\lim_{\epsilon \rightarrow 0} \int_{B(0,\alpha)} \hat{g}_\epsilon = 0$ pour tout α . Il y a seulement un couac avec la norme. Nous calculons $\int_{\mathbb{R}^d} \hat{g}_\epsilon(\xi) d\xi$ avec la forme (30.87c). En utilisant sauvagement Fubini⁶ pour séparer les intégrales et en effectuant le changement de variable $u = t/(2\sqrt{\epsilon})$ nous calculons :

$$\int_{\mathbb{R}^d} e^{-|\xi|^2/4\epsilon} d\xi = \prod_{k=1}^d \int_{\mathbb{R}} e^{-t^2/4\epsilon} dt \quad (30.99a)$$

$$= 2\sqrt{\epsilon} \prod_{k=1}^d \int_{\mathbb{R}} e^{-u^2} du \quad (30.99b)$$

$$= \prod_{k=1}^d 2\sqrt{\epsilon} \sqrt{\pi} \quad (30.99c)$$

$$= 2^d (\pi\epsilon)^{d/2}. \quad (30.99d)$$

Nous avons utilisé l'exemple 15.262 pour le calcul de l'intégrale gaussienne. Avec tout cela nous avons

$$\int_{\mathbb{R}^d} \hat{g}_\epsilon = (2\pi)^d. \quad (30.100)$$

Donc $\frac{1}{(2\pi)^d} \hat{g}_{1/n}$ est une suite régularisante. □

Le corollaire suivant regroupe les résultats à propos des suites régularisantes, leur utilité et leur existence.

Corollaire 30.21.

Si la suite régularisante ρ_n est dans $L^1(\mathbb{R}^d) \cap C^\infty(\mathbb{R}^d)$ alors pour $f \in L^p(\mathbb{R}^d)$ en posant $f_n = \rho_n * f$ nous avons

$$(1) f_n \in C^\infty(\mathbb{R}^d) \cap L^p(\mathbb{R}^d)$$

$$(2) f_n \xrightarrow{L^p} f$$

De plus, de telles suites existent.

Démonstration. Le fait que f_n soit de classe C^∞ est le corollaire 28.58, et la convergence est la proposition 30.19(2).

De telles suites existent, par exemple celle donnée par le lemme 30.20. □

30.3.1 Formule d'inversion

Proposition 30.22 (Formule d'inversion de Fourier[43]).

Si $f \in \mathcal{S}(\mathbb{R})$, alors nous avons la formule d'inversion

$$f(x) = \frac{1}{2\pi} \int_{\mathbb{R}} e^{ikx} \hat{f}(k) dk. \quad (30.101)$$

6. Le pauvre!

Cette formule peut d'écrire de plusieurs autres façons :

$$\mathcal{F}(\mathcal{F}(f))(x) = 2\pi f(-x), \quad (30.102a)$$

$$\mathcal{F}^{-1}(f)(x) = \frac{1}{2\pi} \hat{f}(-x), \quad (30.102b)$$

$$f(x) = \frac{1}{2\pi} \mathcal{F}(\hat{f})(-x). \quad (30.102c)$$

Démonstration. Pour $\epsilon > 0$ nous posons

$$f_\epsilon(k) = e^{-\epsilon k^2} e^{ikx} \hat{f}(k). \quad (30.103)$$

Nous allons calculer

$$\lim_{\epsilon \rightarrow 0} \int_{\mathbb{R}} e^{-\epsilon k^2} e^{ikx} \hat{f}(k) dk \quad (30.104)$$

de deux façons.

D'abord en utilisant directement le théorème de la convergence dominée 15.184. La fonction \hat{f} est dans $\mathcal{S}(\mathbb{R})$ (théorème 30.14) et par conséquent $f_\epsilon \in L^1(\mathbb{R})$ parce que le facteur $e^{-\epsilon k^2}$ ne va certainement pas empêcher de converger. De plus $|f_\epsilon| \leq |\hat{f}|$ et $\hat{f} \in L^1$. Le théorème est de la convergence dominée est applicable et

$$\lim_{\epsilon \rightarrow 0} \int_{\mathbb{R}} e^{-\epsilon k^2} e^{ikx} \hat{f}(k) dk = \int_{\mathbb{R}} e^{ikx} \hat{f}(k) dk. \quad (30.105)$$

Pour le deuxième calcul nous allons faire appel à Fubini⁷ pour la fonction

$$\begin{aligned} u: \mathbb{R} \times \mathbb{R} &\rightarrow \mathbb{R} \\ (k, y) &\mapsto e^{ik(x-y)} e^{-\epsilon k^2} f(y). \end{aligned} \quad (30.106)$$

D'abord nous nous assurons que $u \in L^1(\mathbb{R} \times \mathbb{R})$ par le corollaire 15.258, et ensuite nous utilisons le théorème de Fubini 15.259 pour manipuler les intégrales (et en particulier les inverser). Dans un premier temps nous avons :

$$\int_{\mathbb{R}} \int_{\mathbb{R}} |e^{ik(x-y)} e^{-\epsilon k^2} f(y)| dy dk \leq \int_{\mathbb{R}} e^{-\epsilon k^2} \left[\int_{\mathbb{R}} |f(y)| dy \right] dk < \infty \quad (30.107)$$

parce que f étant dans $\mathcal{S}(\mathbb{R})$, l'intégrale intérieure se réduit à un nombre. Nous savons maintenant que $u \in L^1(\mathbb{R} \times \mathbb{R})$. Nous pouvons alors calculer un peu ...

$$\int_{\mathbb{R}} e^{ikx} e^{-\epsilon k^2} \hat{f}(k) dk = \int_{\mathbb{R}} \int_{\mathbb{R}} e^{ikx} e^{-\epsilon k^2} e^{-iky} f(y) dy dk \quad (30.108a)$$

$$= \int_{\mathbb{R}} \left[\int_{\mathbb{R}} e^{ik(x-y)} e^{-\epsilon k^2} f(y) dk \right] dy \quad (30.108b)$$

$$= \int_{\mathbb{R}} f(y) \left[\int_{\mathbb{R}} e^{ik(x-y)} e^{-\epsilon k^2} dk \right] dy \quad (30.108c)$$

$$= \int_{\mathbb{R}} f(y) \hat{g}_\epsilon(y-x) dy \quad (30.108d)$$

$$= \sqrt{\frac{\pi}{\epsilon}} \int_{\mathbb{R}} f(y) e^{-(y-x)^2/4\epsilon} dy \quad (30.108e)$$

$$= 2\sqrt{\epsilon} \sqrt{\frac{\pi}{\epsilon}} \int_{\mathbb{R}} f(x + 2\sqrt{\epsilon}t) e^{-t^2} dt \quad (30.108f)$$

$$= 2\sqrt{\pi} \int_{\mathbb{R}} f(x + 2\sqrt{\epsilon}t) e^{-t^2} dt \quad (30.108g)$$

$$(30.108h)$$

Justifications :

7. Parce qu'il est toujours plus simple de refiler le boulot aux autres que de le faire soi-même. ... pauvre Fubini!

- La fonction g_ϵ est la gaussienne dont la transformée de Fourier est a été l'objet du lemme 30.17.
- Nous avons effectué le changement de variables $t = (y - x)/(2\sqrt{\epsilon})$ qui donne $dt = dy/2\sqrt{\epsilon}$.

La fonction f étant Schwartz (en particulier bornée), dans la dernière intégrale, nous pouvons effectuer la majoration

$$f(x + 2\sqrt{\epsilon}t)e^{-t^2} \leq \|f\|_\infty e^{-t^2}, \quad (30.109)$$

qui est une fonction intégrable. Nous pouvons donc permuter la limite et l'intégrale. Dans l'égalité

$$\lim_{\epsilon \rightarrow 0} \int_{\mathbb{R}} e^{ikx} e^{-\epsilon k^2} \hat{f}(k) dk = \lim_{\epsilon \rightarrow 0} 2\sqrt{\pi} \int_{\mathbb{R}} f(x + 2\sqrt{\epsilon}t) e^{-t^2} dt \quad (30.110)$$

À gauche nous avons déjà la limite depuis (30.105), et à droite nous obtenons

$$\lim_{\epsilon \rightarrow 0} 2\sqrt{\pi} \int_{\mathbb{R}} f(x + 2\sqrt{\epsilon}t) e^{-t^2} dt = \int_{\mathbb{R}} f(x) e^{-t^2} dt = 2\sqrt{\pi} f(x) \sqrt{\pi} = 2\pi f(x) \quad (30.111)$$

où nous avons utilisé l'intégrale gaussienne faite dans l'exemple 15.262.

En remettant tout ensemble,

$$2\pi f(x) = \lim_{\epsilon \rightarrow 0} \int_{\mathbb{R}} e^{-\epsilon k^2} e^{ikx} \hat{f}(k) dk = \int_{\mathbb{R}} e^{ikx} \hat{f}(k) dk, \quad (30.112)$$

ce qu'il fallait prouver. □

Corollaire 30.23.

Nous avons la formule

$$\int_{\mathbb{R}} \int_{\mathbb{R}} e^{-ikx} f(x) dx dk = 2\pi f(0). \quad (30.113)$$

Démonstration. Poser $x = 0$ dans l'équation (30.101). □

30.24.

Les physiciens qui n'ont que rarement peur écrivent souvent la formule (30.113) sous la forme

$$\int_{\mathbb{R}} e^{-ikx} dk = \delta(x) \quad (30.114)$$

où δ serait la fonction de Dirac qui vaut zéro partout sauf en $x = 0$ où elle vaudrait l'infini, mais pas n'importe quel infini ; juste celui qu'il faut pour que l'intégrale vaille 1.

Lemme 30.25.

Si $\phi \in \mathcal{S}(\mathbb{R} \times \mathbb{R}^n)$, alors

$$\partial_t \hat{\phi} = \widehat{\partial_t \phi} \quad (30.115)$$

où le chapeau dénote la transformée de Fourier par rapport à la variable dans \mathbb{R}^n et non par rapport à celle dans \mathbb{R} . Le t par contre est la variable dans \mathbb{R} .

Démonstration. Par définition de la transformée de Fourier nous avons

$$(\partial_t \hat{\phi})(t, \xi) = \frac{\partial}{\partial t} \int_{\mathbb{R}^n} \phi(t, x) e^{-ix\xi} dx. \quad (30.116)$$

Notre but est de permuter l'intégrale et la dérivée en utilisant le théorème 18.18. Il nous faut une fonction $G: \mathbb{R}^n \rightarrow \mathbb{R}$ qui soit intégrable sur \mathbb{R}^n et telle que

$$\left| \frac{\partial \phi}{\partial t} \phi(t, x) \right| \leq G(x) \quad (30.117)$$

pour tout $t \in B(t_0, \delta)$. Étant donné que la fonction $\partial_t \phi$ est tout autant Schwartz que ϕ elle-même nous pouvons alléger les notations et chercher une fonction G qui convient pour toute fonction $\varphi \in \mathcal{S}(\mathbb{R} \times \mathbb{R}^n)$. Soit la fonction

$$G(x) = \sup_{t \in B(t_0, \delta)} |\varphi(t, x)|. \quad (30.118)$$

Pour tout multiindice α nous avons alors

$$\sup_{x \in \mathbb{R}^n} |x^\alpha G(x)| \leq \sup_{(t, x) \in \mathbb{R} \times \mathbb{R}^n} |x^\alpha \varphi(t, x)| \leq M_\alpha \in \mathbb{R}. \quad (30.119)$$

Grâce à la proposition 28.146, cela signifie que φ décroît plus vite que n'importe quel polynôme; G est donc intégrable sur \mathbb{R}^n . \square

30.4 Transformée de Fourier sur $L^2(\mathbb{R}^d)$

La théorie des transformées de Fourier est intéressante sur $L^2(\mathbb{R}^d)$ parce qu'elle y donne une isométrie. Nous allons la donner avec des fonctions à valeurs dans \mathbb{C} .

Remarque 30.26.

Une remarque qui vaut ce qu'elle vaut, mais si u est une classe de fonction pour la relation $u \sim v$ si et seulement si $u(x) = v(x)$ pour presque tout x alors l'intégrale

$$\hat{u}(\xi) = \int_{\mathbb{R}^d} u(x) e^{ix\xi} dx \quad (30.120)$$

ne dépend pas du choix du représentant. Nous pouvons donc parfaitement parler de transformée de Fourier d'une classe de fonctions.

30.4.1 Le problème

Nous avons défini en général la transformée de Fourier d'une fonction $f: \mathbb{R} \rightarrow \mathbb{C}$ par la formule

$$\hat{f}(\xi) = \int_{\mathbb{R}} e^{i\xi x} f(x) dx \quad (30.121)$$

tant que cette intégrale existe.

Il se fait que cette intégrale n'existe pas toujours pour des fonctions dans $L^2(\mathbb{R})$. Donc nous devons faire mieux pour définir la transformée de Fourier sur L^2 .

Exemple 30.27([1])

Prenons la fonction

$$f(x) = \begin{cases} 0 & \text{si } x < 1 \\ \frac{1}{x} & \text{si } x \geq 1. \end{cases} \quad (30.122)$$

Vu que l'intégrale $\int_1^\infty \frac{1}{x^2} dx$ existe et est finie (proposition 15.248(2)), la fonction f est dans $L^2(\mathbb{R})$.

Cependant l'intégrale (30.121) n'existe pas. Pour nous convaincre de cela, nous pouvons simplement nous souvenir de la définition d'une intégrale à valeurs dans un espace vectoriel (définition 15.176). Nous fixons $\xi \in \mathbb{R}$ et nous posons $g(x) = f(x) e^{i\xi x}$.

Bien évidemment, $|f(x)| = \frac{1}{x}$ sur $]1, \infty[$. Donc $\int_{\mathbb{R}} |g| = \infty$, et la fonction g n'est pas intégrable. Fin de l'histoire.

Nous pouvons toujours essayer de comprendre mieux. Vu que $\int_{\mathbb{R}} |g| = \infty$, la proposition 15.177 nous dit qu'au moins une des intégrales parmi

$$\int f_{re}^+, \int f_{im}^+, \int f_{re}^-, \int f_{im}^- \quad (30.123)$$

est égale à $+\infty$.

Note qu'en travaillant un peu, on se convainc qu'en réalité, elles divergent toutes les quatre.

\triangle

30.4.2 Extension de $L^1 \cap L^2$ vers L^2

Théorème 30.28 (Extention de la transformée de Fourier vers $L^2(\mathbb{R}^d)$ [436]).

Soit $f \in L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d)$. Alors

(1) Nous avons $\mathcal{F}(f) \in L^2$ et $\|\hat{f}\|_{L^2} = (2\pi)^d \|f\|_{L^2}$.

(2) L'application $\mathcal{F}: L^1 \cap L^2 \rightarrow L^2$ peut être étendue en une application $\mathcal{F}: L^2(\mathbb{R}^d) \rightarrow L^2(\mathbb{R}^d)$ vérifiant

$$\|\hat{f}\|_{L^2} = (2\pi)^d \|f\|_{L^2} \tag{30.124}$$

pour tout $f \in L^2(\mathbb{R}^d)$.

Démonstration. Le fait que $f \in L^1$ implique $\|\mathcal{F}(f)\|_\infty \leq \|f\|_1$ (c'est le lemme 30.8). En particulier, $|\mathcal{F}(f)(\xi)|^2$ est majoré et l'intégrale

$$\clubsuit = \int_{\mathbb{R}^d} |\hat{f}|^2 e^{-\epsilon|\xi|^2} d\xi \tag{30.125}$$

existe et est finie.

Découper l'intégrale Dans un premier temps nous développons les intégrales. Dans les égalités suivantes, $x\xi$ est le produit scalaire $x \cdot \xi$ dans \mathbb{R}^d .

$$\clubsuit = \int_{\mathbb{R}^d} \left(\int_{\mathbb{R}^d} \overline{f(x)} e^{ix\xi} dx \right) \left(\int_{\mathbb{R}^d} f(y) e^{-y\xi} dy \right) e^{-\epsilon|\xi|^2} d\xi \tag{30.126a}$$

$$= \int_{\mathbb{R}^d} \left[\int_{\mathbb{R}^d \times \mathbb{R}^d} \overline{f(x)} f(y) e^{i\xi(x-y)} dx dy \right] e^{-\epsilon|\xi|^2} d\xi. \tag{30.126b}$$

Nous avons utilisé le théorème de Fubini pour regrouper les intégrales⁸. Vu que $f \in L^1(\mathbb{R}^d)$, la fonction $(x, y, \xi) \mapsto \overline{f(x)} f(y) e^{i\xi(x-y)}$ est dans $L^1(\mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^d)$ et le théorème de Fubini 15.259 avec $\Omega_1 = \mathbb{R}^d \times \mathbb{R}^d$ et $\Omega_2 = \mathbb{R}^d$ nous permet de permuter les intégrales pour avoir

$$\clubsuit = \int_{\mathbb{R}^d \times \mathbb{R}^d} \overline{f(x)} f(y) \left[\int_{\mathbb{R}^d} e^{i\xi(x-y)} e^{-\epsilon|\xi|^2} d\xi \right] dx dy. \tag{30.127}$$

Discuter de cette gaussienne En posant

$$g(x) = e^{-|x|^2/2} \tag{30.128a}$$

$$g_\epsilon(x) = g(\sqrt{2\epsilon}x) = e^{-\epsilon|x|^2} \tag{30.128b}$$

nous avons $g_\epsilon \in \mathcal{S}(\mathbb{R}^d)$ et le lemme 30.20 nous autorise à écrire

$$\hat{g}(\xi) = (2\pi)^{d/2} g(\xi) \tag{30.129a}$$

$$\hat{g}_\epsilon(\xi) = \left(\frac{\pi}{\epsilon}\right)^{d/2} e^{-|\xi|^2/4\epsilon} \tag{30.129b}$$

Nous voyons que $\hat{g}_\epsilon \in \mathcal{S}(\mathbb{R}^d)$ (c'était gagné d'avance par la proposition 30.14) et que \hat{g}_ϵ est une fonction paire (encore une fois, c'était gagné d'avance parce que la transformée de Fourier d'une fonction paire est paire).

Tout cela pour dire que l'intégrale entre crochet dans (30.127) est $\hat{g}_\epsilon(y-x) = \hat{g}_\epsilon(x-y)$, et donc

$$\clubsuit = \int_{\mathbb{R}^d \times \mathbb{R}^d} \overline{f(x)} f(y) \hat{g}_\epsilon(x-y) dx dy. \tag{30.130}$$

Encore une fois le théorème de Fubini permet de séparer les intégrales et de calculer l'intégrale sur y en premier. Vu que $f \in L^1$ et que $\hat{g}_\epsilon \in \mathcal{S}(\mathbb{R}^d)$, le produit de convolution $f * \hat{g}_\epsilon$ est un élément de $\mathcal{S}(\mathbb{R}^d)$ par la proposition 28.155. Nous avons donc

$$\clubsuit = \int_{\mathbb{R}^d} \overline{f(x)} (f * \hat{g}_\epsilon)(x) dx. \tag{30.131}$$

8. Dans la suite nous allons encore utiliser Fubini quelques fois pour regrouper et dégroupier des intégrales.

Là, nous reconnaissons un produit scalaire dans $L^2(\mathbb{R}^d)$, et donc

$$\int_{\mathbb{R}^d} |\hat{f}|^2 e^{-\epsilon \xi^2} d\xi = \langle f, f * \hat{g}_\epsilon \rangle_{L^2(\mathbb{R}^d)}. \quad (30.132)$$

Notons que tout a un sens : $f \in L^2(\mathbb{R}^d)$ et $f * \hat{g}_\epsilon \in \mathcal{S}(\mathbb{R}^d) \subset L^2(\mathbb{R}^d)$.

Suite régularisante Nous prenons la suite régularisante du lemme 30.20 donnée par

$$\rho_n = \frac{1}{(2\pi)^d} \hat{g}_{1/n}. \quad (30.133)$$

Première conclusion Nous reprenons (30.132)

$$\int_{\mathbb{R}^d} |\hat{f}|^2 e^{-|\xi|^2/n} d\xi = \langle f, f * \hat{g}_{1/n} \rangle_{L^2(\mathbb{R}^d)} = (2\pi)^d \langle f, f * \rho_n \rangle. \quad (30.134)$$

En prenant la limite $n \rightarrow \infty$ nous trouvons

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}^d} |\hat{f}|^2 e^{-\epsilon \xi^2} d\xi = (2\pi)^d \|f\|^2. \quad (30.135)$$

Pour effectuer la limite du membre de gauche nous devons remarquer qu'en posant

$$g_n(\xi) = |\hat{f}(\xi)| e^{-|\xi|^2/n}, \quad (30.136)$$

nous avons une suite décroissante de fonction (c'est-à-dire que à ξ fixé, c'est décroissant en n). Par ailleurs ces fonctions sont toujours à valeurs dans $[0, \infty]$ et nous pouvons utiliser le théorème de la convergence monotone 15.160 pour permuter la limite et l'intégrale. Au final :

$$\|\hat{f}\|_{L^2} = (2\pi)^d \|f\|_{L^2}. \quad (30.137)$$

En ce qui concerne l'extension, soit $f \in L^2(\mathbb{R}^d)$ et une suite (f_n) dans $L^1 \cap L^2$ telle que $f_n \xrightarrow{L^2} f$.

Existence d'une telle suite Si $f \in L^2(\mathbb{R}^d)$, alors nous pouvons poser

$$f_n(x) = f(x) e^{-|x|^2/n^2}. \quad (30.138)$$

Par l'inégalité de Hölder (28.59) nous avons $f_n \in L^1(\mathbb{R}^d)$; de plus $f_n \in L^2(\mathbb{R}^d)$ parce que pour tout x nous avons $|f_n(x)| \leq |f(x)|$. Montrons que $f_n \xrightarrow{L^2} f$. Nous avons

$$\|f_n - f\|_{L^2}^2 = \int_{\mathbb{R}^d} |f(x)(1 - e^{-|x|^2/n^2})|^2 dx. \quad (30.139)$$

Nous voulons prendre la limite $n \rightarrow \infty$. Pour ce faire à droite nous remarquons que $e^{-|x|^2/n^2}$ est majoré par 1; ce qui se trouve dans l'intégrale est donc majoré (uniformément en n) par $|f(x)|^2$, qui est une fonction L^1 parce que f est L^2 . Le théorème de la convergence dominée 15.184 nous permet alors de permuter la limite et l'intégrale, ce qui donne

$$\lim_{n \rightarrow \infty} \|f_n - f\|_2^2 = \int_{\mathbb{R}^d} \lim_{n \rightarrow \infty} |f(x)(1 - e^{-|x|^2/n^2})|^2 dx = 0. \quad (30.140)$$

Définition de $\mathcal{F}: L^2 \rightarrow L^2$ La suite (f_n) est une suite convergence dans L^2 , et elle est donc de Cauchy. De plus pour chaque n, m nous avons

$$\|\hat{f}_n - \hat{f}_m\| = (2\pi)^d \|f_n - f_m\|. \quad (30.141)$$

Nous voyons donc que la suite (\hat{f}_n) est également de Cauchy, dans l'espace $L^2(\mathbb{R}^d)$ qui est complet (lemme 28.73). Nous posons

$$\hat{f} = \lim_{n \rightarrow \infty} \hat{f}_n. \quad (30.142)$$

Indépendance aux choix Nous devons montrer que la définition de \hat{f} ne dépend pas de la suite approximant f dans $L^1 \cap L^2$. Soient dans deux suites $f_n \xrightarrow{L^2} f$ et $g_n \xrightarrow{L^2} f$ telles que $\hat{f}_n \xrightarrow{L^2} F$ et $\hat{g}_n \xrightarrow{L^2} G$. Alors

$$\|\hat{f}_n - \hat{g}_n\| = (2\pi)^d \|f_n - g_n\| \leq (2\pi)^d \|f_n - f\| + (2\pi)^d \|g_n - f\| \rightarrow 0. \quad (30.143)$$

Par conséquent $(\hat{f}_n - \hat{g}_n)_n$ est une suite qui converge vers zéro. Par unicité de la limite, $F = G$.

□

Remarque 30.29.

Une autre suite possible, à la place de (30.138), est

$$f_n(x) = f(x) \mathbb{1}_{|x| < n}. \quad (30.144)$$

C'est-à-dire la fonction f limitée à une boule de rayon n autour de 0.

Chapitre 31

Distributions

Nous donnons ici une partie de la théorie sur les distributions. L'utilisation des distributions dans le cadre des équations différentielles est mise dans le chapitre sur les équations différentielles, section 33.11.

Proposition 31.1 ([439]).

Soient un ouvert Ω de \mathbb{R} et une fonction intégrable $f: (\Omega, \mathcal{B}(\Omega), \lambda) \rightarrow \mathbb{C}$ telle que

$$\int_{\Omega} f(t)\varphi(t)dt = 0 \quad (31.1)$$

pour toute fonction $\varphi \in \mathcal{D}(\Omega)$. Alors $f = 0$ presque partout sur Ω .

Démonstration. Nous commençons par prouver que f est nulle sur tout compact de Ω . Soit un compact K de Ω . Le lemme d'Urysohn 16.121 nous donne une fonction θ à support compact qui vaut 1 sur K .

Nous considérons une suite régularisante (ϕ_k) de fonctions toujours strictement positives (par exemple celle du lemme 30.20). Vu que $f\theta$ est à support compact, elle est dans $L^p(\Omega)$ et le corollaire 30.21 s'applique :

$$\phi_k * (\theta f) \rightarrow \theta f. \quad (31.2)$$

Mais, x et k étant fixés, nous avons

$$(\phi_k * (\theta f))(x) = \int_{\mathbb{R}} \phi_k(x-t)\theta(t)f(t). \quad (31.3)$$

La fonction

$$t \mapsto \phi_k(x-t)\theta(t) \quad (31.4)$$

étant à support compact, l'hypothèse à propos de f fait que l'intégrale (31.3) est nulle :

$$\phi_k * (\theta f) = 0 \quad (31.5)$$

pour tout k . En prenant la limite $k \rightarrow \infty$,

$$\theta f = 0. \quad (31.6)$$

Vu que $\theta(x) = 1$ pour tout $x \in K$, nous avons $f(x) = 0$ pour tout $x \in K$.

Nous avons démontré que f était nulle sur tout compact de Ω .

Nous considérons maintenant une suite exhaustive (K_n) de compacts (lemme 9.59). La fonction f est nulle sur chaque K_n , et comme $\Omega = \bigcup_{n=0}^{\infty} K_n$, la fonction f est nulle sur Ω . \square

31.1 Dérivée faible

31.1.1 Dérivée partielle au sens faible

Lemme-définition 31.2.

Soit $f \in L^p(I)$ où I est l'intervalle ouvert $]a, b[$. Il existe au maximum¹ une fonction g telle que

$$\int_I f\varphi' = - \int_I g\varphi \quad (31.7)$$

pour tout $\varphi \in \mathcal{D}(I)$. Lorsqu'une telle fonction existe, nous la nommons **dérivée faible** de f .

Démonstration. Soient $g, h \in L^2$ tels que

$$\int_I u\varphi' = - \int_I g\varphi = - \int_I h\varphi \quad (31.8)$$

pour tout $\varphi \in C_c^\infty(I)$. Nous avons alors

$$\int_I (g - h)\varphi = 0. \quad (31.9)$$

Cela implique que $g - h = 0$ presque partout par la proposition 28.132². □

Exemple 31.3 (Dérivée faible de $\mathbb{1}_{\mathbb{Q}}$)

Vu que \mathbb{Q} est de mesure nulle dans \mathbb{R} , nous avons

$$\int_{\mathbb{R}} \mathbb{1}_{\mathbb{Q}}\varphi' = 0 \quad (31.10)$$

pour tout $\varphi \in \mathcal{D}(\mathbb{R})$. Pour $g = 0$ nous avons aussi $\int_{\mathbb{R}} g\varphi = 0$. Donc $g = 0$ est la dérivée faible de $\mathbb{1}_{\mathbb{Q}}$.

Cela n'est pas étonnant du fait qu'en théorie de l'intégration, les parties de mesure nulle ne comptent pas. De ce point de vue, $\mathbb{1}_{\mathbb{Q}} = 0$. D'ailleurs cette égalité est vraie dans L^p (les classes et tout ça). △

Exemple 31.4 (La fonction de Heaveside n'a pas de dérivée faible)

Nous montrons que la fonction

$$H(x) = \begin{cases} 0 & \text{si } x < 0 \\ 1 & \text{si } x \geq 0 \end{cases} \quad (31.11)$$

n'a pas dérivée faible. Nous nommons g une hypothétique fonction vérifiant les conditions pour être la dérivée faible de H .

Soit une fonction $\varphi \in \mathcal{D}(\mathbb{R})$, dont le support est contenu dans $]0, \infty[$. Sur le support de φ , et donc aussi de φ' , nous avons $H(x) = 1$ et donc

$$\int_{\mathbb{R}} \varphi' = - \int_{\mathbb{R}} g\varphi. \quad (31.12)$$

Vu que φ est à support compact, $\int_{\mathbb{R}} \varphi' = 0$. En effet, si le support de φ est contenu dans $[-M, M]$, alors en utilisant le théorème fondamental de l'analyse 15.233, nous trouvons $\int_{\mathbb{R}} \varphi = \int_{-M}^M \varphi' = \varphi(M) - \varphi(-M) = 0 - 0 = 0$.

Donc g doit satisfaire

$$\int_{\mathbb{R}} g(x)\varphi(x) = 0 \quad (31.13)$$

1. En réalité, c'est une classe au sens de l'égalité presque partout.

2. Ou alors par le lemme 28.59 qui est moins général mais tout aussi bien pour ici.

pour tout $\varphi \in \mathcal{D}(x > 0)$. La proposition 31.1 nous dit que $g = 0$ presque partout sur $]0, \infty[$.

Le même raisonnement dit que $g = 0$ presque partout sur les négatifs. Que g soit maintenant nulle ou non en $x = 0$ ne change pas le fait que $g = 0$ presque partout sur \mathbb{R} .

Par conséquent, $\int_{\mathbb{R}} g\varphi = 0$ pour toute $\varphi \in \mathcal{D}(\mathbb{R})$. Hélas, nous avons d'autre part

$$\int_{\mathbb{R}} H(x)\varphi'(x)dx = \int_0^{\infty} \varphi'(x)dx, \quad (31.14)$$

qui n'est pas forcément nul. Notons que pour avoir un exemple de φ qui donne $\int_{\mathbb{R}} H\varphi \neq 0$, il faut chercher des fonctions dont le support contient des négatifs et des positifs.

La fonction H n'a donc pas de dérivée faible. Notons cependant que cela ne présume en rien la possibilité d'accepter une dérivée au sens des distributions. \triangle

31.5.

Nous verrons dans la proposition 31.24 que la dérivée de Heaveside au sens des distributions est le delta de Dirac.

Exemple 31.6 (Dérivée faible de la valeur absolue)

L'exemple de base de fonction continue qui n'est pas dérivable est la valeur absolue $f(x) = |x|$ prise en $x = 0$. Nous allons montrer ici que la fonction

$$H(x) = \begin{cases} 1 & \text{si } x < 0 \\ -1 & \text{si } x > 0 \end{cases} \quad (31.15)$$

est la dérivée faible de f .

Commençons par noter que H peut valoir la valeur qu'on veut en zéro; de toutes façons la dérivée faible n'est définie qu'à partie de mesure nulle près.

Soit $\varphi \in \mathcal{D}(\mathbb{R})$ et $M > 0$ tel que le support de φ soit contenu dans $[-M, M]$. Nous avons d'une part

$$\int_{\mathbb{R}} |x|\varphi'(x)dx = - \int_{-M}^0 x\varphi'(x)dx + \int_0^M x\varphi'(x)dx \quad (31.16a)$$

$$= -[\varphi x]_{-M}^0 + \int_{-M}^0 \varphi + [x\varphi]_0^M - \int_0^M \varphi \quad (31.16b)$$

où nous avons utilisé l'intégration par partie de la proposition 21.132 en posant

$$u = x \quad v' = \varphi' \quad (31.17a)$$

$$u' = 1 \quad v = \varphi. \quad (31.17b)$$

Tout cela pour dire que

$$\int_{\mathbb{R}} |x|\varphi'(x)dx = \int_{-M}^0 \varphi - \int_0^M \varphi. \quad (31.18)$$

D'autre part, l'égalité

$$\int_{\mathbb{R}} H(x)\varphi(x)dx = \int_{-M}^0 \varphi - \int_0^M \varphi \quad (31.19)$$

est immédiate.

Nous en déduisons que H est bien la dérivée faible de $x \mapsto |x|$. \triangle

Remarque 31.7.

La dérivée faible ne doit pas être confondue avec la dérivée au sens des distributions qui sera définie en 31.22. Nous avons donc trois notions distinctes de dérivation pour une fonction :

- la dérivée usuelle,

- la dérivée au sens des distributions,
- la dérivée faible.

La dérivée faible d'une fonction reste une fonction, tandis que la dérivée distributionnelle d'une fonction est une distribution.

Je vous mets en garde contre l'idée que l'existence de l'une impliquerait trop facilement l'existence d'une autre³.

31.1.2 Dérivée faible partielle

La notion de dérivée partielle faible est la même que l'autre. Histoire de nous mettre dans le bain, nous écrivons la définition avec les notation du produit scalaire au lieu de l'intégrale.

Définition 31.8.

Si $i = 1, \dots, n$, la **dérivée faible** de v dans la direction e_i est l'application⁴ notée $\partial_i v$ définie par

$$\langle \partial_i v, \phi \rangle = -\langle v, \partial_i \phi \rangle \quad (31.20)$$

pour tout $\phi \in C_c^\infty(\Omega)$.

Lemme 31.9.

Si $v \in L^2$ admet une dérivée faible, alors cette dernière est unique.

Démonstration. Supposons f, g telles que $\langle g, \phi \rangle$ et $\langle f, \phi \rangle$ soient tous deux égaux à $-\langle v, \partial_i \phi \rangle$. En particulier pour tout $\phi \in \mathcal{D}(\Omega)$ nous avons $\langle (f - g), \phi \rangle = 0$.

Cela donne $f - g = 0$ par la proposition 28.132. □

Soit Ω un ouvert de \mathbb{R}^d . Le but de notre histoire est de définir une distribution comme étant un élément de l'espace dual (topologique, voir définition 12.38) de l'espace $\mathcal{D}(\Omega)$ des fonctions C^∞ à support compact dans Ω . Pour ce faire nous devons voir un peu de topologie sur différents espaces de fonctions. Notons que l'espace $\mathcal{D}(\Omega)$ n'est pas réduit à la fonction nulle comme en témoigne l'exemple donné par l'équation (16.458).

Pour chaque K compact dans Ω et multiindice $\alpha \in \mathbb{N}^d$ nous considérons sur $C^\infty(\Omega)$ la semi-norme suivante :

$$p_{K,m}(f) = \sum_{|\mu| \leq m} \|\partial^\mu f\|_{K,\infty}. \quad (31.21)$$

En particulier,

$$p_{K,0}(f) = \sup_{x \in K} |f(x)| = \|f\|_{\infty, K}. \quad (31.22)$$

31.2 Topologie et convergence sur des espaces de fonctions

Définition 31.10.

Les topologies que nous allons considérer sont :

- (1) Sur $C^\infty(\Omega)$, la topologie des semi-normes $p_{K,m}$ (avec K et m comme paramètres).
- (2) Sur $\mathcal{D}(K)$, la topologie des semi-normes $p_{K,m}$ (avec seulement m comme paramètre).
- (3) Sur $\mathcal{D}(\Omega)$, la topologie induite de $C^\infty(\Omega)$.

Cela n'est pas très explicite, mais heureusement nous n'aurons souvent pas besoin de plus que de la notion de convergence dans $\mathcal{D}'(\Omega)$. Rappelons que la topologie d'un espace donne la notion de convergence par la définition 7.25.

3. Wikipédia cite l'exemple de la fonction de Cantor qui est dérivable presque partout au sens usuel, mais qui n'est pas faiblement dérivable. Écrivez-moi si vous connaissez des théorèmes qui lient les trois notions de dérivée.

4. En fait c'est une classe au sens de l'égalité presque partout.

Lemme 31.11 (Convergence dans $\mathcal{D}(K)$).

Si α est un multiindice et si $\varphi_n \xrightarrow{\mathcal{D}(K)} \varphi$, alors nous avons

$$\partial^\alpha \varphi_n \xrightarrow{\text{unif}} \partial^\alpha \varphi. \tag{31.23}$$

Démonstration. Quitte à considérer la suite $\varphi_n - \varphi$ nous pouvons supposer $\varphi_n \xrightarrow{\mathcal{D}(K)} 0$. Nous avons

$$\|\partial^\alpha \varphi_n\| \leq \sum_{\mu \leq \alpha} \|\partial^\mu \varphi_n\|_{K, \infty}. \tag{31.24}$$

Vu que le membre de droite tend vers zéro, nous avons

$$\lim_{n \rightarrow \infty} \|\partial^\alpha \varphi_n\|_{K, \infty} \rightarrow 0, \tag{31.25}$$

ce qui revient à dire que $\partial^\alpha \varphi_n$ converge uniformément sur K vers $\partial^\alpha \varphi$. □

Lemme 31.12.

Si une fonction $f: \mathcal{D}(\Omega) \rightarrow \mathbb{R}$ est continue sur chacun des $\mathcal{D}(K)$ pour tout K compact dans Ω alors est continue sur $\mathcal{D}(\Omega)$.

Démonstration. Soit I ouvert dans \mathbb{R} ; nous devons trouver un ouvert \mathcal{O} dans $C^\infty(\Omega)$ tel que $f^{-1}(I) = \mathcal{D}(\Omega) \cap \mathcal{O}$. Vu que f est continue sur chacun des $\mathcal{D}(K)$ avec K compact dans Ω , pour tout tel compact nous avons un ouvert \mathcal{O}_K dans $\mathcal{D}(K)$ tel que $f^{-1}(I) \cap \mathcal{D}(K) = \mathcal{O}_K$. En tant qu'union d'ouverts⁵, l'ensemble

$$\mathcal{O} = \bigcup_{K \text{ compact de } \Omega} \mathcal{O}_K \tag{31.26}$$

est ouvert dans $C^\infty(\Omega)$. Si $\phi \in \mathcal{D}(\Omega)$, nous avons $\phi \in \mathcal{D}(K)$ pour un certain K compact de Ω , donc $f^{-1}(I) \cap \mathcal{D}(K) = \mathcal{O}_K$. A fortiori nous avons $f^{-1}(I) \cap \mathcal{D}(\Omega) \subset \mathcal{O}$.

Dans l'autre sens, si $\phi \in \mathcal{O}$, alors ϕ est dans un des \mathcal{O}_K et donc dans $f^{-1}(I)$. Nous avons donc bien $f^{-1}(I) = \mathcal{D}(\Omega) \cap \mathcal{O}$. □

Théorème 31.13 (Convergence dans $\mathcal{D}(\Omega)$ [117]).

Soit $(\varphi_n)_{n \in \mathbb{N}}$ une suite dans $\mathcal{D}(\Omega)$ et $\varphi \in \mathcal{D}(\Omega)$. Nous avons $\varphi_n \xrightarrow{\mathcal{D}(\Omega)} \varphi$ si et seulement s'il existe K compact dans Ω tel que $\varphi_n \in \mathcal{D}(K)$ pour tout n et $\varphi_n \xrightarrow{\mathcal{D}(K)} \varphi$.

Démonstration. Supposons que $\varphi_n \xrightarrow{\mathcal{D}(\Omega)} \varphi$ et qu'il n'existe pas de compacts contenant tous les supports des φ_n . Alors pour tout compact de Ω il existe un n tel que le support de φ_n ne soit pas dans K . Nous considérons une suite de compacts (K_n) tels que $\text{Int}(K_n) \subset K_{n+1}$ et $\Omega = \bigcup_n K_n$. Une telle suite existe par le lemme 9.59. Ensuite nous construisons des sous-suites de la façon suivante. D'abord $L_1 = K_1$ et $n_1 \in \mathbb{N}$ est choisi de telle sorte que φ_{n_1} ait un support non contenu dans L_1 . Ensuite L_i est un compact de la suite (K_n) choisi plus loin que L_{i-1} et tel que $\varphi_{n_{i-1}} \in \mathcal{D}(L_i)$. Le nombre n_i est alors choisit plus grand que n_{i-1} de telle sorte que $\varphi_{n_i} \notin \mathcal{D}(L_i)$. Ce faisant, en posant $\phi_i = \varphi_{n_i}$ nous avons

$$\phi_i \in \mathcal{D}(L_{i+1}) \setminus \mathcal{D}(L_i) \tag{31.27}$$

et $\text{Int}(L_n) \subset L_{n+1}$ et $\Omega = \bigcup_n L_n$. Étant donné que (ϕ_i) est une sous-suite de (φ_i) nous avons encore $\phi_i \xrightarrow{\mathcal{D}(\Omega)} \varphi$.

Soit $i \in \mathbb{N}$. Nous allons utiliser le résultat 28.139, aka la seconde forme géométrique du théorème de Hahn-Banach pour séparer les parties $\{\phi_i\}$ (compact) et $\mathcal{D}(L_i)$ (fermé) dans $\mathcal{D}(\Omega)$. Nous avons $f_i \in \mathcal{D}'(\Omega)$ telle que

$$\begin{cases} f_i(\phi_i) > \alpha \end{cases} \tag{31.28a}$$

$$\begin{cases} f_i(\mathcal{D}(L_i)) < \alpha. \end{cases} \tag{31.28b}$$

5. Voir définition 7.1.

Nous redéfinissons immédiatement f_i de façon à avoir

$$\begin{cases} f_i(\phi_i) = 0 & (31.29a) \\ f(\mathcal{D}(L_i)) < 0. & (31.29b) \end{cases}$$

Nous introduisons la fonction définie sur $\mathcal{D}(\Omega)$ par

$$p(\phi) = \sum_{i=1}^{\infty} i \frac{f_i(\phi)}{|f_i(\phi_i)|}. \tag{31.30}$$

Si $\phi \in L_k$, alors $f_k(\phi) = 0$ et même $f_l(\phi) = 0$ pour tout $l \geq k$. Donc pour chaque k , la somme définissant p est finie sur $\mathcal{D}(L_k)$. Nous en déduisons que p est continue sur chacun des $\mathcal{D}(L_k)$ et donc sur $\mathcal{D}(\Omega)$ par le lemme 31.12.

L'image de la suite convergente $\phi_k \xrightarrow{\mathcal{D}(\Omega)} \varphi$ par p doit être bornée parce que p est continue. Mais dans la somme (31.30), tous les termes sont positifs et en particulier le terme $i = k$ vaut k , donc $p(\phi_k) \geq k$, ce qui contredit le fait que l'image de la suite soit bornée. Nous en déduisons donc l'existence d'un compact K tel que $\varphi_n \in \mathcal{D}(K)$ pour tout n .

Nous devons encore prouver que $\varphi_n \xrightarrow{\mathcal{D}(K)} \varphi$ pour ce choix de K . Vu que $\varphi_n \xrightarrow{\mathcal{D}(\Omega)} \varphi$, le lemme 7.28 nous dit que nous avons aussi $\varphi_n \xrightarrow{C^\infty(\Omega)} \varphi$, ce qui signifie que pour tout K et m nous avons

$$p_{K,m}(\varphi_n - \varphi) \rightarrow 0. \tag{31.31}$$

En particulier pour le K fixé plus haut nous avons $p_m(\varphi_n - \varphi) \rightarrow 0$, c'est-à-dire que $\varphi_n \xrightarrow{\mathcal{D}(K)} \varphi$. □

Proposition 31.14.

Si K est compact dans Ω , l'espace $\mathcal{D}(K)$ est métrique et complet.

Démonstration. Nous allons d'abord montrer que $\mathcal{D}(K)$ est complet. Ensuite nous allons montrer que sa topologie peut être donnée par une distance.

Complet Nous considérons une suite de Cauchy (φ_n) dans $\mathcal{D}(K)$ au sens de la définition 9.19.

Soient $\epsilon > 0$ et $i \in \mathbb{N}$; si k et l sont assez grands nous avons

$$\varphi_k - \varphi_l \in B_i(0, \epsilon). \tag{31.32}$$

En particulier pour $i = 0$ nous avons l'inégalité

$$\|\varphi_k - \varphi_l\|_\infty \leq \epsilon, \tag{31.33}$$

La suite (φ_n) est donc de Cauchy dans $(C(K), \|\cdot\|_\infty)$ et y converge donc par complétude, proposition 13.282. Il existe donc une fonction $\varphi \in C(K)$ telle que

$$\varphi_n \xrightarrow{unif} \varphi. \tag{31.34}$$

Notre jeu à présent est de prouver que $\varphi \in \mathcal{D}(K)$, c'est-à-dire qu'elle est de classe C^∞ .

Soit un multiindice $\alpha = \mu_1, \dots, \mu_n, i$. Si k et l sont assez grands nous avons

$$\|\partial^\alpha(\varphi_k - \varphi_l)\|_\infty \leq \epsilon, \tag{31.35}$$

c'est-à-dire que

$$\|\partial_i(\partial^\mu \varphi_k) - \partial_i(\partial^\mu \varphi_l)\|_\infty \leq \epsilon. \tag{31.36}$$

Si nous notons $\psi_k = \partial^\mu \varphi_k$ cela signifie que $(\partial_i \psi_n)$ est une suite de Cauchy dans $(C(K), \|\cdot\|_\infty)$. Elle y converge donc et il existe une fonction $g_i \in C(K)$ telle que

$$\partial_i \psi_n \xrightarrow{unif} g_i. \tag{31.37}$$

Dans ce cas le théorème 13.298 nous indique que ψ est de classe C^1 , c'est-à-dire que $\varphi \in C^{m+1}(K)$.

Métrie La proposition 9.80 nous dit que la topologie donnée par l'écart

$$d(\varphi_1, \varphi_2) = \sup_{k \geq 1} \min\left\{\frac{1}{k}, p_{k-1}(\varphi_1 - \varphi_2)\right\} \quad (31.38)$$

est la même que celle de $\mathcal{D}(K)$. Il reste à montrer que cette formule est bien une distance au sens de la définition 7.87.

- (1) Nous avons bien $d(\varphi_1, \varphi_2) \geq 0$ parce que tous les éléments du supremum et du minimum sont positifs.
- (2) Si $d(\varphi_1, \varphi_2) = 0$ alors pour tout k nous devons avoir $p_{k-1}(\varphi_1 - \varphi_2) = 0$; en particulier pour $k = 1$ cela donne $\varphi_1 = \varphi_2$.
- (3) Nous avons

$$p_k(\varphi_1 - \varphi_2) = p_k(-(\varphi_2 - \varphi_1)) = p_l(\varphi_2 - \varphi_1) \quad (31.39)$$

en utilisant la propriété (2) de la définition 9.74 de semi-norme.

- (4) Nous avons

$$p_k(\varphi_1 - \varphi_2) = p_k(\varphi_1 - \varphi_3 + \varphi_3 - \varphi_2) \leq p_k(\varphi_1 - \varphi_3) + p_k(\varphi_3 - \varphi_2) \quad (31.40)$$

en utilisant la propriété (3) de la définition 9.74.

□

Notons que la proposition 9.80 nous dit que $\mathcal{D}(K)$ est complet tout autant pour la topologie des semi-normes que pour celle de la distance que nous venons de décrire. Ces deux topologies sont les mêmes. Étant métrique et complet, l'espace $\mathcal{D}(\Omega)$ et donc de Baire par le théorème 9.87. Ce qui est bien avec ces deux topologies identiques c'est qu'on peut utiliser la propriété de Baire même en ne parlant que des semi-normes.

31.3 Distributions

Si Ω est un ouvert de \mathbb{R}^d , alors l'ensemble $\mathcal{D}(\Omega)$ est contenu dans $C^\infty(\Omega)$. Nous allons commencer par définir une topologie sur $C^\infty(\Omega)$ et ensuite donner à $\mathcal{D}(\Omega)$ la topologie induite⁶.

Définition 31.15 (Distribution).

Une **distribution** sur un ouvert Ω de \mathbb{R}^d est une forme linéaire continue sur $\mathcal{D}(\Omega) = C_c^\infty(\Omega)$. C'est donc un élément de $\mathcal{D}'(\Omega)$.

Le théorème suivant donne quelques façons de vérifier qu'une forme linéaire soit continue. En particulier il nous dit que pour prouver qu'une forme linéaire est une distribution il suffit de prouver la continuité séquentielle.

Théorème 31.16 ([117, 440]).

Soit T une forme linéaire sur $\mathcal{D}(\Omega)$. Nous avons équivalence entre les points suivants.

- (1) T est continue.
- (2) Pour tout compact $K \subset \Omega$ il existe $m \in \mathbb{N}$ et $C \geq 0$ tel que pour tout $\varphi \in \mathcal{D}(K)$ nous ayons

$$|T(\varphi)| \leq C p_{m,K}(\varphi) \quad (31.41)$$

où $p_{m,K}$ est la semi-norme donnée en (31.21).

- (3) T est séquentiellement continue sur $\mathcal{D}(\Omega)$.
- (4) T est séquentiellement continue en 0.
- (5) Pour tout compact $K \subset \Omega$, la restriction de T à $\mathcal{D}(K)$ est continue.

6. Définition 7.10.

Un certain nombre d'ouvrages prennent le point (2) comme la définition d'une distribution.

Définition 31.17 (Topologie sur $\mathcal{D}'(\Omega)$).

Nous munissons l'espace $\mathcal{D}'(\Omega)$ de la **topologie *-faible**, c'est-à-dire celle de la famille de semi-normes

$$\begin{aligned} p_\varphi: \mathcal{D}'(\Omega) &\rightarrow \mathbb{R} \\ T &\mapsto |T(\varphi)| \end{aligned} \quad (31.42)$$

avec $\varphi \in \mathcal{D}(\Omega)$.

Oui, c'est bien une famille de semi-normes indicée par l'ensemble $\mathcal{D}(\Omega)$. Il n'y en a donc a priori pas du tout une quantité dénombrable.

Proposition 31.18 (Convergence au sens des distributions).

Nous avons $T_n \xrightarrow{\mathcal{D}'(\Omega)} T$ si et seulement si $T_n(\varphi) \rightarrow T(\varphi)$ pour tout $\varphi \in \mathcal{D}(\Omega)$.

Démonstration. La convergence $T_n \xrightarrow{\mathcal{D}'(\Omega)} T$ signifie que l'on ait $p_\varphi(T_n - T) \rightarrow 0$ pour tout $\varphi \in \mathcal{D}(\Omega)$, ce qui en retour signifie que

$$|(T_n - T)(\varphi)| \rightarrow 0. \quad (31.43)$$

□

Cette proposition suppose que l'on ait une distribution T qui vérifie $T_n(\varphi) \rightarrow T(\varphi)$ et conclut qu'on a une convergence dans les distributions. Le théorème suivant est plus fort : il va seulement supposer que $T_n(\varphi)$ converge dans \mathbb{C} et va conclure que $T: \varphi \mapsto \lim_{n \rightarrow \infty} T_n(\varphi)$ est une distribution.

Théorème 31.19 ([441]).

Soit (T_n) une suite dans $\mathcal{D}'(\Omega)$ et nous supposons que pour tout $\varphi \in \mathcal{D}(\Omega)$ la suite $(T_n(\varphi))$ converge dans \mathbb{C} . Alors il existe $T \in \mathcal{D}'(\Omega)$ telle que $T_n \xrightarrow{\mathcal{D}'(\Omega)} T$.

Proposition 31.20.

L'application

$$\begin{aligned} i: L^2(\Omega) &\rightarrow \mathcal{D}'(\Omega) \\ f &\mapsto T_f \end{aligned} \quad (31.44)$$

est une injection continue.

Démonstration. Le fait que ce soit une injection est le fait que si $T_f = T_g$ alors pour tout $\phi \in \mathcal{D}(\Omega)$ nous avons $\langle f - g, \phi \rangle = 0$, et cela implique que $f - g$ est nulle presque partout en tant que fonction et est simplement nulle en tant que classe de fonction dans L^2 .

En ce qui concerne la continuité, il suffit de la prouver en zéro (par linéarité). Soit donc $f_n \xrightarrow{L^2} 0$ et montrons que $T_{f_n} \xrightarrow{\mathcal{D}'(\Omega)} T_0$. Pour prouver cela, la proposition 31.18 nous indique qu'il est suffisant de tester $T_n(\phi) \rightarrow 0$ pour tout $\phi \in \mathcal{D}(\Omega)$.

Notons que si $\phi \in \mathcal{D}$ a fortiori $\phi \in L^2$. Nous avons

$$T_{f_n}(\phi) = \int_{\Omega} f_n \phi \leq \|f_n \phi\|_{L^1} \leq \|f_n\|_2 \|\phi\|_2 \rightarrow 0 \quad (31.45)$$

où nous avons utilisé l'inégalité de Hölder de la proposition 28.33. □

Cette proposition permet de donner un sens à des phrases du type « Soit une distribution T . Si $T \in L^2$, alors ... ». Cela signifie qu'il existe $u \in L^2$ tel que $T = T_u$. Notons que dans ce cas, la distribution est définie sur L^2 et non seulement sur \mathcal{D} .

31.3.1 Multiplication d'une distribution par une fonction

Définition 31.21.

Si $T \in \mathcal{D}'(\Omega)$ et si $f \in C^\infty(\Omega)$ nous définissons la distribution fT par

$$(fT)(\varphi) = T(f\varphi). \quad (31.46)$$

Souvent écrit sous la forme plus compacte $\langle fu, \phi \rangle = \langle u, f\phi \rangle$.

Cela a un sens parce que si $\varphi \in \mathcal{D}(\Omega)$ alors $f\varphi$ est aussi dans $\mathcal{D}(\Omega)$.

Cette définition est motivée par ce que l'on ferait pour une distribution à densité. Si T est une distribution de densité notée également T , nous avons $T(\phi) = \int T(x)\phi(x)$ et donc

$$(fT)(\phi) = \int (fT)(x)\phi(x) = \int T(x)f(x)\phi(x) = \int T(x)(f\phi)(x) = T(f\phi). \quad (31.47)$$

En ce qui concerne les distributions tempérées, nous pouvons définir le produit avec une fonction $f \in \mathcal{S}'(\Omega)$ par la même formule : si $f, \varphi \in \mathcal{S}'(\Omega)$ alors le produit $f\varphi$ est encore Schwartz. Notons toutefois que nous ne pouvons pas définir fT dans $\mathcal{S}'(\Omega)$ si f est seulement dans $C^\infty(\Omega)$.

31.3.2 Dérivée de distribution

Proposition-définition 31.22.

Soit T une distribution sur Ω et $\alpha \in \mathbb{N}^d$. Alors la formule

$$(\partial^\alpha T)(\varphi) = (-1)^{|\alpha|} T(\partial^\alpha \varphi) \quad (31.48)$$

définit une distribution $\partial^\alpha T$.

Cette distribution $\partial^\alpha T$ sera la **dérivée distributionnelle** de T . Notons que le même résultat est encore valide pour des distributions tempérées, et la démonstration est la même.

Démonstration. La forme linéaire $\partial^\alpha T$ sera continue si elle est séquentiellement continue par le théorème 31.16. Nous considérons donc une suite $\varphi_n \xrightarrow{\mathcal{D}(\Omega)} \varphi$ et nous vérifions que

$$\lim_{n \rightarrow \infty} (\partial^\alpha T)(\varphi_n) = (\partial^\alpha T)(\varphi). \quad (31.49)$$

D'abord T étant une distribution (et donc continue) nous pouvons la permuter avec la limite :

$$\lim_{n \rightarrow \infty} (\partial^\alpha T)(\varphi_n) = \lim_{n \rightarrow \infty} (-1)^{|\alpha|} T(\partial^\alpha \varphi_n) = (-1)^{|\alpha|} T\left(\lim_{n \rightarrow \infty} \partial^\alpha \varphi_n\right). \quad (31.50)$$

Notons qu'à gauche la limite est une limite dans \mathbb{R} tandis qu'à droite c'est une limite dans $\mathcal{D}(\Omega)$. Ensuite le lemme 31.11 nous dit que l'hypothèse $\varphi_n \xrightarrow{\mathcal{D}(\Omega)} \varphi$ signifie en particulier que nous avons un compact $K \subset \Omega$ contenant tous les supports des φ_n et que $\partial^\alpha \varphi_n$ converge uniformément (sur K et donc sur Ω) vers $\partial^\alpha \varphi$. Donc

$$\lim_{n \rightarrow \infty} (\partial^\alpha T)(\varphi_n) = (-1)^{|\alpha|} T\left(\lim_{n \rightarrow \infty} \partial^\alpha \varphi_n\right) = (-1)^{|\alpha|} T(\partial^\alpha \varphi) = (\partial^\alpha T)(\varphi), \quad (31.51)$$

ce qui est la relation demandée. □

Le lemme suivant montre une compatibilité entre la dérivée des distributions, la dérivée faible et l'injection de L^2 dans l'espace des distributions.

Lemme 31.23.

Soit Ω un ouvert bornée de \mathbb{R}^n et $f \in L^2(\Omega)$. Alors nous avons

$$\partial_i(T_f) = T_{\partial_i f} \quad (31.52)$$

où la dérivée à droite est la dérivée faible définie en 31.8.

Démonstration. En utilisant la définition de la dérivation de distribution, pour tout $\phi \in \mathcal{D}$ nous avons

$$\partial_i(T_f)\phi = -T_f(\partial_i\phi) = -\langle f, \partial_i\phi \rangle = \langle \partial_i f, \phi \rangle = T_{\partial_i f}(\phi). \tag{31.53}$$

Nous avons utilisé la définition (31.20) de la dérivée faible. □

Nous avons déjà vu dans l'exemple 31.4 que la fonction de Heaveside n'a pas de dérivée faible. Nous allons à présent voir que cette fonction a une dérivée au sens des distributions. Intuitivement, la fonction de Heaveside a une dérivée nulle partout sauf en $x = 0$ où sa dérivée serait infinie ; nous nous attendons à un delta de Dirac.

Proposition 31.24.

La dérivée de la fonction de Heaveside

$$H: \mathbb{R} \rightarrow \mathbb{R}^+ \\ x \mapsto \begin{cases} 0 & \text{si } x \leq 0 \\ 1 & \text{si } x > 0 \end{cases} \tag{31.54}$$

est le delta de Dirac.

Démonstration. Par définition, la dérivée de H au sens des distributions est la distribution H' qui fait

$$H'(\varphi) = -\langle H, \varphi' \rangle \tag{31.55}$$

pour tout élément $\varphi \in \mathcal{D}(\mathbb{R})$. Un petit calcul :

$$-\langle H, \varphi' \rangle = -\int_{\mathbb{R}} H(t)\varphi'(t)dt \tag{31.56a}$$

$$= -\int_0^{\infty} \varphi'(t)dt \tag{31.56b}$$

$$= -\lim_{x \rightarrow \infty} \int_0^x \varphi'(t)dt \tag{31.56c}$$

$$= -\lim_{x \rightarrow \infty} (\varphi(x) - \varphi(0)) \tag{31.56d}$$

$$= \varphi(0). \tag{31.56e}$$

Justifications :

- Pour (31.56b), c'est que $H(t) = 0$ pour $t \in]-\infty, 0[$.
- Pour (31.56c), c'est le lemme 15.226.
- Pour (31.56d), c'est le corollaire 15.247.

□

31.3.3 Ordre et support d'une distribution

Définition 31.25 (support d'une distribution[4]).

*Soit T une distribution. Le **support** de T est le complémentaire de l'union des ouverts \mathcal{O} tels que $T(\varphi) = 0$ pour tout φ à support dans \mathcal{O} .*

Définition 31.26.

*Si T est une distribution sur Ω , nous disons que T est d'**ordre** inférieur ou égal à $p \in \mathbb{N}$ si pour tout compact K de Ω , il existe $C_K \in \mathbb{R}$ tel que pour tout $\varphi \in \mathcal{D}(K)$,*

$$|\langle T, \varphi \rangle| \leq C_K \max_{|\alpha| \leq p} \|\partial^\alpha \varphi\|_\infty. \tag{31.57}$$

Ici α est un multiindice.

La distribution T est d'ordre p si elle est d'ordre inférieur ou égal à p mais pas à $p - 1$.

Pour la proposition suivante, on peut se remémorer la définition 31.10 de la topologie sur $C^\infty(\Omega)$.

Proposition 31.27 ([442]).

Restriction entre C^∞ et \mathcal{D} .

- (1) Si $T \in C^\infty(\Omega)'$, alors la restriction de T à $\mathcal{D}(\Omega)$ est une distribution à support⁷ compact.
- (2) Si T est une distribution à support compact alors elle se prolonge de façon unique en une forme linéaire continue sur $C^\infty(\Omega)$.

Proposition 31.28 ([117]).

Une distribution à support compact est d'ordre fini.

Lemme 31.29 ([168]).

Soit $u \in \mathcal{D}'(\mathbb{R})$ et $\phi \in \mathcal{D}(\mathbb{R})$ tels que $\text{supp}(u) \cap \text{supp}(\phi) = \emptyset$. Alors $\langle u, \phi \rangle = 0$.

Démonstration. Soit $x \notin \text{supp}(u)$. Alors il existe un voisinage V_x de x tel que $\langle u, \psi \rangle = 0$ pour tout $\psi \in \mathcal{D}(V_x)$. En particulier, si $x \in \text{supp}(\phi)$, alors x n'est pas dans le support de u et les ensembles $\{V_x \text{ tel que } x \in \text{supp}(\phi)\}$ recouvrent $\text{supp}(\phi)$. Cependant ϕ est à support compact et nous pouvons extraire un sous-recouvrement fini de $\text{supp}(\phi)$: il existe x_1, \dots, x_p tels que

$$\text{supp}(\phi) \subset \bigcup_{i=1}^p V_{x_i}. \quad (31.58)$$

Nous prenons une partition de l'unité⁸ subordonnée à ce recouvrement. C'est-à-dire des fonctions $\chi_i \in \mathcal{D}(V_{x_i})$ telles que pour tout $x \in \text{supp}(\phi)$,

$$\sum_{i=1}^p \chi_i(x) = 1. \quad (31.59)$$

En particulier nous avons $\sum_i \chi_i(x)\phi(x) = \phi(x)$, et donc

$$\langle u, \phi \rangle = \langle u, \sum \chi_i \phi \rangle = \sum \langle u, \chi_i \phi \rangle = 0 \quad (31.60)$$

parce que $\text{supp}(\chi_i \phi) \subset V_{x_i}$. □

Lemme 31.30 ([168]).

Si u est une distribution d'ordre fini N sur \mathbb{R} , si $\text{supp}(u) = \{x_0\}$ et si

$$\phi(x_0) = \dots = \phi^{(N)}(x_0) = 0 \quad (31.61)$$

alors $\langle u, \phi \rangle = 0$.

Démonstration. Les fonctions plateaux dont nous avons parlé dans la section 16.13.1 nous permettent de considérer une fonction $\chi \in \mathcal{D}(\mathbb{R})$ vérifiant

$$\chi(x) = \begin{cases} 1 & \text{si } x \in \overline{B(0,1)} \\ 0 & \text{si } |x| > 2 \end{cases} \quad (31.62)$$

Ensuite nous posons $\chi_n(x) = \chi(n(x-x_0))$. Par conséquent $\chi_n(x_0) = \chi(0) = 1$ et même $\chi_n(x_0+\epsilon) = \chi(\epsilon) = 1$ tant que ϵ est plus petit que disons $\frac{1}{2}$ pour être sur. Nous en déduisons que la fonction $1 - \chi_n$ s'annule sur un voisinage de x_0 et que donc x_0 n'est pas dans le support de $1 - \chi_n$. Donc $\text{supp}(1 - \chi_n) \cap \text{supp}(u) = \emptyset$ et le lemme 31.29 est utilisable : $\langle u, (1 - \chi_n \phi) \rangle = 0$, ou encore :

$$\langle u, \phi \rangle = \langle u, \chi_n \phi \rangle \quad (31.63)$$

7. Définition 31.25.

8. Lemme 21.20.

pour tout n . Vu que le but est de prouver que $\langle u, \phi \rangle = 0$, nous allons prouver que

$$|\langle u, \chi_n \phi \rangle| \xrightarrow{n \rightarrow \infty} 0. \quad (31.64)$$

Dans ce dessin nous posons

$$\|f\|_n = \sup_{x \in B(0, \frac{2}{n})} \|f(x)\| \quad (31.65)$$

et

$$\|f\|_{(p)} = \sup_{i \leq p} \|\partial^i f\|_\infty. \quad (31.66)$$

La distribution u est d'ordre fini N , et nous en écrivons la définition 31.26 en prenant $\text{supp}(\chi_n \phi)$ en tant que K :

$$|\langle u, \chi_n \phi \rangle| \leq C \max_{k \leq N} \|\partial^k(\chi_n \phi)\|_\infty. \quad (31.67)$$

En remplaçant le maximum par une somme de $k = 0$ à $k = N$, nous majorons. De plus le support de χ_n étant contenu dans $B_n = B(x_0, 2/n)$ nous ne changeons rien en utilisant $\|\cdot\|_n$ au lieu de $\|\cdot\|_\infty$. Donc

$$|\langle u, \chi_n \phi \rangle| \leq C \sum_{k=0}^N \|\partial^k(\chi_n \phi)\|_n \leq C \sum_{k=0}^N \binom{k}{i} \|\partial^i \chi_n\|_n \|\partial^{k-i} \phi\|_n. \quad (31.68)$$

Notons que la seconde inégalité est une inégalité du type $\|fg\| \leq \|f\| \|g\|$. En dérivant un petit peu nous trouvons que

$$(\partial^i \chi_n)(x) = n^i (\partial^i \chi)(n(x - x_0)). \quad (31.69)$$

Donc⁹

$$\|\partial^i \chi_n\|_n = \sup_{x \in B_n} n^i |(\partial^i \chi)(n(x - x_0))| = n^i \sup_{y \in [-2, 2]} |(\partial^i \chi)(y)| = n^i \|\partial^i \chi\|_\infty. \quad (31.70)$$

Nous pouvons donc remplacer $\|\partial^i \chi_n\|_n$ par $n^i \|\partial^i \chi\|_\infty$.

D'autre part nous voulons majorer $\|\partial^{k-i} \phi\|_n$ par quelque chose ne dépendant ni de k ni de i . Nous faisons le théorème des accroissements finis 12.151 : $\|\partial^l \phi\|_n \leq \frac{2}{n} \|\partial^{l+1} \phi\|_n$. Ce n au dénominateur est salutaire parce que nous avons un n^i apparu à cause du remplacement (31.70). Nous faisons donc $i + 1$ fois le théorème des accroissements finis :

$$\|\partial^{k-i} \phi\|_n \leq \left(\frac{2}{n}\right)^{i+1} \|\partial^{k+1} \phi\|_n. \quad (31.71)$$

Toutes ces majorations donnent

$$|\langle u, \chi_n \phi \rangle| \leq C \sum_{k=0}^N \sum_{i=0}^k \binom{i}{k} n^i \underbrace{\|\partial^i \chi\|_\infty}_{\leq \|\chi\|_{(N)}} \left(\frac{2}{n}\right)^{i+1} \underbrace{\|\partial^{k+1} \phi\|_n}_{\leq \|\phi\|_{(N+1)}} \quad (31.72a)$$

$$\leq C \|\chi\|_{(N)} \|\phi\|_{(N+1)} \frac{1}{n} \sum_{k=0}^N \sum_{i=0}^k \binom{k}{i} 2^{i+1} \quad (31.72b)$$

$$= \frac{C'}{n} \quad (31.72c)$$

où C' est une constante qui dépend de χ , de ϕ et de N , mais pas de n . Vu que $\frac{C'}{n} \rightarrow 0$ nous avons bien

$$\langle u, \phi \chi_n \rangle = 0, \quad (31.73)$$

ce qu'il fallait démontrer. \square

Proposition 31.31 ([168]).

Soit $u \in \mathcal{D}(\mathbb{R})'$ avec $\text{supp}(u) = \{x_0\}$. Alors $u = \sum_{i=0}^N a_i \partial^i \delta_{x_0}$ où N est l'ordre de u .

9. Dans [168], la dernière égalité vient avec une inégalité, et je comprends pas pourquoi.

Démonstration. D'abord il faut préciser que l'ordre de u est fini parce que son support est compact (proposition 31.28); nous notons N cet ordre.

Soit $\phi \in \mathcal{D}(\mathbb{R})$. Nous considérons $\chi \in \mathcal{D}(\mathbb{R})$ telle que

$$\chi(x) = \begin{cases} 1 & \text{si } x \in \overline{B(x_0, 1)} \\ 0 & \text{si } |x - x_0| > 2. \end{cases} \quad (31.74)$$

Encore une fois, $1 - \chi$ s'annule sur un voisinage autour de x_0 , ce qui fait que

$$\text{supp}(u) \cap \text{supp}((1 - \chi)\phi) = \emptyset, \quad (31.75)$$

et donc $\langle u, (1 - \chi)\phi \rangle = 0$. Au final,

$$\langle u, \phi \rangle = \langle u, \chi\phi \rangle. \quad (31.76)$$

C'est le moment de poser

$$\psi(x) = \chi(x) \left[\phi(x) - \sum_{k=1}^N \frac{1}{k!} (\partial^k \phi)(x_0) (x - x_0)^k \right] \quad (31.77)$$

La fonction ψ ayant un support disjoint de celui de u , nous avons aussi $\langle u, \psi \rangle = 0$, ce qui donne

$$\langle u, \phi \rangle = \langle u, \chi\phi \rangle = \langle u, \chi \sum_{k=0}^N \frac{1}{k!} (\partial^k \phi)(x_0) (x - x_0)^k \rangle. \quad (31.78)$$

En posant $a_k = \frac{1}{k!} \langle u, x \mapsto \chi(x)(x - x_0)^k \rangle$ nous avons alors

$$\langle u, \phi \rangle = \sum_{k=0}^N a_k (\partial^k \phi)(x_0) = \sum_k (-1)^k a_k (\partial^k \delta_{x_0})(\phi). \quad (31.79)$$

□

31.4 Distributions tempérées

L'espace de Schwartz¹⁰ $\mathcal{S}(\Omega)$ est défini dans la définition 28.143; sa topologie y est discutée.

Définition 31.32.

Une **distribution tempérée** est une forme linéaire continue sur $\mathcal{S}(\mathbb{R}^d)$. L'ensemble des distributions tempérées est noté $\mathcal{S}'(\mathbb{R}^d)$. Si T est une telle distribution, nous notons $\langle T, \varphi \rangle$ l'image de φ par T .

31.33.

Pour rappel, la topologie sur $\mathcal{S}(\mathbb{R}^d)$ est celle de la définition 28.148.

Si f est une fonction sur \mathbb{R}^d telle que $f\varphi \in L^1(\mathbb{R}^d)$ pour tout $\varphi \in \mathcal{S}(\mathbb{R}^d)$, alors nous définissons la distribution $T_f \in \mathcal{S}'(\mathbb{R}^d)$ par

$$\langle T_f, \varphi \rangle = \int_{\mathbb{R}^d} f(x)\varphi(x)dx. \quad (31.80)$$

Cette définition ne fonctionne pas pour toutes les fonctions. Par exemple pour $f(x) = e^{x^2}$, et $\varphi(x) = e^{-x^2} \in \mathcal{S}(\mathbb{R})$ nous avons $f\varphi = 1$ qui n'est pas du tout intégrable sur \mathbb{R} .

10. Attention : ce Schwartz (avec un t) est le Schwartz des distributions dont le prénom est Laurent. À ne pas confondre avec Schwarz (sans t) dont le prénom est Cauchy.

Lemme-définition 31.34.

L'application

$$\begin{aligned} \delta: \text{Fun}(\mathbb{R}^d) &\rightarrow \mathbb{C} \\ \varphi &\mapsto \varphi(0). \end{aligned} \quad (31.81)$$

est

(1) une distribution,

(2) une distribution tempérée.

Cette distribution est nommée **distribution de Dirac**.

Démonstration. Juste pour rappel, $\text{Fun}(X)$ est l'ensemble de toutes les fonctions sur X . Pour prouver que δ est une distribution, nous devons démontrer que $\delta: \mathcal{D}(\mathbb{R}) \rightarrow \mathbb{C}$ est continue. Et pour qu'elle soit une distribution tempérée, il faut démontrer que $\delta: \mathcal{S}(\mathbb{R}) \rightarrow \mathbb{C}$ est continue.

Nous utilisons le théorème 31.16(2) pour démontrer la continuité de δ sur $\mathcal{D}(\mathbb{R})$. Soit un compact $K \subset \mathbb{R}$ et $\varphi \in \mathcal{D}(K)$. En prenant $m = 0$ nous avons la majoration

$$|\delta(\varphi)| = \varphi(0) \leq \|\varphi\|_{K,\infty} = p_{0,K}(\varphi). \quad (31.82)$$

Pour la continuité de δ sur $\mathcal{S}(\mathbb{R}^d)$, nous utilisons les résultats de 28.149. Soit une suite $\varphi_n \xrightarrow{\mathcal{S}} 0$. En particulier, $p_{0,0}(\varphi_n) = \sup_x |\varphi_n(x)| \rightarrow 0$. Donc $\varphi_n(0) \rightarrow 0$ comme il le faut. \square

Exemple 31.35

La **valeur principale** de la fonction $x \mapsto \frac{1}{x}$ est la distribution

$$\begin{aligned} T: \mathcal{S}(\mathbb{R}) &\rightarrow \mathbb{R} \\ \varphi &\mapsto \lim_{\substack{\epsilon \rightarrow 0 \\ \epsilon > 0}} \int_{|x| > \epsilon} \frac{\varphi(x)}{x}. \end{aligned} \quad (31.83)$$

Montrons que cela définit bien une distribution tempérée.

D'abord l'intégrale existe pour tout ϵ , par exemple parce que pour les grands $|x|$ nous avons par exemple $|\varphi(x)| \leq x^3$ et donc $\varphi(x)/x \leq 1/x^2$ dont l'intégrale converge. Nous devons maintenant regarder la limite.

Nous considérons une suite $\epsilon_n \rightarrow 0$ et la suite

$$a_n = \int_{|x| \geq \epsilon_n} \frac{\varphi(x)}{x} dx. \quad (31.84)$$

Nous montrons que cette suite converge dans \mathbb{R} en montrant qu'elle est de Cauchy. Pour cela nous travaillons un peu la forme de φ :

$$\varphi(x) = \varphi(0) + \int_0^x \varphi'(t) dt = \varphi(0) + \int_0^1 x \varphi'(x\theta) d\theta. \quad (31.85)$$

Ce qui est dans l'intégrale est borné par $K = \|M_x \varphi'\|_\infty$ qui est parfaitement fini parce que φ est à décroissance rapide. Lorsque nous calculons $|a_m - a_n|$, le terme $\varphi(0)/x$ donne une intégrale nulle parce que le domaine d'intégration $\epsilon_n \leq |x| \leq \epsilon_m$ est symétrique alors que la fonction $1/x$ est impaire.

$$|a_m - a_n| \leq \left| \int_{\epsilon_m < |x| < \epsilon_n} K \right| = 2|\epsilon_n - \epsilon_m|K \quad (31.86)$$

Tout cela nous dit que T est bien définie. Nous devons encore étudier sa continuité.

Soit χ une fonction dans $C_c^\infty(\mathbb{R})$ telle valant 1 sur $[-1, 1]$, paire et à valeurs dans $[0, 1]$. Pour tout $\epsilon > 0$ nous avons $\int_{|x| > \epsilon} \frac{\chi(x)}{x} dx = 0$.

Nous avons aussi $\varphi = \chi\varphi + (1 - \chi)\varphi$, et donc

$$\int_{|x|>\epsilon} \frac{\varphi(x)}{x} dx = \int_{|\epsilon|>0} \chi(x) \frac{\varphi(x) - \varphi(0)}{x} dx + \int_{|\epsilon|>0} (1 - \chi(x)) \frac{\varphi(x)}{x} dx \quad (31.87a)$$

$$= \int_{|\epsilon|>0} \chi(x) \int_0^1 \underbrace{\varphi'(\theta x)}_{\leq \|\varphi'\|_\infty} d\theta + \int_{|x|\geq 1} (1 - \chi(x)) \frac{\varphi(x)}{x} dx \quad (31.87b)$$

$$\leq \|\varphi'\|_\infty \int_{|x|\geq \epsilon} \chi(x) dx + \|\varphi\|_{L^1} \quad (31.87c)$$

$$= C\|\varphi'\|_\infty + \|\varphi\|_1. \quad (31.87d)$$

Cela est valable pour toute fonction $\varphi \in \mathcal{S}(\mathbb{R})$. Mais nous savons que si $\varphi_n \xrightarrow{\mathcal{S}(\mathbb{R})} 0$, alors $\|\varphi_n\|_\infty \rightarrow 0$, $\|\varphi_n'\|_\infty \rightarrow 0$ et $\|\varphi_n\|_1 \rightarrow 0$; donc si $\varphi_n \xrightarrow{\mathcal{S}(\mathbb{R})} 0$, alors

$$T(\varphi_n) = \lim_{\substack{\epsilon \rightarrow 0 \\ \epsilon > 0}} \int_{|x|>\epsilon} \frac{\varphi(x)}{x} \leq C\|\varphi_n'\|_\infty + \|\varphi_n\|_1 \rightarrow 0. \quad (31.88)$$

△

31.4.1 Topologie

La topologie que nous mettons sur l'espace $\mathcal{S}'(\mathbb{R}^d)$ est le même type que celle que nous mettons sur $\mathcal{D}'(\mathbb{R}^d)$, c'est-à-dire celle des semi-normes $p_\varphi(T) = |T(\varphi)|$. La définition 31.17 et la proposition 31.18 restent.

Proposition 31.36.

Nous avons $T_n \xrightarrow{\mathcal{S}'(\mathbb{R}^d)} T$ si et seulement si pour tout $\varphi \in \mathcal{S}(\mathbb{R}^d)$ nous avons $T_n(\varphi) \rightarrow T(\varphi)$.

31.4.2 Distributions associées à des fonctions

Si $f \in L^1_{loc}(\mathbb{R}^d)$ alors nous lui associons une distribution $T_f \in \mathcal{D}'(\mathbb{R}^d)$ définie par la formule

$$T_f(\varphi) = \int_{\mathbb{R}^d} f(x)\varphi(x)dx. \quad (31.89)$$

Proposition 31.37.

L'application ainsi définie

$$\begin{aligned} L^1_{loc}(\mathbb{R}^d) &\rightarrow \mathcal{D}'(\mathbb{R}^d) \\ f &\mapsto T_f \end{aligned} \quad (31.90)$$

est injective.

Démonstration. Si $T_f = 0$ alors pour tout $\varphi \in \mathcal{D}$ nous avons $\int_{\mathbb{R}^d} f\varphi = 0$. En vertu de la proposition 28.133 cela implique $f = 0$ presque partout. □

31.4.3 Composition avec une fonction

Proposition 31.38 ([441], page 113 et 32).

Soit $T \in \mathcal{S}'(\Omega)$ et $f \in C^k(A \times \Omega)$ où A est ouvert dans \mathbb{R}^d . Nous posons

$$\begin{aligned} F: A &\rightarrow \mathbb{R} \\ \lambda &\mapsto T(f(\lambda, \cdot)). \end{aligned} \quad (31.91)$$

Alors $F \in C^k(A)$.

31.4.4 Transformée de Fourier d'une distribution tempérée

Définition 31.39.

La **transformée de Fourier** de la distribution tempérée $T \in \mathcal{S}'(\mathbb{R}^d)$ est la distribution \hat{T} définie par

$$\hat{T}(\varphi) = T(\hat{\varphi}) \quad (31.92)$$

pour tout $\varphi \in \mathcal{S}(\mathbb{R}^d)$.

Lemme 31.40.

Si $f \in \mathcal{S}(\mathbb{R}^d)$, nous avons $\hat{T}_f = T_{\hat{f}}$.

Démonstration. En utilisant les définitions,

$$\hat{T}_f(\varphi) = T_f(\hat{\varphi}) = \int_{\mathbb{R}^d} f(x)\hat{\varphi}(x)dx = \int_{\mathbb{R}^d} f(x) \left[\int_{\mathbb{R}^d} \varphi(y)e^{-iyx}dy \right] dx \quad (31.93)$$

où nous avons noté xy le produit scalaire $x \cdot y$. Nous permutons les intégrales en utilisant le théorème de Fubini 15.259 avec la fonction

$$(x, y) \mapsto f(x)\varphi(y)e^{-ixy} \quad (31.94)$$

qui est parfaitement dans $L^1(\mathbb{R}^d \times \mathbb{R}^d)$. Nous écrivons alors

$$\hat{T}_f(\varphi) = \int_{\mathbb{R}^d} \left[\int_{\mathbb{R}^d} f(x)\varphi(y)e^{-iyx}dx \right] dy = \int_{\mathbb{R}^d} \varphi(y)\hat{f}(y)dy = T_{\hat{f}}(\varphi). \quad (31.95)$$

□

31.4.5 Convolution d'une distribution par une fonction

Nous savons que si $\psi \in \mathcal{S}(\mathbb{R}^d)$ et si $x \in \mathbb{R}^d$ alors la fonction $y \mapsto \psi(x - y)$ est encore une fonction dans $\mathcal{S}(\mathbb{R}^d)$. Donc si $T \in \mathcal{S}'(\mathbb{R}^d)$ nous pouvons considérer la fonction $T * \psi = \psi * T$ définie par

$$(T * \psi)(x) = T(y \mapsto \psi(x - y)). \quad (31.96)$$

Notons que $T * \psi$ est bien une fonction et non une distribution.

Le but de la définition est d'avoir

$$T_f * \psi = f * \psi. \quad (31.97)$$

En effet

$$(T_f * \psi)(x) = T_f(y \mapsto \psi(x - y)) = \int_{\mathbb{R}^d} f(y)\psi(x - y)dy = (f * \psi)(x). \quad (31.98)$$

Exemple 31.41

La distribution de Dirac est le neutre pour le produit de convolution. En effet

$$(\delta * \psi)(x) = \delta(y \mapsto \psi(x - y)) = \psi(x), \quad (31.99)$$

c'est-à-dire $\delta * \psi = \psi$.

△

Proposition 31.42 ([4]).

Si $T \in \mathcal{S}'(\mathbb{R}^d)$ et $\psi \in \mathcal{S}(\mathbb{R}^d)$, alors la distribution associée à la fonction $T * \psi$ est tempérée.

Démonstration. En agissant sur $\varphi \in \mathcal{D}(\mathbb{R}^d)$ nous avons

$$T_{T*\psi}(\varphi) = \int_{\mathbb{R}^d} T(y \mapsto (t_x\psi)(y))\varphi(x)dx \quad (31.100a)$$

$$= \int_{\mathbb{R}^d} T(y \mapsto \varphi(x)\psi(x-y))dx \quad (31.100b)$$

$$= T\left(y \mapsto \int_{\mathbb{R}^d} \varphi(x)\psi(x-y)dx\right) \quad (31.100c)$$

$$= T(y \mapsto (\varphi * \check{\psi})(y)) \quad (31.100d)$$

$$= T(\varphi * \check{\psi}). \quad (31.100e)$$

Attention : Problèmes et choses à faire

Le passage à la ligne (31.100c) n'est pas justifié. □

31.4.6 Approximation de la distribution de Dirac

Lemme 31.43 ([1]).

Soient des fonctions $j_n: \mathbb{R} \rightarrow \mathbb{R}^+$ de classe C^∞ telles que

- (1) Pour chaque n , la fonction $x \mapsto j_n(|x|)$ est strictement décroissante et converge ponctuellement vers zéro.
- (2) Pour chaque x , la suite $n \mapsto j_n(x)$ est décroissante et converge vers 0.
- (3) Pour tout $M > 0$, la suite j_n converge vers zéro uniformément sur $B(0, M)^c$.
- (4) Pour tout δ et ϵ , il existe un $N \in \mathbb{N}$ tel que $|\int_{B(0, \delta)} j_n(x)dx - 1| \leq \epsilon$.
- (5) Pour tout n , nous avons $\int_{\mathbb{R}} j_n = 1$.

Alors si $u \in \mathcal{S}(\mathbb{R})$ nous avons

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}} u(x)j_n(x)dx = u(0). \quad (31.101)$$

Démonstration. Nous posons

$$I_n = \int_{\mathbb{R}} j_n u \quad (31.102a)$$

$$I_{\delta, n} = \int_{B(0, \delta)} j_n u \quad (31.102b)$$

$$Z_{\delta, n} = \int_{B(0, \delta)} u(0)j_n \quad (31.102c)$$

Nous allons progressivement montrer qu'en prenant δ assez petit et n assez grand, les quantités $|I_n - I_{\delta, n}|$, $|I_{\delta, n} - Z_{\delta, n}|$ et $|Z_{\delta, n} - u(0)|$ peuvent être simultanément majorées par ϵ .

Soient $\delta > 0$ et $\epsilon > 0$; vu que $u \in \mathcal{S}(\mathbb{R})$, il existe M tel que $\int_{|x| > M} |u| < \epsilon$. Soit $N_1 \in \mathbb{N}$ tel que pour tout $n > N_1$ nous avons $|j_n(x)| < 1$ dès que $|x| > M$ (hypothèse (3)). Alors

$$\int_{|x| > M} |j_n(x)u(x)| < \epsilon. \quad (31.103)$$

De plus en posant $s = \max\{|u(x)| \text{ tel que } \delta \leq |x| \leq M\}$ (qui existe parce que u est continue et prise sur un compact) nous pouvons considérer N_2 tel que $j_n(x) < \epsilon/s$ pour tout $|x| > \delta$.

Avec $n > \max\{N_1, N_2\}$ nous avons

$$\left| \int_{B(0, \delta)} j_n u - \int_{\mathbb{R}} j_n u \right| = \left| \int_{B(0, \delta)^c} j_n u \right| \quad (31.104a)$$

$$\leq \int_{\delta \leq |x| \leq M} |j_n u| + \int_{|x| \geq M} |j_n u| \quad (31.104b)$$

$$\leq \epsilon(1 + |M - \delta|). \quad (31.104c)$$

En redéfinissant le ϵ nous avons donc montré que pour tout ϵ et δ , il existe un $N \in \mathbb{N}$ tel que

$$|I_{\delta,n} - I_n| \leq \epsilon \quad (31.105)$$

dès que $n \geq N$.

La fonction u est uniformément continue sur tout $\overline{B(0,\delta)}$, et nous pouvons donc choisir δ tel que $|u(0) - u(x)| \leq \epsilon$ pour tout $x \in \overline{B(0,\delta)}$. Pour ce δ , nous avons déjà trouvé un N tel que $|I_{\delta,n} - I_n| \leq \epsilon$ dès que $n > N$. Nous avons :

$$|I_{\delta,n} - Z_{\delta,n}| \leq \int_{B(0,\delta)} |u(x) - u(0)| j_n(x) dx \quad (31.106a)$$

$$\leq \epsilon \int_{B(0,\delta)} j_n \quad (31.106b)$$

$$\leq \epsilon. \quad (31.106c)$$

Nous avons donc prouvé que pour tout $\epsilon > 0$, il existe un δ et un N tels que

$$\begin{cases} |I_{\delta,n} - I_n| \leq \epsilon \\ |I_{\delta,n} - Z_{\delta,n}| \leq \epsilon \end{cases} \quad (31.107a)$$

$$(31.107b)$$

dès que $n \geq N$.

Enfin nous avons

$$|z_{\delta,n} - u(0)| = u(0) \left(\int_{B(0,\delta)} j_n - 1 \right), \quad (31.108)$$

et par l'hypothèse (4) nous pouvons choisir n assez grand pour que la parenthèse soit plus petite que ϵ .

Pour ϵ donné, nous avons donc trouvé un δ et un N tels que

$$|I_n - u(0)| \leq |I_n - I_{\delta,n}| + |I_{\delta,n} - Z_{\delta,n}| + |Z_{\delta,n} - u(0)| \leq 3\epsilon. \quad (31.109)$$

En passant à la limite nous avons bien $I_n \rightarrow u(0)$ dans \mathbb{R} . □

Il va sans dire que nous connaissons de telles fonctions. Nous en donnons une maintenant.

Exemple 31.44([443])

Nous introduisons la fonction f_ϵ ($\epsilon > 0$) donnée par

$$f_\epsilon(x) = e^{-\epsilon|x|}. \quad (31.110)$$

Nous calculons la transformée de Fourier de f_ϵ en divisant le domaine d'intégration :

$$\hat{f}_\epsilon(k) = \int_{\mathbb{R}} e^{-ikx} e^{-\epsilon|x|} dx = \int_{-\infty}^0 e^{\epsilon x} e^{-ikx} dx + \int_0^{\infty} e^{-\epsilon x} e^{-ikx} dx \quad (31.111)$$

En décomposant les parties imaginaires et réelles, et avec un peu de changement de variables, nous pouvons utiliser les intégrales (18.488) et (18.489) pour obtenir

$$\hat{f}_\epsilon(k) = \frac{2\epsilon}{k^2 + \epsilon^2}. \quad (31.112)$$

Sachant que $\text{arctg}(x)$ est une primitive de $\frac{1}{x^2+1}$ et avec encore un peu de changement de variables, nous avons¹¹

$$\int_{\mathbb{R}} \hat{f}_\epsilon(k) dk = \int_{-\infty}^{\infty} \frac{2\epsilon}{k^2 + \epsilon^2} = 2[\text{arctg}(x/\epsilon)]_{-\infty}^{\infty} = 2\pi. \quad (31.113)$$

11. Et en écrivant correctement l'intégrale sur \mathbb{R} comme une limite, etc.

Cela montre que si nous introduisons la fonction δ_ϵ donnée par

$$\delta_\epsilon(k) = \frac{1}{\pi} \frac{\epsilon}{\epsilon^2 + k^2}, \quad (31.114)$$

alors nous avons une fonction qui tout en même temps ressemble à \hat{f}_ϵ et vérifie

$$\int_{\mathbb{R}} \delta_\epsilon(k) dk = 1 \quad (31.115)$$

pour tout ϵ .

Jusqu'ici nous avons montré que

$$\int_{\mathbb{R}} e^{-ikx} e^{-\epsilon|x|} dx = 2\pi \delta_\epsilon(k). \quad (31.116)$$

Pour chaque $\epsilon > 0$ nous avons $\delta_\epsilon \in L^1(\mathbb{R})$. △

Proposition 31.45 ([1]).

Soit $g \in \mathcal{S}(\mathbb{R})$. Alors nous avons

$$\int_{\mathbb{R}} \int_{\mathbb{R}} g(x) e^{-ixy} dx dy = 2\pi g(0). \quad (31.117)$$

Démonstration. Soit $u \in \mathcal{S}(\mathbb{R})$; nous multiplions l'équation (31.116) par $u(k)$ et nous intégrons par rapport à k :

$$\int_{\mathbb{R}} u(k) \left[\int_{\mathbb{R}} e^{-ikx} e^{-\epsilon|x|} dx \right] dk = 2\pi \int_{\mathbb{R}} u(k) \delta_\epsilon(k) dk. \quad (31.118)$$

Il s'agit de passer à la limite dans l'équation (31.118). Les intégrales à gauche peuvent effectuées séparément parce qu'elles respectent le théorème de Fubini. En effet soit la fonction

$$f(k, x) = u(k) e^{-ikx} e^{-\epsilon|x|} \quad (31.119)$$

qui est dans $L^1(\mathbb{R} \times \mathbb{R})$ en vertu du critère du corollaire 15.258 et du fait que à la fois $k \mapsto |u(k)|$ et $x \mapsto e^{-\epsilon|x|}$ sont dans $L^1(\mathbb{R})$.

Nous pouvons donc grouper et dégroupier les intégrales et en particulier les inverser. Si nous effectuons d'abord l'intégrale sur k nous trouvons

$$\int_{\mathbb{R}} u(k) \left[\int_{\mathbb{R}} e^{-ikx} e^{-\epsilon|x|} dx \right] dk = \int_{\mathbb{R}} e^{-\epsilon|x|} \int_{\mathbb{R}} u(k) e^{-ikx} dk dx = \int_{\mathbb{R}} e^{-\epsilon|x|} \hat{u}(x) dx. \quad (31.120)$$

La fonction $x \mapsto |e^{-\epsilon|x|} \hat{u}(x)|$ est majorée (uniformément en ϵ) par $x \mapsto \hat{u}(x)$ qui est intégrable parce que la transformée de Fourier d'une fonction de \mathcal{S} est dans \mathcal{S} par la proposition 30.14. Le théorème de la convergence dominée de Lebesgue 15.184 nous permet de permuter la limite $\epsilon \rightarrow 0$ avec l'intégrale et obtenir

$$\lim_{\epsilon \rightarrow 0} \int_{\mathbb{R}} u(k) \int_{\mathbb{R}} e^{-ikx} e^{-\epsilon|x|} dx dk = \int_{\mathbb{R}} \hat{u}(x) dx = \int_{\mathbb{R}} \int_{\mathbb{R}} u(k) e^{-ikx} dk dx. \quad (31.121)$$

Notons qu'en passant à la limite nous avons perdu le droit de permuter les intégrales.

Nous devons encore prouver que

$$\lim_{\epsilon \rightarrow 0} \int_{\mathbb{R}} u(k) \delta_\epsilon(k) dk = u(0). \quad (31.122)$$

Cela n'est rien d'autre que le lemme 31.43 appliqué à la suite de fonctions $j_n = \delta_{1/n}$. □

31.46.

Notons que les intégrales dans (31.117) ne peuvent pas être permutées parce que $\int_{\mathbb{R}} e^{-ixy} dy$ n'existe pas. Il faut avouer que, malgré tous les conseils du type « attention : permuter des intégrales doit être fait avec prudence », ce n'est pas tous les jours que nous trouvons des intégrales qui ne peuvent pas être permutées, autrement que dans des exemples fait exprès.

31.4.7 Peigne de Dirac

Proposition 31.47.

La formule

$$\Delta_a = \sum_{k \in \mathbb{Z}} \delta_{ka} \quad (31.123)$$

définit un élément de $\mathcal{D}'(\mathbb{R})$.

La forme linéaire Δ_a est le **peigne de Dirac** de pas a .

Démonstration. Nous utilisons le critère de continuité séquentielle en zéro du théorème 31.16. Soit une suite $\varphi_n \rightarrow 0$ dans $\mathcal{D}(\mathbb{R})$. Par le théorème 31.13 il existe un compact K de \mathbb{R} pour lequel $\varphi_n \in \mathcal{D}(K)$ pour tout n et $\varphi_n \rightarrow 0$ dans $\mathcal{D}(K)$. La somme 31.123 est donc finie et nous pouvons la permuter avec une limite :

$$\lim_{n \rightarrow \infty} \Delta_a(\varphi_n) = \sum_{k \in \mathbb{Z}} \lim_{n \rightarrow \infty} \varphi_n(ka). \quad (31.124)$$

La limite $\varphi_n \rightarrow 0$ dans $\mathcal{D}(K)$ signifie que nous avons convergence uniforme de la fonction et de toutes ses dérivées vers 0. En particulier $\|\varphi_n\|_\infty \rightarrow 0$; disons que la somme (qui est finie) fasse s termes :

$$\sum_{k \in \mathbb{Z}} \varphi_n(ka) \leq s \|\varphi_n\|_\infty. \quad (31.125)$$

Le terme de droite tend vers zéro lorsque n tend vers l'infini. □

Donc Δ_a est bien une distribution au sens de la définition 31.15.

Lemme 31.48 ([417]).

Le peigne de Dirac vérifie la relation

$$\Delta_a = \frac{1}{a} \Delta_1 \circ D_a \quad (31.126)$$

où D_a est l'application $D_a: \mathcal{D}(\mathbb{R}) \rightarrow \mathcal{D}(\mathbb{R})$,

$$(D_a f)(x) = a f(ax). \quad (31.127)$$

Démonstration. Pour $\varphi \in \mathcal{D}(\mathbb{R})$ nous avons

$$\Delta_a(\varphi) = \sum_{k \in \mathbb{Z}} \varphi(ka) = \frac{1}{a} \sum_{k \in \mathbb{Z}} (D_a \varphi)(k) = \frac{1}{a} \Delta_1(D_a \varphi). \quad (31.128)$$

□

Proposition 31.49.

Le peigne de Dirac est une distribution tempérée.

Notez qu'il y a plus de fonctions dans $\mathcal{S}(\mathbb{R})$ que dans $\mathcal{D}(\mathbb{R})$; il est donc plus difficile de rentrer dans $\mathcal{S}'(\mathbb{R})$ que dans $\mathcal{D}'(\mathbb{R})$: il est plus compliqué d'avoir existence de $T(\varphi)$ pour tout $\varphi \in \mathcal{S}(\mathbb{R})$ que pour tout $\varphi \in \mathcal{D}(\mathbb{R})$.

Démonstration. Soit $\varphi \in \mathcal{S}(\mathbb{R})$. Nous avons

$$|\Delta_a(\varphi)| = \left| \sum_k \varphi(ak) \right| = \left| \sum_k \frac{(1 + a^2 k^2) \varphi(ak)}{1 + a^2 k^2} \right| \leq \sup_{x \in \mathbb{R}} |(1 + x^2) \varphi(x)| \sum_k \frac{1}{1 + a^2 k^2}. \quad (31.129)$$

La somme $\sum_k \frac{1}{1 + a^2 k^2}$ est une somme convergente, et son supremum est borné par la proposition 28.146 en prenant $Q(x) = 1 + x^2$. En effet sur $\overline{B(0, r)}$ la fonction $x \mapsto (1 + x^2) \varphi(x)$ est bornée par ce que c'est une fonction continue sur un compact, et à l'extérieur de $B(0, r)$ cette fonction est alors bornée par 1. □

Si aucune ambiguïté n'est à craindre, nous noterons f la distribution T_f .

Exemple 31.50

La transformée de Fourier de la distribution de Dirac est la fonction constante : $\hat{\delta} = 1$. En effet si nous agissons sur une fonction test,

$$\hat{\delta}(\varphi) = \delta(\hat{\varphi}) = \hat{\varphi}(0) = \int_{\mathbb{R}^d} \varphi(x) dx. \quad (31.130)$$

△

31.5 L'espace $C^\infty(\mathbb{R}, \mathcal{D}'(\mathbb{R}^d))$

D'abord parlons un peu de continuité en recopiant la proposition 9.84 dans notre contexte.

Proposition 31.51.

Soient I un intervalle ouvert de \mathbb{R} et $u: I \rightarrow \mathcal{D}'(\mathbb{R}^d)$ une fonction continue. Alors

- (1) Pour tout $\varphi \in \mathcal{D}(\mathbb{R}^d)$, l'application $t \mapsto u_t(\varphi)$ est continue.
- (2) Pour tout $\varphi \in \mathcal{D}(\mathbb{R}^d)$, nous avons la limite

$$\lim_{t \rightarrow t_0} u_t(\varphi) = u_{t_0}(\varphi). \quad (31.131)$$

- (3) Nous avons la limite dans $\mathcal{D}'(\mathbb{R}^d)$

$$\lim_{t \rightarrow t_0} u_t = u_{t_0}. \quad (31.132)$$

En ce qui concerne la définition de l'espace $C^\infty(I, \mathcal{D}'(\mathbb{R}^d))$, c'est la définition 9.85. Grâce au point (1) de la proposition 31.51, nous retenons que la propriété fondamentale d'une application $T \in C^k(I, \mathcal{D}'(\Omega))$ est que pour tout $\varphi \in \mathcal{D}(\Omega)$, l'application

$$\begin{aligned} I &\rightarrow \mathbb{C} \\ t &\mapsto T_t(\varphi) \end{aligned} \quad (31.133)$$

est de classe C^k .

Proposition 31.52.

Soit $(T_t) \in C^0(I, \mathcal{D}'(\Omega))$ et $\psi \in \mathcal{D}(I \times \Omega)$. Alors l'application

$$t \mapsto T_t(\psi(t, \cdot)) \quad (31.134)$$

est continue sur I .

Démonstration. La fonction dont nous voulons prouver la continuité est une fonction $\mathbb{R} \rightarrow \mathbb{R}$; il est donc loisible de se contenter de la continuité séquentielle. Soient $t_0 \in I$ et (t_j) une suite dans I convergeant vers t_0 . Nous posons $U_j = T_{t_j}$ et $\psi_j = \psi(t_j, \cdot)$. Par hypothèse de continuité de (T_t) nous avons $U_j \xrightarrow{\mathcal{D}'(\Omega)} T_{t_0}$. D'autre part le support de ψ étant compact nous avons $\text{supp}(\psi) \subset [c, d] \times K$ où $[c, d] \subset I$ et K est compact dans Ω . Par conséquent nous avons aussi $\text{supp}(\psi_j) \subset K$.

Afin d'alléger les notations notons $\tilde{\psi}(x) = \psi(t_0, x)$. Pour tout multiindice α et pour tout j nous avons

$$p_\alpha(\psi_j - \tilde{\psi}) = \left| \partial^\alpha \psi(t_j, x) - \partial^\alpha \psi(t_0, x) \right| \leq |t_j - t_0| \sup_{\substack{t \in [c, d] \\ x \in K}} |\partial_t \partial^\alpha \psi(t, x)| \rightarrow 0. \quad (31.135)$$

Nous avons donc la convergence

$$\psi_j \xrightarrow{\mathcal{D}'(K)} \psi(t_0, \cdot). \quad (31.136)$$

Étant donné que $U_j \xrightarrow{\mathcal{D}'(\Omega)} T_{t_0}$ et $\psi_j \xrightarrow{\mathcal{D}'(\Omega)} \tilde{\psi}$, le point (3) du corollaire 28.6 nous donne la convergence

$$U_j(\psi_j) \rightarrow T_{t_0}(\tilde{\psi}) \quad (31.137)$$

dans \mathbb{C} . Cela est bien la continuité de la fonction $t \mapsto T_t(\psi(t, \cdot))$. \square

Proposition 31.53 ([441]).

Soit $(T_t) \in C^0(I, \mathcal{D}'(\Omega))$. Nous définissons l'application $T: \mathcal{D}(\Omega) \rightarrow \mathbb{C}$ par la formule

$$T(\psi) = \int_I T_t(\psi(t, \cdot)) dt \quad (31.138)$$

pour tout $\psi \in \mathcal{D}(I \times \Omega)$. Alors $T \in \mathcal{D}'(I \times \Omega)$.

Démonstration. La proposition 31.52 nous indique que la fonction $t \mapsto T_t(\psi(t, \cdot))$ est continue. Étant donné qu'elle est seulement non nulle sur un compact, l'intégrale

$$\int_I T_t(\psi(t, \cdot)) dt \quad (31.139)$$

a un sens et est finie. L'application $T: \mathcal{D}(I \times \Omega) \rightarrow \mathbb{C}$ ainsi définie est linéaire. Il reste à voir qu'elle est continue. Pour cela nous allons utiliser le théorème 31.16(2) qui nous dit que nous pouvons nous fixer un compact $[c, d] \times K \subset I \times \Omega$ et considérer $\psi \in \mathcal{D}([c, d] \times K)$.

Soit, pour commencer, donnée une application $\varphi \in \mathcal{D}(K)$. L'application $t \mapsto T_t(\varphi)$ est continue et non nulle sur le et il existe donc $C_\varphi > 0$ tel que

$$|T_t(\varphi)| \leq C_\varphi \quad (31.140)$$

pour tout $t \in [c, d]$.

Nous voulons utiliser le théorème de Banach-Steinhaus dans sa version 28.7 sur la famille d'applications paramétrée par $u \in [c, d]$:

$$\begin{aligned} U_u: \mathcal{D}([c, d] \times K) &\rightarrow \mathbb{R} \\ \psi &\mapsto T_u(\psi(u, \cdot)). \end{aligned} \quad (31.141)$$

Commençons par prouver que cela est une application continue pour chaque u . Ce sera le cas si la projection

$$\begin{aligned} \mathbf{proj}: \mathcal{D}([c, d] \times K) &\rightarrow \mathcal{D}(K) \\ \psi &\mapsto \psi(u, \cdot) \end{aligned} \quad (31.142)$$

est continue. Pour cela nous notons P_{kl} la semi-norme sur $\mathcal{D}([c, d] \times K)$ donnée par

$$P_{k,l}(\psi) = \sum_{n \leq k} \sum_{|\alpha| \leq l} \sup_{\substack{t \in [c,d] \\ x \in K}} |\partial_t^n \partial^\alpha \psi(t, x)|. \quad (31.143)$$

Nous montrons que \mathbf{proj} est séquentiellement continue ; étant donné que les topologies sur $\mathcal{D}(K)$ et $\mathcal{D}([c, d] \times K)$ sont données par des métriques (proposition 31.14), cela suffit pour assurer la continuité grâce à la proposition 9.14. Montrons que si $\psi_n \xrightarrow{\mathcal{D}([c,d] \times K)} 0$, alors $\mathbf{proj}(\psi_n) \xrightarrow{\mathcal{D}(K)} 0$. Pour cela nous remarquons que

$$p_j(\mathbf{proj}(\psi)) = \sum_{|\alpha| \leq j} \sup_{x \in K} |\partial^\alpha \psi(u, x)| \quad (31.144a)$$

$$\leq \sum_{|\alpha| \leq j} \sup_{t \in [c,d]} \sup_{x \in K} |\partial^\alpha \psi(t, x)| \quad (31.144b)$$

$$= P_{0,j}(\psi). \quad (31.144c)$$

Par conséquent

$$p_j(\mathbf{proj}(\psi_n)) \leq P_{0,j}(\psi_n) \rightarrow 0 \quad (31.145)$$

où nous avons utilisé la proposition 9.77. Utilisant cette même proposition à l'envers, nous déduisons que $\mathbf{proj}(\psi_n) \xrightarrow{\mathcal{D}(K)} 0$. Les applications U_u sont donc continues ; elles sont également bornées parce que si $\psi \in \mathcal{D}([c, d] \times K)$ nous avons

$$\sup_{u \in [c, d]} |U_u(\psi)| = \sup_u |T_u(\psi(u, \cdot))|, \quad (31.146)$$

et la continuité déjà évoquée, sur le compact $[c, d]$, nous dit que cette quantité est finie. Le théorème de Banach-Steinhaus peut maintenant être appliqué et il existe $C > 0$ et $k, l \in \mathbb{N}$ tels que pour tout $\psi \in \mathcal{D}([c, d] \times K)$,

$$|U_u(\psi)| \leq CP_{k,l}(\psi) = C \sum_{|\alpha| \leq k} \sum_{n \leq l} \sup_{t,x} |\partial_t^n \partial^\alpha \psi(t, x)| \leq C \sum_{|\alpha| + n \leq k+l} \sup_{t,x} |\partial_t^n \partial^\alpha \psi(t, x)|. \quad (31.147)$$

Quelques remarques

- Nous n'avons pas mis de maximum devant le supremum (alors que la conclusion (28.5) en demande) parce que dans le cas des semi-normes P_{kl} , c'est toujours celle avec k et l le plus grand possible qui sont les plus grandes parce qu'elles sont des sommes emboîtées les unes les autres.
- La fusion de deux sommes est bien une majoration parce qu'il y a plus de termes dans la seconde que dans la première.
- La quantité la plus à droite est (à part le C) ce que nous pouvons noter $P_{k+l}(\psi)$: c'est bien une des semi-normes associées à l'espace de dimension $d + 1$.

Nous majorons maintenant $T(\psi)$ par

$$|T(\psi)| \leq \int_c^d |T_t(\psi(t, \cdot))| dt = \int_c^d |U_t(\psi)| dt \leq C|d - c|P_{k+l}(\psi). \quad (31.148)$$

Maintenant le théorème 31.16(2) appliqué à l'ouvert $I \times \Omega$ et avec ψ au lieu de φ nous informe que $T \in \mathcal{D}(I \times K)$. \square

31.5.1 Dérivation

Quelques propriétés de dérivation des fonctions $I \rightarrow \mathcal{D}(\Omega)$ seront directement énoncées et démontrées dans le cas des distributions tempérées. Les résultats 31.60 et 31.61 seront a fortiori valables si nous remplaçons \mathcal{S} par \mathcal{D} .

31.6 L'espace $C^\infty(\mathbb{R}, \mathcal{S}'(\mathbb{R}^d))$

Dans cette section nous notons I un ouvert de \mathbb{R} et Ω un ouvert de \mathbb{R}^d ; si ψ est une fonction sur $I \times \Omega$ nous allons noter $\psi_t : \Omega \rightarrow \mathbb{R}$ la fonction $\psi_t(x) = \psi(t, x)$. C'est une notation plus légère que $\psi(t, \cdot)$.

31.6.1 Propriétés générales

La définition de l'espace $C^\infty(I, \mathcal{S}'(\Omega))$ est encore la définition 9.85 et les propriétés énoncées dans la proposition 31.51 sont encore bonnes ici.

D'abord parlons un peu de continuité en recopiant la proposition 9.84 dans notre contexte.

Proposition 31.54.

Soient I un intervalle ouvert de \mathbb{R} et $u : I \rightarrow \mathcal{S}'(\mathbb{R}^d)$ une fonction continue. Alors

- (1) Pour tout $\varphi \in \mathcal{S}(\mathbb{R}^d)$, l'application $t \mapsto u_t(\varphi)$ est continue.

(2) Pour tout $\varphi \in \mathcal{S}(\mathbb{R}^d)$, nous avons la limite

$$\lim_{t \rightarrow t_0} u_t(\varphi) = u_{t_0}(\varphi). \quad (31.149)$$

(3) Nous avons la limite dans $\mathcal{S}'(\mathbb{R}^d)$

$$\lim_{t \rightarrow t_0} u_t = u_{t_0}. \quad (31.150)$$

Lemme 31.55.

Nous avons $C^\infty(I, \mathcal{S}'(\Omega)) \subset C^\infty(I, \mathcal{D}'(\Omega))$.

Démonstration. Soit $(T_t) \in C^\infty(I, \mathcal{S}'(\Omega))$. Pour chaque t nous avons

$$T_t \in \mathcal{S}'(\Omega) \subset \mathcal{D}'(\Omega). \quad (31.151)$$

Ensuite il suffit de dire que pour tout $\varphi \in \mathcal{D}(\Omega)$ la fonction

$$t \mapsto T_t(\varphi) \quad (31.152)$$

est de classe C^∞ parce que c'est le cas pour toute fonction dans $\mathcal{S}(\Omega)$. La proposition 31.51 (en changeant \mathcal{D} en \mathcal{S}) conclut que $(T_t) \in C^\infty(I, \mathcal{D}'(\Omega))$. \square

Proposition 31.56.

L'espace $\mathcal{S}(\Omega)$ est complet et métrisable.

Démonstration. En ce qui concerne le métrisable nous reprenons la formule de l'écart (9.112). Dans notre cas pour l'écrire explicitement il faudrait une énumération de \mathbb{N}^2 à partir de 1 (et non de zéro). Cette formule donne bien une distance parce que si $d(\varphi_1 - \varphi_2) = 0$ alors en particulier $p_{00}(\varphi_1 - \varphi_2) = \|\varphi_1 - \varphi_2\|_\infty = 0$ et donc $\varphi_1 = \varphi_2$.

Nous montrons maintenant que $\mathcal{S}(\Omega)$ est complet en y considérant une suite de Cauchy (φ_n) . Soit $\epsilon > 0$ et $\alpha, \beta \in \mathbb{N}$ ainsi que k, l assez grands pour que $\varphi_k - \varphi_l \in B_{\alpha\beta}(0, \epsilon)$. En particulier pour $\alpha = \beta = 0$ nous avons $\|\varphi_k - \varphi_l\|_\infty \leq \epsilon$, ce qui signifie que nous avons une suite vérifiant le critère de Cauchy uniforme 13.279. Elle converge donc uniformément vers une certaine fonction φ que la proposition 13.280 nous assure être continue. Il existe donc $\varphi \in C(\Omega)$ telle que

$$\varphi_k \xrightarrow{\text{unif}} \varphi. \quad (31.153)$$

Nous devons montrer que $\varphi \in \mathcal{S}(\Omega)$. Le fait que φ soit de classe C^∞ s'obtient en utilisant les semi-normes $p_{0,\alpha}(\varphi) = \|\partial^\alpha \varphi\|_\infty$ de la même façon que dans la preuve que $\mathcal{D}(\Omega)$ était complet (proposition 31.14). Nous obtenons en particulier que

$$\partial^\alpha \varphi_k \xrightarrow{\text{unif}} \partial^\alpha \varphi \quad (31.154)$$

pour tout multiindice α . Montrons encore que φ est à décroissance rapide : nous devons montrer que pour tout α et β nous avons

$$p_{\alpha\beta}(\varphi) = \sup_{x \in \Omega} |x^\beta (\partial^\alpha \varphi)(x)| < \infty. \quad (31.155)$$

Étant donné que (φ_n) est de Cauchy dans $\mathcal{S}(\Omega)$ nous avons (pour ϵ fixé et k, l assez grands) :

$$|x^\beta (\partial^\alpha \varphi_k - \partial^\alpha \varphi_l)(x)| \leq \epsilon \quad (31.156)$$

pour tout $x \in \Omega$. En considérant l fixé et en prenant la limite $k \rightarrow \infty$ et en utilisant la convergence uniforme (31.154) nous trouvons que

$$|x^\beta (\partial^\alpha \varphi - \partial^\alpha \varphi_l)(x)| \leq \epsilon \quad (31.157)$$

Du coup nous pouvons faire la majoration

$$\sup_{x \in \Omega} |x^\beta (\partial^\alpha \varphi)(x)| \leq \sup_x |x^\beta (\partial^\alpha \varphi - \partial^\alpha \varphi_l)(x)| + \sup_x |(\partial^\alpha \varphi_l)(x)| \leq \epsilon + p_{\alpha\beta}(\varphi_l) < \infty \quad (31.158)$$

du fait que $p_{\alpha\beta}(\varphi_l) < \infty$ parce que $\varphi_l \in \mathcal{S}(\Omega)$.

Donc $\varphi \in \mathcal{S}(\Omega)$ et ce dernier est alors complet. \square

Proposition 31.57.

Soit $(T_t) \in C^0(I, \mathcal{S}'(\Omega))$ et $\psi \in \mathcal{S}(I \times \Omega)$. Alors la fonction

$$t \mapsto T_t(\psi_t) \quad (31.159)$$

est continue sur I .

Démonstration. Soient $t_0 \in I$ et une suite convergente vers $t_0 : t_j \rightarrow t_0$ dans \mathbb{R} . Vu que (T_t) est continue en t , elle est en particulier séquentiellement continue et nous avons

$$T_{t_j} \xrightarrow{\mathcal{S}'(\Omega)} T_{t_0}. \quad (31.160)$$

Montrons que nous avons aussi $\psi_{t_j} \xrightarrow{\mathcal{S}(\Omega)} \psi_{t_0}$. Pour cela nous utilisons les semi-normes¹² $p_{\alpha\beta}$ définies en (28.567) :

$$p_{\alpha\beta}(\psi_{t_j} - \psi_{t_0}) = \sum_{x \in \Omega} \left| x^\beta (\partial^\alpha \psi(t_j, x) - \partial^\alpha \psi(t_0, x)) \right| \quad (31.161a)$$

$$\leq \sup_{x \in \Omega} \left| x^\beta |t_0 - t_j| \sup_{t \in [t_0, t_j]} |\partial_t \partial^\alpha \psi(t, x)| \right| \quad (31.161b)$$

$$\leq |t_0 - t_j| \sup_{x \in \Omega} \sup_{t \in I} \left| x^\beta \partial_t \partial^\alpha \psi(t, x) \right| \quad (31.161c)$$

$$\leq |t_0 - t_j| P_{(\alpha t), \beta}(\psi). \quad (31.161d)$$

Pour la première majoration nous avons utilisé le théorème des accroissements finie 13.252. Pour la dernière ligne nous avons noté $P_{\alpha\beta}$ les semi-normes de $\mathcal{S}(I \times \Omega)$ et $(t\alpha)$ est le multiindice qui commence par la variable t et qui continue par α . Étant donné que $P_{(\alpha t)\beta}(\psi) < \infty$ nous avons bien

$$p_{\alpha\beta}(\psi_{t_j} - \psi_{t_0}) \rightarrow 0 \quad (31.162)$$

et donc $\psi_{t_j} \xrightarrow{\mathcal{S}(\Omega)} \psi_{t_0}$.

Étant donné que $\mathcal{S}(\Omega)$ est métrisable et complet, le corollaire 28.6 nous dit que

$$T_{t_j}(\psi_{t_j}) \rightarrow T_{t_0}(\psi_{t_0}), \quad (31.163)$$

ce qui est bien le critère de continuité séquentielle de la fonction (31.159). \square

Remarque 31.58.

La proposition 31.53 nous dit, a fortiori, que si $(T_t) \in C^\infty(I, \mathcal{S}'(\Omega))$ alors la formule

$$\tilde{T}(\psi) = \int_I T_t(\psi_t) \quad (31.164)$$

donne un élément $\tilde{T} \in \mathcal{D}'(I \times \Omega)$. Au cas où aucune confusion n'est à craindre, nous pourrions noter également T l'élément de $\mathcal{D}'(I \times \Omega)$ déduit de $T \in C^\infty(I, \mathcal{S}'(\Omega))$.

Notons que ce T ne sera pas toujours une distribution tempérée comme le montre l'exemple suivant.

Exemple 31.59

En posant $T_t(\varphi) = e^{t^2} \varphi(0)$ avec $I = \mathbb{R}$, l'intégrale

$$T(\psi) = \int_{\mathbb{R}} T_t(\psi_t) = \int_{\mathbb{R}} e^{t^2} \psi(t, 0) dt \quad (31.165)$$

ne converge pas pour tout $\psi \in \mathcal{S}(\mathbb{R} \times \Omega)$. En effet par rapport à t , la fonction $\psi(t, 0)$ décroît rapidement mais pas spécialement assez rapidement pour compenser e^{t^2} . \triangle

12. Pas parce que nous en avons envie, mais bien parce qu'elles font partie de la définition de la convergence et de tous ces trucs.

31.6.2 Dérivation

Proposition 31.60 ([441]).

Soit $T \in C^k(I, \mathcal{S}'(\Omega))$ et $0 \leq l \leq k$. Pour tout $t_0 \in I$ l'application

$$T_{t_0}^{(l)} : \mathcal{S}(\Omega) \rightarrow \mathbb{C} \\ \varphi \mapsto \left(\frac{d^l}{dt} T_t(\varphi) \right) (t_0) \quad (31.166)$$

est bien définie, est une distribution et de plus

$$t \mapsto T_t^{(l)} \in C^{k-l}(I, \mathcal{S}'(\Omega)). \quad (31.167)$$

Attention que la formule (31.166) est bonne si $\varphi \in \mathcal{S}(\Omega)$. Si par contre $\psi \in \mathcal{S}(I \times \Omega)$ et qu'on veut regarder $u_t^{(1)}(\psi_t)$ alors il faut regarder la proposition 31.61 et utiliser la formule (31.172) dans laquelle se trouve $u_t^{(1)}(\psi_t)$.

Démonstration. Pour $k = 0$ nous avons $T_t^{(0)} = T_t$ et c'est bon. Pour le cas $k = 1$ et $l = 0$ c'est encore $T_t^{(0)} = T_t$ qui fonctionne.

Le premier cas non trivial à traiter est $k = 1$ et $l = 1$. Nous considérons $t_0 \in I$; par définition de la dérivée, pour tout $\varphi \in \mathcal{S}(\Omega)$, nous avons (pour peu que les limites existent) :

$$T_{t_0}^{(1)}(\varphi) = \frac{d}{dt} [T_t(\varphi)]_{t=t_0} = \lim_{j \rightarrow \infty} \frac{T_{t_0+\epsilon_j}(\varphi) - T_{t_0}(\varphi)}{\epsilon_j} = \lim_{j \rightarrow \infty} U_j(\varphi) \quad (31.168)$$

où

$$U_j = \frac{1}{\epsilon_j} (T_{t_0+\epsilon_j} - T_{t_0}). \quad (31.169)$$

et (ϵ_j) est une suite de réels tendant vers zéro.

Vu que $(T_t) \in C^k(I, \mathcal{S}'(\Omega))$, l'application $t \mapsto T_t(\varphi)$ est de classe C^k et en particulier l'expression (31.168) a une limite lorsque $j \rightarrow \infty$. Donc $T_{t_0}^{(1)}(\varphi)$ est bien définie. Le point (1) du corollaire 28.6 nous dit que $\lim_{j \rightarrow \infty} U_j(\varphi) = T_{t_0}^{(1)}(\varphi)$ et $T_{t_0}^{(1)}$ est une distribution (linéaire et continue).

Nous devons encore voir que $t \mapsto T_t^{(1)}$ est une application $C^0(I, \mathcal{S}'(\Omega))$. Cela est une conséquence du fait que (T_t) soit de classe C^1 , ce qui se traduit par le fait que l'application

$$t \mapsto \frac{d}{dt} (T_t(\varphi)) \quad (31.170)$$

est continue (définition de la dérivée et point (1) de la proposition 31.54 appliquée à la dérivée).

Les cas $k \geq 1$ se traitent par récurrence. \square

Proposition 31.61 ([441]).

Soit $(T_t) \in C^1(I, \mathcal{S}'(\Omega))$ et $\psi \in \mathcal{S}(I \times \Omega)$. Alors la fonction

$$t \mapsto T_t(\psi(t, \cdot)) \quad (31.171)$$

est de classe C^1 sur I et

$$\frac{d}{dt} (T_t(\psi(t, \cdot))) = T_t^{(1)}(\psi(t, \cdot)) + T_t \left(\frac{\partial \psi}{\partial t}(t, \cdot) \right) \quad (31.172)$$

Démonstration. Soient $t_0 \in I$ et $\epsilon_j \rightarrow 0$ une suite réelle. Le membre de gauche de (31.172), écrit en t_0 , donne

$$\spadesuit = \lim_{j \rightarrow \infty} \frac{T_{t_0+\epsilon_j}(\psi(t_0 + \epsilon_j, \cdot)) - T_{t_0}(\psi(t_0, \cdot))}{\epsilon_j} \quad (31.173)$$

Afin d'alléger les notations nous allons écrire $\psi_t = \psi(t, \cdot)$. Dans le numérateur de (31.173) nous ajoutons et soustrayons la quantité $T_{t_0+\epsilon_j}(\psi_{t_0})$ et nous découpons la limite en deux morceaux :

$$\spadesuit = \lim_{j \rightarrow \infty} \frac{T_{t_0+\epsilon_j}(\psi_{t_0+\epsilon_j} - \psi_{t_0})}{\epsilon_j} + \lim_{j \rightarrow \infty} \frac{(T_{t_0+\epsilon_j} - T_{t_0})(\psi_{t_0})}{\epsilon_j} \quad (31.174)$$

Le second terme vaut

$$\frac{d}{dt} \left(T_t(\psi_{t_0}) \right)_{t=t_0} = T_{t_0}^{(1)}(\psi_{t_0}) \quad (31.175)$$

par la proposition 31.60. Occupons nous de l'autre morceau de \spadesuit . Nous posons $U_j = T_{t_0+\epsilon_j}$ et

$$\varphi_j = \frac{1}{\epsilon_j}(\psi_{t_0+\epsilon_j} - \psi_{t_0}). \quad (31.176)$$

Nous voulons utiliser le corollaire 28.6(3) pour obtenir

$$\lim_{j \rightarrow \infty} U_j(\varphi_j) = T_{t_0} \left(\frac{\partial \psi}{\partial t}(t_0, \cdot) \right). \quad (31.177)$$

D'une part (T_t) est de classe C^∞ en t et nous avons donc la convergence $U_j \xrightarrow{\mathcal{S}'(\Omega)} T_{t_0}$. Reste à prouver que

$$\varphi_j \xrightarrow{\mathcal{S}'(\Omega)} \frac{\partial \psi}{\partial t}(t_0, \cdot). \quad (31.178)$$

Cela en remarquant bien que la variable de dérivation n'est pas celle par rapport à laquelle nous voulons la convergence Schwartz¹³. Soient α et β des naturels et calculons un peu :

$$p_{\alpha\beta} \left(\varphi_j - \frac{\partial \psi}{\partial t}(t_0, \cdot) \right) = \sup_{x \in \Omega} \left| x^\beta \partial^\alpha \left(\frac{1}{\epsilon_j} (\psi(t_0 + \epsilon_j, x) - \psi(t_0, x)) - \frac{\partial \psi}{\partial t}(t_0, x) \right) \right| \quad (31.179)$$

Il est à présent l'heure d'utiliser un développement de Taylor avec le reste de la proposition 13.370 :

$$\psi(t_0 + \epsilon_j, x) = \psi(t_0, x) + \epsilon_j \frac{\partial \psi}{\partial t}(t_0, x) + \frac{\epsilon_j^2}{2} \frac{\partial^2 \psi}{\partial t^2}(\bar{t}, x) \quad (31.180)$$

pour un certain $\bar{t} \in [t_0, t_0 + \epsilon_j]$. En mettant ça dans le calcul (31.179) nous restons avec

$$p_{\alpha\beta} \left(\varphi_j - \frac{\partial \psi}{\partial t}(t_0, \cdot) \right) = \sup_{x \in \Omega} \left| x^\beta \partial^\alpha \left(\epsilon_j \frac{\partial^2 \psi}{\partial t^2}(\bar{t}, x) \right) \right| \leq \epsilon_j P_{\alpha, 2; \beta, 0}(\psi) \quad (31.181)$$

où $P_{\alpha, k; \beta, l}$ sont les semi-normes de $\mathcal{S}(I \times \Omega)$ avec la notation plus ou moins évidente de prendre α dérivations sur x , k sur t puis de multiplier par $x^\beta t^l$. Au final nous avons bien

$$\lim_{j \rightarrow \infty} p_{\alpha\beta} \left(\varphi_j - \frac{\partial \psi}{\partial t}(t_0, \cdot) \right) = 0 \quad (31.182)$$

et donc la convergence $\varphi_j \xrightarrow{\mathcal{S}'(\Omega)} \frac{\partial \psi}{\partial t}(t_0, \cdot)$. □

Lemme 31.62.

Soit $(T_t) \in C^1(I, \mathcal{S}'(\Omega))$ alors si \mathcal{F} dénote la transformée de Fourier nous avons

$$\mathcal{F}(T_t^{(1)}) = (\mathcal{F}T)_t^{(1)} \quad (31.183)$$

où $(\mathcal{F}T)$ est la famille de distributions $(\mathcal{F}T)_t = \mathcal{F}T_t$.

Démonstration. Pour la preuve il suffit de tester l'égalité sur une fonction $\varphi \in \mathcal{S}(\Omega)$:

$$(\mathcal{F}T_t^{(1)})(\varphi) = T_t^{(1)}(\mathcal{F}\varphi) = \frac{d}{dt} \left(T_t(\mathcal{F}\varphi) \right) = \frac{d}{dt} \left((\mathcal{F}T_t)(\varphi) \right) = (\mathcal{F}T)_t^{(1)}(\varphi). \quad (31.184)$$

□

13. Je ne sais pas si je me suis bien fait comprendre là.

31.7 Une équation de distribution

Nous allons étudier l'équation

$$(x - x_0)^\alpha u = 0 \quad (31.185)$$

pour $u \in \mathcal{D}'(\mathbb{R})$ et $\alpha \in \mathbb{N}$ est donné fixé. Notons tout de suite que (31.185) est un petit abus de notation pour dire qu'en vertu de la définition 31.21 du produit d'une distribution par une fonction, pour tout $\phi \in \mathcal{D}(\mathbb{R})$, nous avons $u(x \mapsto (x - x_0)\phi(x)) = 0$.

Lemme 31.63 ([168]).

Soit $\alpha \in \mathbb{N}$. Une solution à l'équation

$$(x - x_0)^\alpha u = 0 \quad (31.186)$$

est une distribution à support dans $\{x_0\}$ et d'ordre fini.

Démonstration. Nous commençons par prouver que u est une solution de (31.186) si et seulement si¹⁴ $\langle u, \phi \rangle = 0$ pour tout $\phi \in \mathcal{D}$ telle que

$$\phi(x_0) = \dots = \partial^{\alpha-1}\phi(x_0) = 0. \quad (31.187)$$

Condition nécessaire Supposons que u soit une solution. Alors le corollaire 18.29 du théorème de Hadamard donne $\psi \in \mathcal{D}(\mathbb{R})$ telle que $\phi(x) = (x - x_0)^\alpha \psi(x)$. Dans ce cas, si u est solution de (31.186), alors

$$0 = \langle (x - x_0)^\alpha u, \psi \rangle = \langle u, (x - x_0)^\alpha \psi(x) \rangle = \langle u, \phi \rangle. \quad (31.188)$$

Nous avons vu que si u est solution, alors $\langle u, \phi \rangle = 0$ pour tout ϕ satisfaisant la condition (31.187).

Condition suffisante Supposons maintenant l'inverse : u est une distribution s'annulant sur toute fonction $\phi \in \mathcal{D}'$ satisfaisant (31.187). Nous allons alors prouver que u est une solution. Soit donc $\psi \in \mathcal{D}$ et calculons

$$\langle (x - x_0)u, \psi \rangle = \langle u, (x - x_0)\psi \rangle = 0 \quad (31.189)$$

parce que la fonction $(x - x_0)\psi(x)$ vérifie la condition (31.187).

Nous passons maintenant au cœur de la preuve : nous supposons que u est une solution. Si le support de ϕ est contenu dans $\mathbb{R} \setminus \{x_0\}$ alors ϕ est nulle dans un voisinage de x_0 (et donc $\partial^k \phi = 0$ pour tout k) et $\langle u, \phi \rangle = 0$. Autrement dit, pour tout $\phi \in \mathcal{D}(\mathbb{R} \setminus \{x_0\})$ nous avons $\langle u, \phi \rangle = 0$, ce qui signifie que $\text{supp}(u) \cap (\mathbb{R} \setminus \{x_0\}) = \emptyset$ ou encore que $\text{supp}(u) = \{x_0\}$.

Maintenant que u a un support compact, la proposition 31.28 nous indique qu'elle est d'ordre fini. □

Théorème 31.64 ([168]).

Soit $\alpha \in \mathbb{N}$ et l'équation

$$(x - x_0)^\alpha u = 0 \quad (31.190)$$

pour $u \in \mathcal{D}'(\mathbb{R})$. Les solutions sont les combinaisons linéaires des dérivées de δ_{x_0} jusqu'à la α^{e} exclue.

Démonstration. D'abord montrons que les $\partial^i \delta_{x_0}$ sont des solutions. Avec les définitions 31.21 et 31.22 des dérivées de distributions et de leur produits avec des fonctions¹⁵,

$$(x - x_0)^\alpha \partial^i \delta_{x_0}(\phi) = \delta_{x_0} \left(\partial^i \left((x - x_0)^\alpha \phi(x) \right) \right) \quad (31.191)$$

14. En réalité nous n'aurons besoin que de la condition nécessaire, en particulier pour le théorème 31.64.

15. Comme souvent, dans l'expression suivante, il y a un abus de notation parce que x est une variable muette : il faudrait écrire « $x \mapsto$ » au début de la grande parenthèse.

Si $i < \alpha$ alors dans chaque terme de Leibnitz, il y aura un facteur $(x - x_0)$, et la prise de δ_{x_0} annulera. Si par contre $i \geq \alpha$ alors il y aura le terme

$$\binom{i}{\alpha} \partial^\alpha ((x - x_0)^\alpha) \partial^{i-\alpha} = \binom{i}{\alpha} \alpha! (\partial^{i-\alpha} \phi)(x_0) \quad (31.192)$$

qui est le seul terme contenant $(\partial^{i-\alpha} \phi)(x_0)$. Il suffit alors de choisir $\phi \in \mathcal{D}(\mathbb{R})$ de sorte que

$$(\partial^k \phi)(x_0) = \begin{cases} 0 & \text{si } k \neq i - \alpha \\ 1 & \text{si } k = i - \alpha \end{cases} \quad (31.193)$$

et alors on est certain que le tout n'est pas nul, et donc que $(x - x_0)^\alpha (\partial^i \delta_{x_0}) \neq 0$.

Jusqu'ici nous avons prouvé que $\partial^i \delta_{x_0}$ est solution si et seulement si $0 \leq i < \alpha$.

Il faut encore prouver que les solutions sont toutes des combinaisons linéaires de dérivées de delta de Dirac centrées en x_0 . Pour cela nous invoquons d'abord le lemme 31.63 qui nous assure que u est d'ordre fini et de support $\{x_0\}$. Ensuite la proposition 31.31 nous indique que u doit alors être une combinaison linéaire de dérivées de Dirac. \square

Chapitre 32

Espaces de Sobolev, équations elliptiques

32.1 Espaces de Sobolev

Rappel : la définition de la dérivée faible est [31.2](#).

32.1.1 Sur un intervalle de \mathbb{R}

Sauf mention du contraire dans cette section I est un intervalle borné ouvert $I =]a, b[$ de \mathbb{R} .

Définition 32.1.

Soit $I =]a, b[$ un ouvert borné de \mathbb{R} . L'espace de Sobolev $H^1(I)$ est l'ensemble

$$H^1(I) = \left\{ u \in L^2(I) \text{ tel que } \exists g \in L^2(I) \text{ tel que } \forall \varphi \in C_c^\infty(I), \int_I u \varphi' = - \int_I g \varphi \right\}. \quad (32.1)$$

L'unique élément g de $L^2(I)$ vérifiant $\int_I u \varphi' = - \int_I g \varphi$ est noté u' et est nommé **dérivée**; nous verrons dans les prochaines pages pourquoi.

L'espace H^1 accepte le produit scalaire suivant :

$$\langle u, v \rangle = \int_I uv + \int_I u'v', \quad (32.2)$$

et nous notons $\|\cdot\|_{H^1}$ la norme correspondante qui n'est autre que

$$\|u\|_{H^1} = \langle u, u \rangle = \|u\|_{L^2}^2 + \|u'\|_{L^2}. \quad (32.3)$$

Nous introduisons l'espace $L^1_{loc}(I)$ des fonctions étant L^1 sur tout compact de I .

Corollaire 32.2.

Si $u \in H^1(I)$ et si $u' = 0$ alors il existe une constante C telle que $u = C$ presque partout.

Démonstration. L'hypothèse $u' = 0$ signifie que pour toute fonction $\varphi \in C_c^\infty(I)$,

$$\int_I u \varphi' = \int_I u' \varphi = 0. \quad (32.4)$$

La proposition [28.133](#) nous dit alors qu'il existe une constante C telle que $u = C$ presque partout. \square

Lemme 32.3.

Tout élément de $H^1(I)$ admet un unique représentant continu.

Nous verrons dans le corollaire [32.5](#) que ce représentant pourra être prolongé par continuité sur \bar{I} .

Démonstration. Soient $y_0 \in I$ et $u \in H^1(I)$. Nous considérons la fonction

$$\bar{u}(x) = \int_{y_0}^x u'(t) dt. \quad (32.5)$$

Notons que par définition, $u' \in L^2$ donc l'intégrale ne pose pas de problèmes. Montrons que \bar{u} est continue sur \bar{I} . Pour cela nous considérons $x \in \bar{I}$ et h tel que $x + h \in \bar{I}$. Alors

$$|\bar{u}(x+h) - \bar{u}(x)| = \left| \int_x^{x+h} u' \right| \leq \int_x^{x+h} |u'|. \quad (32.6)$$

Mais la fonction $|u'|$ est dans $L^1_{loc}(I)$ par le lemme 28.35 ; elle est en particulier intégrable sur un ouvert contenant x et par conséquent la dernière intégrale tend vers zéro lorsque h tend vers 0.

Nous prouvons à présent que \bar{u} est dans $H^1(I)$ et que sa dérivée est égale à u' ; pour cela nous allons montrer que pour tout $\varphi \in C_c^\infty(I)$,

$$\int_I \bar{u} \varphi' = - \int_I u' \varphi. \quad (32.7)$$

Nous avons

$$\int_I \bar{u} \varphi' = \int_I \left(\int_{y_0}^x u'(t) dt \right) \varphi'(x) dx = \int_a^{y_0} \left(\int_{y_0}^x u'(t) dt \right) \varphi'(x) dx + \int_{y_0}^b \left(\int_{y_0}^x u'(t) dt \right) \varphi'(x) dx. \quad (32.8)$$

Pour faire plus court, nous notons $f(t, x) = u'(t) \varphi'(x)$. La première intégrale vaut

$$\int_a^{y_0} \left(\int_{y_0}^x u'(t) \varphi'(x) \right) = \int_a^{y_0} \left(\int_{y_0}^a f(t, x) \mathbb{1}_{t < x}(t, x) dt \right) dx \quad (32.9a)$$

$$= \int_{y_0}^a \int_a^{y_0} f(t, x) \mathbb{1}_{t > x} dx dt \quad (32.9b)$$

$$= \int_{y_0}^a \int_a^t f(t, x) dx dt \quad (32.9c)$$

$$= - \int_a^{y_0} \int_a^t u'(t) \varphi'(x) dx dt \quad (32.9d)$$

La permutation d'intégrales pour obtenir (32.9b) est due au théorème de Fubini 15.259(3). Par le même petit jeu, la seconde intégrale vaut

$$\int_{y_0}^b \int_t^b u'(t) \varphi'(x) dx dt. \quad (32.10)$$

En refaisant la somme,

$$\int_I \bar{u} \varphi' = - \int_a^{y_0} u'(t) \left(\int_a^t \varphi'(x) dx \right) dt + \int_{y_0}^b u'(t) \left(\int_t^b \varphi'(x) dx \right) dt \quad (32.11a)$$

$$= - \int_a^{y_0} u'(t) (\varphi(t) - \varphi(a)) dt + \int_{y_0}^b u'(t) (\varphi(b) - \varphi(t)) \quad (32.11b)$$

$$= - \int_a^b u' \varphi \quad (32.11c)$$

$$= - \int_I u' \varphi. \quad (32.11d)$$

Notons que $\varphi(a) = \varphi(b) = 0$ parce que φ est à support compact dans $]a, b[$. Nous avons donc prouvé que \bar{u} est dans $H^1(I)$ et que $\bar{u}' = u'$. Par le corollaire 32.2, nous avons une constante C telle que $\bar{u} = u + C$ presque partout, c'est-à-dire $u = \bar{u} + C$ dans $H^1(I)$.

En résumé, $\tilde{u} = \bar{u} + C$ est un représentant continu de u dans $L^2(I)$.

L'unicité du représentant continu est simplement le fait que deux fonctions continues égales presque partout sont égales (proposition 21.154). □

Proposition 32.4.

Si $u \in H^1(I)$, alors

$$u(x) - u(y) = \int_y^x u' \quad (32.12)$$

pour tout $x, y \in I$.

Démonstration. Pour fixer les idées, nous supposons $x < y$. Nous considérons une suite $\varphi_n \in C_c^\infty(I)$ convergeant uniformément sur I vers $\mathbb{1}_{[x,y]}$. Nous exigeons de plus que

- φ_n' est positive sur $[a, x + \frac{1}{n}]$
- φ_n' est négative sur $[y - \frac{1}{n}, b]$
- $\varphi_n = 1$ sur $[x + \frac{1}{n}, y - \frac{1}{n}]$.
- $\varphi_n = 0$ sur $[a, x - 1/n]$ et sur $[y + 1/n, b]$.

Pour chaque n , nous découpons l'intégrale comme

$$-\int_I u' \varphi_n = \int_I u \varphi_n' = \int_a^{x-1/n} u \varphi_n' + \int_{x-1/n}^{x+1/n} u \varphi_n' + \int_{x+1/n}^{y-1/n} u \varphi_n' + \int_{y-1/n}^{y+1/n} u \varphi_n' + \int_{y+1/n}^b u \varphi_n'. \quad (32.13)$$

Par construction de φ_n , de ces 5 morceaux, il n'en reste que deux de non nulles :

$$\int_I u \varphi_n' = \underbrace{\int_{x-1/n}^{x+1/n} u(t) \varphi_n'(t) dt}_A + \underbrace{\int_{y-1/n}^{y+1/n} u(t) \varphi_n'(t) dt}_B \quad (32.14)$$

Soit $\epsilon > 0$ et n suffisamment grand pour avoir $u(t) \in B(u(x), \epsilon)$ pour tout $t \in B(x, \frac{1}{n})$ et (en même temps) $u(t) \in B(u(y), \epsilon)$ pour tout $t \in B(y, \frac{1}{n})$. C'est la continuité de u qui permet de trouver un tel n . Pour cette valeur de n , en tenant compte des hypothèses sur la positivité de φ_n' nous avons

$$\int_{x-1/n}^{x+1/n} (u(x) - \epsilon) \varphi_n'(t) dt \leq \int_{x-1/n}^{x+1/n} u(t) \varphi_n'(t) dt \leq \int_{x-1/n}^{x+1/n} (u(x) + \epsilon) \varphi_n'(t) dt, \quad (32.15)$$

mais par hypothèse sur φ_n nous trouvons

$$\int_{x-1/n}^{x+1/n} \varphi_n'(t) dt = \varphi_n(x + \frac{1}{n}) - \varphi_n(x - \frac{1}{n}) = 1. \quad (32.16)$$

donc

$$u(x) - \epsilon \leq \int_{x-1/n}^{x+1/n} u(t) \varphi_n'(t) dt \leq u(x) + \epsilon. \quad (32.17)$$

Pour encadrer la seconde, il faut être plus prudent avec les signes parce que φ_n' y est négative. En posant $\psi_n = -\varphi_n$ nous avons

$$-B = \int_{y-1/n}^{y+1/n} u(t) \psi_n(t) dt, \quad (32.18)$$

et donc

$$u(y) - \epsilon \leq -B \leq u(y) + \epsilon \quad (32.19)$$

ou encore

$$-\epsilon - u(y) \leq B \leq \epsilon - u(y). \quad (32.20)$$

En additionnant avec (32.17) nous voyons que pour tout $\epsilon > 0$ il existe un $N(\epsilon)$ tel que nous avons

$$u(x) - u(y) - 2\epsilon \leq \int_I u' \varphi_n \leq u(x) - u(y) + 2\epsilon \quad (32.21)$$

pour tout $n \geq N$. Nous voulons évidemment prendre la limite $\epsilon \rightarrow 0$, c'est-à-dire $n \rightarrow \infty$. Étant donné que $\varphi_n(t) < 1$ pour tout t et pour tout n , la fonction $t \mapsto u'(t)\varphi_n(t)$ est dominée par u' , qui est dans $L^1(I)$ par le lemme 28.35. Le théorème de la convergence dominée nous permet donc d'affirmer que

$$\lim_{n \rightarrow \infty} \int_I u' \varphi_n = \int_I u' \mathbb{1}_{[x,y]} = \int_x^y u', \quad (32.22)$$

et donc les inégalités (32.21) donnent le résultat, grâce au signe dans (32.13). \square

Corollaire 32.5.

Si $[u] \in H^1(I)$, le représentant continu $u \in C^0(I)$ peut être prolongé par continuité en $u \in C^0(\bar{I})$.

Démonstration. Soit (x_n) une suite strictement croissante dans $]a, b[$ convergeant vers b . Nous voulons montrer que la suite $(u(x_n))$ est de Cauchy dans \mathbb{R} , ce qui nous permettra de définir

$$u(b) = \lim_{n \rightarrow \infty} u(x_n). \quad (32.23)$$

qui sera évidemment continue. Cette construction ne dépendra pas du choix de la suite (x_n) parce que deux fonctions continues sur \bar{I} et égales sur I sont égales sur \bar{I} .

En notant u' la dérivée de u dans H^1 , nous avons par construction du représentant continu : $u(x) = \int_{y_0}^x u'(t) dt$. Et donc

$$|u(x_n) - u(x_{n+p})| = \left| \int_{y_0}^{x_n} u' - \int_{y_0}^{x_{n+p}} u' \right| = \left| \int_{x_n}^{x_{n+p}} u' \right|. \quad (32.24)$$

Vu que la suite (x_n) est de Cauchy et que u' est intégrable (même sur \bar{I}), la limite $n \rightarrow \infty$ de cela est zéro, quelle que soit la valeur de p . Donc $(u(x_n))$ est de Cauchy dans \mathbb{R} et est donc convergente. \square

Proposition 32.6 ([43, 1]).

Quelques propriétés de l'espace de Sobolev $H^1(I)$ où $I =]a, b[$ est un ouvert borné de \mathbb{R} .

- (1) $H^1(I)$ est un espace de Hilbert.
- (2) $H^1(I)$ s'injecte de façon compacte dans $C^0(\bar{I})$.
- (3) $H^1(I)$ s'injecte de façon continue dans $L^2(I)$.

Démonstration. Nous prouvons point par point.

- (1) Le seul critère à vérifier est la complétude. Pour cela nous considérons une suite de Cauchy (u_n) dans $H^1(I)$. Si $\epsilon > 0$, alors il existe $N > 0$ tel que pour tout $p \geq 0$ nous avons $\|u_{n+p} - u_n\|_{H^1}^2 \leq \epsilon$, c'est-à-dire

$$\|u_{n+p} - u_n\|_{L^2}^2 + \|u'_{n+p} - u'_n\|_{L^2}^2 + \quad (32.25)$$

En particulier les suites (u_n) et (u'_n) sont de Cauchy dans L^2 qui est complet par le théorème de Fischer-Riesz 28.41. Nous notons donc

$$u_n \xrightarrow{L^2} u \quad (32.26a)$$

$$u'_n \xrightarrow{L^2} v. \quad (32.26b)$$

Nous allons maintenant montrer quelques limites.

$u_n \varphi \xrightarrow{L^2} u \varphi$ Si M est une constante qui majore φ alors $\|u_n \varphi - u \varphi\|_2 \leq M \|u_n - u\|_2 \rightarrow 0$.

$u'_n \varphi \xrightarrow{L^2} v \varphi$ C'est la même chose avec $\|u'_n \varphi - v \varphi\|_2 \leq M \|u'_n - v\|_2 \rightarrow 0$.

$u \in H^1(I)$ avec $u' = v$ Attendu le corollaire 28.36 qui permet de permuter intégrale et limite dans $L^2(I)$ et les limites que nous venons de prouver,

$$\int_I u\varphi' = \lim_{n \rightarrow \infty} \int_I u_n\varphi' = - \lim_{n \rightarrow \infty} \int_I u'_n\varphi = - \int_I v\varphi. \tag{32.27}$$

Cela signifie que v est la dérivée faible de $u : u' = v$.

$u_n \xrightarrow{H^1} u$ Nous pouvons alors prouver que $u_n \rightarrow u$ dans $H^1(I)$:

$$\|u_n - u\|_{H^1(I)}^2 = \|u_n - u\|_{L^2}^2 + \|u'_n - u'\|_{L^2}^2. \tag{32.28}$$

Mais nous savons déjà que $u_n \rightarrow u$ dans L^2 (d'ailleurs c'est la définition de u) et que $u' = v$ alors que par définition de v , nous avons $u'_n \rightarrow v$ dans L^2 .

Tout cela donne que $u_n \rightarrow u$ dans $H^1(I)$ et donc que $H^1(I)$ est un espace complet.

(2) L'application que nous allons prouver être compacte entre $H^1(I)$ et $C^0(\bar{I})$ est

$$\begin{aligned} \psi: H^1(I) &\rightarrow C^0(\bar{I}) \\ [u] &\mapsto \tilde{u} \end{aligned} \tag{32.29}$$

où $[u]$ désigne une classe de fonction dans $H^1(I)$ et \tilde{u} est son représentant continu prolongé par continuité à \bar{I}^1 , qui existe par le lemme 32.3 et le corollaire 32.5. Cette application est une injection par l'unicité du représentant continu. Nous allons prouver que c'est une application compacte en utilisant le critère (2) de la proposition 28.3. Pour cela nous allons commencer par utiliser le théorème d'Ascoli sur l'ensemble $\tilde{\mathcal{B}}$ des représentants continus des éléments de \mathcal{B} , prolongés par continuité sur \bar{I} ; c'est-à-dire $\tilde{\mathcal{B}} \subset C^0(\bar{I})$.

Soit $u \in \tilde{\mathcal{B}}$; par la proposition 32.4, nous avons

$$|u(x) - u(y)| = \left| \int_y^x u'(t) dt \right| \tag{32.30a}$$

$$= \left| \int_I \mathbb{1}_{[x,y]}(t) u'(t) dt \right| \tag{32.30b}$$

$$\leq \| \mathbb{1}_{[x,y]} \|_{L^2} \| u' \|_{L^2} \tag{32.30c}$$

$$\leq \sqrt{|x - y|} \| u' \|_{H^1} \tag{32.30d}$$

$$\leq \sqrt{|x - y|}. \tag{32.30e}$$

où nous insistons sur le fait que la continuité n'impliquant pas la dérivabilité, le u' ici est la dérivé au sens de H^1 , et non la dérivée usuelle. Quoi qu'il en soit, l'ensemble $\tilde{\mathcal{B}}$ est équicontinu². Nous montrons à présent qu'il est également borné pour la norme uniforme. Soit $u \in \tilde{\mathcal{B}}$; vu la construction du représentant continu au lemme 32.3, nous avons

$$|u(x)| = \left| \frac{1}{b - a} \int_a^b u(x) dy \right| \tag{32.31a}$$

$$= \left| \frac{1}{b - a} \int_a^b \left(\int_y^x u'(t) dt - u(y) \right) dy \right| \tag{32.31b}$$

$$= \left| \frac{1}{b - a} \int_a^b \int_y^x u'(t) dt dy - \frac{1}{b - a} \int_a^b u(y) dy \right| \tag{32.31c}$$

$$\leq \frac{1}{b - a} \int_a^b \int_a^b |u'(t)| dt dy + \frac{1}{b - a} \int_a^b |u(y)| dy. \tag{32.31d}$$

1. Encore que par soucis d'économie d'encre nous n'allons pas écrire toujours les tildes et noter u le représentant continu prolongé à \bar{I} par le corollaire 32.5.

2. Définition 9.69.

À ce niveau, il faut remarquer que dans la première intégrale, le passage de la valeur absolue à l'intérieur de l'intégrale en même temps que l'élargissement des bornes n'a rien d'innocent. Si $x < y$, les bornes ne sont pas « dans le bon ordre » et nous ne pouvons pas faire la majoration usuelle en entrant simplement la valeur absolue. Ici nous tenons compte de cela en élargissant les bornes, et en les mettant dans le bon ordre. Le passage exact est le suivant : si $x, y \in]a, b[$, nous avons

$$\left| \int_y^x f(t) dt \right| \leq \left| \int_y^x |f(t)| dt \right| \leq \left| \int_a^b |f(t)| dt \right| = \int_a^b |f(t)| dt. \quad (32.32)$$

Notons en particulier que dans le cas du passage vers l'équation (32.31d), le nombre x est fixé alors que y est une variable d'intégration. Donc l'ordre des deux est certainement de temps en temps le « mauvais ».

Quoi qu'il en soit, la première intégrale se réduit à une multiplication par $b - a$ et le calcul continue :

$$|u(x)| \leq \int_I |u'(t)| dt + \frac{1}{b-a} \int_I |u| \quad (32.33a)$$

$$\leq \sqrt{b-a} \|u'\|_{L^2} + \frac{1}{\sqrt{b-a}} \|u\|_{L^2} \quad (32.33b)$$

$$\leq \left(\sqrt{b-a} + \frac{1}{\sqrt{b-a}} \right) (\|u'\|_{L^2} + \|u\|_{L^2}) \quad (32.33c)$$

$$\leq \left(\sqrt{b-a} + \frac{1}{\sqrt{b-a}} \right) \|u\|_{H^1} \quad (32.33d)$$

$$= \sqrt{b-a} + \frac{1}{\sqrt{b-a}}. \quad (32.33e)$$

Donc $\tilde{\mathcal{B}}$ est borné pour la norme L^∞ . Et c'est même borné par un nombre facilement calculable connaissant I . En particulier l'ensemble

$$\{u(x) \text{ tel que } u \in H^1\} \quad (32.34)$$

est pour, tout x , contenu dans la boule de rayon $\sqrt{a-b} + \frac{1}{\sqrt{a-b}}$ et donc est relativement compact dans \mathbb{R} . Par conséquent le théorème d'Ascoli 28.5 nous dit que l'ensemble $\tilde{\mathcal{B}}$ est relativement compact dans $C^0(I)$.

Par conséquent nous avons montré que l'image par ψ de la boule unité fermée \mathcal{B} de $H^1(I)$ est relativement compacte dans $C^0(\bar{I})$, ce qui signifie que ψ est une application compacte.

- (3) Les éléments de $H^1(I)$ sont des éléments de $L^2(I)$; donc l'identité est une injection. Nous devons seulement étudier la continuité. Si (u_n) est une suite dans H^1 convergeant dans H^1 vers u , alors

$$\|u_n - u\|_{L^2} \leq \|u_n - u\|_{L^2} + \|u'_n - u'\|_{L^2} = \|u_n - u\|_{H^1} \rightarrow 0. \quad (32.35)$$

Donc la suite des images (par l'identité) converge dans L^2 . L'identité est donc continue. \square

32.1.2 Sur un ouvert de \mathbb{R}^n

Soit Ω , un ouvert de \mathbb{R}^n et $v \in L^2(\Omega)$ (voir 28.72). Les fonctions considérées sont à valeurs réelles.

32.1.2.1 Définition

Définition 32.7 (Espace de Sobolev $H^1(\Omega)$).

Soit Ω une partie de \mathbb{R}^n . L'espace de **Sobolev** $H^1(\Omega)$ est :

$$H^1(\Omega) = \{v \in L^2(\Omega) \text{ tel que } \forall i = 1, \dots, n, \partial_i v \in L^2(\Omega)\}. \quad (32.36)$$

Nous munissons cet espace d'un produit scalaire

$$(u, v)_{H^1} = \langle u, v \rangle_{L^2} + \langle \nabla u, \nabla v \rangle_{L^2}, \quad (32.37)$$

où $\nabla u = \sum_i \partial_i u \in L^2$.

L'existence des intégrales dans le produit scalaire est assurée par le fait que $u, v, \nabla u$ et ∇v sont dans $L^2(\Omega)$. La définition du produit scalaire dans L^2 est la définition 28.200 (mais sans la conjugaison complexe).

Pour la même raison, $(u, u)_{H^1} = 0$ demande que chacun des deux termes est séparément nul, et nous avons $u = 0$ dans L^2 , et donc aussi dans H^1 .

Théorème 32.8 ([444]).

L'espace $H^1(\Omega)$ est un espace de Hilbert³.

Démonstration. Nous devons nous assurer que l'espace H^1 est complet. Pour cela nous considérons une suite de Cauchy (u_n) dans H^1 . Soit $\epsilon > 0$; il existe $N > 0$ tel que si $n, m > N$ alors $\|u_n - u_m\|_{H^1} < \epsilon$. Dans ce cas nous avons en particulier

$$\|u_m - u_n\|_{H^1}^2 = (u, u)_{H^1} = \langle u, u \rangle + \langle \nabla u, \nabla u \rangle = \|u\|_{L^2}^2 + \|\nabla u\|_{L^2}^2, \quad (32.38)$$

et en particulier les suites (u_n) et (∇u_n) sont de Cauchy dans L^2 . Vu que L^2 , lui, est complet (théorème 28.40), il existe $u \in L^2$ et $v_i \in L^2$ tels que

$$u_n \xrightarrow{L^2} u \quad (32.39a)$$

$$\partial_i u_n \xrightarrow{L^2} v_i. \quad (32.39b)$$

Nous savons que l'injection $i: L^2 \rightarrow \mathcal{D}'$ est continue par la proposition 31.20. Nous avons donc aussi les limites

$$T_{u_n} \xrightarrow{\mathcal{D}'} T_u \quad (32.40a)$$

$$T_{\partial_i u_n} \xrightarrow{\mathcal{D}'} T_{v_i}. \quad (32.40b)$$

La dérivée étant une opération continue sur \mathcal{D}' nous avons de plus

$$\partial_i(T_{u_n}) \xrightarrow{\mathcal{D}'} \partial_i(T_u) \quad (32.41)$$

En utilisant le lemme 31.23 nous avons alors

$$T_{\partial_i u_n} = \partial_i(T_{u_n}) \xrightarrow{\mathcal{D}'} \partial_i(T_u) = T_{\partial_i u}. \quad (32.42)$$

En comparant avec (32.40b) et par l'unicité de la limite, nous avons $T_{v_i} = T_{\partial_i u}$. Cela implique $v_i = \partial_i u$.

Vu que $v_i \in L^2$ nous avons aussi $\partial_i u \in L^2$. Par conséquent $u \in H^1(\Omega)$ parce que ses dérivées sont dans L^2 .

Nous devons maintenant prouver que $u_n \xrightarrow{H^1} u$. Nous avons

$$\|u_n - u\|_{H^1} = \|u_n - u\|_{L^2} + \|\nabla u_n - \nabla u\|_{L^2} \quad (32.43)$$

Le premier terme tend vers zéro parce que $u_n \xrightarrow{L^2} u$ et le second parce que $\partial_i u_n \xrightarrow{L^2} \partial_i u$. \square

3. Définition 26.1.

32.1.3 Espace de Sobolev fractionnaire

Définition 32.9.

Pour $m \in \mathbb{N}$ et un ouvert Ω de \mathbb{R}^d nous définissons l'espace de Sobolev

$$H^m(\Omega) = \{u \in L^2(\Omega) \text{ tel que } \partial^\alpha u \in L^2(\Omega) \forall |\alpha| \leq m\}. \quad (32.44)$$

Nous définissons également un produit scalaire sur H^m par

$$(u, v)_{H^m} = \sum_{|\alpha| \leq m} \langle \partial^\alpha u, \partial^\alpha v \rangle_{L^2}. \quad (32.45)$$

En particulier la topologie est celle de la norme dérivée du produit scalaire :

$$\|u\|_{H^m(\Omega)} = \sum_{|\alpha| \leq m} \|\partial^\alpha u\|_{L^2(\Omega)}. \quad (32.46)$$

Le lemme suivant montre que la proposition 30.13 fonctionne encore avec L^2 au lieu de \mathcal{S} .

Lemme 32.10 (Lemme de transfert[445], thème 67).

Soit $f \in H^m(\mathbb{R}^d)$. Alors pour tout multiindice α avec $|\alpha| \leq m$ nous avons

$$\mathcal{F}(\partial^\alpha f) = [\xi \mapsto i^{|\alpha|} \xi^\alpha \hat{f}(\xi)]. \quad (32.47)$$

Lemme 32.11.

Il existe des constantes c_1 et c_2 telles que pour tout $x \in \mathbb{R}^d$,

$$c_1(1 + \|x\|^2)^m \leq \sum_{|\alpha| \leq m} (x^\alpha)^2 \leq c_2(1 + \|x\|^2)^m. \quad (32.48)$$

Lemme 32.12.

Soit $u \in L^2(\mathbb{R}^d)$. Nous avons $u \in H^m(\mathbb{R}^d)$ si et seulement si l'application

$$\xi \mapsto (1 + |\xi|^2)^{k/2} \hat{u} \quad (32.49)$$

est dans $L^2(\mathbb{R}^d)$ pour tout $k \leq m$. Ici $|\xi|$ est la norme euclidienne de ξ dans \mathbb{R}^d .

Démonstration. Vu le lemme 32.10, il suffit de montrer que

$$(1 + |\xi|^2)^{k/2} \hat{u} \quad (32.50)$$

est dans L^2 pour tout $k \leq m$ si et seulement si

$$\xi^\alpha \hat{u} \quad (32.51)$$

l'est pour tout α avec $|\alpha| \leq m$.

L'expression (32.50) est une somme d'expressions du type (32.51). Donc l'implication dans un sens est montrée. Pour l'autre sens, nous savons que

$$\xi^\alpha = \xi_1^{\alpha_1} \dots \xi_n^{\alpha_n}, \quad (32.52)$$

et donc

$$|\xi^\alpha| \leq |\xi_1|^{\alpha_1} \dots |\xi_n|^{\alpha_n}. \quad (32.53)$$

Or $|\xi|^{|\alpha|} = |\xi|^{\sum_i \alpha_i} = |\xi|^{\alpha_1} \dots |\xi|^{\alpha_n}$ et $|\xi| \geq |\xi_i|$ pour tout i , donc

$$|\xi^\alpha| \leq |\xi|^{|\alpha|}. \quad (32.54)$$

D'autre part pour tout $x \in \mathbb{R}^+$ et tout k positif nous avons

$$(1 + x^2)^{k/2} \geq x^k \quad (32.55)$$

qui est facile à vérifier en prenant le carré des deux membres.

En remettant tout ensemble,

$$|\xi^\alpha \hat{u}| \leq |\xi^\alpha| |\hat{u}| \leq |\xi|^{|\alpha|} |\hat{u}| \leq (1 + |\xi|^2)^{|\alpha|/2} |\hat{u}|. \quad (32.56)$$

Donc si le membre de droite est de carré intégrable, celui de gauche l'est également. \square

Définition 32.13 (Espace de Sobolev H^s [446]).

Pour $s > 0$ nous définissons l'espace de Sobolev $H^s(\mathbb{R}^d)$ par

$$H^s(\mathbb{R}^d) = \{u \in L^2(\mathbb{R}^d) \text{ tel que } (1 + \|\xi\|^2)^{s/2} \hat{u} \in L^2(\mathbb{R}^d)\}. \quad (32.57)$$

Nous y mettons le produit scalaire

$$(u, v)_{H^s} = \int_{\mathbb{R}^d} \hat{u}(\xi) \overline{\hat{v}(\xi)} (1 + \|\xi\|^2)^s d\xi. \quad (32.58)$$

32.14.

Vu que $\mathcal{D}(\mathbb{R}^d)$ est dense dans $L^2(\mathbb{R}^d)$ (théorème 28.47), on pourrait croire à la densité a fortiori dans $H^s(\mathbb{R}^d)$. Mais attention : $\mathcal{D}(\mathbb{R}^d)$ est dense dans L^2 pour la norme L^2 . Nous n'avons encore rien dit pour la norme $H^s(\mathbb{R}^d)$.

Proposition 32.15 ([447]).

La partie $\mathcal{S}(\mathbb{R}^d)$ est dense dans $H^s(\mathbb{R}^d)$.

Démonstration. Soit $u \in H^s(\mathbb{R}^d)$. Par définition l'application

$$\xi \mapsto (1 + \|\xi\|^2)^{s/2} \hat{u} \quad (32.59)$$

est dans $L^2(\mathbb{R}^d)$. Elle peut donc être approximée au sens L^2 par des fonctions dans $\mathcal{D}(\mathbb{R}^d)$ (théorème 28.47(5)), c'est-à-dire qu'il existe des fonctions $\phi_n \in \mathcal{D}(\mathbb{R}^d)$ telles que

$$\phi_n \xrightarrow{L^2(\mathbb{R}^d)} (1 + \|\xi\|^2)^{s/2} \hat{u}. \quad (32.60)$$

Nous posons

$$\psi_n = \frac{\phi_n}{(1 + \xi^2)^{s/2}} \quad (32.61)$$

Cela est encore une fonction de $\mathcal{D}(\mathbb{R}^d)$, et donc de $\mathcal{S}(\mathbb{R}^d)$. Vu que la transformée de Fourier est une bijection de $\mathcal{D}(\mathbb{R}^d)$ (proposition 30.14), nous pouvons considérer une suite $\varphi_n \in \mathcal{D}(\mathbb{R}^d)$ telle que $\hat{\varphi}_n = \psi_n$, et nous allons montrer que $\varphi_n \xrightarrow{H^s(\mathbb{R}^d)} u$.

Nous avons :

$$\|\varphi_n - u\|_{H^s}^2 = \int_{\mathbb{R}^d} |\hat{\varphi}_n - \hat{u}|^2 (1 + \xi^2)^s d\xi \quad (32.62a)$$

$$= \int_{\mathbb{R}^d} \left| \frac{\phi_n(\xi)}{(1 + \xi^2)^{s/2}} - \hat{u}(\xi) \right|^2 (1 + \xi^2)^s d\xi \quad (32.62b)$$

$$= \int_{\mathbb{R}^d} |\phi_n(\xi) - \hat{u}(\xi)(1 + \xi^2)^{s/2}|^2 d\xi \quad (32.62c)$$

$$= \|\phi_n - (1 + \xi^2)^{s/2} \hat{u}\|_{L^2}^2. \quad (32.62d)$$

Par définition de la suite ϕ_n nous avons donc bien

$$\|\varphi_n - u\|_{H^s}^2 = \|\phi_n - (1 + \xi^2)^{s/2} \hat{u}\|_{L^2}^2 \rightarrow 0. \quad (32.63)$$

Notons que même si ϕ_n est dans $\mathcal{D}(\mathbb{R}^d)$, nous n'avons pas prouvé la convergence $\phi_n \xrightarrow{H^s} u$, mais bien $\varphi_n \xrightarrow{H^s} u$. Or les fonctions φ_n sont dans $\mathcal{S}(\mathbb{R}^d)$, et rien n'assure qu'elles soient à support compact. Nous avons donc bien prouvé la densité de \mathcal{S} et non celle de \mathcal{D} . \square

Remarque 32.16.

Pour qui a tout compris, cela peut sembler une évidence, mais nous précisons que nous parlons de densité de $\mathcal{S}(\mathbb{R}^d)$ dans $H^s(\mathbb{R}^d)$, à aucun moment la topologie de $\mathcal{S}(\mathbb{R}^d)$ n'entre en compte.

Un peu moins évident : ce que nous avons réellement montré est la densité de $\iota(\mathcal{D}(\mathbb{R}^d))$ dans $H^s(\mathbb{R}^d)$ où ι est l'application « prise de classe ». Nous n'avons pas insisté là-dessus, mais il faut

dire que dans la preuve de la proposition 32.15, u est un représentant d'un élément choisi dans $H^s(\mathbb{R}^d)$.

Nous avons ensuite prouvé la convergence $\|\varphi_n - u\|_{H^s(\mathbb{R}^d)} \rightarrow 0$ qui est une convergence d'une suite dans \mathbb{R} , et dans laquelle l'opération $\|\cdot\|_{H^s}$ est définie sur un espace de fonctions et n'est pas une norme (c'est pour que cela devienne une norme que l'on prend les classes).

Nous en avons déduit la convergence $\varphi_n \xrightarrow{H^s(\mathbb{R}^d)} u$ où maintenant φ_n et u sont des classes dans $H^s(\mathbb{R}^d)$.

Proposition 32.17.

La partie $\mathcal{D}(\mathbb{R}^d)$ est dense dans $(H^s(\mathbb{R}^d), \|\cdot\|_{H^s(\mathbb{R}^d)})$.

Démonstration. Nous savons déjà que $\mathcal{D}(\mathbb{R}^d)$ est dense dans $H^s(\mathbb{R}^d)$ par la proposition 32.15. Nous devons seulement prouver que $\mathcal{D}(\mathbb{R}^d)$ est dense dans $(\mathcal{S}(\mathbb{R}^d), \|\cdot\|_{H^2(\mathbb{R}^d)})$. Pour cela nous utilisons la densité de $\mathcal{D}(\mathbb{R}^d)$ dans $\mathcal{S}(\mathbb{R}^d)$ de la proposition 28.153. Soit donc $f \in \mathcal{D}(\mathbb{R}^d)$ et une suite f_k dans $\mathcal{D}(\mathbb{R}^d)$ telle que

$$f_k \xrightarrow{\mathcal{S}(\mathbb{R}^d)} f. \quad (32.64)$$

Vu que la transformée de Fourier est continue sur $\mathcal{S}(\mathbb{R}^d)$ (proposition 30.14) nous avons aussi

$$\hat{f}_k \xrightarrow{\mathcal{S}(\mathbb{R}^d)} \hat{f}, \quad (32.65)$$

et en particulier pour tout polynôme P nous avons la convergence uniforme

$$P\hat{f}_k \xrightarrow{unif} P\hat{f}. \quad (32.66)$$

D'autre part la fonction $\xi \mapsto |\hat{f}_k(\xi) - \hat{f}(\xi)|^2(1 + \xi^2)^s$ est Schwartz et en tout point décroissante en k . Soient $\epsilon > 0$ et $r > 0$ choisis de telle sorte à avoir

$$\int_{\|\xi\|>r} |\hat{f}_k(\xi) - \hat{f}(\xi)|(1 + \xi^2)^s d\xi < \epsilon. \quad (32.67)$$

pour tout k . La convergence uniforme (32.66) permet de considérer k_0 tel que pour tout $k > k_0$,

$$|\hat{f}_k - \hat{f}|(1 + \xi^2)^s < \frac{\epsilon}{\text{Vol}(B(0, r))} \quad (32.68)$$

dans $B(0, r)$. Avec tout cela, dès que $k > k_0$ nous avons

$$\|f_k - f\|_{H^s(\mathbb{R}^d)} = \int_{\mathbb{R}} |\hat{f}_k - \hat{f}|(1 + \xi^2)^s d\xi = \int_{B(0, r)} \dots + \int_{\|\xi\|>r} \dots \leq 2\epsilon. \quad (32.69)$$

Donc nous avons bien $\|f_k - f\|_{H^s(\mathbb{R}^d)} \rightarrow 0$ et convergence de f_k vers f dans $H^s(\mathbb{R}^d)$. \square

32.2 Trace

Définition 32.18 ([446]).

Nous définissons la **trace** d'une fonction par

$$\begin{aligned} \gamma_0: \mathcal{D}(\mathbb{R}^d) &\rightarrow \mathcal{D}(\mathbb{R}^{d-1}) \\ (\gamma_0 v)(x_1, \dots, x_{d-1}) &= v(x_1, \dots, x_{d-1}, 0). \end{aligned} \quad (32.70)$$

Théorème 32.19 ([448, 446]).

Si $s > \frac{1}{2}$, alors γ_0 accepte une unique extension en opérateur linéaire borné

$$\gamma_0: H^s(\mathbb{R}^d) \rightarrow H^{s-\frac{1}{2}}(\mathbb{R}^{d-1}). \quad (32.71)$$

Démonstration. Nous subdivisons la preuve en plusieurs pas.

Une inégalité pour $\varphi \in \mathcal{D}(\mathbb{R}^d)$ Nous commençons par considérer $v \in \mathcal{D}(\mathbb{R}^d)$ (fonction C^∞ à support compact). Nous allons alors prouver que

$$\|\gamma_0 \varphi\|_{H^{s-\frac{1}{2}}(\mathbb{R}^{d-1})} \leq K \|\varphi\|_{H^s(\mathbb{R}^d)} \tag{32.72}$$

pour une certaine constante K (qui ne dépend en particulier pas de φ).

Nous avons

$$\|\gamma_0 \varphi\|_{H^{s-\frac{1}{2}}(\mathbb{R}^{d-1})}^2 = (\gamma_0 \varphi, \gamma_0 \varphi)_{H^{s-\frac{1}{2}}} \tag{32.73a}$$

$$= \int_{\mathbb{R}^{d-1}} |\widehat{\gamma_0 \varphi}(\xi)|^2 (1 + \|\xi\|^2)^{s-\frac{1}{2}} d\xi \tag{32.73b}$$

$$= \int_{\mathbb{R}^{d-1}} \left| \int_{\mathbb{R}^{d-1}} (\gamma_0 \varphi)(x) e^{-i\xi x} dx \right|^2 (1 + \|\xi\|^2)^{s-\frac{1}{2}} d\xi \tag{32.73c}$$

$$\tag{32.73d}$$

Nous appliquons la trace en appliquant la formule du corollaire 30.23,

$$(\gamma_0 \varphi)(x) = \varphi(x, 0) = \frac{1}{2\pi} \int_{\mathbb{R}} \int_{\mathbb{R}} e^{-iky} \varphi(x, y) dy dk \tag{32.74}$$

En remplaçant dans (32.73c) nous avons

$$\|\gamma_0 \varphi\|_{H^{s-\frac{1}{2}}(\mathbb{R}^{d-1})}^2 = \frac{1}{2\pi} \int_{\mathbb{R}^{d-1}} \left| \int_{\mathbb{R}^{d-1}} \int_{\mathbb{R}} e^{-iky} e^{-i\xi x} \varphi(x, y) dy dk dx \right|^2 (1 + \|\xi\|^2)^{s-\frac{1}{2}} d\xi. \tag{32.75}$$

Nous voudrions permuter les intégrales en k et en x . Pour cela nous étudions la fonction $u: \mathbb{R} \times \mathbb{R}^{d-1} \rightarrow \mathbb{C}$ donnée par

$$u(k, x) = e^{-ikx} \int_{\mathbb{R}} e^{-i\xi x} \varphi(x, y) dy \tag{32.76}$$

Effectuer l'intégrale par rapport à y revient à calculer la transformée de Fourier partielle dont nous parlons dans la proposition 30.15⁴. Elle est donc une fonction Schwartz de k et de x (conjointement et non seulement séparément) et est donc dans $L^1(\mathbb{R} \times \mathbb{R}^{d-1})$. Les intégrales sur k et sur x peuvent donc être réunies et permutées par le théorème de Fubini 15.259 (n'oubliez tout de même pas de vous convaincre que la condition (2) est remplie).

Nous avons donc

$$\|\gamma_0 \varphi\|_{H^{s-\frac{1}{2}}(\mathbb{R}^{d-1})}^2 = \frac{1}{2\pi} \int_{\mathbb{R}^{d-1}} \left| \int_{\mathbb{R}} \int_{\mathbb{R}^{d-1}} \int_{\mathbb{R}} e^{-iky} e^{-i\xi x} \varphi(x, y) dy dx dk \right|^2 (1 + \|\xi\|^2)^{s-\frac{1}{2}} d\xi. \tag{32.77}$$

Étant donné que φ est à support compact, les intégrales sur x et sur y peuvent se réunir en utilisant encore le théorème de Fubini ; ces intégrales donnent :

$$\int_{\mathbb{R}^{d-1} \times \mathbb{R}} e^{-iky} e^{-i\xi x} \varphi(x, y) dx \otimes dy = \int_{\mathbb{R}^{d-1} \times \mathbb{R}} e^{-i(\xi, k) \cdot (x, y)} \varphi(x, y) dx \otimes dy = \hat{\varphi}(\xi, k). \tag{32.78}$$

Nous restons avec

$$\|\gamma_0 \varphi\|_{H^{s-\frac{1}{2}}(\mathbb{R}^{d-1})}^2 = \frac{1}{2\pi} \int_{\mathbb{R}^{d-1}} \left| \int_{\mathbb{R}} \hat{\varphi}(\xi, k) dk \right|^2 (1 + \|\xi\|^2)^{s-\frac{1}{2}} d\xi. \tag{32.79}$$

Nous allons maintenant traiter la partie du milieu :

$$\clubsuit = \left| \int_{\mathbb{R}} \hat{\varphi}(\xi, k) dk \right| = \left| \int_{\mathbb{R}} \hat{\varphi}(\xi, k) (1 + \xi^2 + k^2)^{s/2} \frac{1}{(1 + \xi^2 + k^2)^{s/2}} dk \right| = |\langle f_1, f_2 \rangle_{L^2(\mathbb{R}^d)}| \tag{32.80}$$

4. Dont une relecture de la preuve ne serait vraiment pas de trop, ainsi que la preuve de 28.147.

Ici ξ est vu comme une constante et les fonctions f_1 et f_2 sont

$$f_1: k \rightarrow \hat{\varphi}(\xi, k)(1 + \xi^2 + k^2)^{s/2} \quad (32.81a)$$

$$f_2: k \rightarrow \frac{1}{(1 + \xi^2 + k^2)^{s/2}} \quad (32.81b)$$

Nous pouvons utiliser l'inégalité de Cauchy-Schwarz 11.10 :

$$\clubsuit \leq \left(\int_{\mathbb{R}} |\hat{\varphi}(\xi, k)|^2 (1 + \xi^2 + k^2)^s dk \right)^{1/2} \left(\int_{\mathbb{R}} \frac{1}{(1 + \xi^2 + k^2)^s} dk \right)^{1/2} \quad (32.82)$$

Nous notons $g(\xi)$ ce qui se trouve dans la seconde parenthèse (après intégration sur k). Avec cela nous continuons :

$$\|\gamma_0 \varphi\|_{H^{s-\frac{1}{2}}(\mathbb{R}^{d-1})}^2 \leq \frac{1}{2\pi} \int_{\mathbb{R}^{d-1}} \int_{\mathbb{R}} |g(\xi)| |\hat{\varphi}(\xi, k)|^2 (1 + \xi^2 + k^2)^s (1 + \|\xi\|^2)^{s-\frac{1}{2}} dk d\xi. \quad (32.83)$$

Vu que $\hat{\varphi}$ est Schwartz, la fonction qui est à l'intérieur des deux intégrales est dans $L^1(\mathbb{R}^{d-1} \times \mathbb{R})$ et nous pouvons réunir les deux intégrales :

$$\|\gamma_0 \varphi\|_{H^{s-\frac{1}{2}}(\mathbb{R}^{d-1})}^2 \leq \frac{1}{2\pi} \int_{\mathbb{R} \times \mathbb{R}^{d-1}} |g(\xi)| |\hat{\varphi}(\xi, k)|^2 (1 + \xi^2 + k^2)^s (1 + \|\xi\|^2)^{s-\frac{1}{2}} dk \otimes d\xi. \quad (32.84)$$

À ce point nous démontrons qu'en réalité la combinaison $g(\xi)(1 + \xi^2)^{s-\frac{1}{2}}$ ne dépend pas de ξ . En effet

$$g(\xi)(1 + \xi^2)^{s-\frac{1}{2}} = (1 + \xi^2)^{s-\frac{1}{2}} \int_{\mathbb{R}} \frac{1}{(1 + \xi^2 + k^2)} dk \quad (32.85a)$$

$$= \frac{1}{(1 + \xi^2)^{1/2}} \int_{\mathbb{R}} \left(\frac{1 + \xi^2}{1 + \xi^2 + k^2} \right)^s dk \quad (32.85b)$$

$$= \frac{1}{(1 + \xi^2)^{1/2}} \int \left(\frac{1}{1 + \frac{k^2}{1+\xi^2}} \right)^s dk. \quad (32.85c)$$

Nous effectuons le changement de variables $t = \frac{k}{\sqrt{1+\xi^2}}$, $dk = (1 + \xi^2)^{1/2} dt$, et le tout vaut

$$\int_{\mathbb{R}} \left(\frac{1}{1 + t^2} \right)^s dt, \quad (32.86)$$

qui est effectivement indépendant de ξ . Nous nommons cela K (auquel nous ajoutons le $\frac{1}{2\pi}$) :

$$\|\gamma_0 \varphi\|_{H^{s-\frac{1}{2}}(\mathbb{R}^{d-1})}^2 \leq K \int_{\mathbb{R} \times \mathbb{R}^{d-1}} |\hat{\varphi}(\xi, k)|^2 (1 + \xi^2 + k^2)^s dk \otimes d\xi = K \|\varphi\|_{H^s(\mathbb{R}^d)}^2. \quad (32.87)$$

Nous avons donc prouvé pour tout $\varphi \in \mathcal{D}(\mathbb{R}^d)$ (avec redéfinition du K) :

$$\|\gamma_0 \varphi\|_{H^{s-\frac{1}{2}}(\mathbb{R}^{d-1})} \leq K \|\varphi\|_{H^s(\mathbb{R}^d)}. \quad (32.88)$$

À propos de classes Il serait tentant de conclure en disant que $\mathcal{D}(\mathbb{R}^d)$ est dense dans $H^s(\mathbb{R}^d)$.

Hélas, **techniquement**, l'ensemble $\mathcal{D}(\mathbb{R}^d)$ n'est même pas un sous-ensemble de $H^s(\mathbb{R}^d)$ parce que ce dernier est un ensemble de *classes* de fonctions. Ce petit détail a ici son importance parce que γ_0 n'est pas une application qui descend aux classes. En effet, \mathbb{R}^{d-1} étant de mesure nulle dans \mathbb{R}^d , deux fonctions de la même classe peuvent différer en *tous* les points de \mathbb{R}^{d-1} en même temps.

Si nous notons ι l'application qui consiste à prendre la classe ce qui est dense dans $H^s(\mathbb{R}^d)$, c'est $\iota(\mathcal{D}(\mathbb{R}^d))$. Or chaque classe contient au maximum une seule fonction continue (qui sera même de classe C^∞ à support compact pour les éléments de $\iota(\mathcal{D})$).

L'application γ_0 considérée est l'application composée entre le γ_0 classique et le choix du représentant continu dans la classe. La formule (32.88) que nous venons de prouver est valide pour l'application γ_0 vue comme

$$\gamma_0: \iota(\mathcal{D}(\mathbb{R}^d)) \rightarrow \mathcal{D}(\mathbb{R}^{d-1}). \quad (32.89)$$

Densité et conclusion Ce que la majoration (32.88) prouve est la continuité de l'application

$$\gamma_0: (\mathcal{D}(\mathbb{R}^d), \|\cdot\|_{H^s(\mathbb{R}^d)}) \rightarrow (H^{s-\frac{1}{2}}(\mathbb{R}^{d-1}), \|\cdot\|_{H^{s-\frac{1}{2}}(\mathbb{R}^{d-1})}). \quad (32.90)$$

Mais la proposition 32.17 nous donne la densité de la partie $\mathcal{D}(\mathbb{R}^d)$ dans $H^s(\mathbb{R}^d)$. La proposition 18.121 nous donne alors une extension

$$\gamma_0: (H^s(\mathbb{R}^d), \|\cdot\|_{H^s(\mathbb{R}^d)}) \rightarrow (H^{s-\frac{1}{2}}(\mathbb{R}^{d-1}), \|\cdot\|_{H^{s-\frac{1}{2}}(\mathbb{R}^{d-1})}). \quad (32.91)$$

□

Remarque 32.20.

L'extension n'est pas évidente parce que les éléments de $H^s(\mathbb{R}^d)$ sont en général des classes de fonctions dont les valeurs sur le bord ne sont pas du tout fixées du fait que le bord soit de mesure nulle.

32.3 Théorème de plongement

L'objet des théorèmes de plongement de Sobolev est de montrer que si $s > \frac{d}{2} + k$ alors les éléments de $H^s(\mathbb{R}^d)$ possèdent des représentants de classe C^k . Avant de démontrer le théorème, pour alléger, nous allons donner deux lemmes.

Lemme 32.21.

Soit (u_j) une suite dans $\mathcal{S}(\mathbb{R}^d)$ telle que

$$u_j \xrightarrow{H^s(\mathbb{R}^d)} u \quad (32.92)$$

avec $s > 0$. Alors nous avons aussi la convergence

$$u_j \xrightarrow{L^2(\mathbb{R}^d)} u. \quad (32.93)$$

Démonstration. Vu que $s > 0$ nous avons $(1 + k^2)^s > 1$ (ici nous écrivons k^2 pour $\|k\|^2$). Par conséquent

$$(u, v)_{H^s(\mathbb{R}^d)} = \int_{\mathbb{R}^d} \hat{u} \bar{\hat{v}} (1 + k^2)^s dk \geq \int_{\mathbb{R}^d} \hat{u} \bar{\hat{v}} dk = \langle \hat{u}, \hat{v} \rangle_{L^2(\mathbb{R}^d)}. \quad (32.94)$$

Nous avons alors

$$\|u_j - u\|_{L^2} = \frac{1}{(2\pi)^d} \|\hat{u}_j - \hat{u}\|_{L^2} \quad (32.95a)$$

$$= \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} |\hat{u}_j - \hat{u}|^2 \quad (32.95b)$$

$$\leq \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} |\hat{u}_j - \hat{u}|^2 (1 + k^2)^s dk \quad (32.95c)$$

$$= \frac{1}{(2\pi)^d} \|u_j - u\|_{H^s(\mathbb{R}^d)}. \quad (32.95d)$$

□

Lemme 32.22.

Soient des fonctions $u_j \in \mathcal{S}(\mathbb{R}^d)$ telles que

$$u_j \xrightarrow{(C_0^0(\mathbb{R}^d), \|\cdot\|_\infty)} v. \quad (32.96)$$

Alors nous avons la convergence

$$\int_{\mathbb{R}^d} u_j \varphi \rightarrow \int_{\mathbb{R}^d} v \varphi \quad (32.97)$$

pour tout $\varphi \in \mathcal{S}(\mathbb{R}^d)$.

Démonstration. La suite (u_j) est équibornée. En effet il existe une queue de suite pour laquelle $\|u_j - v\|_\infty < \epsilon$; cette queue de suite est alors équibornée par $\|v\|_\infty + \epsilon$. Le début de la suite est un nombre fini de fonctions, toutes bornées. Le maximum des bornes donne alors une borne.

Soit donc $M > 0$ tel que $|u_j(x)| < M$ pour tout $x \in \mathbb{R}^d$ et pour tout $j \in \mathbb{N}$. Nous avons alors $|u_j\varphi| < M|\varphi|$ pour tout j et les fonctions $|u_j\varphi|$ sont majorées par la fonction $M|\varphi|$ qui est intégrable. Nous pouvons donc utiliser le théorème de la convergence dominée de Lebesgue 15.184 nous donne

$$\lim_{j \rightarrow \infty} \int_{\mathbb{R}^d} u_j \varphi = \int_{\mathbb{R}^d} v \varphi. \quad (32.98)$$

□

Nous pouvons écrire la conclusion du lemme 32.22 sous la forme

$$\langle u_j, \varphi \rangle_{L^2(\mathbb{R}^d)} \rightarrow \langle v, \varphi \rangle_{L^2(\mathbb{R}^d)} \quad (32.99)$$

pour tout $\varphi \in \mathcal{S}(\mathbb{R}^r)$ (et non pour tout $\varphi \in L^2(\mathbb{R}^d)$).

Théorème 32.23 (Théorème de Sobolev avec $k = 0$ [449]).

Soit $s > \frac{d}{2}$ et $u \in H^s(\mathbb{R}^d)$. Alors u possède un représentant dans $C_0^0(\mathbb{R}^d)$ (les fonctions continues et qui s'annulent à l'infini). Nous écrivons cela $H^s(\mathbb{R}^d) \subset C_0^0(\mathbb{R}^d)$.

Démonstration. Nous commençons par supposer que $u \in H^s(\mathbb{R}^d) \cap \mathcal{S}(\mathbb{R}^d)$, et dans ce cas nous notons u le représentant dans $\mathcal{S}(\mathbb{R}^d)$. Nous allons prouver l'inégalité

$$\|u\|_\infty \leq c \|u\|_{H^s(\mathbb{R}^d)}. \quad (32.100)$$

La formule d'inversion de Fourier 30.22 appliquée à u_j donne

$$u(x) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} e^{ikx} \hat{u}(k) dk, \quad (32.101)$$

nous avons alors

$$(2\pi)^d |u(x)| \leq \int_{\mathbb{R}^d} |\hat{u}(k)| dk \quad (32.102a)$$

$$= \int_{\mathbb{R}^d} (1+k)^{s/2} |\hat{u}(k)| (1+k^2)^{-s+2} dk \quad (32.102b)$$

$$= \int_{\mathbb{R}^d} \underbrace{(|\hat{u}(k)|^2 (1+k^2)^s)^{1/2}}_f \underbrace{((1+k^2)^{-s})^{1/2}}_g dk \quad (32.102c)$$

$$= \langle f, g \rangle_{L^2(\mathbb{R}^d)}. \quad (32.102d)$$

Ici il convient nous arrêter un instant pour nous convaincre que f et g sont réellement des éléments de L^2 . En ce qui concerne f , c'est facile : \hat{u} est une fonction Schwartz. En ce qui concerne g il faut l'intégrabilité de $|g|^2$, c'est-à-dire de $k \mapsto (1+k^2)^{-s}$. Cela a lieu si et seulement si $2s > n$ et donc a lieu dans les hypothèses du théorème. Nous utilisons le théorème de Cauchy-Schwarz⁵ pour continuer :

$$(2\pi)^d |u(x)| \leq \|f\|_{L^2} \|g\|_{L^2} \quad (32.103a)$$

$$= c \left(\int_{\mathbb{R}^d} |\hat{u}(k)|^2 (1+k^2)^s dk \right)^{1/2} \quad (32.103b)$$

$$= c \|u\|_{H^s(\mathbb{R}^d)}. \quad (32.103c)$$

Donc en introduisant le facteur $(2\pi)^d$ dans la constante c nous avons

$$\|u\|_\infty \leq c \|u\|_{H^s(\mathbb{R}^d)}. \quad (32.104)$$

5. Formule 11.9.

Cela est tout ce que nous voulions faire avec $u \in \mathcal{S}(\mathbb{R}^d)$.

Nous considérons maintenant $u \in H^s(\mathbb{R}^d)$. Vu que la densité des fonctions Schwartz dans H^s (proposition 32.15) nous considérons une suite (u_j) dans $\mathcal{S}(\mathbb{R}^d)$ telle que

$$u_j \xrightarrow{H^s(\mathbb{R}^d)} u \quad (32.105)$$

Ici u est une classe, mais nous identifions u_j avec sa classe (parce qu'il ne faut pas exagérer non plus). La suite (u_j) est de Cauchy dans H^s , donc si $\epsilon > 0$ est donné, il existe N tel que si $n, m > N$, $\|u_m - u_n\| \leq \epsilon$. Nous avons alors aussi

$$\|u_m - u_n\|_\infty \leq c\epsilon, \quad (32.106)$$

ce qui signifie que (u_j) est également une suite de Cauchy dans $(C^0(\mathbb{R}^d), \|\cdot\|_\infty)$ qui est un espace complet par la proposition 13.282.

Il existe donc une fonction $v \in C_0^0(\mathbb{R}^d)$ telle que

$$u_j \xrightarrow{(C_0^0(\mathbb{R}^d), \|\cdot\|_\infty)} v. \quad (32.107)$$

La question est de savoir si nous pouvons déduire que v est un représentant de u .

Par le lemme 32.21 nous avons également la convergence

$$u_j \xrightarrow{L^2(\mathbb{R}^d)} u. \quad (32.108)$$

Pour tout $\varphi \in \mathcal{S}(\mathbb{R}^d)$ nous avons alors

$$\langle u_j, \varphi \rangle_{L^2} \rightarrow \langle u, \varphi \rangle_{L^2}. \quad (32.109)$$

Mais en même temps, la convergence (32.107) couplée au lemme 32.22 donne également

$$\langle u_j, \varphi \rangle_{L^2} \rightarrow \langle v, \varphi \rangle_{L^2}. \quad (32.110)$$

Par unicité de la limite (dans \mathbb{R}) nous avons

$$\langle v, \varphi \rangle_{L^2} = \langle u, \varphi \rangle_{L^2} \quad (32.111)$$

pour tout $\varphi \in \mathcal{S}(\mathbb{R}^d)$. La proposition 31.1 appliquée à $u - v$ montre alors que $u - v = 0$ presque partout, c'est-à-dire que v est bien un représentant de u .

Le représentant v de u est non seulement continu (comme limite uniforme de fonctions continues), mais également bornée, comme limite uniforme de fonctions Schwartz. \square

Proposition 32.24 ([449]).

Si $u \in H^s(\mathbb{R}^d)$ ($s \in \mathbb{R}$) alors

$$(1) \partial^\alpha u \in H^{s-|\alpha|}(\mathbb{R}^d),$$

(2) l'application

$$\partial^\alpha : H^s(\mathbb{R}^d) \rightarrow H^{s-|\alpha|}(\mathbb{R}^d) \quad (32.112)$$

est continue.

Note : ici ∂ est l'opération de dérivée faible.

Démonstration. Nous allons seulement prouver que $\partial_j : H^s(\mathbb{R}^d) \rightarrow H^{s-1}(\mathbb{R}^d)$ est bien définie⁶ et continue. Par composition, la thèse suivra.

Soit $u \in H^s(\mathbb{R}^d)$ par le lemme 32.10 nous avons

$$\widehat{\partial_j u} = i\xi_j \hat{u}. \quad (32.113)$$

6. Au sens où l'espace d'arrivée est bien celui-là.

D'autre part, la fonction

$$\begin{aligned} f: \mathbb{R}^n &\rightarrow \mathbb{R} \\ x &\mapsto \frac{x_i}{1 + \|x\|^2} \end{aligned} \quad (32.114)$$

est bornée (et même indépendamment de i) par une constante K . Donc nous avons pour tout ⁷ s :

$$k_i(1 + \|k\|^2)^{-s} < K(1 + \|k\|^2)^{-s+1}. \quad (32.115)$$

Avec cela nous pouvons calculer un peu : si $u \in H^s(\mathbb{R}^d)$, nous avons

$$\|\partial_j u\|_{H^{s-1}(\mathbb{R}^d)} = \int_{\mathbb{R}^d} |\widehat{\partial_j u}| (1 + k^2)^{s-1} dk \quad (32.116a)$$

$$= \int_{\mathbb{R}^d} k_j |\hat{u}| (1 + k^2)^{s-1} dk \quad (32.116b)$$

$$\leq \int_{\mathbb{R}^d} K |\hat{u}| (1 + k^2)^s dk \quad (32.116c)$$

$$= K \|u\|_{H^s(\mathbb{R}^d)}. \quad (32.116d)$$

Nous avons donc que $\|\partial_j u\|_{H^{s-1}(\mathbb{R}^d)}$ est fini lorsque $u \in H^s(\mathbb{R}^d)$.

La majoration $\|\partial_j u\| \leq K \|u\|$ donne la majoration suivante pour la norme de l'opérateur ∂_j :

$$\|\partial_j\| = \sup_{\|u\|_{H^s}=1} \|\partial_j u\|_{H^{s-1}} \leq K. \quad (32.117)$$

Le fait d'être borné implique d'être continu par la proposition 12.25. □

Théorème 32.25 (Théorème de plongement de Sobolev [449]).

Soient $k \in \mathbb{N}$ et $m > \frac{d}{2} + k$. Alors

$$H^s(\mathbb{R}^d) \subset C_0^k(\mathbb{R}^d). \quad (32.118)$$

Remarques :

- L'espace $C_0^k(\mathbb{R}^d)$ est l'ensemble des fonctions de classe C^k qui s'annulent à l'infini.
- L'inclusion (32.118) signifie que tout élément dans H^s possède un représentant dans $C_0^k(\mathbb{R}^d)$.

Démonstration. Pour $k = 0$, c'est le théorème 32.23. Si $|\alpha| < k$ nous savons que $\partial^\alpha u \in H^{s-k} \subset C_0^0(\mathbb{R}^d)$. Cela signifie que les dérivées faibles sont continues, mais pas qu'il existe un représentant qui est réellement k fois continument dérivable.

Soit $u \in H^s(\mathbb{R}^d)$ et une suite (u_j) dans $\mathcal{S}(\mathbb{R}^d)$ telle que

$$u_j \xrightarrow{H^s(\mathbb{R}^d)} u. \quad (32.119)$$

Vu que l'espace topologique $(C_0^k(\mathbb{R}^d), \|\cdot\|_\infty)$ est complet il existe $v \in C_0^k$ tel que

$$u_j \xrightarrow{C_0^k} v. \quad (32.120)$$

Il reste à montrer que v est un représentant de u . Cela se fait comme plus haut en montrant que $u_j \xrightarrow{L^2} u$. □

7. Question : dans [449], il faut dépendre cette constante de s . Je ne comprends pas pourquoi.

Chapitre 33

Équations différentielles ordinaires

Une équation différentielle ordinaire est la recherche de toutes les fonctions définie sur une partie de \mathbb{R} satisfaisant à une certaine égalité, faisant intervenir les dérivées de la fonction recherchée.

Dans la suite, I désignera un intervalle de \mathbb{R} . Une fonction sera **dérivable sur I** si elle est dérivable au sens usuel sur l'intérieur de I , et si elle est dérivable à droite (resp. à gauche) sur l'éventuel bord gauche (resp. droit) de I .

Définition 33.1.

Une **équation différentielle ordinaire d'ordre n sur I** est la recherche d'une fonction $y : I \rightarrow \mathbb{R}$ dérivable n fois, satisfaisant à une équation du type

$$F(t, y(t), y'(t), \dots, y^{(n)}(t)) = 0 \quad \text{pour tout } t \in I \quad (33.1)$$

où I est un intervalle de \mathbb{R} et $F : (I \times D) \subset (\mathbb{R} \times \mathbb{R}^{n+1}) \rightarrow \mathbb{R}$ est une fonction donnée.

Remarque 33.2.

L'équation différentielle (33.1) sera raccourcie sous la forme

$$F(t, y, y', \dots, y^{(n)}) = 0 \quad (33.2)$$

où la dépendance en t est sous-entendue.

Exemple 33.3

Soit $f : I \rightarrow \mathbb{R}$ une fonction continue fixée. L'équation différentielle

$$y' = f(t) \quad (33.3)$$

se ramène à la recherche des primitives de f sur l'intervalle I . △

Le lemme suivant sert de temps en temps.

Lemme 33.4 (Lemme de Grönwall).

Soient ϕ et ψ deux fonctions telles que pour tout $t \in [t_0, t_1]$, $\phi(t) \geq 0$, $\psi(t) \geq 0$ et

$$\phi(t) \leq K + L \int_{t_0}^t \psi(s)\phi(s)ds \quad (33.4)$$

où K et L sont des constantes positives. Alors

$$\phi(t) \leq K \exp\left(L \int_{t_0}^t \psi\right). \quad (33.5)$$

Lemme 33.5 (Lemme de Grönwall[450]).

Si $u, a, b \in C^0([0, T], \mathbb{R}^+)$ sont telles que

$$u(t) \leq b(t) + \int_0^t a(s)u(s)ds \quad (33.6)$$

pour tout $t \in [0, T]$ alors pour tout $t \in [0, T]$ nous avons aussi

$$u(t) \leq b(t) + \int_0^t b(s)a(s)e^{\int_s^t a(u)du} ds. \quad (33.7)$$

33.1 Équation homogène, solution particulière

Voici un petit morceau d'algèbre linéaire. Soient des espaces vectoriels V et W ainsi qu'une application linéaire $D: V \rightarrow W$. Nous voulons résoudre $D(u) = v$, c'est-à-dire déterminer l'ensemble

$$D^{-1}(v) = \{u \in V \text{ tel que } Du = v\}. \quad (33.8)$$

Lemme 33.6.

Soient des espaces vectoriels V et W ainsi qu'une application linéaire $D: V \rightarrow W$. Si $u_P \in V$ satisfait à $Du_P = v$ alors

$$D^{-1}(v) = \ker(D) + u_P. \quad (33.9)$$

Démonstration. Si $u \in \ker(D) + u_P$ alors $u = k + u_P$ avec $Dk = 0$, ce qui donne tout de suite $Du = Dk + Du_P = v$. Donc $u \in D^{-1}(v)$.

Dans l'autre sens, si $u \in D^{-1}(v)$ alors nous pouvons écrire $u = (u - u_P) + u_P$. Vu que $u - u_P \in \ker(D)$ nous avons bien $u \in \ker(D) + u_P$. \square

Ce petit lemme explique pourquoi la résolution d'équation différentielles passe par le principe « générale de l'homogène plus particulière de la non-homogène ». Cela marche autant pour les équations différentielles ordinaires que pour celles aux dérivées partielles.

Exemple 33.7

Considérons l'équation différentielle ordinaire

$$y' - y = 4. \quad (33.10)$$

L'opérateur dont nous parlons est par exemple

$$D: C^\infty(\mathbb{R}) \rightarrow C^\infty(\mathbb{R})y \mapsto y' - y. \quad (33.11)$$

Nous devons résoudre $Dy = 4$ où « 4 » est l'élément fonction constante égale à 4 dans $C^\infty(\mathbb{R})$. L'ensemble $\ker(D)$ sont les éléments $y \in C^\infty(\mathbb{R})$ tels que $y' = y$:

$$\ker(D) = \{t \mapsto Ke^t \text{ tel que } K \in \mathbb{R}\}. \quad (33.12)$$

Nous devons trouver un élément quelconque y_P de $D^{-1}(4)$. Facile : $y_P(t) = -4$.

Au final,

$$D^{-1}(4) = \{t \mapsto Ke^t - 4 \text{ tel que } K \in \mathbb{R}\}. \quad (33.13)$$

\triangle

Dans cet exemple nous avons pris $V = W = C^\infty(\mathbb{R})$. Mais souvent nous sommes amenés à considérer des espaces plus subtils, parce qu'il existe simplement pas de solutions dans C^∞ , ou alors parce que beaucoup de solutions n'y sont pas.

33.2 Que faire avec $f(z)dz = g(t)dt$?

Dans de nombreux exercices d'équations différentielles, nous tombons sur $u' = f(t)$, et nous faisons formellement

$$\frac{du}{dt} = f(t) \Rightarrow du = f(t)dt, \quad (33.14)$$

et ensuite, il y a la formule un peu magique

$$u - u_0 = \int_{t_0}^t f(t)dt. \quad (33.15)$$

Voyons ce qu'il en est. Tout d'abord, il faut comprendre ce que signifie la formule

$$f(z)dz = g(t)dt. \quad (33.16)$$

Il s'agit d'une égalité entre deux formes différentielles sur \mathbb{R} où z est une fonction de t . Étant donné que z est une fonction de t , il faut voir dz comme la différentielle de cette fonction. La différentielle d'une fonction à une variable est donné par la dérivée :

$$dz_t = z'(t)dt \quad (33.17)$$

Écrire l'équation (33.16) pour chaque t revient donc à écrire

$$f(z(t))z'(t)dt = g(t)dt \quad (33.18)$$

Cela est une égalité entre deux formes différentielles. Nous avons donc égalité entre les intégrales des formes sur un chemin. Prenons un chemin tout simple de t_0 vers t :

$$\int_{t_0}^t f(z(t))z'(t)dt = \int_{t_0}^t g(t)dt. \quad (33.19)$$

Dans le premier membre, nous faisons un changement de variable $\xi = z(t)$, $d\xi = z'(t)dt$, et nous obtenons

$$\int_{z_0}^{z(t)} f(\xi)d\xi = \int_{t_0}^t g(t)dt. \quad (33.20)$$

où nous avons remplacé la constante $z(t_0)$ par z_0 dans la borne d'intégration. Si F est une primitive de f et G une primitive de g , nous avons

$$F(z) - F(z_0) = G(t) - G(t_0). \quad (33.21)$$

Si aucun problème de Cauchy n'est donné, les constantes $F(z_0)$ et $G(t_0)$ sont mises en une seule et nous écrivons la solution

$$F(z(t)) = G(t) + C, \quad (33.22)$$

qui est une équation implicite pour $z(t)$.

Nous trouvons assez souvent le cas simple

$$f(z)dz = dt. \quad (33.23)$$

En remplaçant $g(t) = 1$ dans (33.20), nous trouvons la fameuse

$$t - t_0 = \int_{z_0}^z f(z)dz, \quad (33.24)$$

dans laquelle il y a un abus de notation terrible entre le z de la borne (que les étudiants oublient souvent) et la variable d'intégration z !!

Le passage de (33.23) à (33.24) sera très souvent utilisé dans le cours de mécanique par exemple.

33.3 Équations linéaires du premier ordre

Une **équation différentielle linéaire** est une équation de la forme

$$y' + u(t)y = v(t). \quad (33.25)$$

Exemple 33.8

Tant qu'il n'y a pas de second membre, c'est facile. Prenons l'exemple suivant :

$$y' + 2ty = 0. \quad (33.26)$$

Nous mettons tous les t d'un côté et tous les y et y' de l'autre :

$$\frac{y'}{y} = -2t, \quad (33.27)$$

et puis on intègre sans oublier la constante d'intégration :

$$\ln(y) = -t^2 + C, \quad (33.28)$$

et donc $y(t) = Ke^{-t^2}$. △

Exemple 33.9

Lorsqu'il y a un second membre, il y a une astuce. Prenons par exemple

$$y' + 2ty = 4t. \quad (33.29)$$

L'astuce est de commencer par résoudre l'équation sans le second membre (l'équation homogène associée). Nous notons y_H la solution. Ici, la réponse est

$$y_H(t) = Ke^{-t^2}. \quad (33.30)$$

Ensuite le truc est d'essayer de trouver la solution de l'équation (33.29) sous la forme

$$y(t) = K(t)e^{t^2}. \quad (33.31)$$

L'idée est de prendre la même que la solution de l'équation homogène (sans second membre), mais en disant que K est une fonction. Afin de trouver la fonction K qui donne la solution, il suffit de remettre l'essai (33.31) dans l'équation (33.29) :

$$\underbrace{K'e^{-t^2} - 2tKe^{-t^2}}_{y'(t)} + \underbrace{2tKe^{-t^2}}_{2ty(t)} = 4t \quad (33.32)$$

Les deux termes avec K se simplifient et il reste

$$K'(t) = 4te^{t^2}, \quad (33.33)$$

ce qui signifie $K(t) = 2e^{t^2+C}$. Nous avons donc déterminé la fonction qui fait fonctionner l'essai, et la solution à l'équation est

$$y(t) = (2e^{t^2} + C)e^{-t^2} = 2 + Ce^{-t^2}. \quad (33.34)$$

△

La technique pour résoudre cette équation est de commencer par résoudre l'équation homogène associée. Si $U(t)$ est une primitive de $u(t)$, nous avons

$$\begin{aligned} y_H'(t) + u(t)y_H(t) &= 0 \\ \frac{y_H'}{y_H} &= -u(t) \\ \ln(y_H) &= -U(t) + C \\ y_H(t) &= e^{-U(t)+C} = Ke^{-U(t)} \end{aligned} \quad (33.35)$$

où $K = e^C$.

Cela fournit la solution générale de l'équation homogène. Il existe un truc génial qui permet d'en tirer la solution générale du système non homogène. Lorsque nous avons trouvé $y_H(t) = Ke^{-U(t)}$, le symbole K désigne une constante. La méthode de **variation des constantes** consiste à essayer la solution

$$y(t) = K(t)e^{-U(t)}, \quad (33.36)$$

c'est-à-dire à dire que la constante est en réalité une fonction. Afin de trouver quelle fonction $K(t)$ fait en sorte que l'essai (33.36) soit une solution, nous la remplaçons dans l'équation de départ $y' + uy = v$. Maintenant,

$$y'(t) = K'(t)e^{-U(t)} - K(t)u(t)e^{-U(t)}. \quad (33.37)$$

En remettant dans l'équation,

$$y' + uy = K'e^{-U} - Kue^{-U} + uKe^{-U} = K'e^{-U} = v. \quad (33.38)$$

Notez que les termes en K se sont miraculeusement simplifiés. Cela est directement dû au fait que e^{-U} est solution de l'équation homogène. Nous restons avec l'équation

$$K' = \frac{v}{e^{-U}} \quad (33.39)$$

pour $K(t)$. La solution générale du problème non homogène est donc finalement donnée par

$$y(t) = (W(t) + C)e^{-U(t)} \quad (33.40)$$

si $W(t)$ est une primitive de $v(t)e^{U(t)}$.

Tout ceci est un peu heuristique. La proposition suivante dit dans quels cas ça fonctionne.

Proposition 33.10.

Soient u et v continues sur I et U , une primitive de u sur I et W une primitive de ve^{-U} sur I . Une fonction $y: I \rightarrow \mathbb{R}$ est solution de $y' + u(t)y = v(t)$ si et seulement s'il existe une constante $C \in \mathbb{R}$ telle que

$$y(t) = (W(t) + C)e^{U(t)} \quad (33.41)$$

pour tout $t \in I$.

33.3.1 Pourquoi la variation des constantes fonctionne toujours ?

Prenons une équation non homogène

$$z'(t) = f(t)z(t) + g(t), \quad (33.42)$$

et supposons avoir une solution de l'homogène associée sous la forme $z_H(t) = Ch(t)$. Le coup de la variation des constantes consiste à essayer une solution pour l'équation non homogène sous la forme¹

$$z(t) = K(t)h(t). \quad (33.43)$$

1. Je ne sais plus qui a eu l'idée de changer le nom de la constante de C vers K au moment de la transformer en fonction, mais c'est une bonne idée.

Nous injectons cette solution dans l'équation de départ en utilisant le fait que $z'(t) = K'(t)h(t) + K(t)h'(t)$:

$$K'(t)h(t) + K(t)h'(t) = f(t)K(t)h(t) + g(t). \quad (33.44)$$

Le terme $K(t)h'(t)$ se réécrit en utilisant la propriété de définition de h , c'est-à-dire que $h'(t) = f(t)h(t)$. Nous voyons que les termes ne contenant pas de K' se simplifient ; il reste

$$K'h = g. \quad (33.45)$$

Cette équation a comme solution

$$K = \int \frac{f}{h} + C. \quad (33.46)$$

J'insiste sur la constante d'intégration ! En réalité, celles et ceux qui auront compris l'équation (33.24) sauront que K est donné par

$$K(t) = \int_{\xi_0}^t \frac{f(\xi)}{g(\xi)} d\xi \quad (33.47)$$

où ξ_0 joue le rôle de la constante d'intégration.

Quoi qu'il en soit, la solution générale de l'équation non homogène est

$$z(t) = K(t)h(t) = \left(\int \frac{g}{h} + C \right) h. \quad (33.48)$$

Cette solution comprend deux termes : Ch qui est solution de l'homogène, et $(\int \frac{g}{h}) h$ qui est une particulière de l'équation non homogène.

Quelques conclusions :

- (1) Si vous avez encore du K (et pas que du K') dans votre équation qui donne K , c'est que vous n'êtes pas dans le cadre d'une équation de type (33.42). Le plus souvent, c'est que vous avez fait une faute de calcul quelque part.
- (2) La méthode des variations des constantes n'est pas en contradiction avec le principe de « SGEH+SPENH ». En effet, la SGEP et la SPENH sont toutes deux dans la solution (33.48).
- (3) La variation des constantes peut être vue comme une façon cool de trouver une solution particulière de l'équation non homogène.
- (4) La simplification ne se fait que après avoir remplacé Kh' par Kfh , c'est-à-dire après avoir utilisé le fait que z_H est solution de l'homogène. Sinon, la simplification n'est pas du tout évidente a priori. Il se peut même que, visuellement, les termes Kh' et Kfh ne se ressemblent pas du tout. Un exemple de cela arrivera par exemple dans l'exemple 33.13, pour arriver à l'équation (33.64).

33.4 Équations à variables séparées

Une **équation à variables séparées** est une équation de la forme

$$y' = u(t)f(y) \quad (33.49)$$

où $u: I \rightarrow \mathbb{R}$ et $f: J \rightarrow \mathbb{R}$ sont deux fonctions continues données. Les propositions 33.11 et 33.12 résolvent ce cas, mais avant de voir cela, nous allons donner quelques indications « pratiques ».

33.4.1 La méthode rapide

On peut évidemment mettre tous les y et y' d'un côté :

$$\frac{y'}{f(y)} = u(x). \quad (33.50)$$

Une fois que cela est fait, on écrit $y' = \frac{dy}{dx}$, et on envoie le dx du côté des x :

$$\frac{dy}{f(y)} = u(x)dx. \quad (33.51)$$

Maintenant il suffit de prendre l'intégrale des deux côtés : comme la position des dx et dy l'indiquent, il faut intégrer par rapport à y d'un côté et par rapport à dx de l'autre côté.

L'intégrale à gauche est facile : c'est $\ln(y)$. À droite, par contre, ça dépend tout à fait de u .

33.4.2 La méthode plus propre

$$y'(t) = u(t)f(y(t)). \quad (33.52)$$

Nous considérons U , une primitive de u sur I et G , une primitive de $1/f$ sur J . Si $I' \subseteq I$ et $y: I' \rightarrow J$, alors y est solution de (33.49) si et seulement s'il existe une constante C telle que

$$G(y(t)) = U(t) + C. \quad (33.53)$$

La recherche des solutions de l'équation différentielle se ramène donc à la recherche de primitives et de solutions d'une équation algébrique (il faut isoler $y(t)$ dans (33.53)). Réciproquement toute solution régulière de cette dernière relation est solution de l'équation différentielle.

Remarque : lorsque nous cherchons U et G , nous ne cherchons que *une* primitive. Il ne faut pas considérer des constantes d'intégration à ce niveau.

33.4.3 Les théorèmes

Proposition 33.11.

Nous considérons l'équation (33.49) avec $u(t)$ continue sur I et f continue sur J avec $f(\eta) \neq 0$ pour tout $\eta \in J$. Soit U , une primitive de u sur I , et G , une primitive de $1/f$ sur J .

Si $y: I' \rightarrow J$ est une fonction sur un intervalle $I' \subset I$, alors y est solution de l'équation (33.49) si et seulement s'il existe $C \in \mathbb{R}$ tel que

$$G(y(t)) = U(t) + C. \quad (33.54)$$

Cette proposition dit que toutes les solutions qui ne s'annulent jamais sur un intervalle ont la forme $G(y(t)) = U(t) + C$ et peuvent donc être trouvées en calculant des primitives.

La formule (33.54) peut être obtenue de la façon heuristique suivante, en écrivant $y' = dy/dt$, et en passant le dt à droite. Nous trouvons successivement

$$\begin{aligned} y' &= u(t)f(y) \\ dy &= u(t)f(y)dt \\ \frac{dy}{f(y)} &= u(t)dt \\ \int \frac{dy}{f(y)} &= \int u(t)dt \\ G(y) &= U(t) + C. \end{aligned} \quad (33.55)$$

Proposition 33.12.

Soient u continue sur I et f continue sur J , et $f(\eta) \neq 0$ sur J . Soient $t_0 \in I$ et $y_0 \in J$. Alors il existe $I' \subset I$ avec $t_0 \in I'$ et $f \in C^1(I' \rightarrow J)$ tels que

- (1) y est solution de (33.49) sur I' et vérifie $y(t_0) = y_0$,
- (2) si z est une solution de (33.49) sur $I'' \subset I'$ avec $t_0 \in I''$ et $z(t_0) = y_0$, alors $I'' \subset I'$ et $z(t) = y(t)$ pour tout $t \in I''$.

Exemple 33.13

Résoudre l'équation différentielle

$$y - \cos(t)y' = \cos(t)(1 - \sin(t))y^2. \quad (33.56)$$

La fonction $y = 0$ est solution. En posant $z = 1/y$, nous trouvons l'équation

$$z + \cos(t)z' = \cos(t)(1 - \sin(t)) \quad (33.57)$$

à laquelle z doit satisfaire. L'équation homogène est

$$z'_H = -\frac{z_H}{\cos(t)}. \quad (33.58)$$

Ceci est une équation à variables séparées que nous résolvons en suivant les méthodes données plus haut : nous posons

$$\begin{aligned} u(t) &= \frac{1}{\cos(t)}, \\ f(z) &= -z, \\ U(t) &= \ln \left[\tan \left(\frac{\pi}{4} + \frac{t}{2} \right) \right] \quad (\text{voir formulaire}), \\ G(z) &= \ln \left(\frac{1}{z} \right). \end{aligned} \quad (33.59)$$

La solution z_H est donnée par l'équation

$$\ln \left(\frac{1}{z} \right) = \ln \left[K \tan \left(\frac{\pi}{4} + \frac{t}{2} \right) \right], \quad (33.60)$$

c'est-à-dire

$$z_H(t) = \frac{K}{\tan \left(\frac{\pi}{4} + \frac{t}{2} \right)}. \quad (33.61)$$

Nous appliquons maintenant la méthode de variation des constantes sur cette solution afin de trouver la solution générale de l'équation (33.57). En utilisant la règle de Leibnitz, $z' = K'z_H + Kz'_H$, nous trouvons

$$\frac{K}{\tan \left(\frac{\pi}{4} + \frac{t}{2} \right)} + \cos(t) \left(\frac{K'}{\tan \left(\frac{\pi}{4} + \frac{t}{2} \right)} - \frac{K}{2 \sin^2 \left(\frac{\pi}{4} + \frac{t}{2} \right)} \right) = \cos(t)(1 - \sin(t)). \quad (33.62)$$

Malgré leurs apparences, les deux termes en K se simplifient. En effet, en vertu de l'équation $z'_H = \frac{-z_H}{\cos(t)}$, nous avons

$$\frac{-K}{2 \sin^2 \left(\frac{\pi}{4} + \frac{t}{2} \right)} = \frac{-K}{\cos(t) \tan \left(\frac{\pi}{4} + \frac{t}{2} \right)}. \quad (33.63)$$

Le travail de voir quel est le lien entre $\sin^2 \left(\frac{\pi}{4} + \frac{t}{2} \right)$, $\tan \left(\frac{\pi}{4} + \frac{t}{2} \right)$ et $\cos(t)$ est en réalité fait dans votre formulaire au moment où vous l'avez utilisé pour intégrer u pour obtenir le $U(t)$ de (33.59).

Après cette simplification durement méritée, nous trouvons l'équation suivante pour $K(t)$:

$$\frac{K'}{\tan \left(\frac{\pi}{4} + \frac{t}{2} \right)} = 1 - \sin(t). \quad (33.64)$$

Résoudre cela revient à trouver la primitive de

$$(1 - \sin(t)) \tan \left(\frac{\pi}{4} + \frac{t}{2} \right), \quad (33.65)$$

ce qui est relativement compliqué. La réponse est

$$\begin{aligned} K(t) &= \ln \left(\sin \left(\frac{2x + \pi}{4} \right) + 1 \right) + \ln \left(\sin \left(\frac{2x + \pi}{4} \right) - 1 \right) \\ &\quad + 2 \ln \sec \left(\frac{2x + \pi}{4} \right) + 2 \sin^2 \left(\frac{2x + \pi}{4} \right) \end{aligned} \quad (33.66)$$

Nous pouvons un peu simplifier en utilisant le fait que $\ln(a+b) + \ln(a-b) = \ln(a^2 - b^2)$:

$$K(t) = \ln \left(-\cos^2 \left(\frac{2x + \pi}{4} \right) \right) + 2 \ln \sec \left(\frac{2x + \pi}{4} \right) + 2 \sin^2 \left(\frac{2x + \pi}{4} \right). \quad (33.67)$$

Il me semble toutefois qu'il faudrait prendre des valeurs absolues pour les logarithmes.

△

33.5 Équations linéaires d'ordre supérieur

33.5.1 Équations et systèmes linéaire à coefficients constants

Nous regardons l'équation

$$y^{(n)} + a_1 y^{(n-1)} + \dots + a_{n-1} y' + a_n y = v(t) \quad (33.68)$$

où les coefficients a_k sont maintenant des constantes. Il faut commencer par résoudre le polynôme caractéristique

$$r^n + a_1 r^{n-1} + \dots + a_n = 0. \quad (33.69)$$

Si $\lambda_1, \dots, \lambda_k$ sont les solutions avec multiplicité μ_1, \dots, μ_k , alors le **système fondamental** de solutions linéairement indépendantes est l'ensemble suivant de solutions à l'équation homogène :

$$\begin{aligned} e^{\lambda_1 t}, t e^{\lambda_1 t}, \dots, t^{\mu_1-1} e^{\lambda_1 t} \\ \vdots \\ e^{\lambda_k t}, t e^{\lambda_k t}, \dots, t^{\mu_k-1} e^{\lambda_k t}. \end{aligned} \quad (33.70)$$

Nous notons y_i ces solutions. La solution générale de l'équation homogène est donc donnée par

$$y_H = \sum_i c_i y_i. \quad (33.71)$$

Afin de trouver la solution générale de l'équation non homogène, nous appliquons la méthode de variation des constantes, en imposant les $n - 1$ conditions

$$\sum_{i=1}^n c_i'(t) y_i^{(l)}(t) = 0 \quad (33.72)$$

avec $l = 0, \dots, n - 2$. Ces conditions plus l'équation de départ (33.68) forment un système de n équations différentielles pour les n fonctions inconnues $c_i(t)$.

Cette condition peut paraître mystérieuse. Il est cependant encore possible de travailler sans poser la condition (33.72) en suivant la recette, en calculant des déterminants de Wronskien. Des exemples sont donnés dans les exercices sur le second ordre.

33.5.2 Si les coefficients ne sont pas constants ?

Une équation différentielle linéaire d'ordre n sur I est une équation de la forme

$$y^{(n)} + u_1(t) y^{(n-1)} + \dots + u_{n-1}(t) y' + u_n(t) y = v(t) \quad (33.73)$$

où v et u_k sont des fonctions continues fixées de I vers \mathbb{R} .

Pour résoudre cette équation, il faut commencer par résoudre l'équation homogène correspondante (c'est-à-dire celle que l'on obtient en posant $v(t) = 0$). Ensuite, nous trouvons la solution de l'équation (33.73) en appliquant la méthode de la **variation des constantes**.

Donnons un exemple du pourquoi la méthode de variations des constantes est efficace. Soit l'équation

$$u' + f(t)u = g(t), \quad (33.74)$$

et disons que u_H est une solution de l'équation homogène. La méthode de variations des constantes consiste à poser $u(t) = K(t)u_H(t)$, et donc $u'(t) = K'u_H + Ku'_H$. En remettant dans l'équation de départ,

$$K'u_H + Ku'_H + fKu_H = g. \quad (33.75)$$

La somme $Ku'_H + fKu_H$ est nulle, par définition de u_H . Par conséquent, il ne reste que

$$K' = \frac{g(t)}{u_H}. \quad (33.76)$$

Lorsqu'on utilise la méthode de variation des constantes, nous trouvons toujours une simplification « miraculeuse ».

Dans l'immédiat, nous ne considérons que le cas où les u_i sont des constantes. Le cas où les u_i deviennent des fonctions de t sera vu plus tard.

33.6 Système d'équations linéaires

33.6.1 La magie de l'exponentielle...

Prenons l'équation différentielle très simple

$$y' = ay. \quad (33.77)$$

La solution est $y(t) = Ae^{at}$. Et si on a la donnée de Cauchy $y(t_0) = y_0$, alors

$$y(t) = Ae^{at}e^{-at_0}e^{at_0} = e^{a(t-t_0)}y(t_0). \quad (33.78)$$

Donc on a le facteur multiplicatif $e^{a(t-t_0)}$ qui sert à faire passer de $y(0)$ à $y(t)$. C'est un peu un opérateur d'évolution. Ce qui fait la magie de l'exponentielle, c'est son développement en série

$$e^x = 1 + x + \frac{x^2}{2} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots \quad (33.79)$$

qui est tel que chaque terme est la dérivée du terme suivant.

Maintenant, si on a un système

$$\bar{y}' = A\bar{y}, \quad (33.80)$$

il n'est pas du tout étonnant d'avoir comme solution $\bar{y}(t) = e^{At}$ où l'exponentielle de la matrice est définie exactement par la série (33.79). C'est un peu longuet, mais dans le cours, c'est effectivement ce qui est prouvé. La matrice résolvante $R(t, t_0): \bar{y}_0 \rightarrow \bar{y}(t; t_0, y_0)$ est donné par

$$R(t, t_0) = e^{(t-t_0)A}, \quad (33.81)$$

exactement comme dans l'équation (33.78).

33.6.2 ... mais la difficulté

Maintenant, il est suffisant de calculer des exponentielles de matrices pour résoudre des systèmes. Hélas, il est en général très difficile de calculer des exponentielles. Tu peux essayer de prouver les deux suivantes :

$$\begin{aligned} A &= \begin{pmatrix} 0 & a \\ -a & 0 \end{pmatrix} \rightsquigarrow e^A = \begin{pmatrix} \cos(a) & \sin(a) \\ -\sin(a) & \cos(a) \end{pmatrix} \\ S &= \begin{pmatrix} 0 & a \\ a & 0 \end{pmatrix} \rightsquigarrow e^S = \begin{pmatrix} \cosh(a) & \sinh(a) \\ \sinh(a) & \cosh(a) \end{pmatrix}. \end{aligned} \quad (33.82)$$

La première, tu vas la revoir si tu fais de la géométrie différentielle ou de la mécanique quantique : l'algèbre de Lie du groupe des matrices orthogonales de déterminant 1 est l'algèbre des matrices antisymétriques.

La seconde se retrouve en relativité parce que e^S est la matrice qui préserve $x^2 - y^2$, tout comme e^A préserve $x^2 + y^2$. Quelques mots sur l'utilisation des fonctions hyperboliques en relativité dans [43.4.8.1](#).

33.6.3 La recette

Afin d'éviter de devoir calculer explicitement des exponentielles de matrices, nous faisons appel à toutes sortes de trucs, dont la forme de Jordan. Le résultat final est la méthode suivante. Soit le système homogène

$$\bar{y}' = A\bar{y}. \quad (33.83)$$

- (1) D'abord, nous calculons les valeurs propres de A .
- (2) Ensuite les vecteurs propres.
- (3) Une bonne valeur propre, c'est une valeur propre dont l'espace propre a une dimension égale à sa multiplicité. C'est-à-dire que si λ est de multiplicité m , alors on a, dans les bons cas, m vecteurs propres linéairement indépendants.

Dans ce cas, si v_1, \dots, v_m sont les vecteurs, alors on a les solutions linéairement indépendantes suivantes :

$$\begin{pmatrix} \vdots \\ v_1 \\ \vdots \end{pmatrix} e^{\lambda t}, \dots, \begin{pmatrix} \vdots \\ v_m \\ \vdots \end{pmatrix} e^{\lambda t}. \quad (33.84)$$

Pour chaque bonne valeur propre, ça nous fait un tel paquet de solutions linéairement indépendantes.

- (4) Si λ n'est pas une bonne valeur propre, alors les choses se compliquent. Mettons que λ ait k vecteurs propres en moins que sa multiplicité. Dans ce cas, il faut chercher des solutions sous la forme

$$\begin{pmatrix} a_1^{(k)} t^k + \dots + a_1^{(0)} \\ \vdots \\ a_n^{(k)} t^k + \dots + a_n^{(0)} \end{pmatrix} e^{\lambda t}. \quad (33.85)$$

C'est-à-dire qu'on prend comme coefficient de $e^{\lambda t}$, un vecteur de polynômes de degré k . Il faut mettre cela dans l'équation de départ pour voir quelles sont les contraintes sur les constantes $a_i^{(j)}$ introduites.

- (5) Nous avons un cas particulier du cas précédent. Si λ est une valeur propre de multiplicité m qui n'a que un seul vecteur propre v , alors il faut chercher des polynômes de degré $m - 1$, et on peut directement fixer le coefficient de t^{m-1} , ce sera l'unique vecteur propre :

$$\left[\begin{pmatrix} \vdots \\ v \\ \vdots \end{pmatrix} + \begin{pmatrix} a_1^{(m-2)} \\ \vdots \\ a_n^{(m-2)} \end{pmatrix} t^{m-2} + \dots \right] e^{\lambda t}. \quad (33.86)$$

Cela économise quelques calculs par rapport à poser brutalement [\(33.85\)](#).

33.6.4 Système d'équations linéaires avec matrice constante

Nous considérons l'équation différentielle

$$y'(t) = Ay(t) \quad (33.87)$$

pour la fonction $y: \mathbb{R} \rightarrow \mathbb{R}^n$ et A est une matrice ne dépendant pas de t . Nous supposons que A est diagonalisable pour les vecteurs propres v_i et les valeurs propres λ_i correspondantes.

La matrice

$$R(t) = [e^{\lambda_1 t} v_1 \dots e^{\lambda_n t} v_n] \quad (33.88)$$

est la **matrice résolvente** du système. Alors la solution du système (33.87) pour la condition initiale $y(0) = y_0$ est

$$y(t) = R(t)y_0. \quad (33.89)$$

En effet

$$AR(t) = \left[A \begin{pmatrix} \uparrow \\ e^{\lambda_1 t} v_1 \\ \downarrow \end{pmatrix} \dots A \begin{pmatrix} \uparrow \\ e^{\lambda_n t} v_n \\ \downarrow \end{pmatrix} \right] = R'(t). \quad (33.90)$$

Par conséquent $y'(t) = R'(t)y_0 = AR(t)y_0 = Ay(t)$.

33.6.5 Système d'équations linéaires avec matrice non constante

Théorème 33.14 ([451]).

Soient I un intervalle de \mathbb{R} et $M: \mathbb{R} \rightarrow \mathcal{L}(\mathbb{R}^n, \mathbb{R}^n)$ une fonction. Si les composantes M_{ij} sont des fonctions continues sur I alors :

(1) pour tout $t_0 \in I$ et pour tout $y_0 \in \mathbb{R}^n$ le système

$$y'(t) = M(t)y(t) \quad (33.91)$$

admet une unique solution maximale définie sur I telle que $y(t_0) = y_0$;

(2) l'ensemble des solutions de l'équation (33.91) sur I est un espace vectoriel de dimension n .

33.7 Réduction de l'ordre

Afin de diminuer l'ordre d'une équation dans laquelle le paramètre n'apparaît pas, il y a deux changements de variables très utiles. Le premier, le plus simple, est simplement de poser $z(t) = y'(t)$, ce qui donne $z'(t) = y''(t)$. Le second, qui n'est pas le même, est $z(y(t)) = y'(t)$, qui entraîne $y''(t) = z'(y(t))z(t)$. Dans ce second cas, il faut également changer de variable, et utiliser $y(t)$ comme variable au lieu de t .

Si ça ne marche pas, il faut suivre la procédure ci-après.

Nous supposons avoir une équation différentielle d'ordre p dans laquelle $y^{(p)}$ est isolée des autres dérivées :

$$y^{(p)}(t) = f(t, y(t), y'(t), \dots, y^{(p-1)}(t)) \quad (33.92)$$

où f est une fonction $f: \mathbb{R} \times \mathbb{R}^p \rightarrow \mathbb{R}$, et la fonction cherchée est $y: \mathbb{R} \rightarrow \mathbb{R}$.

La méthode proposée ici consiste à transformer cette équation d'ordre p en un système d'équations d'ordre 1. Pour cela nous posons

$$F: \mathbb{R} \times \mathbb{R}^p \rightarrow \mathbb{R}^p$$

$$(t, x) \mapsto \begin{pmatrix} x_2 \\ \vdots \\ x_p \\ f(t, x_1, \dots, x_p) \end{pmatrix}. \quad (33.93)$$

Nous considérons alors l'équation différentielle

$$Y'(t) = F(t, Y(t)). \quad (33.94)$$

pour $Y: \mathbb{R} \rightarrow \mathbb{R}^p$. La fonction $y = Y_1$ résout l'équation (33.92) si et seulement si la fonction Y résout l'équation (33.94).

De plus si l'équation (33.92) est donnée avec les conditions initiales $y^{(k)} = a_k$ ($k = 0, \dots, p-1$) alors l'équation (33.94) vient avec les conditions initiales

$$Y(t_0) = \begin{pmatrix} a_0 \\ \vdots \\ a_{p-1} \end{pmatrix}, \quad (33.95)$$

c'est-à-dire $Y(t_0) = A_0$ avec $A_0 \in \mathbb{R}^p$.

Le théorème de Cauchy-Lipschitz 18.40 nous donne existence et unicité locale de la solution au système 33.94. Lorsque le système est linéaire, c'est-à-dire sous la forme $Y'(t) = M(t)Y(t)$, alors il y a mieux : le théorème 33.14.

Exemple 33.15

Nous considérons l'équation différentielle

$$\begin{cases} -u''(t) - u(t) = 1 & (33.96a) \\ u(0) = a_0 \in \mathbb{R} & (33.96b) \\ u'(0) = a_1 \in \mathbb{R}, & (33.96c) \end{cases}$$

et nous voulons montrer que ce système accepte une unique solution. Vu que l'équation différentielle se présente sous la forme $u'' = f(t, u)$ avec $f(t, u) = -1 - u$ nous posons

$$\begin{aligned} F: \mathbb{R} \times \mathbb{R}^2 &\rightarrow \mathbb{R}^2 \\ (t, x) &\mapsto \begin{pmatrix} x_2 \\ -1 - x_1 \end{pmatrix}, \end{aligned} \quad (33.97)$$

et nous considérons l'équation différentielle

$$Y' = F(t, Y) \quad (33.98)$$

pour la fonction $Y: \mathbb{R} \rightarrow \mathbb{R}^2$. La fonction F est Lipschitz (et même globalement) par rapport à Y . En effet,

$$\|F(x) - F(y)\| = (x_2 - y_2)^2 + (x_1 - y_1)^2 = \|x - y\|^2. \quad (33.99)$$

Le théorème de Cauchy-Lipschitz 18.40 s'applique et en posant $Y_0 = (a_0, a_1)$, il existe une unique solution à l'équation $Y' = F(t, Y)$ vérifiant $Y(0) = Y_0$. Nous notons $t \mapsto Y(t)$ cette solution.

En quoi cela nous aide ? Nous posons $u(t) = Y_1(t)$. Alors

$$u'(t) = Y_1'(t) = F_1'(t, Y) = Y_2(t). \quad (33.100)$$

En dérivant encore,

$$u''(t) = Y_2'(t) = F_2'(t, Y) = -1 - Y_1(t) = -1 - u(t), \quad (33.101)$$

ce qu'il fallait. La fonction $t \mapsto Y_1(t)$ est solution de notre équation de départ. Quid des conditions initiales ? Vu que $u = Y_1$ et $u' = Y_2$ nous avons

$$u(0) = Y_1(0) = a \quad (33.102)$$

et

$$u'(0) = Y_2(0) = b. \quad (33.103)$$

Toutes les prescriptions sont respectées.

Si vous voulez vraiment résoudre cette équation, il faudra plus de travail. D'abord résoudre l'équation homogène associée, c'est-à-dire l'équation caractéristique $r^2 + 1 = 0$, ce qui va donner

$$u_H(t) = Ae^{it} + Be^{-it}, \quad (33.104)$$

et ensuite faire le coup de la variation des constantes pour déterminer la solution générale du problème non homogène. \triangle

Exemple 33.16

Nous reprenons l'équation différentielle

$$-u''(t) - u(t) = 1. \quad (33.105)$$

Nous avons déjà vu dans l'exemple 33.15 que cette équation avait une solution unique pour toute condition initiale. Cette fois nous voulons étudier les solutions lorsque nous imposons les conditions aux limites $u(0) = u(\pi) = 0$. Nous allons voir qu'il n'y a pas de telles solutions.

Pour ce faire, soit une solution u . D'abord u'' existant, la fonction u est de classe au moins C^1 . Mais $u'' = -1 - u$, donc u'' est également C^1 , ce qui donne la régularité C^3 pour u . En continuant ainsi nous trouvons que u est de classe C^∞ .

Le truc est de considérer la fonction $v(t) = \sin(t)$ qui vérifie l'équation différentielle

$$\begin{cases} -v'' - v = 0 & (33.106a) \\ v(0) = 0 & (33.106b) \\ v'(0) = 0. & (33.106c) \end{cases}$$

Nous calculons le produit scalaire sur $L^2(]0, \pi[)$ de (33.105) avec v :

$$\langle u'', v \rangle + \langle u, v \rangle = -\langle v, 1 \rangle. \quad (33.107)$$

Le calcul de $\langle v, 1 \rangle$ est simplement l'intégrale de $\sin(t)$ pour t allant de 0 à π , c'est-à-dire $\langle v, 1 \rangle = 2$. Vu que u et v sont toutes deux des fonctions qui s'annulent en 0 et en π nous pouvons faire des intégrations par partie les yeux fermés et exprimer $\langle u'', v \rangle$ sans dérivées sur u :

$$\langle u'', v \rangle = -\langle u', v' \rangle = \langle u, v'' \rangle = -\langle u, v \rangle \quad (33.108)$$

où la dernière égalité n'est autre que le fait que $v = \sin$, donc $v'' = -v$. Le membre de gauche de (33.107) vaut donc zéro alors que celui de droite vaut -2 .

Nous concluons que le problème aux limites posé n'admet pas de solutions. \triangle

33.8 Autour de Cauchy-Lipschitz

Dans cette section nous étudions les équations différentielles du type

$$\begin{cases} y'(t) = f(t, y(t)) & (33.109a) \\ y(t_0) = y_0. & (33.109b) \end{cases}$$

33.8.1 Fuite des compacts et explosion en temps fini

Théorème 33.17 (Fuite des compacts[452, 285]).

Nous considérons l'équation différentielle

$$\begin{cases} y'(t) = f(t, y(t)) & (33.110a) \\ y(t_0) = y_0, & (33.110b) \end{cases}$$

où $f: I \times \Omega \rightarrow \mathbb{R}^n$ est continue et Ω ouvert dans \mathbb{R}^n . Soit la solution maximale $y_M: J_M =]t_{\min}, t_{\max}[\rightarrow \Omega$. Si $t_{\max} < \sup(I)$ alors $y_M(t)$ sort de tout compact de Ω lorsque $t \rightarrow t_{\max}$.

Démonstration. Soit K un compact de Ω et nous considérons une suite (t_m) dans $]t_{min}, t_{max}[$ telle que $t_m \rightarrow t_{max}$. Si nous supposons que $y_M(t)$ ne sort pas de K alors nous avons $y_M(t_m) \in K$, c'est-à-dire une suite dans un compact. Quitte à passer à une sous-suite, nous supposons qu'elle est convergente. Soit $x_1 \in K$ la limite $\lim_{m \rightarrow \infty} y_M(t_m) = x_1$.

Vu que $t_{max} \in I$, la condition initiale $y(t_{max}) = x_1$ est valide et le théorème de Cauchy-Lipschitz 18.40 nous donne une unique solution maximale y_P définie sur un ouvert J_P autour de t_{max} .

Nous allons maintenant construire une solution au problème initial qui contredit la maximalité de y_M . Attention : il n'est pas évident a priori que $y_P(t) = y_M(t)$ sur l'intersection des domaines. Si c'était évident, la proposition serait démontrée.

Soit $\tilde{J} = J_M \cup J_P \cap]t_{min}, +\infty[$ et la fonction

$$\tilde{y}(t) = \begin{cases} y_M(t) & \text{si } t < t_{max} \\ y_P(t) & \text{si } t \geq t_{max}. \end{cases} \tag{33.111}$$

La fonction \tilde{y} est continue par construction parce que

$$\lim_{t \rightarrow t_{max}} y_M(t) = x_1 = y_P(t_{max}). \tag{33.112}$$

Nous vérifions à présent que \tilde{y} est une solution : $\tilde{y}'(t_{max}) = f(t_{max}, y(t_{max}))$:

$$\lim_{\epsilon \rightarrow 0} \frac{\tilde{y}(t_{max} - \epsilon) - \tilde{y}(t_{max})}{\epsilon} = \lim_{\epsilon \rightarrow 0} \frac{y_M(t_{max} - \epsilon) - y_P(t_{max})}{\epsilon} \tag{33.113a}$$

$$= \lim_{\epsilon \rightarrow 0} \frac{y_M(t_{max} - \epsilon) - y_P(t_{max} - \epsilon) + y_P(t_{max} - \epsilon) - y_P(t_{max})}{\epsilon} \tag{33.113b}$$

$$= \lim_{\epsilon \rightarrow 0} \frac{y_P(t_{max} - \epsilon) - y_P(t_{max})}{\epsilon} \tag{33.113c}$$

$$= y_P'(t_{max}). \tag{33.113d}$$

Donc \tilde{y} est solution pour la condition initiale $\tilde{y}(t_{max}) = x_1$ et coïncide avec y_P en t_{max} et avec y_M avant t_{max} . Donc en réalité y_P, y_M et \tilde{y} sont identiques et cela contredit la maximalité de y_M . \square

Corollaire 33.18 (Explosion en temps fini).

Soit (y_m, J) la solution maximale du problème de Cauchy (18.146) :

$$\begin{cases} y' = f(t, y) \\ y(t_0) = y_0, \end{cases} \tag{33.114a}$$

$$\tag{33.114b}$$

avec $f: U = I \times \Omega \rightarrow \mathbb{R}^n$ où I est ouvert dans \mathbb{R} et Ω ouvert dans \mathbb{R}^n . Nous supposons que f est continue sur U et localement Lipschitz par rapport à y .

Si la solution maximale est définie sur $J =]t_{min}, t_{max}[$ alors nous avons l'alternative suivante :

- (1) Soit $t_{max} = \sup(I)$,
- (2) soit $t_{max} < \sup(I)$ et $\lim_{t \rightarrow t_{max}} \|y(t)\| = \infty$.

Le résultat tient aussi mutatis mutandis pour t_{min} .

Remarque 33.19.

Attention : ceci n'est pas une simple paraphrase de la fuite des compacts. L'information supplémentaire que ce corollaire donne est que la solution sort de tout compact pour ne plus y retourner.

Démonstration. L'hypothèse $t_{max} < \sup(I)$ signifie que la solution finit d'exister avant que les hypothèses sur f cessent d'être vraies. C'est-à-dire que la solution maximale est moindre que ce que nous aurions pu espérer.

Soit un compact K . Supposons que pour tout $t_0 < t_{max}$ il existe $t \in]t, t_{max}[$ tel que $y_M(t) \in K$. Alors cela crée une suite t_k dans J telle que $y_M(t_k)$ est dans K . Comme dans le théorème de la fuite des compacts nous concluons l'impossibilité de la chose.

Donc pour tout compact K de Ω , il existe $T < t_{max}$ tel que $y_M(t) \in \Omega \setminus K$ pour tout $t \in [T, t_{max}[$. En prenant des boules fermées de plus en plus grandes en guise de compacts nous concluons que

$$\lim_{t \rightarrow t_{max}} \|y_M(t)\| = \infty. \quad (33.115)$$

□

33.20.

Notons que si $t_{max} < \infty$, si nous sommes dans l'alternative 33.18(2) et si la solution maximale y est de classe C^1 (ce qui est le cas lorsqu'on utilise Cauchy-Lipschitz 18.40) alors la dérivée de y est également non bornée dans un voisinage de t_{max} .

Mais si f est globalement bornée, alors dans l'équation $y' = f(t, y)$, la dérivée y' sera globalement bornée. Dans ce cas, la solution ne peut pas exploser en temps fini et existe donc globalement.

Problèmes et choses à faire

Êtes-vous d'accord avec 33.20 ?

Exemple 33.21

Soit l'équation différentielle

$$\begin{cases} y' = y(y-1) \sin(yt) & (33.116a) \\ y(0) = \frac{1}{2}. & (33.116b) \end{cases}$$

La fonction $f(t, y) = y(y-1) \sin(yt)$ ayant une dérivée bornée partout, elle est localement Lipschitz et le théorème de Cauchy-Lipschitz 18.40 s'applique. Pour toute condition initiale, une solution maximale unique existe.

Si nous oublions la condition initiale, il est facile de trouver des solutions constantes : $y' = 0$ avec $y(t) = k$ donne l'équation

$$0 = k(k-1) \sin(kt). \quad (33.117)$$

Les solutions $y_1(t) = 0$ et $y_2(t) = 1$ sont des solutions existant pour tout t .

Le graphe de la solution correspondante à la condition initiale $y(0) = \frac{1}{2}$ ne pouvant pas croiser les graphes de y_1 et y_2 , elle est obligée d'exister pour tout t parce qu'elle ne peut pas exploser en temps fini. \triangle

33.8.2 Écart entre deux conditions initiales

Proposition 33.22 ([450, 453]).

Soit une fonction $f: I \times \Omega \rightarrow \mathbb{R}^n$ continue et globalement Lipschitz en sa seconde variable (Ω est un ouvert de \mathbb{R}^n). Soient deux solutions $y_1: I_1 \rightarrow \mathbb{R}^n$ et $y_2: I_2 \rightarrow \mathbb{R}^n$ aux problèmes

$$\begin{cases} y'_i(t) = f(t, y_i(t)) & (33.118a) \\ y_i(t_0) = a_i. & (33.118b) \end{cases}$$

Alors pour tout $t \in I_1 \cap I_2$ nous pouvons estimer l'écart entre y_1 et y_2 par la formule

$$\|y_1(t) - y_2(t)\| \leq e^{L|t-t_0|} \|a_1 - a_2\|, \quad (33.119)$$

où L est la constante de Lipschitz de f .

Démonstration. Nous avons d'abord les majorations suivantes, qui semblent juste jouer avec les notations, mais qui utilisent le fait (contenu dans le théorème de Cauchy-Lipschitz) que y_i soit de

classe C^1 :

$$\|y_1(t) - y_2(t)\| = \left\| \int_{t_0}^t (y_1'(s) - y_2'(s)) \right\| \quad (33.120a)$$

$$\leq L \int_{t_0}^t \|f(s, y_1(s)) - f(s, y_2(s))\| ds \quad (33.120b)$$

$$= L \int_{t_0}^t \|y_1(s) - y_2(s)\| ds. \quad (33.120c)$$

C'est à ce moment que nous utilisons le lemme de Grönwall. Vu que

$$\|y_1(t) - y_2(t)\| \leq L \int_{t_0}^t \|y_1(s) - y_2(s)\| ds, \quad (33.121)$$

nous sommes dans les hypothèses de Grönwall 33.5 en posant

$$u(t) = \|y_1(t) - y_2(t)\| \quad (33.122a)$$

$$b(t) = \|y_1(0) - y_2(0)\| \quad (33.122b)$$

$$a(t) = L. \quad (33.122c)$$

Nous avons la majoration

$$\|y_1(t) - y_2(t)\| \leq \|y_1(0) - y_2(0)\| + L \int_0^t \|y_1(0) - y_2(0)\| e^{L(t-s)} ds. \quad (33.123)$$

Le calcul de l'intégrale intérieure donne

$$\int_0^t e^{L(t-s)} ds = -\frac{1}{L}(e^{-Lt} - 1). \quad (33.124)$$

Avec ça, nous avons

$$\|y_1(t) - y_2(t)\| \leq e^{Lt} \|y_1(0) - y_2(0)\|. \quad (33.125)$$

□

33.23.

Notons que la proposition 33.22 est plutôt une mauvaise nouvelle parce que les solutions restent seulement linéairement proches l'une de l'autre lorsqu'on rapproche les conditions initiales, mais elle divergent exponentiellement vite avec le temps. Donc deux trajectoires arbitrairement proches au départ finissent assez vite par être bien séparées.

Cette proposition est cependant cruciale parce qu'elle explique que pour des petits t , les solutions ne s'écartent pas beaucoup, c'est-à-dire que pour t fixé, l'application qui à une donnée initiale fait correspondre la solution en t est continue. C'est le premier pas pour parler de régularité du flot.

33.8.3 Flot d'un champ de vecteurs

Nous reprenons l'équation différentielle du théorème de Cauchy-Lipschitz 18.40. En ce qui concerne les notations, I est un intervalle ouvert de \mathbb{R} contenant 0 et l'application $f: I \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ est continue et localement Lipschitz en sa seconde variable. Pour $a \in \mathbb{R}^n$, nous notons (J_a, y_a) la solution maximale (donc $y_a: J_a \rightarrow \mathbb{R}^n$) du problème

$$\begin{cases} y_a(t) = f(t, y_a(t)) & (33.126a) \\ y_a(0) = a. & (33.126b) \end{cases}$$

Nous noterons aussi de temps en temps $\varphi(t, a) = y_a(t)$.

Nous savons que $t \mapsto y_a(t)$ est de classe C^1 , et cela est directement dans le théorème de Cauchy-Lipschitz. Une question d'une toute autre difficulté est la régularité de $a \mapsto y_a(t)$ pour t fixé, et encore pire : celle de $(t, a) \mapsto y_a(t)$.

Il se fait que l'application $(t, a) \mapsto y_a(t)$ a la même régularité que celle de f , mais cela va être un peu long à prouver. En ce qui concerne la régularité C^1 , ce sera le théorème 33.28 dont la démonstration, comme vous pouvez le voir sera copieuse et demandera des propositions intermédiaires pas simples.

Définition 33.24.

Si t est fixé, l'application

$$\begin{aligned} \varphi_t: \mathbb{R}^n &\rightarrow \mathbb{R}^n \\ x &\mapsto \varphi(t, x) = y_x(t), \end{aligned} \quad (33.127)$$

est le **flot** du problème de Cauchy (33.126).

L'application $t \mapsto \varphi_t$ est ce qui est appelé le groupe à un paramètre de flot, pour des raisons qui arriveront plus tard².

Le but est d'étudier les propriétés du flot : est-il continu, un difféomorphisme, existe, pour quels t ? Où se cache le champ de vecteurs du titre dans l'équation différentielle?

Nous posons

$$\mathcal{D} = \bigcup_{x \in \Omega} (J_x \times \{x\}). \quad (33.128)$$

Comme tout produit d'espaces métrique, l'ensemble \mathcal{D} est muni d'une métrique via la définition 9.65.

Proposition 33.25 ([450]).

Soit un intervalle I ouvert de \mathbb{R} contenant 0 et Ω un ouvert connexe de \mathbb{R}^n . Soit une application $f: I \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ continue et localement Lipschitz en sa seconde variable. Pour $a \in \mathbb{R}^n$, nous notons (J_a, y_a) la solution maximale (donc $y_a: J_a \rightarrow \mathbb{R}^n$) du problème

$$\begin{cases} y_a(t) = f(t, y_a(t)) \\ y_a(0) = a. \end{cases} \quad \begin{array}{l} (33.129a) \\ (33.129b) \end{array}$$

Nous posons

$$\mathcal{D} = \bigcup_{x \in \Omega} (J_x \times \{x\}). \quad (33.130)$$

Nous définissons la fonction φ par $\varphi(t, x) = y_x(t)$ là où ça existe.

L'ensemble \mathcal{D} est ouvert. L'application $\varphi: \mathcal{D} \rightarrow \Omega$ est localement Lipschitz.

Démonstration. Soit $(s, a) \in \mathcal{D}$ et (J_a, y_a) la solution maximale passant par a en $t = 0$. Par définition de \mathcal{D} nous avons $s \in J_a$. Nous considérons J , un compact inclus dans J_a et contenant 0 et s en son intérieur. Nous posons

$$K = J \times y_a(J). \quad (33.131)$$

Vu que y_a est continue, cela est un compact. Chaque point de K possède un voisinage ouvert sur lequel f est Lipschitz³; nous considérons un sous recouvrement fini et le maximum des constantes de Lipschitz. Cela nous crée un voisinage V de K dans $I \times \Omega$ dans lequel f est Lipschitz.

Vu que V est ouvert et K est compact avec $K \subset V$, nous pouvons trouver un ouvert V' et un compact K' tels que

$$K \subset V' \subset K' \subset V. \quad (33.132)$$

Sur ce V' , la fonction f est de plus bornée parce que continue sur le compact K' . Nous renommons V' en V . Sur V nous avons :

$$\text{— } \|f\|_{\infty, V} \leq M,$$

2. ou pas...

3. Cela est à peu près la définition d'être localement Lipschitz : 13.257, voir aussi 13.258.

— f est Lipschitz en sa seconde variable, de constante de Lipschitz L .

En tant qu'espace produit, nous avons une distance sur $I \times \Omega$ donnée en 9.65 :

$$d((t, y), (t', y')) = \max\{|t - t'|, \|y - y'\|\}. \quad (33.133)$$

Nous posons

$$V_\epsilon(K) = \{z \in I \times \Omega \text{ tel que } d(z, K) < \epsilon\} \quad (33.134a)$$

$$W_\epsilon = \{(t, y) \in J \times \Omega \text{ tel que } \|y - y_a(t)\| < \epsilon\}. \quad (33.134b)$$

$\overline{W_\epsilon} \subset \overline{V_\epsilon(K)}$ Soit $(t, y) \in W_\epsilon$. Nous avons :

$$d((t, y), K) = \inf_{(t', y') \in K} d((t, y), (t', y')) \quad (33.135a)$$

$$= \inf_{(t', y') \in K} \max\{|t - t'|, \|y - y'\|\}. \quad (33.135b)$$

Mais demander $(t, y) \in W_\epsilon$ signifie que $t \in J$ et $\|y - y_a(t)\| \leq \epsilon$. Dans K nous avons l'élément $(t, y_a(t))$ qui vérifie

$$d((t, y), (t, y_a(t))) = \|y - y_a(t)\| \leq \epsilon. \quad (33.136)$$

Donc l'infimum de (33.135b) est majoré par ϵ . Nous avons prouvé que $W_\epsilon \subset V_\epsilon(K)$ et donc même inclusion pour les fermetures.

Il existe $\epsilon > 0$ tel que $\overline{V_\epsilon(K)} \subset V$ Supposons que $\overline{V_\epsilon(K)}$ ne soit inclus dans V pour aucun ϵ .

Alors nous considérons

$$z_n \in \overline{V_{1/n}(K)} \setminus V. \quad (33.137)$$

Nous avons par définition $d(z_n, K) \leq \frac{1}{n}$. Vu que K est compact, il comprend (au moins) un élément réalisant la distance : soit $z'_n \in K$ tel que

$$d(z_n, z'_n) = d(z_n, K). \quad (33.138)$$

Nous avons $d(z_n, z'_n) \rightarrow 0$, de telle sorte que les valeurs d'adhérence de (z_n) et (z'_n) sont les mêmes. Et comme (z'_n) est une suite dans un compact, elle a des valeurs d'adhérence. Soit z_∞ l'une d'elles. Vu que c'est une valeur d'adhérence d'une suite contenue dans le compact K , elle est également dans K : $z_\infty \in K$. Mais en même temps, z_n est hors de l'ouvert V , et donc dans le fermé V^c . Les valeurs d'adhérences restent dans le fermé, c'est-à-dire $z_\infty \notin V$. Vu que $K \subset V$, il y a contradiction.

Donc il existe $\epsilon > 0$ tel que $\overline{V_\epsilon(K)} \subset V$.

Il existe ϵ tel que $\overline{W_\epsilon} \subset V$ Il suffit de prendre le ϵ dont nous venons de parler pour avoir

$$\overline{W_\epsilon} \subset \overline{V_\epsilon(K)} \subset V. \quad (33.139)$$

Soit le ϵ en question, et $T > 0$ tel que $J \subset [-T, T]$. Nous posons $r = \epsilon e^{-LT}$. Soit $b \in \overline{B(a, r)}$ et

$$X = \{\tau \in J_+ \text{ tel que }]0, \tau] \subset J_b \text{ et } (t, y_b(t)) \in \overline{W_\epsilon} \forall t \in [0, \tau]\}. \quad (33.140)$$

Nous allons prouver que $X = J_+$ en prouvant qu'il est ouvert, fermé et non vide dans $J_+ = J \cap]0, \infty[$. Nous parlons bien de la topologie de J_+ , celle induite⁴ de \mathbb{R} . Vu que $0 \in J$, l'ensemble J_+ est ouvert à gauche, mais comme il est compact, il ne va certainement pas jusqu'à $+\infty$, de telle sorte qu'il est fermé à gauche. Les ouverts de J_+ sont les ensembles de la forme $\mathcal{O} \cap J_+$ où \mathcal{O} est ouvert de \mathbb{R} . Il y en a de la forme $]0, m]$.

X est fermé C'est parce que $\overline{W_\epsilon}$ et J_b sont fermés.

4. Définition 7.10.

X est ouvert Soit $\tau \in X$. Si $\tau = \sup J_+$ alors $X = J_+$ est un ouvert de J_+ . Supposons donc que $0 < \tau < \sup J_+$. Dans ce cas nous avons

$$(\tau, y_b(\tau)) \in \overline{W}_\epsilon \subset V, \quad (33.141)$$

et nous pouvons résoudre localement le problème de Cauchy

$$\begin{cases} y'(t) = f(t, y(t)) \end{cases} \quad (33.142a)$$

$$\begin{cases} y(\tau) = \varphi(\tau, b) = y_b(\tau). \end{cases} \quad (33.142b)$$

Ce y existe jusqu'à $\tau + \eta$ (pour au moins un petit η), et par l'unicité de la solution, $y = y_b$ sur $[\tau, \tau + \eta[$. Ceci pour dire que le flot $\varphi(\cdot, b)$ existe au moins jusqu'à $\tau + \eta$.

Grâce à la proposition 33.22 nous pouvons évaluer

$$\|\varphi(\tau, b) - \varphi(\tau, a)\| = \|y_a(\tau) - y_b(\tau)\| \leq e^{L\tau} \|b - a\|. \quad (33.143)$$

Comme nous avons choisi $r = \epsilon e^{-LT}$ et $b \in \overline{B}(a, r)$ nous avons aussi $\|b - a\| \leq \epsilon e^{-LT}$ et donc

$$\|\varphi(\tau, b) - \varphi(\tau, a)\| \leq \epsilon e^{L(\tau-T)} < \epsilon \quad (33.144)$$

parce que nous avons $\tau < \sup J_+ \leq T$, ce qui garantit que $e^{L(\tau-T)} < 1$.

Est-ce que ceci nous garantit que $\tau + \eta \in X$? Il faudrait $(\tau + \eta, y_b(\tau + \eta)) \in \overline{W}_\epsilon$, c'est-à-dire $\|y_b(\tau + \eta) - y_a(\tau + \eta)\| \leq \epsilon$. L'ensemble J_+ étant fermé dans l'ouvert J_a , ce dernier déborde certainement. Prenons donc η assez petit pour que y_a existe jusqu'en $\tau + \eta$.

Vu que y_a et y_b sont continues, et qu'en τ elles sont distantes de moins de ϵ , en $\tau + \eta$, elles restent distantes de moins de ϵ (quitte à prendre encore η plus petit).

Ceci nous permet de conclure que X est ouvert.

X est non vide La solution y_b au problème

$$\begin{cases} y'_b(t) = f(t, y_b(t)) \end{cases} \quad (33.145a)$$

$$\begin{cases} y_b(0) = b \end{cases} \quad (33.145b)$$

existe au moins localement et vérifie $\|y_b(0) - y_a(0)\| = \|b - a\| \leq \epsilon e^{-LT} < \epsilon$. Par continuité nous avons

$$\|y_b(t) - y_a(t)\| < \epsilon \quad (33.146)$$

pour tout t dans un voisinage de 0. Donc X est non vide.

Conclusion pour X La partie X est ouverte, fermée et non vide dans J_+ qui est connexe. Donc $X = J_+$ par la proposition 7.39(3).

La conclusion $X = J_+$ nous enseigne que pour tout $t \in J_+$ nous avons $]0, t] \in J_b$ et $(t, y_b(t)) \in \overline{W}_\epsilon$. Nous pouvons faire la même chose pour J_- et au final nous avons que pour tout $\tau \in J$ nous avons d'abord $\tau \in J_b$, ce qui prouve $J \subset J_b$. De plus pour tout $t \in J$ nous avons aussi

$$(t, y_b(t)) \in \overline{W}_\epsilon \subset V. \quad (33.147)$$

Nous en concluons que

$$J \times \overline{B}(a, r) \subset V. \quad (33.148)$$

Nous savons de plus que pour tout $b \in \overline{B}(a, r)$, $J \subset J_b$. Cela signifie que

$$J \times \overline{B}(a, r) \subset \mathcal{D}. \quad (33.149)$$

Mais $J \times \overline{B}(a, r)$ est un voisinage de (s, a) qui était au début de la preuve un point générique choisi dans \mathcal{D} . Donc \mathcal{D} est ouvert parce qu'il contient un voisinage de chacun de ses points.

Il nous reste à voir que $\varphi: \mathcal{D} \rightarrow \Omega$ est localement Lipschitz. Soit donc le point générique (s, a) dans \mathcal{D} et l'ensemble V qui avait été construit plus haut. Nous allons montrer que φ est Lipschitz sur $J \times \overline{B}(a, r) \subset V$. D'abord sur V , l'application f est Lipschitz, donc

$$\|\varphi(t, b_1) - \varphi(t, b_2)\| \leq e^{Lt} \|b_1 - b_2\| \leq e^{LT} \|b_1 - b_2\| \quad (33.150)$$

pour tout $t \in J$ et $b_1, b_2 \in \overline{B(a, r)}$.

Ensuite, f est bornée, majorée par M sur V , donc

$$\|\varphi(t_1, b) - \varphi(t_2, b)\| = \left| \int_{[t_1, t_2]} y'(s) ds \right| \quad (33.151a)$$

$$= \left| \int_{[t_1, t_2]} f(s, y_b(s)) ds \right| \quad (33.151b)$$

$$\leq \int_{[t_1, t_2]} |f(s, y_b(s))| ds \quad (33.151c)$$

$$\leq M|t_1 - t_2|. \quad (33.151d)$$

Et enfin nous prouvons que φ est localement Lipschitz. En posant $k = \max\{e^{LT}, M\}$ nous avons

$$\|\varphi(t_1, b_1) - \varphi(t_2, b_2)\| \leq \|\varphi(t_1, b_1) - \varphi(t_1, b_2)\| + \|\varphi(t_1, b_2) - \varphi(t_2, b_2)\| \quad (33.152a)$$

$$\leq e^{LT}\|b_1 - b_2\| + M|t_1 - t_2| \quad (33.152b)$$

$$\leq k(\|b_1 - b_2\| + |t_1 - t_2|) \quad (33.152c)$$

$$\leq 2k \max\{\|b_1 - b_2\|, |t_1 - t_2|\} \quad (33.152d)$$

$$= 2kd((b_1, t_1), (b_2, t_2)). \quad (33.152e)$$

Le flot φ est donc Lipschitz de constante $2k$. □

Exemple 33.26

Problèmes et choses à faire

Cet exemple doit être lu attentivement. Il me semble prouver que le flot n'est pas dérivable en la condition initiale sans que f le soit. Le document [454] semble dire le contraire. Je ne suis pas assez sûr de mon coup pour contredire.

Il n'y a pas de raisons de penser que $a \mapsto y_a(t)$ soit mieux que continue en sans hypothèses supplémentaires sur f . Pour illustrer cela nous considérons l'équation différentielle

$$\begin{cases} \frac{\partial X}{\partial s} = f(X(s), s) \\ X(t) = x \end{cases} \quad (33.153a)$$

$$\quad (33.153b)$$

où t et x sont des paramètres fixés. Nous allons étudier la dérivabilité de X en x lorsque

$$f(x, t) = |x|. \quad (33.154)$$

Cela est un exemple typique de fonction autant Lipschitz que l'on veut sans être dérivable. L'équation différentielle est

$$\frac{\partial X}{\partial s}(s) = |X(s)|. \quad (33.155)$$

Si $x > 0$ alors $X(s) > 0$ dans un voisinage de $s = t$ et nous avons $X(s) = Ke^s$. La constante K se fixe par la condition initiale $X(t) = x$:

$$X(s) = xe^{s-t}. \quad (33.156)$$

Et cette solution tient en réalité pour tout s parce que $X(s)$ est alors toujours positif.

Si au contraire $x < 0$ nous avons la solution

$$X(s) = xe^{t-s}. \quad (33.157)$$

Au final,

$$X(s; x, t) = \begin{cases} xe^{t-s} & \text{si } x < 0 \\ 0 & \text{si } x = 0 \\ xe^{s-t} & \text{si } x > 0 \end{cases} \quad (33.158)$$

L'application $(s, x, t) \mapsto X(s; x, t)$ est continue. En ce qui concerne la dérivée partielle $\partial_x X$ en $x = 0$ nous avons :

$$\frac{\partial X}{\partial x}(s, 0, t) = \lim_{\epsilon \rightarrow 0} \frac{X(s, \epsilon, t) - X(s, 0, t)}{\epsilon} = \lim_{\epsilon \rightarrow 0} \frac{X(s, \epsilon, t)}{\epsilon}. \quad (33.159)$$

La limite à droite donne :

$$\lim_{\epsilon \rightarrow 0^+} \frac{X(s, \epsilon, t)}{\epsilon} = \frac{\epsilon e^{s-t}}{\epsilon} = e^{s-t}. \quad (33.160)$$

La limite à gauche donne :

$$\lim_{\epsilon \rightarrow 0^-} \frac{X(s, \epsilon, t)}{\epsilon} = e^{t-s}. \quad (33.161)$$

Les deux limites n'étant pas égales, la limite (33.159) n'existe pas⁵ et l'application $(s, x, t) \mapsto X(s, x, t)$ n'est pas dérivable par rapport à x . \triangle

Lemme 33.27 ([455]).

Soit un application $A: \overline{B(t_0, \tau)} \times \overline{B(a, R)} \rightarrow \mathcal{L}(\mathbb{R}^n)$ continue par rapport à sa première variable ($t_0 \in \mathbb{R}$ et $a \in \mathbb{R}^n$). Alors en posant l'équation

$$\begin{cases} \frac{\partial \psi}{\partial t}(t, b) = A(t, b)\psi(t, b) \\ \psi(t_0, b) = \psi_0. \end{cases} \quad (33.162a)$$

$$\psi(t_0, b) = \psi_0. \quad (33.162b)$$

Nous avons l'estimation

$$\begin{aligned} \|\psi(t, v) - \psi(t, w)\| &\leq \|\psi_0\| \tau \max_{s \in \overline{B(t_0, \tau)}} \|A(s, v) - A(s, w)\| \times \\ &\quad \times \exp \left(\tau \max_{s \in \overline{B(t_0, \tau)}} \max\{\|A(s, v)\|, \|A(s, w)\|\} \right) \end{aligned} \quad (33.163)$$

pour tout $t \in \overline{B(t_0, \tau)}$ et $v, w \in V$.

Théorème 33.28 (Régularité C^1 du flot [455]).

Soit un intervalle ouvert I de \mathbb{R} et un ouvert connexe Ω de \mathbb{R}^n . Soit une fonction $f \in C^1(I \times \Omega, \mathbb{R})$, $a \in \Omega$ et $t_0 \in I$.

Il existe un voisinage $W \times V = \overline{B(t_0, \tau)} \times \overline{B(a, r)}$ de (t_0, a) dans $I \times \Omega$ et une unique application $\varphi: W \times V \rightarrow \Omega$ telle que

$$\begin{cases} \frac{\partial \varphi}{\partial t}(t, x) = f(t, \varphi(t, x)) \\ \varphi(t_0, x) = x \end{cases} \quad (33.164a)$$

$$\varphi(t_0, x) = x \quad (33.164b)$$

pour tout $x \in V$.

L'application $(t, x) \mapsto \varphi(t, x)$ est de classe C^1 .

Problèmes et choses à faire

La preuve qui suit doit être lue avec beaucoup d'attention, en particulier sur les incohérences possibles de notations, et sur les oublis possibles de précautions oratoires type « quitte à encore réduire les voisinages V et W ».

Démonstration. En termes de notations, pour $x \in \Omega$ fixé nous écrivons $y_x(t)$ pour $\varphi(t, x)$ et pour $t \in I$ fixé nous notons $\varphi_t(x)$ pour $\varphi(t, x)$.

De plus lorsque nous écrirons des choses comme $g: \mathbb{R} \rightarrow \mathbb{R}$, nous n'entendons pas que g est effectivement définie sur tout \mathbb{R} . La notation « $g: \mathbb{R} \rightarrow \mathbb{R}$ » indiquera seulement que la variable de g est réelle, et que nous comptons préciser le domaine plus tard. Cette remarque s'applique seulement à cette démonstration et non à l'ensemble du livre.

5. Si vous comptez donner ça à manger au jury d'un concours, soyez prudent et n'écrivez pas l'équation (33.159) au tableau. Réfléchissez comment rédiger cela correctement.

Nous considérons $R > 0$ tel que $\overline{B(a, 2R)} \subset \Omega$ et ensuite nous posons $V = \overline{B(a, R)}$. La fonction y_x , solution pour la condition initiale $y_x(t_0) = x$ est définie sur $W = [t_0 - \tau, t_0 + \tau]$ et prend ses valeurs dans $\overline{B(x, R)}$. Ceci est parce que y_x est continue, alors en prenant τ assez petit, la valeur de $y_x(t)$ ne va pas s'éloigner de x lorsque t ne s'éloigne pas de t_0 .

Nous savons déjà de la proposition 33.25 que φ est C^1 en t et localement Lipschitz en sa seconde variable, avec une constante Lipschitz uniforme sur $W \times V$. Elle est donc continue en tant que fonction

$$\varphi: V \times W \rightarrow \mathbb{R}^d. \tag{33.165}$$

La différentielle partielle Df Pour t fixé nous notons $Df_{(t,x)}$ la différentielle de f par rapport à x . C'est-à-dire que

$$Df_{(t,x)}: \mathbb{R}^n \rightarrow \mathbb{R}^n$$

$$u \mapsto \frac{d}{ds} \left[f(t, x + su) \right]_{s=0}. \tag{33.166}$$

C'est un élément de $\mathcal{L}(\mathbb{R}^n)$, l'ensemble des applications linéaires de \mathbb{R}^n vers \mathbb{R}^n . Nous allons montrer que

$$(t, x) \mapsto Df_{(t,x)} \tag{33.167}$$

est continue en tant qu'application $\mathbb{R} \times \mathbb{R}^n \rightarrow \mathcal{L}(\mathbb{R}^n)$. Pour cela nous introduisons l'application d'inclusion $i: \mathbb{R}^n \rightarrow \mathbb{R} \times \mathbb{R}^n$, $i(u) = (0, u)$. Elle donne

$$Df_{(t,x)}(u) = \frac{d}{ds} \left[f((t, x) + s(0, u)) \right]_{s=0} = df_{(t,x)} \circ i(u). \tag{33.168}$$

Autrement dit

$$Df_{(t,x)} = df_{(t,x)} \circ i. \tag{33.169}$$

Or l'application $(t, x) \mapsto df_{(t,x)}$ est continue par hypothèse (f est de classe C^1) et l'application

$$\begin{aligned} \mathcal{L}(\mathbb{R} \times \mathbb{R}^n, \mathbb{R}^n) &\rightarrow \mathcal{L}(\mathbb{R}^n, \mathbb{R}^n) \\ A &\mapsto A \circ i \end{aligned} \tag{33.170}$$

est également continue. Donc $(t, x) \mapsto Df_{(t,x)}$ est continue⁶.

L'équation aux variations Soit $x \in \Omega$. Nous introduisons l'opérateur

$$S_x: \mathbb{R} \times \mathcal{L}(\mathbb{R}^n) \rightarrow \mathcal{L}(\mathbb{R}^n)$$

$$S_x(t, \psi) = Df_{(t, y_x(t))} \circ \psi. \tag{33.171}$$

Par ce que nous avons raconté, cela est une fonction continue en sa première variable et Lipschitz en sa seconde variable. Nous identifions $\mathcal{L}(\mathbb{R}^n)$ à \mathbb{R}^{2^n} .

Toujours pour chaque x considéré nous posons l'équation différentielle ordinaire

$$\begin{cases} \frac{\partial \psi}{\partial t}(t, x) = S_x(t, \psi(t, x)) \\ \psi(t_0, x) = \text{Id}. \end{cases} \tag{33.172a}$$

$$\tag{33.172b}$$

qui est une équation différentielle ordinaire pour $\psi: \mathbb{R} \times \mathbb{R}^n \rightarrow \mathcal{L}(\mathbb{R}^n)$ rentrant dans le cadre de Cauchy-Lipschitz.

Quel est le domaine de définition de ψ pour sa première variable? C'est un ouvert autour de t_0 . Nous réduisons W de telle sorte que la solution ψ soit définie sur W . Idem pour la variable x qui est dans un voisinage de a .

L'équation (33.172) s'appelle l'**équation aux variations**. Nous allons montrer dans la douleur que ψ est continue et est la différentielle de φ_t , c'est-à-dire que

$$(d\varphi_t)_b = \psi(t, b). \tag{33.173}$$

6. Si quelqu'un peut prouver ça de façon moins verbeuse, je suis preneur. Il me semble que quel que soit la façon dont on s'y prend, sous le capot, on passe par la continuité de l'application (33.170).

ψ est continue en (t, x) (début) Il s'agit de majorer les deux termes de

$$\|\psi(t_1, a_1) - \psi(t_2, a_2)\| \leq \|\psi(t_1, a_1) - \psi(t_2, a_1)\| + \|\psi(t_2, a_1) - \psi(t_2, a_2)\|. \quad (33.174)$$

Premier terme Nous avons

$$\|\psi(t_1, b) - \psi(t_2, b)\| = \left\| \int_{[t_1, t_2]} \frac{\partial \psi}{\partial t}(s, b) ds \right\| \quad (33.175a)$$

$$\leq \int_{[t_1, t_2]} \|Df_{(s, y_b(s))} \circ \psi(s, b)\| ds \quad (33.175b)$$

$$\leq \int_{[t_1, t_2]} \|Df_{(s, t_b(s))}\| \|\psi(s, b)\| ds \quad (33.175c)$$

$$\leq |t_1 - t_2| \max_{s \in [t_1, t_2]} \|Df_{(s, y_b(s))}\| \max_{s \in [t_1, t_2]} \|\psi(s, b)\|. \quad (33.175d)$$

Nous allons majorer le second maximum. Prenons $t \in [0, \tau]$; et posons $A(u, b) = Df_{(u, y_b(u))}$ pour alléger les notations. Par l'équation de définition de ψ nous avons

$$\psi(t, b) = \psi(0, b) + \int_{[0, t]} A(u, b) \psi(u, b) du, \quad (33.176)$$

et donc

$$\|\psi(t, b)\| \leq \|\psi_0\| + \int_{[0, t]} \|A(u, b)\| \|\psi(u, b)\| du. \quad (33.177)$$

En y appliquant le lemme de Grönwall dans sa version 33.4 nous trouvons

$$\|\psi(s, b)\| \leq \|\psi_0\| \exp \left(\int_{[0, s]} \|A(u, b)\| du \right) \quad (33.178a)$$

$$\leq \|\psi_0\| \exp \left(s \max_{u \in [0, s]} \|A(u, b)\| \right). \quad (33.178b)$$

En retournant à (33.175d) nous avons $\psi_0 = \text{Id}$ et donc $\|\psi_0\| = 1$ et

$$\max_{s \in [t_1, t_2]} \|\psi(s, b)\| \leq \max_{s \in [t_1, t_2]} \exp \left(s \max_{u \in [0, t]} \|Df_{(u, y_b(u))}\| \right) \quad (33.179)$$

Là dedans nous pouvons remplacer t par $\max\{|t_1|, |t_2|\}$. Posons enfin, pour alléger les expressions

$$a(t_1, t_2, b) = \max_{s \in [t_1, t_2]} \|Df_{(s, y_b(s))}\|. \quad (33.180)$$

La majoration que nous retenons est :

$$\|\psi(t_1, b) - \psi(t_2, b)\| \leq |t_1 - t_2| a(t_1, t_2, b) \exp \left(\max\{|t_1|, |t_2|\} a(0, t, b) \right). \quad (33.181)$$

Cela tend vers zéro lorsque $t_1 \rightarrow t_2$.

Deuxième terme En ce qui concerne le second terme,

$$\|\psi(t, b_1) - \psi(t, b_2)\| \quad (33.182)$$

nous utilisons le lemme 33.27 qui donne, pour $t \in [t_0 - \tau, t_0 + \tau]$,

$$\begin{aligned} \|\psi(t, b_1) - \psi(t, b_2)\| &\leq \tau \max_{s \in B(0, \tau)} \|Df_{s, y_{b_1}(s)} - Df_{s, y_{b_2}(s)}\| \times \\ &\quad \times \exp \left(\tau \max\{\|Df_{s, y_{b_1}(s)}\|, \|Df_{s, y_{b_2}(s)}\|\} \right). \end{aligned} \quad (33.183)$$

Dans notre cas, $t_0 = 0$, donc $t \in [-\tau, \tau]$. Vu la continuité de Df , nous avons

$$\max_{s \in B(0, \tau)} \|Df_{s, y_{b_1}(s)} - Df_{s, y_{b_2}(s)}\| \rightarrow 0 \quad (33.184)$$

lorsque $b_1 \rightarrow b_2$.

ψ est continue en (t, x) (fin) Les deux bons calculs faits, nous avons, en repartant de (33.174),

$$\lim_{(t_1, b_1) \rightarrow (t_2, b_2)} \|\psi(t_1, b_1) - \psi(t_2, b_2)\| = 0, \quad (33.185)$$

ce qui signifie que ψ est une fonction continue de ses deux variables en même temps.

Différentiabilité de φ (début) Nous montrons maintenant que $D\varphi(t, x)$ existe. Pour rappel, D est la différentielle par rapport à la seconde variable. Nous sommes à étudier l'existence de $D\varphi_{(t,b)} = d(\varphi_t)_b$. Nous posons

$$\theta(t, h) = \varphi(t, b + h) - \varphi(t, b) = y_{b+h}(t) - y_b(t) \quad (33.186)$$

où b est le point où nous étudions la différentiabilité. Il est dans un voisinage du point a fixé depuis le début et autour duquel il existe un voisinage qui donne un sens à tout ce que nous avons fait jusqu'à présent. La dépendance de θ en b est implicite. Vu que φ est Lipschitz en sa seconde variable, nous avons la majoration

$$\|\theta(t, h)\| \leq C\|h\| \quad (33.187)$$

dès que $t \in V$ et $b, b + h \in W$.

De plus, parce que t_0 est le temps de la condition initiale nous avons

$$\theta(t_0, h) = y_{b+h}(t_0) - y_b(t_0) = a + h - a = h. \quad (33.188)$$

Et aussi, par définition de ψ :

$$\psi(t, b) = \psi_0 + \int_{t_0}^t \frac{\partial \psi}{\partial t}(s, b) ds = \psi_0 + \int_{t_0}^t Df_{(s, y_b(s))} \circ \psi(s, b) \quad (33.189)$$

En appliquant à h et en se souvenant que $\psi_0 = \text{Id}$,

$$\psi(t, b)h = h + \int_{t_0}^t \left(Df_{s, y_b(s)} \circ \psi(s, b) \right) h ds. \quad (33.190)$$

Puis on peut faire un calcul assez classique en se souvenant que $\theta(t_0, h) = h$:

$$\theta(t, h) = \theta(t_0, h) + \int_{t_0}^t \left[\frac{\partial \varphi}{\partial t}(s, b + h) - \frac{\partial \varphi}{\partial t}(s, b) \right] ds \quad (33.191a)$$

$$= h + \int_{t_0}^t [f(s, y_{b+h}(s)) - f(s, y_b(s))] ds. \quad (33.191b)$$

On fait la différence entre les deux :

$$\theta(t, h) - \psi(t, b)h = - \int_{t_0}^t [Df_{s, y_b(s)} \circ \psi(s, b)h - f(s, y_{b+h}(s)) + f(s, y_b(s))] ds. \quad (33.192)$$

Nous y ajoutons et soustrayons $Df_{s, y_b(s)}\theta(s, h)$ et nous retenons la majoration suivante :

$$\begin{aligned} \|\theta(t, h) - \psi(t, b)h\| &\leq \int_{t_0}^t \|Df_{(s, y_b(s))}\psi(s, b) - Df_{(s, y_b(s))}\theta(s, h)\| ds \\ &\quad + \int_{t_0}^t \|f(s, y_{b+h}(s)) - f(s, y_b(s)) + Df_{(s, y_b(s))}\theta(s, h)\| ds. \end{aligned} \quad (33.193)$$

Nous allons encore majorer ces deux termes séparément. Soit $\epsilon > 0$.

Premier terme Ce qui est dans la norme à majorer est

$$Df_{(s,y_b(s))}(\psi(s,b)h - \theta(s,h)). \quad (33.194)$$

Vu que Df est continue et que y_b est continue⁷, l'application $s \mapsto Df_{s,y_b(s)}$ est continue et donc de norme majorée sur le compact $[t_0, t]$. Nous rapellons la notation

$$a(t_0, t, b) = \max_{s \in [t_0, t]} \|Df_{s,y_b(s)}\|, \quad (33.195)$$

et nous majorons encore et toujours. D'abord

$$\int_{t_0}^t \|Df_{(s,y_b(s))}(\psi(s,b)h - \theta(s,h))\| ds \leq a(t_0, t, b) \int_{t_0}^t \|\psi(s,b) - \theta(s,h)\| ds. \quad (33.196)$$

Deuxième terme Pour traiter le deuxième terme, nous allons provisoirement noter $x = y_b(s)$ et $y = y_{b+h}(s)$; entre autres, $y - x = \theta(s,h)$. Ce qui est écrit dans le second terme de (33.193) est

$$f(s,y) - f(s,x) + Df_{(s,x)}\theta(s,h) = f(s,y) - f(s,x) + Df_{(s,x)}(y - x) \quad (33.197)$$

Comme D ne s'applique pas à la variable s , nous pouvons alléger la notation et déduire de la différentiabilité de f qu'il existe un $\eta > 0$ tel que $x, y \in W$ avec $\|y - x\| \leq \eta$ implique

$$\|f(y) - f(x) - Df_x(y - x)\| \leq \epsilon \|y - x\|. \quad (33.198)$$

Prenons $\|h\| \leq \eta/C$ (le C de (33.187)); en déballant les notations,

$$\|f(s, y_{b+h}(s)) - f(s, y_b(s)) - Df_{(s,y_b(s))}\theta(s,h)\| \leq \epsilon \|\theta(s,h)\| \leq \epsilon C \|h\|. \quad (33.199)$$

Les deux termes ensemble En remettant les deux dans (33.193) nous trouvons la majoration

$$\|\theta(t,h) - \psi(t,b)h\| \leq |t - t_0| \epsilon C \|h\| + a(t_0, t, v) \int_{t_0}^t \|\psi(s,b)h - \theta(s,h)\| ds \quad (33.200)$$

qui est encore de la graine à Grönwall avec

$$\begin{cases} u(t) = \|\theta(t,h) - \psi(t,b)h\| & (33.201a) \\ b(t) = |t - t_0| \epsilon C \|h\| & (33.201b) \\ a(s) = a(t_0, t, v), & (33.201c) \end{cases}$$

la troisième étant une fonction constante. Cela donne, pour $t \in [t_0 - \tau, t_0 + \tau]$,

$$\|\theta(t,h) - \psi(t,b)h\| \leq |t - t_0| \epsilon C \|h\| + \int_{t_0}^t (s - t_0) \epsilon C \|h\| a(t_0, t, b) \exp\left(\int_s^t a(t_0, t, b) du\right) ds. \quad (33.202)$$

En valeur absolue, la différence $s - t_0$ est majorée par τ , l'intégrale dans l'exponentielle vaut $(t - s)a(t_0, t, b)$, et restons avec

$$\|\theta(t,h) - \psi(t,b)h\| \leq \tau \epsilon C \|h\| + \tau \int_{t_0}^t \epsilon C \|h\| a(t_0, t, b) e^{a(t_0, t, b)(t-s)} ds. \quad (33.203)$$

En supposant $t > t_0$ nous pouvons calculer l'intégrale. Si vous m'avez suivi jusqu'ici, vous devriez avoir de tels maux de tête que je vous donne la réponse :

$$\int_{t_0}^t a(t_0, t, b) e^{(t-s)a(t_0, t, b)} ds = e^{(t_0-t)a(t_0, t, b)} - 1. \quad (33.204)$$

7. Il faut encore réduire les voisinages V et W pour que ceci ait un sens.

En remettant dans l'expression,

$$\|\theta(t, h) - \psi(t, b)h\| \leq \tau \epsilon C \|h\| + \epsilon C \|h\| a(t_0, t, b) \tau (e^{(t-t_0)} - 1) = \tau \epsilon C \|h\| e^{(t-t_0)a(t, t_0, b)}. \tag{33.205}$$

Nous pouvons majorer $t - t_0$ par τ et $a(t, t_0, b)$ par $a(t_0 - \tau, t_0 + \tau, b)$ pour avoir la majoration

$$\|\theta(t, h) - \psi(t, b)h\| \leq \tau \epsilon C \|h\| e^{\tau a(t_0 - \tau, t_0 + \tau, b)}. \tag{33.206}$$

Différentiabilité de $\varphi(t, b)$ (fin) Nous écrivons la définition 12.133 de la différentiabilité : nous voulons vérifier que

$$\lim_{h \rightarrow 0} \frac{\varphi(t, b+h) - \varphi(t, b) - \psi(t, b)h}{\|h\|} = 0. \tag{33.207}$$

Nous remplaçons $\varphi(t, b+h) - \varphi(t, b)$ par $\theta(t, h)$ et prenons la norme avec les majorations données :

$$\lim_{h \rightarrow 0} \frac{\|\varphi(t, b+h) - \varphi(t, b) - \psi(t, b)h\|}{\|h\|} \leq \lim_{h \rightarrow 0} \tau \epsilon C e^{\tau a(t_0 - \tau, t_0 + \tau, b)}. \tag{33.208}$$

Cela étant valable pour tout ϵ , nous en déduisons la nullité de la limite.

Nous avons démontré que φ était différentiable par rapport à sa deuxième variable et que

$$D\varphi_{(t,b)} = \psi(t, b). \tag{33.209}$$

Conclusion : φ est de classe C^1 Nous avons déjà prouvé que $(t, b) \mapsto \psi(t, b)$ est continue. Donc de (33.209) nous déduisons que les dérivées partielles $(t, b) \mapsto \frac{\partial \varphi}{\partial x_i}(t, b)$ sont continues. Mais comme φ est Lipschitz en t , la dérivée partielle $(t, b) \mapsto \frac{\partial \varphi}{\partial t}(t, b)$ est également continue. La continuité de toutes les dérivées partielles de φ nous donne la classe C^1 pour φ par la proposition 13.239.

□

Proposition 33.29 (Régularité C^p du flot[456]).

Soit un intervalle ouvert I de \mathbb{R} et un ouvert connexe Ω de \mathbb{R}^n . Soit une fonction $f \in C^p(I \times \Omega, \mathbb{R})$, ainsi que $a \in \Omega$ et $t_0 \in I$.

Il existe un voisinage $W \times V = \overline{B}(t_0, \tau) \times \overline{B}(a, r)$ de (t_0, a) dans $I \times \Omega$ et une unique application $\varphi : W \times V \rightarrow \Omega$ telle que

$$\begin{cases} \frac{\partial \varphi}{\partial t}(t, x) = f(t, \varphi(t, x)) & (33.210a) \\ \varphi(t_0, x) = x & (33.210b) \end{cases}$$

pour tout $x \in V$.

L'application $(t, x) \mapsto \varphi(t, x)$ est de classe C^p .

Démonstration. Nous savons déjà par le théorème 33.28 que $(t, x) \mapsto \varphi(t, x)$ est de classe C^1 . Nous supposons que f est de classe C^p avec $p \geq 2$.

Vu que φ et f sont de classe C^1 , nous avons aussi que l'application $(t, x) \mapsto f(t, \varphi(t, x))$ est de classe C^1 . L'équation donne alors immédiatement le fait que

$$(t, x) \mapsto \frac{\partial \varphi}{\partial t}(t, x) \tag{33.211}$$

est de classe C^1 .

En ce qui concerne la régularité par rapport aux autres variables, il faudra travailler un peu plus.

Une équation différentielle pour le flot Nous allons commencer par un habile jeu d'écriture :
la formule

$$\varphi(t, x) = x + \int_{t_0}^t f(s, \varphi(s, x)) ds \quad (33.212)$$

devient

$$\varphi_t(x) = x + \int_{t_0}^t f(s, \varphi_s(x)) ds. \quad (33.213)$$

Dans le même ordre d'idée nous notons $f_s(x) = f(s, x)$, et ce qui se trouve dans l'intégrale (33.213) n'est autre que la fonction

$$g_s(x) = (f_s \circ \varphi_s)(x). \quad (33.214)$$

Tout cela pour différentier l'égalité (33.213) par la proposition 18.26 :

$$(d\varphi_t)_x = \text{Id} + \int_{t_0}^t (dg_s)_x ds \quad (33.215a)$$

$$= \text{Id} + \int_{t_0}^t (df_s)_{\varphi_s(x)} \circ (d\varphi_s)_s ds. \quad (33.215b)$$

Nous dérivons ensuite cela par rapport à t :

$$\frac{\partial}{\partial t} \left((d\varphi_t)_x \right) = (df_t)_{\varphi_t(x)} \circ (d\varphi_t)_x. \quad (33.216)$$

Cela est une égalité dans $\mathcal{L}(\mathbb{R}^n)$.

Nous introduisons la fonction

$$F: I \times \mathcal{L}(\mathbb{R}^n) \times \Omega \rightarrow \mathcal{L}(\mathbb{R}^n) \quad (33.217)$$

$$(t, A, x) \mapsto (df_t)_{\varphi_t(x)} \circ A.$$

En fait, à la place de I et Ω il faut prendre des petits voisinages dans lesquels les choses ont un sens. Ce que dit l'équation 33.216 est que l'application

$$A: I \times \Omega \rightarrow \mathcal{L}(\mathbb{R}^n) \quad (33.218)$$

$$(t, x) \mapsto (d\varphi_t)_x$$

vérifie l'équation différentielle

$$\begin{cases} \frac{\partial A}{\partial t}(t, x) = F(t, A(t, x), x) & (33.219a) \\ A(t_0, x) = \text{Id}. & (33.219b) \end{cases}$$

Une autre équation différentielle Nous n'oublions pas l'équation différentielle pour la dérivée par rapport à t :

$$\frac{\partial \varphi}{\partial t}(t, x) = f(t, \varphi(t, x)). \quad (33.220)$$

Réécriture pour la différentielle Nous allons récrire l'équation (33.219) de façon à ce que le paramètre x soit inclus dans la condition initiale. De cette manière, la solution pourra profiter de la régularité C^1 du flot déjà prouvée dans le théorème 33.28.

Soit

$$g: I \times (\mathcal{L}(\mathbb{R}^n) \times \Omega) \rightarrow \mathcal{L}(\mathbb{R}^n) \times \Omega \quad (33.221)$$

$$(t, (A, x)) \mapsto (F(t, A, x), 0).$$

Nous posons $E = \mathcal{L}(\mathbb{R}^n) \times \Omega$; c'est cet espace qui va jouer le rôle de Ω . Nous considérons à présent l'équation différentielle suivante pour $z_x: I \rightarrow E$:

$$\begin{cases} \begin{pmatrix} z'_1(t) \\ z'_2(t) \end{pmatrix} = z'(t) = g(t, z(t)) = \begin{pmatrix} F(t, z_1(t), z_2(t)) \\ 0 \end{pmatrix} & (33.222a) \\ z(t_0) = (\text{Id}, x). & (33.222b) \end{cases}$$

Il devrait y avoir un indice (Id, x) à z parce que c'est sa condition initiale. La fonction g est de classe C^p , donc cette équation admet une unique solution dont le flot est de classe C^1 . Autrement dit, si S est dans un voisinage de Id , l'application

$$(t, x) \mapsto z(t) \tag{33.223}$$

est de classe C^1 . Nous allons montrer qu'en posant $A(t, x) = z_1(t)$, nous avons une solution de (33.219) (l'unicité de la solution impose que cette solution est effectivement la différentielle de $d\varphi_t$). D'abord, la seconde ligne de l'équation différentielle est $z'_2(t) = 0$, c'est-à-dire $z_2(t) = x$ pour tout t .

Sachant cela, la première équation devient

$$\begin{cases} z'_1(t) = F(t, z_1(t), x) \\ z_1(t_0) = \text{Id}, \end{cases} \tag{33.224a}$$

$$\tag{33.224b}$$

qui est l'équation différentielle pour A . Rappel : il y a partout une dépendance de z en sa condition initiale x que nous n'avons pas écrite pour des raisons de légèreté notionnelle. Il n'en reste pas moins que le flot de l'équation différentielle pour z est C^1 , c'est-à-dire que $(t, x) \mapsto z_1(t)$ est de classe C^1 .

Par conséquent, $(t, x) \mapsto A(t, x)$ est également C^1 .

Régularité C^2 du flot Le fait que A soit C^1 n'implique pas que le flot le soit parce que le flot suit la même équation différentielle que A ne signifie pas que il soit égal. Il y a un raisonnement à faire.

Le fait est que si A est une solution de (33.219), alors $z(t) = (A(t, x), x)$ est solution de (33.222). C'est l'unicité de cette dernière qui permet de déduire l'unicité de la solution pour A .

Nous avons donc que l'unique solution A du système (33.219) est égale à $A(t, x) = (d\varphi_t)_x$ et est de classe C^1 par rapport à (t, x) .

Donc $(t, x) \mapsto \varphi(t, x)$ est de classe C^2 .

Régularité C^p Nous avons vu que le flot de $y' = f(t, y)$ est de classe C^2 dès que f est de classe C^2 . Supposons que f soit de classe C^p et montrons que si le flot est de classe C^k ($k < p$) alors il est de classe C^{k+1} .

Vu que le flot d'une équation différentielle de classe C^p est de classe C^k , en particulier celui de (33.222) est de classe C^k . Donc aussi la solution pour $A(t, x) = (d\varphi_t)_x$ est de classe C^k . Et vu que $(t, x) \mapsto (d\varphi_t)_x$ est de classe C^k , l'application φ est de classe C^{k+1} .

□

33.30.

Le théorème d'inversion locale 18.48 nous permet de dire que, pour t fixé, le flot $x \mapsto \varphi_t(x)$ est un C^p -difféomorphisme local.

Proposition 33.31 (Cauchy-Lipschitz avec paramètre, régularité C^p [457, 458, 459]).

Soit un intervalle ouvert I de \mathbb{R} , un connexe ouvert Ω de \mathbb{R}^n et un intervalle ouvert Λ de \mathbb{R}^d . Soit une fonction $f \in C^p(I \times \Omega \times \Lambda, \mathbb{R}^n)$ localement Lipschitz en Ω . Soient $t_0 \in I$, $y_0 \in \Omega$ et $\lambda_0 \in \Lambda$. Il existe un voisinage compact de (t_0, y_0, λ_0) sur lequel le problème

$$\begin{cases} y'_\lambda(t) = f(t, y_\lambda(t), \lambda) \\ y_\lambda(t_0) = y_0 \end{cases} \tag{33.225a}$$

$$\tag{33.225b}$$

possède une unique solution. De plus $(t, \lambda) \mapsto y_\lambda(t)$ est de classe C^p par rapport à ses deux variables.

Démonstration. Nous récrivons immédiatement le problème pour la fonction $y: I \times \Lambda \rightarrow \mathbb{R}^n$ donné par $y(t, \lambda) = y_\lambda(t)$:

$$\begin{cases} \frac{\partial y}{\partial t}(t, \lambda) = f(t, y(t, \lambda), \lambda) \\ y(t_0) = y_0. \end{cases} \tag{33.226a}$$

$$\tag{33.226b}$$

Nous allons montrer que ce problème est en réalité équivalent à un problème sans paramètre. Nous posons $E = \Omega \times \Lambda$ et

$$\begin{aligned} g: I \times E &\rightarrow E \\ (t, x) &\mapsto (f(t, x_1, x_2), 0) \end{aligned} \quad (33.227)$$

où x_1 est la composante Ω de x et x_2 est la composante Λ de x . Pour une valeur $\mu \in \Lambda$ donnée nous considérons le problème au condition initiales

$$\begin{cases} x'(t) = g(t, x(t)) \\ x(t_0) = (y_0, \mu). \end{cases} \quad (33.228a)$$

$$(33.228b)$$

Le théorème de Cauchy-Lipschitz que nous prenons sous la forme 33.29 nous indique que ce problème admet une unique solution maximale et que le flot $(t, (y_0, \mu)) \mapsto x_{(y_0, \mu)}(t)$ est de classe C^p .

Nous passons maintenant à la résolution du problème (33.226)

Existence d'une solution C^1 Nous montrons à présent que la fonction y donnée par

$$y(t, \mu) = x_{(t_0, \mu)}(t)_1 \quad (33.229)$$

est solution de (33.226). Vu que x a deux composantes, nous pouvons un peu débâiller l'équation. Afin d'éviter les notations laborieuses nous allons noter x pour $x_{(t_0, \mu)}$ et donc $x_1(t)$ pour $x_{(t_0, \mu)}(t)$. Nous avons l'équation différentielle

$$\begin{pmatrix} x'_1(t) \\ x'_2(t) \end{pmatrix} = \begin{pmatrix} f(t, x_1(t), x_2(t)) \\ 0 \end{pmatrix} \quad (33.230)$$

avec la condition initiale

$$\begin{pmatrix} x_1(t_0) \\ x_2(t_0) \end{pmatrix} = \begin{pmatrix} y_0 \\ \mu \end{pmatrix}. \quad (33.231)$$

La seconde ligne de l'équation donne immédiatement $x_2(t) = \mu$ pour tout t . En injectant dans la première ligne :

$$x'_1(t) = f(t, x_1(t), \mu). \quad (33.232)$$

Or vue la définition de y , le nombre $x'_1(t)$ n'est autre que $\frac{\partial y}{\partial t}(t, \mu)$. La fonction y que nous avons définie vérifie donc

$$\frac{\partial y}{\partial t}(t, \mu) = f(t, y(t, \mu), \mu) \quad (33.233)$$

et la condition initiale $y(t_0) = x_1(t_0) = y_0$. Elle est donc bien solution du problème initial.

De plus l'application $(t, \mu) \mapsto y(t, \mu) = x_{(t_0, \mu)}(t)_1$ est de classe C^p .

Unicité Pour l'unicité, soit on invoque la proposition 18.43 qui donne l'unicité dans les fonctions continues et a fortiori dans les fonctions C^1 . Soit on fait le jeu inverse : on prouve qu'à chaque solution de (33.226) correspond une solution de (33.228), et l'unicité de la solution x donne l'unicité du côté de y .

□

Lemme 33.32.

Soit le problème

$$\begin{cases} \frac{\partial y}{\partial s}(s) = f(y(s), s) \\ y(t) = x \end{cases} \quad (33.234a)$$

$$(33.234b)$$

avec t et x fixés. Nous supposons que f est de classe C^p .

Alors l'application $t \mapsto y_x(s)$ est de classe C^p .

Démonstration. Soit t fixé, et l'équation différentielle

$$\begin{cases} \frac{\partial z}{\partial s}(s) = f(z(s), t-s) & (33.235a) \\ z(0) = x. & (33.235b) \end{cases}$$

Par le théorème 33.29, La solution z est de classe C^p en (s, x) . En posant $y(s) = z(t-s)$ il est vite vérifié que y est solution de (33.234). C'est alors bien de classe C^p en t . \square

33.8.4 Stabilité de Lyapunov

Définition 33.33.

Dans le cas de l'équation différentielle $y'(t) = f(y(t), t)$ pour $y: \mathbb{R} \rightarrow \mathbb{R}^n$, un point $a \in \mathbb{R}^n$ est un **point d'équilibre** lorsque la fonction constante $y(t) = a$ est une solution.

Le point d'équilibre $a \in \mathbb{R}^n$ est **stable** si pour tout $\epsilon > 0$, il existe $\delta > 0$ tel que $\|y(0) - a\| < \delta$ implique $\|y(t) - a\| < \epsilon$ pour tout t .

Théorème 33.34 (Théorème de stabilité de Lyapunov[1, 460, 168, 452]).

Soit l'équation différentielle

$$\begin{cases} y'(t) = f(y) & (33.236a) \\ y(0) = y_0 & (33.236b) \end{cases}$$

avec une fonction $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ de classe C^1 vérifiant $f(0) = 0$ et $y_0 \in \mathbb{R}^n$. Nous supposons que l'application linéaire df_0 n'a que des valeurs propres dont la partie réelle est strictement négative.

Alors

- (1) Il existe $k > 0$ tel que si $\|y_0\| < k$ alors la solution maximale est définie sur \mathbb{R} ,
- (2) pour le même nombre $k > 0$, si $\|y_0\| < k$ alors $y(t) \xrightarrow{t \rightarrow \infty} 0$ exponentiellement vite,
- (3) la solution $y = 0$ est un point d'équilibre attractif.

Démonstration. Placer ici une phrase intelligente ⁸.

Prolégomène Le théorème de Cauchy-Lipschitz 18.40 nous enseigne que l'équation différentielle considérée possède une unique solution maximale (entre autres parce qu'une fonction de classe C^1 est localement Lipschitz) et nous nommons J l'intervalle sur lequel elle est définie.

Système linéarisé Nous posons $A = df_0$. La fonction $y_L(t) = e^{tA}y_0$ est solution du système linéarisé

$$\begin{cases} y'(t) = Ay(t) & (33.237a) \\ y(0) = y_0. & (33.237b) \end{cases}$$

Pour évaluer la norme de y_L nous utilisons le lemme 16.117 : il existe un polynôme P tel que

$$\|y_L(t)\| \leq P(|t|) \sum_{i=1}^r e^{\operatorname{Re} \lambda_i t} \|y_0\|. \quad (33.238)$$

Mais par hypothèse, $\operatorname{Re}(\lambda_i) < 0$ et si nous posons $\lambda = \max\{\operatorname{Re}(\lambda_i)\}$ nous avons $\lambda < 0$ et

$$\|y_L(t)\| \leq P(|t|) e^{\lambda t} \|y_0\|. \quad (33.239)$$

Donc quel que soit y_0 nous avons $\lim_{t \rightarrow \infty} \|y_L(t)\| = 0$ c'est-à-dire $\lim_{t \rightarrow \infty} y_L(t) = 0$.

8. Parce que sinon l'environnement description qui suit donne un mauvais effet.

Une forme linéaire Nous définissons la forme bilinéaire suivante sur \mathbb{R}^n :

$$b(x, y) = \int_0^\infty \langle e^{tA}x, e^{tA}y \rangle dt. \quad (33.240)$$

D'abord cela est bien défini pour tout $x, y \in \mathbb{R}^n$ parce que

$$|\langle e^{tA}x, e^{tA}y \rangle| \leq \|e^{tA}x\| \|e^{tA}y\| \leq P_1(|t|)P_2(|t|)e^{2\lambda t} \|x\| \|y\|, \quad (33.241)$$

qui est intégrable entre 0 et ∞ à cause de la décroissance exponentielle⁹. Montrons que b est définie positive. Soit donc $x \neq 0$ et calculons

$$b(x, x) = \int_0^\infty \|e^{tA}x\|^2 dt. \quad (33.242)$$

Ce qui est dans l'intégrale est forcément (pas strictement) positif pour tout t . Mais si $x \neq 0$ alors $\|x\|^2$ est strictement positif et sur un voisinage de $t = 0$ nous avons aussi $\|e^{tA}x\|^2$ qui est strictement positif. Ergo $b(x, x) > 0$ dès que $x \neq 0$, ce qui signifie que b est strictement définie positive (lemme 11.196).

Nous notons $q: V \rightarrow \mathbb{R}$ la forme quadratique associée à b et aussi la norme qui va avec : $\|x\|_q = \sqrt{q(x)}$. En ce qui concerne le gradient $\nabla q: V \rightarrow V$, nous avons le petit calcul suivant[460] qui se base sur une des nombreuses formules du lemme 13.195¹⁰ :

$$\nabla q(x) \cdot y = \frac{d}{dt} \left[q(x + ty) \right]_{t=0} \quad (33.243a)$$

$$= \frac{d}{dt} \left[q(x) + t^2 q(y) + 2tb(x, y) \right]_{t=0} \quad (33.243b)$$

$$= 2b(x, y). \quad (33.243c)$$

Nous avons aussi

$$\nabla q(x) \cdot Ax = 2b(x, Ax) \quad (33.244a)$$

$$= 2 \int_0^\infty \langle e^{tA}x, e^{tA}Ax \rangle \quad (33.244b)$$

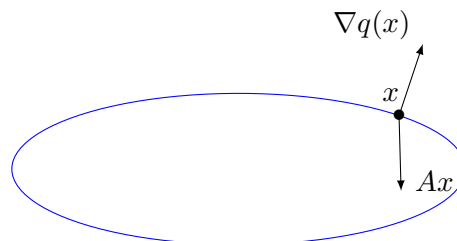
$$= \int_0^\infty \frac{\partial}{\partial t} \left(\langle e^{tA}x, e^{tA}x \rangle \right) (t) dt \quad (33.244c)$$

$$= \lim_{T \rightarrow \infty} \left[\langle e^{tA}x, e^{tA}x \rangle \right]_{t=0}^{t=T}. \quad (33.244d)$$

Mais vu que $\|e^{tA}x\| \rightarrow 0$, pour $t \rightarrow \infty$ il ne reste que terme $t = 0$ de la différence, c'est-à-dire

$$\nabla q(x) \cdot Ax = 2b(x, Ax) = -\|x\|^2. \quad (33.245)$$

Étant donné que $\nabla q(x)$ est le vecteur dirigé vers l'extérieur de l'ellipsoïde de la courbe de niveau de q au point x , le vecteur Ax est dirigé vers l'intérieur.



9. Proposition 16.81.

10. Le fait que q soit différentiable est simplement le fait que b soit bilinéaire.

Majoration de $q(y(t))'$ Nous posons

$$\begin{aligned} r: \mathbb{R}^n &\rightarrow \mathbb{R}^n \\ x &\mapsto f(x) - Ax. \end{aligned} \quad (33.246)$$

Soit y la solution maximale au problème (33.236) que nous pouvons aussi écrire sous la forme

$$y'(t) = r(y(t)) + Ay(t). \quad (33.247)$$

Calculons un peu ...

$$q(y(t))' = b(y(t), y(t))' \quad (33.248a)$$

$$= 2b(y, y') \quad (33.248b)$$

$$= 2b(y, Ay) + 2b(y, r(y)) \quad (33.248c)$$

$$= -\|y\|^2 + 2b(y, r(y)) \quad (33.245) \text{ avec } x = y(t) \quad (33.248d)$$

$$\leq -\|y\|^2 + 2\|y(t)\|_q \|r(y(t))\|_q \quad \text{Cauchy-Schwarz : } |b(a, b)| \leq \|a\|_q \|b\|_q. \quad (33.248e)$$

Chacun des deux termes peut encore être majoré. En ce qui concerne le premier, par équivalence des normes¹¹, il existe une constante C telle que $\|y\| \geq C\|y\|_q$. En renommant immédiatement C^2 en C , $\|y\|^2 \geq C\|y\|_q^2 = Cq(y)$.

Pour le second, nous allons utiliser la différentiabilité de r et le théorème des accroissements finis. Vu que $df_0 = A$ nous avons $dr_0 = df_0 - A = 0$ et de plus r est de classe C^1 parce que f l'est. Toutes les normes étant équivalentes¹² sur \mathbb{R}^n nous pouvons exprimer la continuité de dr pour la norme $\|\cdot\|_q$: si $\epsilon > 0$ est fixé alors il existe $\alpha > 0$ tel que $\|x\| < \alpha$ implique $\|dr_x\|_q < \epsilon$. Nous pouvons écrire les accroissements finis¹³ pour la fonction r :

$$\|r(x) - r(0)\|_q \leq \sup_{a \in [0, x]} \|df_a\| \|x\|_q. \quad (33.249)$$

La chose facile à remarquer est que $r(0) = f(0) = 0$. En ce qui concerne les choses difficiles, vu que dr est continue (parce que r est C^1) il existe un $\delta > 0$ tel que $\|dr_a\|_q < \epsilon$ dès que $a \in B_q(0, \delta)$. Si nous prenons $\|x\|_q < \delta$ alors cette majoration est valable pour tous les éléments sur lequel est pris le supremum dans la formule (33.249). Donc

$$\|r(x)\|_q \leq \epsilon \|x\|_q \quad (33.250)$$

tant que $\|x\|_q \leq \delta$. Par conséquent, tant que $\|y(t)\|_q \leq \delta$ nous avons $\|r(y(t))\| \leq \epsilon \|y(t)\|_q$. Nous continuons le calcul (33.248) :

$$q(y(t))' \leq Cq(y) + 2\epsilon \|y(t)\|_q^2 \quad (33.251a)$$

$$= -(C - 2\epsilon)q(y). \quad (33.251b)$$

Si ϵ est petit on a $C - 2\epsilon > 0$ et on pose $\beta = C - 2\epsilon$ pour écrire

$$q(y(t))' \leq -\beta q(y(t)) \quad (33.252)$$

tant que $\|y(t)\|_q < \delta$.

Si $q(y_0) < \delta$ alors $q(y(t)) < \delta$ Nous posons¹⁴

$$t_1 = \min\{t > 0 \text{ tel que } q(y(t)) = \delta\} \quad (33.253a)$$

$$t_2 = \max\{t < 0 \text{ tel que } q(y(t)) = \delta\}. \quad (33.253b)$$

11. Définition 12.3 et théorème 12.6.

12. Théorème 12.6.

13. Théorème 12.151.

14. t_1 est bien défini et est bien un minimum. J'en veux pour preuve¹⁵ que si $q(y(t_s)) = \delta$, on peut prendre le minimum seulement sur les $t \in [0, t_s]$; or par continuité $q(y(t)) = \delta$ définit un fermé. Bref t_1 est un infimum sur un compact (fermé borné) et donc bien un minimum atteint.

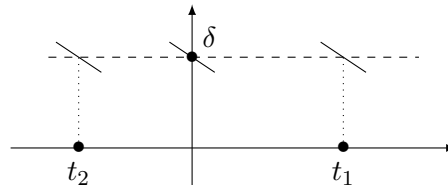
L'inégalité (33.252) est valable pour $t = 0$, $t = t_1$ et $t = t_2$; nous l'écrivons pour t_1 :

$$q(y(t))'_{t=t_1} \leq -\beta q(y(t_1)) \leq -\beta\delta < 0 \quad (33.254)$$

Nous avons donc $q(y(t_1)) = \delta$ et $q(y(t))'_{t=t_1} < 0$. Par conséquent pour tout t proche de t_1 avec $0 < t < t_1$ il y a $q(y(t)) > \delta$.

Pour la même raison, prise en $t = 0$ nous avons pour tout t proche de 0 avec $t > 0$ que $q(y(t)) < \delta$. Par continuité de $t \mapsto q(y(t))$ cette fonction doit passer par la valeur δ dans $]t_2, 0[$ et $]0, t_1[$, ce qui contredit la maximalité de t_2 et la minimalité de t_1 .

Ci-dessous, une partie de ce à quoi ressemble le graphe de $t \mapsto q(y(t))$:



Deux conclusions :

- Vu que $q(y(t))$ est borné pour tout $t \in \mathbb{R}$, nous sommes dans le cas (1) de l'alternative du théorème d'explosion en temps fini 33.18. Donc la solution $y(t)$ existe sur tout \mathbb{R} pourvu que $\|y_0\|$ soit assez petit. Plus précisément par équivalence des normes, il existe un nombre $D > 0$ tel que $\|x\| \geq D\|x\|_q$ pour tout x . Si $\|y_0\| \leq D\delta$ alors

$$D\|y_0\|_q \leq \|y_0\| \leq D\delta, \quad (33.255)$$

qui donne immédiatement $\|y_0\|_q \leq \delta$, ce qui faut pour faire fonctionner l'existence de $y(t)$ pour tout t .

- Nous pouvons maintenant d'utiliser l'inégalité (33.252) pour tout $t \in \mathbb{R}$ sous la seule hypothèse que $q(y_0) < \delta$ au lieu de $q(y(t)) < \delta$.

La partie (1) de ce théorème est prouvée; nous passons au reste à la partie (2). Pour cela nous supposons que $q(y_0) < \delta$.

À propos de $e^{\beta t}q(y)$ En sous-entendant la dépendance en t dans y nous avons

$$\left(e^{\beta t}q(y)\right)' = \beta e^{\beta t}q(y) + e^{\beta t}q(y)' = e^{\beta t}(\beta q(y) + q(y)'), \quad (33.256)$$

mais nous avons déjà prouvé que $q(y)' \leq -\beta q(y)$ (équation (33.252)), donc

$$\left(e^{\beta t}q(y)\right)' \leq 0 \quad (33.257)$$

Décroissance exponentielle Si $t \geq 0$, l'inégalité (33.257) donne

$$e^{\beta t}q(y(t)) \leq q(y_0), \quad (33.258)$$

c'est-à-dire

$$q(y(t)) \leq e^{-\beta t}q(y_0) \quad (33.259)$$

lorsque $t \geq 0$. Par équivalence des normes, nous avons des nombres D_1 et D_2 tels que

$$D_1\|x\|_q \leq \|x\| \leq D_2\|x\|_q \quad (33.260)$$

pour tout $x \in \mathbb{R}^n$. Nous avons donc pour tout $t \geq 0$ que

$$\|y(t)\| \leq D_2\|y(t)\|_q \leq D_2\|y_0\|_q e^{-\beta t}. \quad (33.261)$$

Pour rappel, $\beta > 0$, ce qui prouve la partie (2) du théorème.

Point d'équilibre Le point $y = 0$ est point d'équilibre (définition 33.33) parce que $f(0) = 0$, donc $y(t) = 0$ fonctionne. Dans ce cas, $y_0 = 0$.

Stabilité La stabilité est le fait que $\|y(t)\|_q \leq \delta$ dès que $\|y_0\|_q \leq \delta$.

□

33.8.5 Système proies-prédateurs de Lotka-Volterra

Le système de **Lotka-Volterra** est l'équation différentielle suivante :

$$\begin{cases} x' = ax - bxy & (33.262a) \\ y' = -cy + dxy & (33.262b) \end{cases}$$

où a, b, c, d sont des constantes positives, et avec la condition $x(t_0) > 0, y(t_0) > 0$.

En ce qui concerne l'interprétation des équations[461],

- (1) $x(t)$ est le nombre de proies,
- (2) $y(t)$ est le nombre de prédateurs,
- (3) Les proies ont une reproduction rapide qui mène à une croissance exponentielle en absence de prédation (d'où le terme ax).
- (4) Au contraire, les prédateurs meurent (ou migrent) rapidement lorsqu'ils n'ont pas de proies et nous supposons une décroissance exponentielle du nombre de prédateurs en l'absence de proies. D'où le terme $-cy$ avec le signe négatif.
- (5) Les termes $-bxy$ et dxy sont les termes d'interaction entre les proies et les prédateurs. Ils sont proportionnels à la fréquence de leurs rencontres, lesquelles sont avantageuses pour les prédateurs et problématiques pour les proies.

Théorème 33.35 (Lotka-Volterra[168]).

Soient des constantes positives a, b, c, d et le système équations différentielles

$$\begin{cases} x' = ax - bxy & (33.263a) \\ y' = -cy + dxy & (33.263b) \\ x(t_0) > 0, y(t_0) > 0. & (33.263c) \end{cases}$$

Alors

- (1) Les solutions sont positives sur leur domaines.
- (2) Les solutions existent sur \mathbb{R} .
- (3) Les solutions sont périodiques.

Démonstration. Nous divisons la preuve.

Comment théorème de Cauchy-Lipschitz s'applique Tel quel, le théorème de Cauchy-Lipschitz 18.40

ne s'applique pas parce qu'il demande une condition initiale pour avoir unicité. En ce qui concerne les notations, ce qui est noté « y » dans le théorème est ici le couple x, y et la fonction f est alors

$$f\left(t, \begin{pmatrix} x \\ y \end{pmatrix}\right) = \begin{pmatrix} ax - bxy \\ -cy + dxy \end{pmatrix}. \quad (33.264)$$

C'est une fonction continue localement Lipschitz partout par le lemme 13.261 et la proposition 13.260.

Nous savons cependant que les solutions sont de classe C^1 et que moyennant la donnée d'une condition initiale, la solution est unique.

Les solutions restent positives Supposons $x(s) = 0$ pour un certain $s > t_0$. Alors le solution

$$\begin{cases} x(t) = 0 & (33.265a) \\ y(t) = \exp(-ct) & (33.265b) \end{cases}$$

est une solution pour $[t_0, s + \epsilon]$. Par unicité de la solution avec condition initiale $s(s) = 0$, nous avons aussi $x(t_0) = 0$ pour toutes les solutions, ce qui contredit notre condition.

De la même façon, avoir $y(s) = 0$ donne une solution avec $y(t) = 0$ pour tout t et donc une contradiction.

Solutions sur \mathbb{R} Nous montrons maintenant que les solutions sont définies sur \mathbb{R} .

Nous avons $x' < ax$, donc pour tout t où la solution est définie,

$$0 < x(t) < x(t_0)e^{a(t-t_0)}, \quad (33.266)$$

c'est-à-dire que la solution ne peut pas exploser en temps fini¹⁶ : elle est bornée par le haut et le bas. Elle doit donc exister pour tout $t \in \mathbb{R}$. Par ailleurs, $y' < dxy$ donc

$$0 < y(t) < y(t_0)e^{d \int_{t_0}^t x(s) ds} \quad (33.267)$$

qui est également contraire à l'explosion en temps fini.

4 zones : monotonie Nous divisons \mathbb{R}^2 en quatre zones d'après les signes de $a - by$ et $c - dx$. Nous montrons que dans chacune de ces zones, les solutions sont monotones. Prenons par exemple la partie

$$\{(x, y) \in \mathbb{R}^2 \text{ tel que } a - by > 0\} \times \{c - dx < 0\}. \quad (33.268)$$

Vu l'équation $x' = x(a - by)$, tant que $(x(t), y(t))$ est dans cette zone, la fonction x' a le signe de x et est donc positive. Donc x est croissante dans cette zone.

De la même façon, $y' = -y(c - dx)$, et y' a un signe constant dans la zone.

4 zones : on bouge Nous prouvons à présent qu'une solution ne reste pas dans une zone.

(1) Supposons que $(x(t_0), y(t_0))$ soit dans la zone

$$\{a - by > 0\} \quad \times \quad \{c - dx > 0\} \quad (33.269a)$$

$$x' > 0 \quad \quad \quad y' < 0 \quad (33.269b)$$

et que la solution reste dans cette zone (pour les $t > t_0$). Nous avons en particulier $x' > 0$, donc x est croissante tout en ayant la borne supérieure $x < c/d$. Par conséquent x a une limite que nous appelons $x_1 \in [0, \frac{c}{d}]$.

De la même façon, y est décroissante et bornée vers le bas par zéro. Donc y a une limite que nous notons $y_1 \in [0, y(t_0)]$.

Vu que x est bornée et de classe C^1 nous avons forcément $\lim_{t \rightarrow \infty} x'(t) = 0$. Mais vu que $x' = ax - bxy$ nous devons avoir

$$ax_1 - bx_1y_1 = 0. \quad (33.270)$$

Mais ni $x_1 > 0$ donc $a - by_1 = 0$, ce qui donne $y_1 = \frac{a}{b}$ et aussi $x_1 = \frac{c}{d}$. Bref, y est décroissante et tend vers a/b ; donc $y(t_0) > a/b$, ce qui contredit que $y(t_0)$ soit dans la zone considérée.

Étant donné que $x' > 0$ et $y' < 0$, la solution sort de la zone pour entrer dans la zone ...

(2) Supposons que $(x(t_0), y(t_0))$ soit dans la zone

$$\{a - by > 0\} \quad \times \quad \{c - dx < 0\} \quad (33.271a)$$

$$x' < 0 \quad \quad \quad y' > 0 \quad (33.271b)$$

et que la solution reste dans cette zone (pour les $t > t_0$). Les fonctions x et y sont convergentes. Par conséquent $\ln(y)$ converge aussi et vu que x est croissante,

$$\frac{y'}{y} = -c + dx \geq -x + dx(t_0) > 0 \quad (33.272)$$

Cela signifie que $\ln(y)'$ est toujours positive et bornée par le bas. Cela est impossible si y est borné.

Donc on sort de la zone pour entrer dans ...

16. Voir le corollaire 33.18.

(3) Supposons que $(x(t_0), y(t_0))$ soit dans la zone

$$\{a - by < 0\} \quad \times \quad \{c - dx < 0\} \quad (33.273a)$$

$$x' < 0 \quad \quad \quad y' > 0 \quad (33.273b)$$

et que la solution reste dans cette zone (pour les $t > t_0$).

Le même type de raisonnement fait passer à la zone...

(4) Supposons que $(x(t_0), y(t_0))$ soit dans la zone

$$\{a - by < 0\} \quad \times \quad \{c - dx > 0\} \quad (33.274a)$$

$$x' < 0 \quad \quad \quad y' < 0 \quad (33.274b)$$

et que la solution reste dans cette zone (pour les $t > t_0$). Encore une fois, cela nous fait sortir de la zone et retourne vers la première zone.

À ce moment nous voyons déjà que la relation entre proies et prédateurs, c'est un peu le mythe de Sisyphe...

Une intégrale première Posons la fonction

$$H(x, y) = by + dx - a \ln(y) - c \ln(x). \quad (33.275)$$

Une simple dérivation montre que $x \mapsto H(x(t), y(t))$ est constante. Nous considérons la fonction

$$f: \mathbb{R} \rightarrow \mathbb{R} \quad (33.276)$$

$$s \mapsto H\left(\frac{c}{d}, s\right)$$

dont la dérivée n'est autre que $f'(s) = b - \frac{a}{s}$. La fonction f est donc décroissante sur l'intervalle $[\frac{a}{b}, \infty[$ et donc injective. Sur les changements de zones, il existe un t_0 tel que

$$x(t_0) = \frac{d}{c} \quad (33.277a)$$

$$y(t_0) > 0. \quad (33.277b)$$

Pour cette valeur t_0 nous avons alors $H(x(t_0), y(t_0)) = f(y(t_0))$. En posant $s_0 = y(t_0) > 0$ nous avons

$$H(x_0, y_0) = f(s_0) \quad (33.278)$$

et f étant injective, ce s_0 est la seule valeur de s à vérifier $H(x_0, y_0) = f(s)$.

Conclusion La fonction x passant d'une zone à l'autre, il existe un $t_1 > t_0$ tel que $x(t_1) = a/b$. Nous avons évidemment

$$H(x(t_1), y(t_1)) = H(x_0, y_0) \quad (33.279)$$

parce que H est constante le long du mouvement. Cela se traduit par

$$H\left(\frac{a}{b}, y(t_1)\right) = f(s_0), \quad (33.280)$$

et donc $y(t_1) = f(s_0) = y(t_0)$. Avec tout cela nous avons

$$\left\{ \begin{array}{l} y(t_1) = y(t_0) \\ x(t_1) = x(t_2) = \frac{a}{b} \end{array} \right. \quad (33.281a)$$

$$\left\{ \begin{array}{l} y(t_1) = y(t_0) \\ x(t_1) = x(t_2) = \frac{a}{b} \end{array} \right. \quad (33.281b)$$

Cela est donc un point par lequel la solution repasse. Par unicité de la solution, elle est donc périodique.

□

33.9 Équation du second ordre

33.9.1 Wronskien

Nous considérons ici une équation différentielle de la forme

$$y''(t) + q(t)y(t) = 0 \quad (33.282)$$

Dans ce point nous allons considérer la fonction q sans hypothèse de périodicité. L'équation de Hill (sous-section 33.9.4) sera la même équation, mais en supposant que q est périodique.

Nous commençons par argumenter que si q est continue, alors l'ensemble des solutions de l'équation (33.282) est un espace vectoriel de dimension deux. Pour cela il suffit d'appliquer la méthode de réduction de l'ordre (section 33.7) puis le théorème de dimension pour les systèmes linéaires (théorème 33.14). En effet si la fonction y_1 est solution de (33.282) si et seulement si le vecteur $Y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$ est solution du système linéaire

$$Y'(t) = \begin{pmatrix} 0 & 1 \\ -q(t) & 0 \end{pmatrix} Y(t). \quad (33.283)$$

Soient deux solutions y_1 et y_2 de l'équation différentielle. Le **Wronskien** de ces deux solutions est le déterminant

$$W(t) = \begin{vmatrix} y_1 & y_2 \\ y_1' & y_2' \end{vmatrix}. \quad (33.284)$$

Si nous considérons l'équation différentielle

$$y'' + py' + qy = 0, \quad (33.285)$$

le Wronskien peut être déterminé sans savoir explicitement y_1 et y_2 parce que $W = y_1 y_2' - y_1' y_2$, et en dérivant,

$$W' = y_1 y_2'' + y_1' y_2' - y_1'' y_2 - y_1' y_2' \quad (33.286a)$$

$$= y_1(-p y_2' - q y_2) - (-p y_1' - q y_1) y_2 \quad (33.286b)$$

$$= -p \begin{vmatrix} y_1 & y_2 \\ y_1' & y_2' \end{vmatrix}, \quad (33.286c)$$

c'est-à-dire

$$W' = -pW. \quad (33.287)$$

Il suffit donc de savoir une condition initiale pour obtenir une équation différentielle pour W .

33.9.2 Avec second membre

Une équation différentielle du second ordre avec un second membre se présente sous la forme

$$ay''(t) + by'(t) + cy(t) = v(t) \quad (33.288)$$

où $v(t)$ est une fonction donnée. Le truc est de commencer par résoudre l'équation différentielle sans second membre, c'est-à-dire trouver la fonction $y_H(t)$ telle que

$$ay_H''(t) + by_H'(t) + cy_H(t) = 0. \quad (33.289)$$

Cela se fait en utilisant la méthode du polynôme caractéristique.

Ensuite, il faut trouver une solution particulière $y_P(t)$ de l'équation avec le second membre. Une seule. Pour y parvenir, il faut du doigté et un peu de technique. Il faut faire des essais en fonction de ce à quoi ressemble le $v(t)$:

- (1) Si $v(t)$ est un polynôme, alors il faut essayer un polynôme,
- (2) Si $v(t) = \cos(\omega t)$ ou bien $v(t) = \sin(\omega t)$, alors essayer $y_P(t) = A \cos(t) + B \sin(\omega t)$,
- (3) Si $v(t) = e^{\omega t}$, alors essayer $y_P(t) = Ae^{\omega t}$.

33.9.3 Équation $y'' + q(t)y = 0$

Nous allons donner quelques propriétés des solutions de l'équation

$$y'' + qy = 0 \quad (33.290)$$

en fonction de telle ou telle hypothèse sur q .

Proposition 33.36.

Si $q: \mathbb{R}^+ \rightarrow \mathbb{R}$ est continue et si

$$\int_0^\infty |q(t)| dt \quad (33.291)$$

converge, alors

(1) toute solution bornée de $y'' + qy = 0$ vérifie $\lim_{t \rightarrow \infty} y'(t) = 0$,

(2) l'équation $y'' + qy = 0$ admet des solutions non bornées.

Démonstration. Soit y une solution bornée, et intégrons l'équation différentielle entre 0 et ∞ :

$$\int_0^\infty y''(t) dt = - \int_0^\infty q(t)y(t) dt. \quad (33.292)$$

La fonction y étant bornée, l'hypothèse sur q permet de dire que l'intégrale de droite existe. Par ailleurs,

$$\int_0^\infty y'' = \lim_{a \rightarrow \infty} \int_0^a y'' = \lim_{a \rightarrow \infty} y'(a) - y'(0). \quad (33.293)$$

Cela justifie que la limite $\lim_{t \rightarrow \infty} y'(t)$ existe. Posons $\alpha = \lim_{t \rightarrow \infty} y'(t)$ et supposons par l'absurde que $\alpha \neq 0$. Soit $\epsilon > 0$ et λ assez grand pour que

$$\|y' - \alpha\|_{[\lambda, \infty[} < \epsilon. \quad (33.294)$$

Soit aussi $x > \lambda$. Nous avons

$$y(x) = y(\lambda) + \int_\lambda^x y'(t) dt \quad (33.295a)$$

$$\geq y(\lambda) + \int_\lambda^x (\alpha - \epsilon) dt \quad (33.295b)$$

$$= y(\lambda) + (\alpha - \epsilon)(x - \lambda). \quad (33.295c)$$

En prenant la limite des deux côtés on voit que $y(x) \rightarrow \infty$ dès que $\alpha \neq 0$, ce qui est contraire aux hypothèses. Donc $\alpha = 0$.

Pour la seconde partie de la proposition, nous devons prouver que l'équation $y'' + qy = 0$ possède des solutions non bornées. Si l'équation a seulement des solutions bornées et si $\{u, v\}$ est une base de solutions, alors nous avons $u', v' \rightarrow 0$. Si nous reprenons l'équation (33.287) avec $p = 0$ nous savons que dans notre cas le Wronskien satisfait à $W' = 0$, c'est-à-dire qu'il est constant. Mais vu que u et v sont bornées et que les dérivées tendent vers zéro, nous avons $W(t) \rightarrow 0$ et donc $W(t) = 0$.

Or l'annulation identique du Wronskien contredit que $\{u, v\}$ serait une base de solutions. Donc il existe des solutions non bornées. \square

Proposition 33.37.

Soit l'équation différentielle $y'' + qy = 0$. Si q est C^1 , strictement positive et croissante, alors toutes les solutions sont bornées.

Démonstration. Soit y une solution et multiplions l'équation par $2y'$ (qui est non nulle par hypothèse) :

$$2y'y'' + 2qy'y = 0. \quad (33.296)$$

Nous allons intégrer cela en nous souvenant que $2y'y''$ est la dérivée de $(y')^2$. Pour tout $t > 0$ nous avons

$$0 = y'(t)^2 - y'(0)^2 + 2 \underbrace{\int_0^t q(t)y'(t)y(t)dt}_{\text{par partie}} \quad (33.297a)$$

$$= y'(t)^2 - y'(0)^2 + 2 \left([qy^2]_0^t - \int_0^t q'y^2 \right) \quad (33.297b)$$

$$(33.297c)$$

Le terme qui nous intéresse est celui qui contient $y(t)$:

$$2q(t)y(t)^2 = -y'(t)^2 + y'(0)^2 + 2q(0)y(0)^2 + 2 \int_0^t q'y^2 \quad (33.298)$$

Nous pouvons majorer $-y'(t)^2$ par zéro et remplacer toutes les constantes par K :

$$q(t)y(t)^2 \leq \int_0^t q'y^2 + K = \int_0^t \frac{q'}{q} qy^2. \quad (33.299)$$

C'est le moment d'utiliser le lemme de Grönwall 33.4 avec $\phi = qy^2$ et $\psi = q'/q$. Les hypothèses de croissance et de positivité ont été posées exprès. Bref, on a

$$qy^2 \leq K \exp \left(\int_0^t \frac{q'(s)}{q(s)} ds \right) \quad (33.300a)$$

$$= K \exp \left(\ln \frac{q(t)}{q(0)} \right) \quad (33.300b)$$

$$= K \frac{q(t)}{q(0)}. \quad (33.300c)$$

Notons que $q(0)$ est strictement positif. Nous déduisons que

$$y^2 \leq \frac{K}{q(0)} \quad (33.301)$$

et donc y est bornée. □

33.9.4 Équation de Hill

L'équation de Hill est une équation différentielle de la forme

$$y'' + qy = 0 \quad (33.302)$$

où

- (1) $q \in C^1(\mathbb{R}, \mathbb{R})$,
- (2) q est paire et π -périodique

Nous nous intéressons aux solutions complexes de cette équation différentielle.

Nous nommons $W \subset C^2(\mathbb{R}, \mathbb{C})$ l'espace des solutions complexes de l'équation (33.302). Nous savons par ce qui a été dit en 33.9.3 que cet espace est de dimension deux. De plus avec les hypothèses faites ici sur q , nous savons que les solutions sont de classe C^3 parce que si y est une solution, alors l'équation $y'' = -qy$ nous indique que y est C^1 parce que y'' existe (y' est dérivable et donc continue). Mais si y est de classe C^1 , alors le membre de droite $-qy$ est C^1 et donc y'' est C^1 , ce qui prouve que y est de classe C^3 . La récurrence ne va pas plus loin parce que q est seulement de classe C^1 .

Nous considérons l'application de translation

$$\begin{aligned} T: C^2(\mathbb{R}, \mathbb{C}) &\rightarrow C^2(\mathbb{R}, \mathbb{C}) \\ (Ty)(x) &= y(x + \pi). \end{aligned} \quad (33.303)$$

En utilisant la règle de dérivation de fonctions composées, $(Ty)' = Ty'$ et $(Ty)'' = Ty''$, de telle sorte que si u est solution de l'équation (33.302), alors Tu est également solution. Donc W est un espace stable par T .

Le théorème 33.14 nous permet de choisir une base de W en imposant des conditions. Nous choisissons une base $\{y_1, y_2\}$ telles que

$$\begin{aligned} y_1(0) &= 1 & y_2(0) &= 0 \\ y_1'(0) &= 0 & y_2'(0) &= 1. \end{aligned} \quad (33.304)$$

Le théorème 33.14 nous assure que deux telles solutions existent et qu'elles forment une base de W parce que W est de dimension 2.

Lemme 33.38 ([451]).

Avec ce choix de base $\{y_1, y_2\}$ la matrice de T est donnée par

$$T = \begin{pmatrix} y_1(\pi) & y_2(\pi) \\ y_1'(\pi) & y_2'(\pi) \end{pmatrix}. \quad (33.305)$$

De plus la fonction y_1 est paire et la fonction y_2 est impaire.

Démonstration. Cherchons la matrice de T dans cette base en associant $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$ à y_1 et $\begin{pmatrix} 0 \\ 1 \end{pmatrix}$ à y_2 . Si

$$T = \begin{pmatrix} a & b \\ c & d \end{pmatrix}, \text{ alors}$$

$$Ty_1 = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} a \\ c \end{pmatrix} = ay_1 + cy_2. \quad (33.306)$$

En évaluant cela en $t = 0$,

$$(Ty_1)(0) = ay_1(0) + cy_2(0) = a, \quad (33.307)$$

donc $a = (Ty_1)(0) = y_1(\pi)$. En dérivant (33.306), en tenant compte du fait que $(Ty_1)' = Ty_1'$ et en évaluant en $t = 0$, nous trouvons de même $c = y_1'(\pi)$. Puis le même cinéma avec y_2 donne

$$T = \begin{pmatrix} y_1(\pi) & y_2(\pi) \\ y_1'(\pi) & y_2'(\pi) \end{pmatrix}. \quad (33.308)$$

Passons maintenant à la parité de y_1 et y_2 . Nous posons $\psi(t) = y_1(-t)$. Alors $\psi'(t) = -y_1'(-t)$ et $\psi''(t) = y_1''(t)$, tant et si bien que

$$\psi''(t) + q(t)\psi(t) = y_1''(-t) + q(t)y_1(-t) = 0. \quad (33.309)$$

donc ψ est une solution de l'équation. Mais

$$\begin{cases} \psi(0) = y_1(0) & (33.310a) \\ \psi'(0) = -y_1'(0) = 0, & (33.310b) \end{cases}$$

donc ψ a les mêmes conditions initiales que y_1 . Par conséquent $\psi = y_1$ (par l'unicité donnée dans le théorème de Cauchy-Lipschitz 18.40) et y_1 est paire. Nous procédons de même en partant de $\varphi(t) = -y_2(-t)$ pour trouver que $\varphi = y_2$ et que donc que y_2 est impaire. \square

Remémorons nous toutefois, pour calmer tout enthousiasme excessif, que T dépend de deux solutions et donc de la fonction q donnée dans l'équation.

Proposition 33.39 ([43]).

Nous considérons l'équation $y'' + qy = 0$ et sa base de solutions $\{y_1, y_2\}$ en suivant les notations données plus haut.

- (1) Si $|\operatorname{Tr}(T)| < 2$, alors toutes les solutions de l'équation sont bornées.
- (2) Si $|\operatorname{Tr}(T)| = 2$ alors nous avons une solution non bornée.
- (3) Si $|\operatorname{Tr}(T)| > 2$ alors toutes les solutions de l'équation sont non bornées.
- (4) Le cas $|\operatorname{Tr}(T)| = 2$ se présente si et seulement si $y_1'(\pi)y_2(\pi) = 0$.

Démonstration. Remarquons que le déterminant de la matrice T est égal au Wronskien des solutions y_1 et y_2 calculé en $t = \pi$. Calculons sa valeur :

$$W(y_1, y_2) = \det \begin{pmatrix} y_1 & y_2 \\ y_1' & y_2' \end{pmatrix} = y_1 y_2' - y_1' y_2. \quad (33.311)$$

En dérivant et en remplaçant y_i'' par $-qy_i$, nous trouvons tout de suite $W(y_1, y_2)' = 0$. Donc le Wronskien est constant et il est facile de le calculer en $t = 0$:

$$W(y_1, y_2)(0) = 1 - 0 = 1. \quad (33.312)$$

Donc pour tout t nous avons $W(y_1, y_2)(t) = 1$. En particulier

$$\det(T) = W(y_1, y_2)(\pi) = 1, \quad (33.313)$$

et notons au passage que T est inversible.

Nous écrivons le polynôme caractéristique de T sous la forme $\chi_T = X^2 - \operatorname{Tr}(T)X + \det(T)$, c'est-à-dire

$$\chi_T = X^2 - \operatorname{Tr}(T)X + 1, \quad (33.314)$$

dont le discriminant est $\Delta = \operatorname{Tr}(T)^2 - 4$.

Nous passons à présent aux différents points de la proposition.

- (1) Si $|\operatorname{Tr}(T)| < 2$, alors $\Delta < 0$ et χ_T a deux racines complexes conjuguées que nous notons ρ et $\bar{\rho}$. De plus le produit des racines étant le terme indépendant, $\rho\bar{\rho} = 1$; en particulier $|\rho| = |\bar{\rho}| = 1$. Notons $\{u, v\}$ une base de vecteurs propres : $Tu = \rho u$ et $Tv = \bar{\rho}v$. Il est vite vu que la fonction $|u|$ est π -périodique :

$$|u|(t + \pi) = |u(t + \pi)| = |(Tu)(t)| = |(\rho u)(t)| = |\rho||u|(t) = |u|(t). \quad (33.315)$$

La fonction $|u|$ est continue¹⁷ et périodique ergo bornée. La fonction $|v|$ est bornée pour la même raison et par linéarité, toutes les fonctions de W sont bornées.

- (2) Si $\operatorname{Tr}(T) = \pm 2$, alors $\Delta = 0$ et χ_T a une racine réelle double¹⁸ qui doit être ± 1 . Soit u un vecteur propre de T pour la valeur propre ± 1 . Nous avons

$$|u|(t + \pi) = |Tu(t)| = |\pm u(t)|, \quad (33.316)$$

ce qui prouve encore que $|u|$ est périodique et donc bornée.

Notons que nous n'avons pas d'informations sur le fait qu'une autre solution soit ou non bornée.

- (3) Si $|\operatorname{Tr}(T)| > 2$, alors χ_T a deux racines réelles distinctes r et r' avec $rr' = 1$ (toujours les relations coefficients-racines). En raison de quoi $r' = r^{-1}$ et quitte à échanger r et r' nous supposons $|r| > 1$. L'opérateur est maintenant diagonalisable et nous considérons $\{u, v\}$ une base de vecteurs propres pour les valeurs propres r et r' . Une solution non nulle de l'équation s'écrit donc sous la forme

$$y = \alpha u + \beta v \quad (33.317)$$

avec $(\alpha, \beta) \neq (0, 0)$.

17. La fonction u elle-même n'est cependant pas garantie d'être périodique.

18. Ce qui n'implique pas le fait d'avoir deux vecteurs propres pour cette valeur propre, mais tout de même au moins un, voir l'exemple 11.151.

— Si $\alpha = 0$, alors $\beta \neq 0$ et nous choisissons une valeur t telle que $v(t) \neq 0$. Dans ce cas,

$$y(t + n\pi) = \beta v(t + n\pi) = \beta(T^n v)(t) = \beta(r')^n v(t), \quad (33.318)$$

et en faisant $n \rightarrow -\infty$ nous obtenons $\pm\infty$ suivant le signe de β .

— Si $\alpha \neq 0$, alors nous fixons¹⁹ t tel que $u(t) \neq 0$. Alors

$$y(t + n\pi) = \alpha r'^n u(t) + \beta(r')^n(t). \quad (33.319)$$

En faisant $n \rightarrow \infty$, nous avons $(r')^n \rightarrow 0$ tandis que le premier terme tend vers $\pm\infty$ suivant le signe de α .

(4) D'abord le théorème de Cayley-Hamilton 11.154 nous indique que $\chi_T(T) = 0$, c'est-à-dire que

$$T^2 - \text{Tr}(T)T + 1 = 0. \quad (33.320)$$

Nous avons déjà mentionné le fait que T était inversible. Multiplions donc (33.320) par T^{-1} :

$$T + T^{-1} = \text{Tr}(T)\mathbb{1}_2. \quad (33.321)$$

Vu que T^{-1} est l'endomorphisme $T^{-1}u(t) = u(t - \pi)$, sa matrice est donnée par

$$T^{-1} = \begin{pmatrix} y_1(-\pi) & y_2(-\pi) \\ y_1'(-\pi) & y_2'(-\pi) \end{pmatrix} = \begin{pmatrix} y_1(\pi) & -y_2(\pi) \\ -y_1'(\pi) & y_2'(\pi) \end{pmatrix} \quad (33.322)$$

où nous avons utilisé le fait que y_1 était paire et y_2 impaire (lemme 33.38). Si nous notons $T = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$, alors $T^{-1} = \begin{pmatrix} a & -b \\ -c & d \end{pmatrix}$ et

$$T + T^{-1} = \begin{pmatrix} 2a & 0 \\ 0 & 2b \end{pmatrix}. \quad (33.323)$$

L'équation (33.321) donne alors, vu que $\text{Tr}(T) = a + d$,

$$\begin{pmatrix} 2a & 0 \\ 0 & 2b \end{pmatrix} = \begin{pmatrix} a + d & 0 \\ 0 & a + d \end{pmatrix}, \quad (33.324)$$

ce qui donne immédiatement $a = d$. La matrice de T a donc comme forme $T = \begin{pmatrix} a & b \\ c & a \end{pmatrix}$ et $\text{Tr}(T) = 2a$.

Donc $\text{Tr}(T) = \pm 2$ si et seulement si $a = \pm 1$ et vu que $1 = \det(T) = a^2 - bc$, nous avons $a = \pm 1$ si et seulement si $bc = 0$, ce qui signifie exactement $y_1'(\pi)y_2(\pi) = 0$.

□

33.10 Différents types d'équations différentielles

33.10.1 Équation homogène

Une équation différentielle **homogène** est une équation de la forme

$$y' = f(t, y) \quad (33.325)$$

où $f(\lambda t, \lambda y) = f(t, y)$ pour tout $\lambda \neq 0$.

Elle se présente sous la forme

$$y' = \frac{\text{degré } n \text{ en } t, y}{\text{degré } n \text{ en } t, y}, \quad (33.326)$$

avec pas de y' à droite : juste du y et du t .

19. Mais pas trop hein ; nous aurons encore besoin d'assigner à t d'autres valeurs dans d'autres théorèmes.

Lemme 33.40.

L'équation $y' = f(t, y)$ est homogène si et seulement si $f(t, y)$ est une fonction de y/t seulement.

Pour résoudre l'équation homogène, on pose

$$z(t) = \frac{y(t)}{t}, \quad (33.327)$$

donc $tz = y$, et

$$y'(t) = tv'(t) + v(t), \quad (33.328)$$

à remettre dans l'équation de départ.

33.10.2 Équation de Bernoulli

C'est une équation du type

$$y' = a(t)y + b(t)y^\alpha \quad (33.329)$$

où $\alpha \neq 0$ ou 1 . Pour la résoudre, on divise l'équation par y^α , et on pose $u = y^{1-\alpha}$, et on tombe sur une équation linéaire

$$u' = (1 - \alpha)(a(t)u + b(t)). \quad (33.330)$$

33.10.3 Équation de Riccati

C'est une équation de la forme

$$y' = a(t)y^2 + b(t)y + c(t). \quad (33.331)$$

En général, on ne peut pas la résoudre, mais si on en connaît *a priori* des solutions particulières, alors on peut s'en sortir.

(1) Si on sait que $y_1(t)$ est une solution, alors on pose

$$y(t) = y_1(t) + \frac{1}{u(t)}, \quad (33.332)$$

et on obtient une équation linéaire

$$u' = -(2y_1(t)a(t) + b(t))u - a(t). \quad (33.333)$$

(2) Si y_1 et y_2 sont solutions, alors nous avons y sous forme implicite

$$\frac{y - y_1}{y - y_2} = Ke^{\int a(t)(y_1(t) - y_2(t))dt}. \quad (33.334)$$

Pour résoudre une équation de Riccati, il faut donc d'abord deviner une ou deux solutions.

33.10.4 Équation différentielle exacte**33.10.4.1 Résolution lorsque tout va bien**

Avant de vous lancer dans les équations différentielles exacte, vous devez lire la section sur les formes différentielles 13.19. Une équation différentielle exacte est de la forme $P(t, y) + Q(t, y)y' = 0$ que nous allons écrire sous la forme

$$P(t, y)dt + Q(t, y)dy = 0. \quad (33.335)$$

Nous savons que si $\partial_y P = \partial_t Q$, alors il existe une fonction $f(t, y)$ telle que $Pdt + Qdy = df$. Pour trouver une telle fonction, nous pouvons simplement intégrer la forme $Pdt + Qdy$. En effet, si $\gamma: [0, 1] \rightarrow \mathbb{R}^2$ est un chemin tel que $\gamma(0) = (0, 0)$ et $\gamma(1) = (t, y)$, alors en définissant

$$f(t, y) = \int_\gamma [Pdt + Qdy] = \int_0^1 [(P \circ \gamma)(u)dt + (Q \circ \gamma)(u)](\gamma'(u))du, \quad (33.336)$$

nous avons $df = Pdt + Qdy$. N'importe quel chemin fait l'affaire. Calculons avec $\gamma(u) = (tu, yu)$. La dérivée de ce chemin est donnée par

$$\gamma'(u) = t \begin{pmatrix} 1 \\ 0 \end{pmatrix} + y \begin{pmatrix} 0 \\ 1 \end{pmatrix}. \quad (33.337)$$

Étant donné que $dt \begin{pmatrix} a \\ b \end{pmatrix} = a$ et $dy \begin{pmatrix} a \\ b \end{pmatrix} = b$, nous avons

$$\begin{aligned} f(t, y) &= \int_0^1 [Pdt + Qdy](\gamma(u)) \left(t \begin{pmatrix} 1 \\ 0 \end{pmatrix} + y \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right) du \\ &= \int_0^1 P(\gamma(t))tdu + \int_0^1 Q(\gamma(t))ydu \\ &= \int_0^1 [tP(tu, yu) + yQ(tu, yu)]du. \end{aligned} \quad (33.338)$$

Nous retrouvons exactement la formule (21.303). Si ça t'étonne, c'est que tu n'as pas compris ;)
Dans le cas où nous avons la fonction f qui vérifie $P = \partial_t f$ et $Q = \partial_y f$, l'équation (33.335) devient

$$\frac{\partial f}{\partial t} + \frac{\partial f}{\partial y} \frac{dy}{dt} = 0, \quad (33.339)$$

c'est-à-dire

$$\frac{d}{dt} [f(t, y(t))] = 0, \quad (33.340)$$

dont la solution

$$f(t, y(t)) = C \quad (33.341)$$

donne la solution $y(t)$ sous forme implicite.

33.10.4.2 Facteur intégrant (quand tout ne va pas bien)

Si la forme $Pdt + Qdy$ n'est pas exacte, il n'existe pas de fonction f qui résolve l'affaire. Nous pouvons toutefois essayer de trouver un **facteur intégrant**. Nous cherchons une fonction M telle que

$$(MP)dt + (MQ)dy \quad (33.342)$$

soit exacte. Nous cherchons donc $M(t, y)$ telle que $\partial_y(MP) = \partial_t(MQ)$. En utilisant la règle de Leibnitz, nous trouvons l'équation suivante pour M :

$$M(\partial_y P - \partial_t Q) = Q(\partial_t M) - P(\partial_y M). \quad (33.343)$$

Cette équation est en générale extrêmement difficile à résoudre, mais dans certains cas particuliers, il est possible d'en trouver une solution à tâtons.

33.11 Distributions pour les équations différentielles

Nous commençons par définir l'espace $C^\infty(\mathbb{R}, \mathcal{S}'(\mathbb{R}^d))$ en disant que $t \mapsto u_t$ est dans cet espace si

- (1) pour tout $t \in \mathbb{R}$ nous avons $u_t \in \mathcal{S}'(\mathbb{R}^d)$,
- (2) l'application $t \mapsto u_t$ est de classe C^∞ .

Pour définir ce que nous entendons par une fonction de classe C^k à valeurs dans $\mathcal{S}'(\mathbb{R}^d)$ nous nous souvenons de la proposition 31.36.

33.11.1 Équation de Schrödinger

Théorème 33.41 (Équation de Schrödinger[43]).

Soit $g \in \mathcal{S}'(\mathbb{R}^d)$ et le problème

$$\begin{cases} \partial_t \tilde{u} - i\Delta \tilde{u} = 0 \\ u_0 = g \end{cases} \quad (33.344a)$$

$$(33.344b)$$

où $\tilde{u} \in C^\infty(\mathbb{R}, \mathcal{S}'(\mathbb{R}^d))$ est lié à u par la remarque 31.58. Alors

(1) Il existe une unique solution dans $C^\infty(\mathbb{R}, \mathcal{S}'(\mathbb{R}^d))$.

(2) Cette solution u vérifie de plus $\tilde{u} \in \mathcal{S}'(\mathbb{R} \times \mathbb{R}^d)$.

Démonstration. Nous allons donner explicitement une fonction $u \in C^\infty(\mathbb{R}, \mathcal{S}'(\mathbb{R}^d))$ et nous allons vérifier l'équation (33.344a) en testant sur une fonction $\psi \in \mathcal{S}'(\mathbb{R} \times \mathbb{R}^d)$. Cela prouvera le point (2) ainsi que la partie existence de (1). Dans ce qui suit toutes les transformées de Fourier seront par rapport à la variable $x \in \mathbb{R}^d$ ou par rapport à ξ . Jamais par rapport à $t \in \mathbb{R}$.

Existence Pour $t \in \mathbb{R}$ nous posons²⁰

$$u_t = \mathcal{F}^{-1}(f_t \hat{g}) \quad (33.345)$$

où $f_t \in \mathcal{S}'(\mathbb{R}^d)$ est la fonction $f_t(x) = e^{-it\|x\|^2}$. Pour toute fonction $\varphi \in \mathcal{S}'(\mathbb{R}^d)$ nous avons

$$u_t(\varphi) = (f_t \hat{g})(\mathcal{F}^{-1}(\varphi)) = \hat{g}(f \mathcal{F}^{-1}(\varphi)) = g\left(\mathcal{F}(f \mathcal{F}^{-1}(\varphi))\right). \quad (33.346)$$

Le fait que $\mathcal{F}^{-1}(\varphi)$ soit une fonction Schwartz fait partie de la proposition 30.14. Pour chaque t nous avons bien $u_t \in \mathcal{S}'(\Omega)$.

De plus la fonction $h(t, x) = e^{-it\|x\|^2}(\mathcal{F}^{-1}\varphi)(x)$ est dans $C^\infty(\mathbb{R} \times \mathbb{R}^d)$, et par conséquent l'application

$$t \mapsto \hat{g}(h(t, \cdot)) \quad (33.347)$$

est également C^∞ par la proposition 31.38. Ceci pour dire que $u \in C^\infty(\mathbb{R}, \mathcal{S}'(\mathbb{R}^d))$. Il faut encore vérifier que cette fonction est bien une solution de notre problème. Nous testons cette équation sur $\psi \in \mathcal{S}'(\mathbb{R} \times \mathbb{R}^d)$. Pour alléger les notations nous posons $\psi_t : x \mapsto \psi(t, x)$ et par conséquent aussi $(\partial_t \psi_t)(x) = (\partial_t \psi)(t, x)$. Nous avons :

$$\heartsuit = (\partial_t \tilde{u} - i\Delta \tilde{u})(\psi) \quad (33.348a)$$

$$= -\tilde{u}(\partial_t \psi) - i\tilde{u}(\Delta \psi) \quad (33.348b)$$

$$= -\int_{\mathbb{R}} u_t((\partial_t \psi_t) + i(\Delta \psi_t)) dt \quad (33.348c)$$

Ici nous nous souvenons du lemme 30.25 qui nous dit que nous pouvons permuter \mathcal{F}^{-1} et ∂_t . Et pour l'autre terme il faut utiliser le lemme 30.13 avec $|\alpha| = 2$ et une somme pour obtenir que

$$\widehat{\Delta \varphi}(x) = -\|x\|^2 \hat{\varphi}(x), \quad (33.349)$$

qui dans notre cas s'écrit sous la forme

$$\mathcal{F}^{-1}\left((\Delta \psi_t)\right)(x) = -\|x\|^2 \mathcal{F}^{-1}\psi(t, x). \quad (33.350)$$

En remettant bout à bout,

$$\heartsuit = -\int_{\mathbb{R}} (f_t \hat{g})\left((\partial_t - i\|\cdot\|^2)\mathcal{F}^{-1}\psi_t\right) dt \quad (33.351a)$$

$$= -\int_{\mathbb{R}} \hat{g}\left(x \mapsto e^{-it\|x\|^2}(\partial_t - i\|x\|^2)(\mathcal{F}^{-1}\psi)(t, x)\right) dt \quad (33.351b)$$

20. En utilisant la définition (31.46) du produit d'une distribution par une fonction.

Pour alléger les notations nous notons $\check{\psi}_t(x) = (\mathcal{F}^{-1}\psi)(t, x)$. Nous avons

$$\partial_t \left(e^{-it\|x\|^2} \check{\psi}_t(x) \right) = -i\|x\|^2 e^{-it\|x\|^2} \check{\psi}_t(x) + e^{-it\|x\|^2} (\partial_t \check{\psi}_t), x = e^{-it\|x\|^2} (\partial_t - i\|x\|^2) \check{\psi}_t(x); \tag{33.352}$$

cela nous permet d'un peu factoriser une dérivée dans \heartsuit :

$$\heartsuit = - \int_{\mathbb{R}} \hat{g} \left(\partial_t \left(e^{-it\|\cdot\|^2} \check{\psi}_t(\cdot) \right) \right) dt \tag{33.353a}$$

$$= - \int_{\mathbb{R}} \partial_t \hat{g} \left(e^{-it\|\cdot\|^2} \check{\psi}_t(\cdot) \right) dt \tag{33.353b}$$

$$= - \lim_{N \rightarrow \infty} \left[\hat{g} \left(e^{-i\|\cdot\|^2} \check{\psi}_t(\cdot) \right) \right]_{t=-N}^{t=N}. \tag{33.353c}$$

Histoire de bien comprendre les notations, il ne s'agit pas de calculer $\hat{g}(e^{-it\|\cdot\|^2} \check{\psi}_t)$ pour un t général et de remplacer ensuite t par N et $-N$. En effet la valeur de $\hat{g}(e^{-it\|\cdot\|^2} \check{\psi}_t)$ pour un t donné est celle qu'on obtient en calculant $\hat{g}(\dots)$ après avoir remplacé t par ce que l'on veut. Par conséquent, en posant $\varphi(t, \xi) = e^{-i\|\xi\|^2} \check{\psi}_t(\xi)$ nous avons :

$$\heartsuit = \lim_{N \rightarrow \infty} \left[g \left(x \mapsto \int_{\mathbb{R}} e^{-ix\xi} \varphi(t, \xi) d\xi \right) \right]_{t=-N}^{t=N} \tag{33.354a}$$

$$= \lim_{N \rightarrow \infty} g \left(x \mapsto \int_{\mathbb{R}} e^{-ix\xi} \varphi(N, \xi) d\xi \right) - \lim_{N \rightarrow \infty} g \left(x \mapsto \int_{\mathbb{R}} e^{-ix\xi} \varphi(-N, \xi) d\xi \right) \tag{33.354b}$$

La limite commute avec g parce que cette dernière est une distribution (continue). De plus la limite commute avec l'intégrale parce que ce qui est dedans est Schwartz. La fonction φ étant Schwartz, la limite est nulle. Donc

$$\heartsuit = 0. \tag{33.355}$$

Cela signifie que la fonction u proposée est bien une solution de l'équation de Schrödinger dans $C^\infty(\mathbb{R}, \mathcal{S}'(\mathbb{R}^d))$.

Unicité Nous considérons deux solutions $u_1, u_2 \in C^\infty(\mathbb{R}, \mathcal{S}'(\mathbb{R}^d))$ et la fonction $u = u_1 - u_2$ doit satisfaire au problème

$$\begin{cases} (\partial_t \tilde{u} - i\Delta \tilde{u})(\psi) = 0 \\ u_0 = 0. \end{cases} \tag{33.356a}$$

$$\tag{33.356b}$$

Nous allons montrer que seule la fonction $u_t = 0$ peut satisfaire à cela pour tout $\psi \in \mathcal{S}(\mathbb{R} \times \mathbb{R}^d)$. Nous allons même montrer qu'en imposant ces équations seulement sur la partie de $\mathcal{S}(\mathbb{R} \times \mathbb{R}^d)$ qui est à support compact par rapport à \mathbb{R} , la seule solution est $u_t = 0$. Soit donc $\psi \in \mathcal{S}(\mathbb{R} \times \mathbb{R}^d)$ à support compact vis-à-vis de sa variable t . Alors

$$0 = -\tilde{u}(\partial_t \psi + i\Delta \psi) = - \int_{\mathbb{R}} u_t \left((\partial_t \psi_t) + i(\Delta \psi_t) \right) dt \tag{33.357}$$

où encore une fois $\partial_t \psi_t$ est la fonction $x \mapsto (\partial_t \psi)(t, x)$. Maintenant nous utilisons la proposition 31.61 pour dire que

$$\frac{d}{dt} \left(u_t(\psi_t) \right) = u_t^{(1)}(\psi_t) + u_t \left(\frac{\partial \psi}{\partial t}(t, \cdot) \right) \tag{33.358}$$

pour écrire

$$0 = - \int_{\mathbb{R}} \frac{d}{dt} \left(u_t(\psi_t) \right) - u_t^{(1)}(\psi_t) + u_t(i(\Delta \psi)(t, \cdot)) dt \tag{33.359}$$

Le premier terme est facile :

$$\int_{\mathbb{R}} \frac{d}{dt} \left(u_t(\psi_t) \right) dt = \lim_{N \rightarrow \infty} \left[u_t(\psi_t) \right]_{t=-N}^{t=N} = 0 \tag{33.360}$$

parce que ψ est à support compact par rapport à t . Nous restons donc avec

$$\int_{\mathbb{R}} u_t^{(1)}(\psi_t) - iu_t((\Delta\psi)(t, \cdot)) dt = 0 \quad (33.361)$$

Nous traitons le terme en $u_t^{(1)}$ en utilisant le fait évident $T(\varphi) = (\mathcal{F}T)(\mathcal{F}^{-1}\varphi)$ et en remarquant le lemme 31.62 :

$$u_t^{(1)}(\psi_t) = (\mathcal{F}u_t^{(1)})(\mathcal{F}^{-1}\psi_t) = (\mathcal{F}u)_t^{(1)}(\mathcal{F}^{-1}\psi_t). \quad (33.362)$$

Pour l'autre terme on fait un peu la même chose en nous souvenant ce que fait la transformée de Fourier en traversant le laplacien :

$$u_t(\Delta\psi_t) = (\mathcal{F}u_t)(\mathcal{F}^{-1}\Delta\psi_t) = (\mathcal{F}u_t)(x \mapsto -\|x\|^2(\mathcal{F}^{-1}\psi_t)(x)). \quad (33.363)$$

En recollant encore :

$$\int_{\mathbb{R}} (\mathcal{F}u)_t^{(1)}(\mathcal{F}^{-1}\psi_t) + i(\mathcal{F}u_t)(\|\cdot\|^2\mathcal{F}^{-1}\psi_t) dt = 0. \quad (33.364)$$

Cette équation est valable tant que $\psi \in \mathcal{S}(\mathbb{R} \times \mathbb{R}^d)$ avec support compact en t . Nous allons nous en créer une super cool. D'abord nous choisissons $\varphi \in \mathcal{S}(\mathbb{R}^d)$ et $\chi \in \mathcal{D}(\mathbb{R})$ et nous considérons²¹

$$\psi(t, x) = \mathcal{F}\left(\xi \mapsto e^{it\|\xi\|^2}\varphi(\xi)\chi(t)\right)(x). \quad (33.365)$$

Notons que la transformée de Fourier conserve le fait qu'une fonction soit Schwartz²², mais pas le fait d'avoir support compact. Cependant nous ne prenons que la transformée de Fourier par rapport à x . Le résultat est donc une fonction ψ qui est Schwartz par rapport à x et support compact par rapport à t . Nous pouvons donc écrire (33.364) en utilisant la fonction (33.365) :

$$0 = \int_{\mathbb{R}} (\mathcal{F}u)_t^{(1)}\left(x \mapsto e^{it\|x\|^2}\varphi(x)\chi(t)\right) + i(\mathcal{F}u_t)\left(x \mapsto \|x\|^2 e^{it\|x\|^2}\varphi(x)\chi(t)\right) dt. \quad (33.366)$$

Là dedans, $\chi(t)$ peut sortir à la fois de la transformée de Fourier et de l'application des distributions; il doit seulement rester dans l'intégrale. Dans le second terme nous allons utiliser l'égalité (due entre autre à la proposition 31.61) :

$$\frac{d}{dt}(\hat{u}_t(e^{it\|\cdot\|^2}\varphi)) = \frac{d}{dt}\left(u_t(\mathcal{F}e^{it\|\cdot\|^2}\varphi)\right) \quad (33.367a)$$

$$= u_t^{(1)}(\mathcal{F}e^{it\|\cdot\|^2}\varphi) + u_t\left(\frac{\partial}{\partial t}\mathcal{F}e^{it\|\cdot\|^2}\varphi\right) \quad (33.367b)$$

$$= (\mathcal{F}u_t^{(1)})(x \mapsto e^{it\|x\|^2}\varphi(x)) + (\mathcal{F}u_t)(x \mapsto i\|x\|^2 e^{it\|x\|^2}\varphi(x)) \quad (33.367c)$$

$$= (\mathcal{F}u)_t^{(1)}(x \mapsto e^{it\|x\|^2}\varphi(x)) + (\mathcal{F}u_t)(x \mapsto i\|x\|^2 e^{it\|x\|^2}\varphi(x)). \quad (33.367d)$$

Et là, magie c'est exactement ce qui est dans (33.366). Donc

$$\int_{\mathbb{R}} \frac{d}{dt}\hat{u}_t(x \mapsto e^{it\|x\|^2}\varphi(x))\chi(t) dt = 0 \quad (33.368)$$

pour toute fonctions à support compact χ . Donc la proposition 31.1 nous dit que

$$\partial_t \hat{u}_t(x \mapsto e^{it\|x\|^2}\varphi(x)) = 0. \quad (33.369)$$

C'est zéro partout et non seulement presque partout parce qu'en plus nous avons la continuité. Par conséquent pour tout $t \in \mathbb{R}$ nous avons

$$\hat{u}_t(x \mapsto e^{it\|x\|^2}\varphi(x)) = \hat{u}_0(x \mapsto \varphi(x)) = 0. \quad (33.370)$$

21. Le candidat qui parvient à effectivement présenter ça comme développement, il est fort.

22. Proposition 30.14.

Et cela est vrai pour toute fonction $\varphi \in \mathcal{S}(\mathbb{R}^d)$. Nous considérons donc $t_0 \in \mathbb{R}$ et une fonction $\theta \in \mathcal{S}(\mathbb{R}^d)$ pour construire

$$\varphi(x) = e^{-it_0\|x\|^2}\theta(x). \quad (33.371)$$

Nous avons alors $\hat{u}_{t_0}(x \mapsto \theta(x)) = 0$, ce qui signifie que $\hat{u}_{t_0} = 0$. Du coup pour tout $\theta \in \mathcal{S}(\mathbb{R}^d)$ nous avons $u_{t_0}(\mathcal{F}\theta) = 0$, mais comme la transformée de Fourier est une bijection de $\mathcal{S}(\mathbb{R}^d)$ (proposition 30.14) nous avons en fait $u_{t_0}(\theta) = 0$ pour tout $\theta \in \mathcal{S}(\mathbb{R}^d)$, c'est-à-dire $u_{t_0} = 0$ pour tout $t_0 \in \mathbb{R}$ et au final $u = 0$.

□

33.12 Équations différentielles du premier ordre

Définition 33.42 (Équation différentielle du premier ordre).

Une **équation différentielle du premier ordre** est une équation qui, sur un intervalle donné, I , décrit la relation entre une variable réelle, notée x ou t dans I , une fonction $y : I \rightarrow \mathbb{R}$, et la dérivée première de y qui on note y' .

Souvent on écrit « $y'(x) =$ une formule contenant x et $y(x)$ », c'est à dire

$$y'(x) = f(x, y(x)), \quad \text{pour } x \in I, \quad (33.372)$$

où f est une fonction de deux variables réelles.

Remarque 33.43.

La théorie des fonctions de deux variables ne sera pas abordée dans ce cours, nous allons nous contenter de prendre f dans (33.372) comme une simple notation.

On peut presque toujours omettre d'écrire la dépendance de y en x et écrire simplement (33.372) sous la forme $y' = f(x, y)$.

Définition 33.44 (Solution particulière d'une équation différentielle du premier ordre).

Une **solution particulière** de l'équation (33.372) sur l'intervalle I est une fonction $z : I \rightarrow \mathbb{R}$ telle que :

- (1) z est dérivable sur I ;
- (2) $z'(x) = f(x, z(x))$, pour tout $x \in I$.

Définition 33.45 (Solution générale d'une équation différentielle du premier ordre).

Résoudre une équation différentielle veut dire trouver l'ensemble qui contient toutes ses solutions particulières. Cet ensemble s'appelle **solution générale** de l'équation.

Exemple 33.46

- (1) Résoudre une équation du type $y'(x) = f(x)$ revient à trouver l'ensemble des primitives de la fonction f , qui est donc la solution générale de cette équation. Il y a donc une infinité de solutions particulières, déterminées par une constante additive.

Si $f(x) = \sin(x)$ alors la solution générale sera $\mathcal{Y} = \{-\cos(x) + C : C \in \mathbb{R}\}$.

- (2) L'équation

$$y' = y, \quad x \in \mathbb{R}, \quad (33.373)$$

a peut-être été abordée dans votre cours de terminale lors de la définition de la fonction exponentielle. Sa solution générale est $\mathcal{Y} = \{Ce^x : C \in \mathbb{R}\}$. Ici aussi il y a une infinité de solutions particulières.

△

Remarque 33.47.

La solution générale d'une équation différentielle du premier ordre est une famille à un paramètre de fonctions.

Définition 33.48 (Équation différentielle du second ordre).

Une *équation différentielle du second ordre* est une équation qui, sur un intervalle donne, I , décrit la relation entre une variable réelle, notée x ou t dans I , une fonction $y : I \rightarrow \mathbb{R}$, et les dérivées première et seconde de y qui on note y' et y'' respectivement.

On utilise la forme générale

$$y'' = f(x, y, y'), \quad \text{pour } x \in I. \quad (33.374)$$

où f est une fonction de trois variables réelles.

On peut définir de manière analogue les équations différentielles d'ordre supérieur. Les définitions de solution particulière et de solution générale se généralisent aux équations différentielles d'ordre supérieur à un.

Définition 33.49 (Trajectoire).

La *trajectoire* tracée par une solution particulière y de l'équation (33.372) est le graphe de y en tant que fonction de x .

Exemple 33.50

Nous allons regarder de plus près l'équation (33.373), $y' = y$, pour tout $x \in \mathbb{R}$. Soient y_1 et y_2 deux solutions distinctes de cette équation. S'il existe un point \bar{x} tel que $y_1(\bar{x}) = y_2(\bar{x})$ alors forcément $y_1(\bar{x})/y_2(\bar{x}) = 1$. Or, la solution générale de l'équation est $\mathcal{Y} = \{Ce^x : C \in \mathbb{R}\}$, donc $y_i(x) = C_i e^x$, $i = 1, 2$, où les C_i sont des constantes. Le rapport $y_1(\bar{x})/y_2(\bar{x})$ vaut C_1/C_2 et par conséquent $C_1 = C_2$. Ce résultat contredit l'hypothèse que les deux solutions soient distinctes. On a donc montré que *deux trajectoires distinctes de cette équations ne se croisent jamais*.

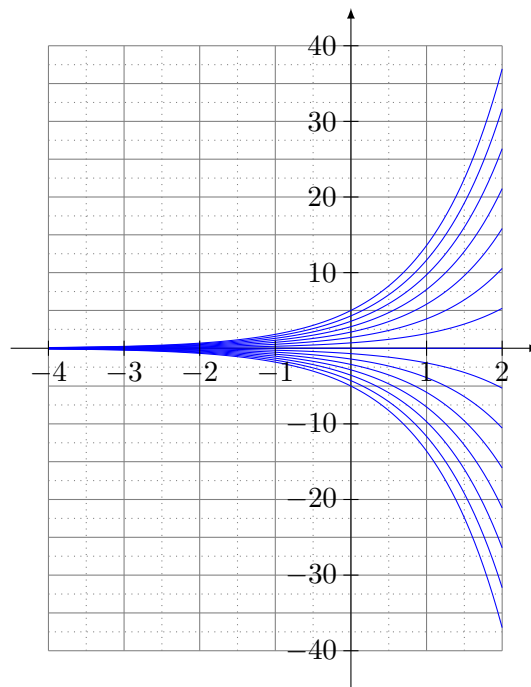


FIGURE 33.1 – Quelques trajectoires de l'équation $y' = y$.

La figure 33.1 représente quelques trajectoires de l'équation. Si on les avait tracées toutes elles recouvriraient tout le plan x - y . Cela veut dire que *par tout point (x, y) passe une et une seule trajectoire de l'équation (33.373)*.

△

Définition 33.51 (Condition initiale).

Une **condition initiale** pour l'équation (33.372) sur l'intervalle I est un point $(\bar{x}, \bar{y}) \in I \times \mathbb{R}$.

On dit que la solution particulière z de (33.372) satisfait la condition initiale $(\bar{x}, \bar{y}) \in I \times \mathbb{R}$ si $z(\bar{x}) = \bar{y}$.

Définition 33.52 (Problème de Cauchy).

L'association d'une équation différentielle et d'une condition initiale est appelée **problème de Cauchy**

$$\begin{cases} y' = f(x, y), & x \in I, \\ y(\bar{x}) = \bar{y}. \end{cases} \quad (33.375)$$

Remarque 33.53.

Sous des conditions assez générales qui seront toujours vérifiées dans ce cours, tout problème de Cauchy admet une et une seule solution.

Pour passer de la solution générale d'une équation différentielle de premier ordre à une solution particulière il faut choisir une valeur du paramètre. Comme il y a un seul paramètre une seule condition (la trajectoire de la solution doit passer par un point fixe du plan) peut suffire. Pour une équation différentielle de second ordre comme (33.374), nous aurons besoin de plus de conditions. Sans rentrer dans les détails, nous allons constater ce fait dans l'exemple suivant.

Exemple 33.54

La solution générale de l'équation

$$y'' = -y, \quad (33.376)$$

est $\mathcal{Y} = \{C_1 \cos(x) + C_2 \sin(x) : C_1, C_2 \in \mathbb{R}\}$. Remarquez que l'équation est du second ordre et que sa solution générale est une famille d'équations à deux paramètres réels. Ce sera toujours le cas pour les équations abordées dans la section 33.15. Pour déterminer une solution particulière de (33.376) il faut fixer les valeurs des deux paramètres et donc, en général, il sera nécessaire de donner deux conditions. △

Remarque 33.55.

Une condition comme $y(0) = 4$ nous dit que la constante $C_1 = 4$ mais elle ne nous permet pas de trouver C_2 . Il y a donc une infinité de solutions de (33.376) qui satisfont à la condition $y(0) = 4$.

On peut fixer les deux conditions de deux manières différentes.

- (1) Problème de Cauchy : on fixe une terne de valeurs réels $\bar{x}, \bar{y}, \bar{y}'$ et on cherche la solution telle que $y(\bar{x}) = \bar{y}$, $y'(\bar{x}) = \bar{y}'$.

Exemple 33.56

Les conditions $y(0) = 4$, $y'(0) = 15$ permettent de trouver la solution $z(x) = 4 \cos(x) + 15 \sin(x)$. △

- (2) Problème aux bords : on fixe deux points dans le plan x - y , $A = (\bar{x}, \bar{y})$ et $B = (\tilde{x}, \tilde{y})$, et on cherche la solution dont la trajectoire passe par A et B , c'est à dire, on impose $y(\bar{x}) = \bar{y}$, $y(\tilde{x}) = \tilde{y}$.

Exemple 33.57

Les conditions $y(0) = 4$, $y(\pi/2) = 15$ permettent de trouver la solution $z(x) = 4 \cos(x) + 15 \sin(x)$. △

33.13 Premier ordre, variables séparables

Pour certaines équations différentielles la recherche d'une solution particulière se réduit à une recherche de primitive moyennant un changement de variables.

Définition 33.58 (Équation différentielle du premier ordre à variables séparables).

Une *équation différentielle du premier ordre à variables séparables* est une équation qui, pour tout les x dans un intervalle donné, I , peut se mettre sous la forme

$$f(y)y' = g(x), \quad (33.377)$$

où f et g sont deux fonctions de \mathbb{R} dans \mathbb{R} .

Nous pouvons intégrer les deux côtés de l'égalité par rapport à x et obtenir

$$\int f(y(x))y'(x) dx = G(x) + C,$$

où G est une primitive de g et C une constante réelle. Il est facile à ce point d'effectuer un changement de variable dans le membre de gauche de l'équation en posant (sans surprise) $y = y(x)$ et donc $y'(x) dx = dy$.

$$\int f(y(x))y'(x) dx = \int f(y) dy = F(y(x)) + C,$$

où F est une primitive de f et C une constante réelle. En somme nous avons

$$F(y(x)) = G(x) + C,$$

et, si F admet une fonction réciproque, alors

$$y(x) = F^{-1}(G(x) + C). \quad (33.378)$$

Remarque 33.59.

L'expression de F^{-1} peut être difficile à calculer. Il sera alors préférable de garder y dans la forme implicite.

Exemple 33.60

L'équation

$$3y^2y' = x, \quad \text{pour tout } x \in \mathbb{R}, \quad (33.379)$$

est une équation à variables séparables. Pour reprendre les notations du début du chapitre, ici $f(y) = 3y^2$ et $g(x) = x$. En intégrant de deux côtés on trouve

$$y^3 = \frac{x^2}{2} + C.$$

La fonction $F(y) = y^3$ est une bijection de \mathbb{R} dans \mathbb{R} , donc nous pouvons écrire la solution générale de l'équation (33.379) dans la forme

$$\mathcal{Y} = \left\{ \left(\frac{x^2}{2} + C \right)^{1/3} \text{ tel que } C \in \mathbb{R} \right\}.$$

△

Exemple 33.61

En intégrant de deux côtés l'équation à variables séparables

$$2yy' = x, \text{ pour tout } x \in \mathbb{R}, \quad (33.380)$$

on trouve

$$y^2 = \frac{x^2}{2} + C.$$

La fonction $F(y) = y^2$ est *n'est pas inversible* sur tout \mathbb{R} , et on sait que $\sqrt{y^2} = |y|$. Au moment de rendre y explicite on doit choisir entre

$$y = \left(\frac{x^2}{2} + C\right)^{1/2} \quad \text{ou} \quad y = -\left(\frac{x^2}{2} + C\right)^{1/2}.$$

Ce choix se fait suivant la condition initiale, si elle est donnée. S'il n'y a pas de condition initiale nous pouvons écrire que la solution générale est l'ensemble

$$\mathcal{Y} = \left\{ y : \mathbb{R} \rightarrow \mathbb{R} \text{ tels que } y^2 = \frac{x^2}{2} + C \text{ et } C \in \mathbb{R} \right\}.$$

△

Exemple 33.62

On considère le problème de Cauchy

$$\begin{cases} e^y y' = \frac{1}{x+3}, & x \in]-\infty, -3[, \\ y(-4) = 0. \end{cases} \quad (33.381)$$

En intégrant des deux côtés nous trouvons

$$e^y = \ln(|x+3|) + C.$$

Nous pouvons alors imposer la condition initiale et obtenir $e^0 = \ln(|-4+3|) + C$, c'est à dire $C = 1 - \ln(1) = 1$.

Remarque 33.63.

L'énoncé du problème de Cauchy dit que x peut varier dans $]-\infty, -3[$, mais nous voyons maintenant que la solution n'est pas définie sur toute la demi-droite, parce que e^y est toujours positif et $\ln(|x+3|) + 1$ est positif seulement pour $x < -(1/e + 3) \approx -3,3679$.

Donc la solution du problème de Cauchy est $y(x) = \ln(|x+3|) + 1$ pour tout $x \in]-\infty, -(1/e+3)[$.
△

Exemple 33.64

Attention, cet exemple est le plus important de la section !

On considère l'équation à variables séparables

$$y' = \sin(x)y, \quad x \in \mathbb{R}. \quad (33.382)$$

Dans ce cas, pour pouvoir écrire l'équation dans la forme (33.377) il faut pouvoir multiplier les deux côtés par $1/y$. Il faut donc éliminer tout de suite le cas où $y = 0$.

Si $y = 0$ alors $y' = 0$ et on a une solution constante (on dit souvent : une solution stationnaire) de l'équation. Par ailleurs les trajectoires des solutions ne peuvent pas se croiser ; donc si y_G est une solution non nulle de l'équation (33.382) alors $y_G(x) \neq 0$ pour tout x ²³. Il n'y a donc aucun danger à diviser par y dans la recherche d'une solution non identiquement nulle.

Supposons maintenant que $y \neq 0$ et écrivons $y'/y = \sin(x)$. En intégrant des deux côtés on trouve

$$\ln(|y|) = -\cos(x) + C,$$

d'où

$$|y| = e^{-\cos(x)+C} = e^C e^{-\cos(x)}.$$

23. Ça vaut la peine de prendre un peu de temps pour bien comprendre cela.

Si on avait imposé une condition initiale alors on pourrait déterminer une solution particulière de l'équation en choisissant une valeur de la constante C . Nous pouvons observer cependant que la fonction exponentielle est bijective de \mathbb{R} dans $\mathbb{R}^{+,*}$ et par conséquent il n'y a pas de perte de généralité en disant que la solution générale de l'équation est

$$\mathcal{Y} = \left\{ y : |y| = Ke^{-\cos(x)}, \text{ pour } K \in \mathbb{R}^{+,*} \right\} \cup \{y \equiv 0\}.$$

Il n'empêche qu'il serait plus élégant d'écrire la solution générale de l'équation sous une forme plus explicite, sans valeur absolue. Nous pouvons le faire en nous rappelant que

$$|x| = \begin{cases} x & \text{si } x \geq 0, \\ -x & \text{si } x < 0, \end{cases}$$

Il suffit alors d'autoriser K dans \mathbb{R}^* pour éliminer la valeur absolue.

Pour écrire la solution générale de façon encore plus compacte nous observons que si $K = 0$ alors $y \equiv 0$, c'est à dire, on retrouve la solution constante nulle.

Finalement, la solution générale de cette équation sera toujours écrite sous la forme suivante

$$\mathcal{Y} = \left\{ y = Ke^{-\cos(x)}, \text{ pour } K \in \mathbb{R} \right\}. \quad (33.383)$$

△

33.14 Équations différentielles linéaires du premier ordre

Définition 33.65 (Équation différentielle linéaire du premier ordre).

Soit $I \subset \mathbb{R}$ un intervalle .

Une *équation différentielle linéaire du premier ordre* est une équation différentielle de la forme

$$a(x)y' + b(x)y = c(x), \quad \text{pour } x \in I, \quad (33.384)$$

où a, b, c sont des fonctions de \mathbb{R} dans \mathbb{R} et $a \neq 0$ pour tout $x \in I$.

On dit que a, b, c sont les coefficients de l'équation (33.384).

Remarque 33.66.

Une fonction $f : \mathbb{R} \rightarrow \mathbb{R}$ est dite *linéaire* si pour tout x_1, x_2 dans \mathbb{R} et pour tout couple de constantes λ et μ on a

$$f(\lambda x_1 + \mu x_2) = \lambda f(x_1) + \mu f(x_2). \quad (33.385)$$

Ces équations différentielles sont dites linéaires parce que la partie de l'équation qui contient y (le membre de gauche) satisfait la propriété (33.385) par rapport à y . En effet par les propriétés de la dérivée nous avons que

$$a(x)(\lambda y_1 + \mu y_2)' + b(x)(\lambda y_1 + \mu y_2) = \lambda(a(x)y_1' + b(x)y_1) + \mu(a(x)y_2' + b(x)y_2).$$

Définition 33.67.

L'équation (33.384) est dite *homogène* quand c est la fonction nulle. Si (33.384) n'est pas homogène on dit que l'équation

$$a(x)y' + b(x)y = 0, \quad (33.386)$$

est son *équation homogène associée*.

Toute équation linéaire du premier ordre homogène est une équation du premier ordre à variables séparables, comme nous en avons vu l'exemple 33.64. Nous n'allons pas répéter les détails du procédé pour trouver sa solution générale, qui aura la forme suivante

À retenir 33.68

$$\mathcal{Y}_h = \left\{ K e^{-\int \frac{b(x)}{a(x)} dx} : K \in \mathbb{R} \right\}. \quad (33.387)$$

Proposition 33.69. (1) Soit y_p une solution particulière de l'équation (33.384) et y_h une solution particulière de l'équation homogène associée (33.386). Alors la fonction somme $z = y_p + y_h$ est encore une solution particulière de l'équation (33.384).

(2) Soient y_1 et y_2 deux solutions particulières de (33.384). Alors la fonction différence $w = y_1 - y_2$ est une solution particulière de (33.386).

Démonstration. (1)

$$a(x)(y_p + y_h)' + b(x)(y_p + y_h) - c(x) = (a(x)y_p' + b(x)y_p - c(x)) + (a(x)y_h' + b(x)y_h) = 0. \quad (33.388)$$

(2)

$$a(x)(y_1 - y_2)' + b(x)(y_1 - y_2) = (a(x)y_1' + b(x)y_1 - c(x)) - (a(x)y_2' + b(x)y_2 - c(x)) = 0. \quad (33.389)$$

□

Cette proposition permet de démontrer le théorème suivant, qui est le plus important de cette section.

Théorème 33.70.

Soit y_p une solution particulière de l'équation (33.384) et \mathcal{Y}_h la solution générale de l'équation (33.386), alors la solution générale de l'équation (33.384) est l'ensemble

$$\mathcal{Y} = \mathcal{Y}_h + y_p = \{z = y_h + y_p : y = h \in \mathcal{Y}_h\}. \quad (33.390)$$

À retenir 33.71

La résolution d'une équation différentielle linéaire du premier ordre comporte trois étapes :

- (1) résolution de l'équation homogène associée ;
- (2) recherche d'une solution particulière de l'équation non homogène ;
- (3) somme de la solution générale de l'équation homogène et de la solution particulière trouvée au point précédent.

La partie qui nous manque encore est de savoir comment trouver une solution particulière de l'équation non homogène (33.384). Si la fonction c dans (33.384) est une constante ou un polynôme simple, ou une exponentielle alors on peut essayer de deviner. Cette méthode cependant n'est pas la plus sûre pour des débutants.

Exemple 33.72

On considère l'équation

$$y' - 5y = 10, \quad x \in \mathbb{R}. \quad (33.391)$$

Comme tous les coefficients de l'équation sont constants on peut essayer de trouver une solution constante.

Toutes les fonctions constantes ont une dérivée nulle, par conséquent, si une solution constante existe elle doit satisfaire $-5y = 10$, ce qui veut dire que la solution constante est $y(x) \equiv -2$. \triangle

Exemple 33.73

On considère l'équation

$$xy' + y = x + 1, \quad x \in \mathbb{R}^{+,*}. \quad (33.392)$$

Comme le membre de droite de l'équation est un polynôme de degré un on cherche une solution de la forme $y(x) = Ax + B$ avec A et B dans \mathbb{R} .

Par substitution on obtient $Ax + (Ax + B) = x + 1$, c'est-à-dire que une solution particulière de l'équation est $y(x) = x/2 + 1$. \triangle

Exemple 33.74

L'équation

$$xy' - y = x + 1, \quad x \in \mathbb{R}^{+,*}. \quad (33.393)$$

ressemble beaucoup à celle de l'exemple précédent, cependant il n'existe pas un polynôme de degré un qui en soit solution.

Dans un cas comme celui-ci, il faut rapidement abandonner la divination et replier sur la méthode, plus technique mais plus sûre, dite *variation de la constante*. \triangle

33.14.1 Méthode de variation de la constante

- Soit \mathcal{Y}_h la solution générale de l'équation homogène associé à (33.384). Il s'agit d'une famille à un paramètre de fonctions. La première étape de cette méthode consiste à construire un candidat solution particulière y_p en remplaçant le paramètre dans \mathcal{Y}_h par une fonction $C : \mathbb{R} \rightarrow \mathbb{R}$ à déterminer.

Exemple 33.75

L'équation homogène associée à $y' - y = \cos(x)$ est $y' - y = 0$, dont la solution générale est $\mathcal{Y}_h = \{Ce^x : C \in \mathbb{R}\}$. Le candidat solution sera alors $y_p = C(x)e^x$, avec C fonction à déterminer. \triangle

- La deuxième étape de cette méthode consiste à injecter y_p dans l'équation. Cela permet de trouver une équation différentielle à variables séparables pour C , en principe plus facile à résoudre que l'équation de départ.

Exemple 33.76

On continue avec l'exemple précédent. On a $y_p' = C'(x)e^x + C(x)e^x$, d'où

$$(C'(x)e^x + C(x)e^x) - C(x)e^x = \cos(x),$$

c'est-à-dire

$$C'(x) = \cos(x)e^{-x}.$$

\triangle

- La troisième étape de la méthode consiste à trouver une solution particulière de l'équation différentielle pour C et, par conséquent déterminer une y_p .

Exemple 33.77

La solution générale de

$$C'(x) = \cos(x)e^{-x}.$$

est $\mathcal{C} = \left\{ e^{-1} \frac{(\sin(x) - \cos(x))}{2} + K : K \in \mathbb{R} \right\}$. Il nous suffit une solution particulière, nous pouvons donc choisir $K = 0$ et alors la solution particulière de (33.384) sera $y_p(x) = \frac{\sin(x) - \cos(x)}{2}$. \triangle

Remarque 33.78.

Le plus souvent en intégrant l'équation pour C on en trouvera la solution générale. Dans ce cas on peut remplacer C par cette solution générale et obtenir d'un seul coup la solution générale de l'équation (33.384), c'est-à-dire sans faire la somme entre la solution générale de l'homogène associée et la solution particulière.

Exemple 33.79

Dans l'exemple qu'on vient de voir la solution générale de (33.384) est

$$\mathcal{Y} = \mathcal{Y}_h + y_p = \left\{ C e^x + \frac{(\sin(x) - \cos(x))}{2} : C \in \mathbb{R} \right\}. \quad (33.394)$$

On obtient le même résultat en écrivant $\mathcal{Y} = \left\{ e^{-x} \left(e^{-1} \frac{(\sin(x) - \cos(x))}{2} + K \right) : K \in \mathbb{R} \right\}$. Notez qu'on a changé le nom du paramètre de C à K seulement pour souligner qu'on obtient de même résultat par deux chemins différents, sinon les deux expressions sont équivalentes! \triangle

33.15 Équations différentielles linéaires du second ordre

Définition 33.80 (Équation différentielle linéaire du second ordre).

Une *équation différentielle linéaire du second ordre* est une équation différentielle de la forme

$$a(x)y'' + b(x)y' + c(x)y = d(x), \quad \text{pour } x \in I, \quad (33.395)$$

où a, b, c et d sont des fonctions de \mathbb{R} dans \mathbb{R} et $a \neq 0$ pour tout $x \in I$.

On dit que a, b, c et d sont les coefficients de l'équation (33.395).

Dans ce cours nous allons étudier exclusivement le cas où a, b et c sont des fonctions constantes.

Définition 33.81 (Équation différentielle linéaire du second ordre homogène).

Une *équation différentielle linéaire du second ordre homogène* est une équation différentielle de la forme (33.395), telle que le coefficient d est nul.

À toute équation de la forme (33.395) on peut associer une équation homogène exactement comme on a fait dans la section précédente pour les équations linéaires du premier ordre.

33.15.1 Équations différentielles linéaires du second ordre homogènes à coefficients constants

Remarque 33.82.

L'application qui à la fonction y fait correspondre $a(x)y'' + b(x)y' + c(x)y$ est linéaire, au sens de la remarque 33.66.

Cela nous dit en particulier, que si y_1 et y_2 sont deux solutions de l'équation homogène alors toute leur combinaison de la forme $z = \lambda y_1 + \mu y_2$, avec λ et μ dans \mathbb{R} , est encore une solution.

Jusqu'ici nous avons toujours travaillé avec des fonctions définies sur \mathbb{R} et à valeurs dans \mathbb{R} . Dans cette section nous nous autorisons à passer par des fonctions définies sur \mathbb{R} et à valeurs dans \mathbb{C} , mais cela sera uniquement une étape dans nos calculs. Au final toutes les solutions que nous allons considérer sont des fonctions à valeurs dans \mathbb{R} .

La solution générale à **valeurs dans les complexes** d'une équation de ce type a la forme

$$\mathcal{Y}_h^{\mathbb{C}} = \{ C_1 e^{r_1 x} + C_2 e^{r_2 x} : C_1, C_2 \in \mathbb{C}, x \in I \}, \quad (33.396)$$

où r_1 et r_2 sont aussi des nombres complexes. Remarquez que la solution générale est une famille à deux paramètres. Il faut aussi observer que en tout cas l'intervalle I dans lequel varie x est un intervalle dans \mathbb{R} , parce que I est une des données du problème.

À partir de cette information nous pouvons, pour toute équation donnée, chercher la solution générale **complexe** par substitution. Il suffit de remplacer y dans l'équation par e^{rx} et chercher les valeurs de r qui nous conviennent.

Si notre équation de départ est

$$ay'' + by' + cy = 0, \quad \text{pour } x \in I, \quad (33.397)$$

alors la substitution nous donne

$$e^{rx} (ar^2 + br + c) = 0.$$

Il est connu que la fonction exponentielle ne prend pas la valeur 0, par conséquent ce qui s'annule est le polynôme de degré deux $ar^2 + br + c$. Il est donc très facile de trouver les valeurs de r qu'on pourra utiliser comme r_1 et r_2 dans la solution générale **complexe**.

Si $b^2 - 4ac > 0$: le polynôme admet deux solutions réelles et distinctes, r_1 et r_2 ;

Si $b^2 - 4ac < 0$: le polynôme admet deux solutions complexes conjuguées, $r_1 = \alpha + i\beta$ et $r_2 = \alpha - i\beta$;

Si $b^2 - 4ac = 0$: le polynôme admet une solution réelle double $r = r_1 = r_2$.

Il faut maintenant écrire la solution générale **réelle** de l'équation, qui est celle que nous intéresse vraiment. La façon de l'obtenir est différente dans les trois cas.

Si $b^2 - 4ac > 0$: la solution générale réelle a la même forme que la solution complexe, (33.396), il suffit de prendre les paramètres C_1 et C_2 dans \mathbb{R} plutôt que dans \mathbb{C} .

$$\mathcal{Y}_h = \{C_1 e^{r_1 x} + C_2 e^{r_2 x} : C_1, C_2 \in \mathbb{R}, x \in I\}, \quad (33.398)$$

Si $b^2 - 4ac < 0$: le polynôme admet deux solutions complexes conjuguées, $r_1 = \alpha + i\beta$ et $r_2 = \alpha - i\beta$; Il faut alors utiliser les formules suivantes

$$\begin{aligned} e^{\alpha+i\beta} &= e^\alpha (\cos(\beta) + i \sin(\beta)) \\ e^{\alpha-i\beta} &= e^\alpha (\cos(\beta) - i \sin(\beta)). \end{aligned} \quad (33.399)$$

La somme $e^{r_1 x} + e^{r_2 x}$, où x est dans $I \in \mathbb{R}$, vaut

$$e^{(\alpha+i\beta)x} + e^{(\alpha-i\beta)x} = e^{\alpha x} (\cos(\beta x) + i \sin(\beta x)) + e^{\alpha x} (\cos(\beta x) - i \sin(\beta x)) = 2e^{\alpha x} \cos(\beta x)$$

et la différence $e^{r_1 x} - e^{r_2 x}$ vaut

$$e^{(\alpha+i\beta)x} - e^{(\alpha-i\beta)x} = e^{\alpha x} (\cos(\beta x) + i \sin(\beta x)) - e^{\alpha x} (\cos(\beta x) - i \sin(\beta x)) = 2e^{\alpha x} \sin(\beta x).$$

Par ces deux calculs élémentaires nous avons trouvé deux fonctions à valeurs dans \mathbb{R} qui n'ont pas de zéros en commun. Elles sont les génératrices de la famille des solutions réelles de l'équation différentielle (la solution générale)

$$\mathcal{Y}_h = \{e^{\alpha x} (C_1 \cos(\beta x) + C_2 \sin(\beta x)) : C_1, C_2 \in \mathbb{R}, x \in I\}, \quad (33.400)$$

Si $b^2 - 4ac = 0$: le polynôme admet une solution réelle double $r = r_1 = r_2$. Dans ce cas la solution générale de l'équation est la famille

$$\mathcal{Y}_h = \{(C_1 + C_2 x)e^{rx} : C_1, C_2 \in \mathbb{R}, x \in I\}. \quad (33.401)$$

Pour justifier cette formule nous observons d'abord que toute fonction $x \mapsto Ce^{rx}$, pour $C \in \mathbb{R}$ est une solution de l'équation différentielle (par construction). Ensuite nous utilisons la méthode de variation de la constante. On trouve rapidement que si une fonction de la forme $x \mapsto C(x)e^{rx}$ est une solution alors $C(x)$ est un polynôme de degré au plus 1, c'est-à-dire $C(x) = C_1 + C_2 x$ avec C_1 et C_2 dans \mathbb{R} .

33.15.2 Linéaires du second ordre à coefficients constants, non homogènes

Nous ne présentons pas une méthode générale pour la résolution de ces équations. Comme dans le cas des équations différentielles linéaires du premier ordre non homogènes, la solution générale de (33.395) est donnée par la somme d'une solution particulière et de la solution générale de l'équation homogène associée. La recherche d'une solution particulière est facilitée par le fait que les coefficients de (33.395) sont supposés constants, c'est-à-dire que a , b et c sont des fonctions constantes. Il faut essayer de deviner la forme d'une solution particulière à partir de la forme du second membre de l'équation, la fonction d . Si d est un polynôme il faut essayer avec un polynôme du même degré, si d est une exponentielle, par exemple $d(x) = e^{5x}$, on pourra essayer avec un multiple de la même fonction exponentielle, dans l'exemple $f(x) = ke^{5x}$, avec k à déterminer. Si d est une combinaison linéaire de sinus et cosinus, comme par exemple $12 \cos(x) + 2 \sin(x)$, on peut essayer avec $k_1 \cos(x) + k_2 \sin(x)$.

Exemple 33.83

On considère l'équation différentielle

$$y'' + 12y' + 36y = -192e^{2x}, \quad x \in \mathbb{R}. \quad (33.402)$$

Son équation homogène associée est

$$y'' + 12y' + 36y = 0, \quad (33.403)$$

dont le polynôme caractéristique est $r^2 + 12r + 36$. Ce polynôme admet une racine double, qui est -6 , par conséquent la solution générale de (33.403) est

$$\mathcal{Y}_h = \{(C_1 + C_2x)e^{-6x} : C_1, C_2 \in \mathbb{R}, x \in \mathbb{R}\}.$$

Le membre de droite de (33.402) est une fonction exponentielle, nous allons donc chercher une solution particulière de (33.402) de la forme $f(x) = ke^{2x}$. Par substitution nous trouvons

$$ke^{2x}(4 + 12 \times 2 + 36) = -192e^{2x},$$

ce qui veut dire que k doit être -3 .

La solution générale de l'équation (33.402) est donc

$$\mathcal{Y} = \{(C_1 + C_2x)e^{-6x} - 3e^{2x} : C_1, C_2 \in \mathbb{R}, x \in \mathbb{R}\}.$$

△

Exemple 33.84

Nous allons résoudre l'équation

$$y'' + 12y' + 36y = 12 \cos(x) + 2 \sin(x), \quad x \in \mathbb{R}. \quad (33.404)$$

Cette équation a comme homogène associée l'équation (33.403), comme dans l'exemple précédent. Il nous suffit donc de trouver une solution particulière de (33.402).

Nous pouvons essayer avec $f(x) = k_1 \cos(x) + k_2 \sin(x)$. Par substitution on trouve

$$\begin{aligned} & -(k_1 \cos(x) + k_2 \sin(x)) + 12(-k_1 \sin(x) + k_2 \cos(x)) + 36(k_1 \cos(x) + k_2 \sin(x)) \\ & = 12 \cos(x) + 2 \sin(x) \end{aligned}$$

Cette équation doit être satisfaite pour toute valeur de x , en particulier pour $x = 0$ et $x = \pi/2$. Cela revient à considérer séparément les coefficients des fonctions sinus et cosinus. Il faut alors que k_1 et k_2 soient solutions du système

$$\begin{cases} -k_1 + 12k_2 + 36k_1 & = 12, \\ -k_2 - 12k_1 + 36k_2 & = 2. \end{cases}$$

On trouve $k_1 = 396/1369$ et $k_2 = 214/1369$, et la solution générale de notre équation est

$$\mathcal{Y} = \left\{ (C_1 + C_2 x)e^{-6x} + \frac{396}{1369} \cos(x) + \frac{214}{1369} \sin(x) : C_1, C_2 \in \mathbb{R}, x \in \mathbb{R} \right\}.$$

△

Exemple 33.85

Nous allons résoudre l'équation

$$y'' + 12y' + 36y = 10x^2 + 3, \quad x \in \mathbb{R}. \quad (33.405)$$

Cette équation a comme homogène associée l'équation (33.403), comme dans l'exemple précédent. Il nous suffit donc de trouver une solution particulière de (33.402).

Nous pouvons essayer avec $f(x) = k_1 x^2 + k_2 x + k_3$. Par substitution on trouve

$$(2k_1) + 12(2k_1 x + k_2) + 36(k_1 x^2 + k_2 x + k_3) = 10x^2 + 3.$$

Pour trouver les bonnes valeurs des coefficients nous devons résoudre le système

$$\begin{cases} 36k_1 & = 10, \\ 24k_1 + 36k_2 & = 0, \\ 2k_1 + 12k_2 + 36k_3 & = 3, \end{cases}$$

ce qui donne $k_1 = 5/18$, $k_2 = -5/27$ et $k_3 = 7/54$. La solution générale de notre équation est

$$\mathcal{Y} = \left\{ (C_1 + C_2 x)e^{-6x} + \frac{5}{18}x^2 - \frac{5}{27}x + \frac{7}{54} : C_1, C_2 \in \mathbb{R}, x \in \mathbb{R} \right\}.$$

△

33.16 Fonction de Green

Soit l'équation différentielle

$$\begin{cases} y''(x) = g(x) & (33.406a) \\ y(0) = y(1) = 0 & (33.406b) \end{cases}$$

pour $x \in]0, 1[$ et où g est continue sur $]0, 1[$.

Nous définissons la fonction de Green

$$G(x, t) = \begin{cases} t(x-1) & \text{si } 0 \leq t \leq x \leq 1 \\ x(t-1) & \text{si } 0 \leq x \leq t \leq 1, \end{cases} \quad (33.407)$$

et nous allons montrer que

$$y(x) = \int_0^1 G(x, t)g(t)dt \quad (33.408)$$

est l'unique solution.

Unicité Si y_1 et y_2 sont des solutions, alors $y_1'' = y_2''$ et donc $y_1(x) = y_2(x) + ax + b$. Les conditions aux bords donnent alors $0 = y_1(0) = y_2(0) + b = b$. D'où $b = 0$. En imposant $y_1(1) = 0$ nous trouvons alors immédiatement $a = 0$, ce qui donne $y_1 = y_2$.

Existence Il est vite vérifié qu'avec (33.408) nous avons $y(0) = y(1) = 0$ parce que $G(0, t) = G(1, t) = 0$ pour tout t . Nous fixons une valeur pour $x \in]0, 1[$ et nous découpons l'intégrale :

$$y(x) = \int_0^x G(x, t)g(t)dt + \int_x^1 G(x, t)g(t)dt. \quad (33.409)$$

Pour calculer $y'(x)$, il faut dériver à la fois à travers l'intégrale et dans la borne. Si vous connaissez une formule pour faire cela, c'est bien pour vous. Nous allons faire ça à la main et poser

$$I(x, y) = \int_0^y t(x-1)g(t)dt. \quad (33.410)$$

La dérivation de I par rapport à x se fait en utilisant le théorème 18.25 :

$$\frac{\partial I}{\partial x}(x, y) = \int_0^y tg(t)dt. \quad (33.411)$$

Pour la dérivation par rapport à y , il s'agit du théorème fondamental de l'analyse, plus précisément le lien primitive et intégrale de la proposition 15.232 :

$$\frac{\partial I}{\partial y}(x, y) = y(x-1)g(y). \quad (33.412)$$

Maintenant nous considérons la fonction $\varphi_I(x) = I(x, x)$. Elle satisfait à

$$\varphi_I'(x) = \frac{\partial I}{\partial x}(x, x) + \frac{\partial I}{\partial y}(x, x) = \int_0^x tg(t)dt + x(x-1)g(x). \quad (33.413)$$

Le même jeu avec $J(x, y) = \int_y^1 x(t-1)g(t)dt$ donne

$$\varphi_J'(x) = \int_0^x fg(t)dt + x(x-1)g(x). \quad (33.414)$$

En remettant les bouts ensemble,

$$y(x) = \int_0^x tg(t)dt + \int_1^x (1-t)g(t)dt. \quad (33.415)$$

Le calcul de la dérivée seconde donne alors

$$y''(x) = xg(x) + (1-x)g(x) = g(x). \quad (33.416)$$

Nous pouvons aussi, sur cette équation, estimer la variation de la solution en termes d'une variation de g . Soit donc une fonction continue δ_g sur $[0, 1]$ et $\tilde{g} = g + \delta_g$. Nous considérons l'équation différentielle

$$\begin{cases} \tilde{y}''(x) = \tilde{g}(x) \\ \tilde{y}(0) = \tilde{y}(1) = 0. \end{cases} \quad (33.417a)$$

$$\quad (33.417b)$$

Par ce que nous venons de faire, l'unique solution est

$$\tilde{y}(x) = \int_0^1 G(x, t)\tilde{g}(t)dt = \int_0^1 G(x, t)g(t)dt + \int_0^1 G(x, t)\delta_g(t)dt = y(x) + \delta_y(x) \quad (33.418)$$

où δ_y est une fonction continue ainsi définie :

$$\delta_y(x) = \int_0^1 G(x, t)\delta_g(t)dt. \quad (33.419)$$

Supposons que $\|\delta_g\|_\infty = \epsilon$. Alors des majorations donnent

$$|\delta_y(x)| \leq \epsilon \int_0^1 |G(x, t)|dt = \epsilon(1-x) \int_0^x tdt + \epsilon x \int_x^1 (1-t)dt = \frac{\epsilon}{2}x(1-x). \quad (33.420)$$

Mais la fonction $x \mapsto x(1-x)$ a son maximum en $x = \frac{1}{2}$, donc nous pouvons donner une majoration indépendante de x :

$$\|\delta_y\|_\infty \leq \frac{1}{8}\|\delta_g\|_\infty. \quad (33.421)$$

Notons que la majoration (33.421) en norme uniforme a l'air plus impressionnante, mais la majoration (33.420) donnant une majoration séparée pour chaque x est en réalité plus précise.

Chapitre 34

Équations aux dérivées partielles

34.1 Symbole principal, équation des caractéristiques

Soit l'équation différentielle semi-linéaire d'ordre k

$$\sum_{|\alpha|=k} a_\alpha(x)(\partial^\alpha u)(x) + F\left(x, u(x), (Du)(x), \dots, (D^{k-1}u)(x)\right) = 0 \quad (34.1)$$

pour la fonction $u: \mathbb{R}^d \rightarrow \mathbb{R}$.

Définition 34.1.

Le *symbole principal* de l'équation (34.1) est l'application

$$\begin{aligned} \sigma: \mathbb{R}^d \times \mathbb{R}^d &\rightarrow \mathbb{R} \\ (x, \xi) &\mapsto \sum_{|\alpha|=k} a_\alpha(x)\xi^\alpha \end{aligned} \quad (34.2)$$

où si $\alpha = (\alpha_1, \dots, \alpha_d)$ et $\xi = (\xi_1, \dots, \xi_d)$ alors $\xi^\alpha = \xi_1^{\alpha_1} \dots \xi_d^{\alpha_d}$.

Définition 34.2.

Les *caractéristiques* de l'équation (34.1) est une surface S de \mathbb{R}^d donné par une équation de la forme $\phi(x) = 0$ où ϕ satisfait à

$$\sigma(x, \nabla\phi(x)) = 0 \quad (34.3)$$

et $\nabla\phi(x) \neq 0$ pour tout $x \in S$.

34.2 Méthode des caractéristiques pour l'ordre 1

Nous [462, 463] voulons étudier l'équation d'ordre 1

$$a(x, y)\frac{\partial u}{\partial x}(x, y) + b(x, y)\frac{\partial u}{\partial y}(x, y) + c(x, y)u(x, y) = f(x, y) \quad (34.4)$$

Le champ de vecteurs associé à cette équation est

$$v = \begin{pmatrix} a \\ b \end{pmatrix}, \quad (34.5)$$

et l'équation peut être écrite sous la forme

$$(v \cdot \nabla) + cu = f. \quad (34.6)$$

Définition 34.3.

Le *flot* de ce champ de vecteurs sont les courbes paramétriques $\gamma(t) = (x(t), y(t))$ vérifiant $\gamma'(t) = v(\gamma(t))$.

Les équations du flot pour l'équation (34.4) sont

$$\begin{cases} x'(t) = a(x(y), y(t)) & (34.7a) \\ y'(t) = b(x(y), y(t)). & (34.7b) \end{cases}$$

Ce sont des équations différentielles ordinaires. Un système de deux équations couplées du premier ordre.

Quel est l'intérêt du flot ? Nous allons voir que sur la ligne $t \mapsto \gamma(t)$, la fonction u est constante. Or des solutions γ au système (34.7), il y en aura plusieurs : une pour chaque valeur des constantes d'intégration. Pour peu que ces lignes recouvrent tout le plan, nous pourrions résoudre l'équation de départ ligne par ligne.

Nous posons

$$\tilde{u}(t) = u(x(t), y(t)) \quad (34.8a)$$

$$\tilde{c}(t) = c(x(t), y(t)) \quad (34.8b)$$

$$\tilde{f}(t) = f(x(t), y(t)). \quad (34.8c)$$

La fonction \tilde{u} est une fonction $\mathbb{R} \rightarrow \mathbb{R}$ normale qui se dérive normalement, en suivant la règle de dérivation des fonctions composées :

$$\tilde{u}'(t) = \frac{\partial u}{\partial x}(x(t), y(t))x'(t) + \frac{\partial u}{\partial y}(x(t), y(t))y'(t) \quad (34.9a)$$

$$= a \frac{\partial u}{\partial x} + b \frac{\partial u}{\partial y} \quad (34.9b)$$

$$= f(x(t), y(t)) - c(x(t), y(t))u(x(t), y(t)) \quad (34.9c)$$

$$= \tilde{f}(t) - \tilde{c}(t)\tilde{u}(t). \quad (34.9d)$$

Nous avons pour \tilde{u} l'équation différentielle ordinaire

$$\tilde{u}' + \tilde{c}\tilde{u} = \tilde{f} \quad (34.10)$$

qui est résolue par la proposition 33.10.

34.2.1 Un exemple complet un peu minimal

Nous considérons l'équation différentielle[462]

$$\frac{\partial u}{\partial x} - \frac{\partial u}{\partial y} - (x - y)u = 0. \quad (34.11)$$

Et nous allons la résoudre.

Les équations du flot, sont simples parce que les coefficients sont des constantes : $x'(t) = 1$, $y'(t) = -1$. Donc

$$x(t) = t + C_1 \quad (34.12a)$$

$$y(t) = -t + C_2. \quad (34.12b)$$

A priori nous avons une caractéristique pour chaque choix de (C_1, C_2) et nous espérons que le tout recouvre le plan \mathbb{R}^2 . En fait seule une des deux constantes doit être laissée libre, l'autre consiste seulement en décaler le paramètre t . Nous posons donc $C_1 = 0$ et nous considérons les courbes caractéristiques

$$\gamma_C(t) = \begin{pmatrix} t \\ -t + C \end{pmatrix}. \quad (34.13)$$

Ces courbes recouvrent bien tout le plan. Pour savoir les valeurs de u sur la courbe γ_C , nous devons résoudre l'équation différentielle ordinaire

$$\tilde{u}'_C + \tilde{c}\tilde{u}_C = \tilde{f}, \quad (34.14)$$

en sachant que $\tilde{c}(t) = c(x(t), y(t)) = -(x(t) - y(t)) = 2t - C$. Cela se fait en suivant la méthode décrite dans l'exemple 33.8 et résumée dans la proposition 33.10.

En termes de notations, $\tilde{u}_C(t) = u(\gamma_C(t))$. Récrivons l'équation :

$$\tilde{u}'(t) - (2t - C)\tilde{u}(t) = 0. \quad (34.15)$$

La méthode pour la résoudre est de mettre les \tilde{u} d'un côté et les t de l'autre :

$$\frac{\tilde{u}'}{\tilde{u}} = 2t - C. \quad (34.16)$$

En intégrant par rapport à t des deux côtés,

$$\ln(\tilde{u}) = t^2 - Ct + K_C, \quad (34.17)$$

c'est-à-dire (avec redéfinition de K_C)

$$\tilde{u}(t) = K_C e^{t^2 - Ct} \quad (34.18)$$

ou encore

$$u(\gamma_C(t)) = K_C e^{t^2 - Ct} \quad (34.19)$$

où C est le paramètre que nous déterminons en sachant sur quelle caractéristique se trouve le point (x, y) où nous voulons calculer $u(x, y)$ et K est une constante (strictement positive parce que si vous avez suivi le mouvement, c'est une exponentielle) qui doit être déterminée par les conditions initiales. Dès que K est fixé pour un des points de la courbe γ_C , alors il est fixé pour tous les points.

Ce que nous avons obtenu est qu'il existe un K_C tel que pour tout t nous avons

$$u(\gamma_C(t)) = K_C e^{t^2 - Ct}. \quad (34.20)$$

Soit donc un point $(x_0, y_0) \in \mathbb{R}^2$. Nous devons d'abord déterminer où ce point se trouve par rapport aux caractéristiques, c'est-à-dire quelle est la valeur de C pour laquelle (x_0, y_0) est sur la courbe γ_C , et ensuite déterminer pour quelle valeur de t nous aurons $\gamma_C(t) = (x_0, y_0)$. À résoudre :

$$\gamma_C(t_0) = \begin{pmatrix} t_0 \\ -t_0 + C \end{pmatrix} = \begin{pmatrix} x_0 \\ y_0 \end{pmatrix}. \quad (34.21)$$

Donc $t_0 = x_0$ et $C = x_0 + y_0$. En reprenant (34.19) nous avons

$$u(\gamma_C(t_0)) = K e^{x_0^2 - (x_0 + y_0)x_0} = K e^{-x_0 y_0}. \quad (34.22)$$

Pour peu que des conditions soient donnée sur chaque caractéristique, nous pouvons déterminer K . Attention : ce K est une constante d'intégration de l'équation différentielle ordinaire pour \tilde{u} . Donc elle n'est valable que sur chaque caractéristique séparément. Cela n'est donc pas du tout une constante sur \mathbb{R}^2 .

Nous pouvons maintenant écrire la solution générale de l'équation de départ. L'équation cartésienne de la courbe γ_C est

$$x + y = C. \quad (34.23)$$

Donc K est une fonction de $x + y$, pas de x et y séparément. Cela est important à comprendre. A priori nous avons

$$u(x, y) = K(x, y) e^{-xy} \quad (34.24)$$

où $K(x, y)$ est constante sur la courbe γ_C contenant (x, y) . Nous avons

- si $x_1 + y_1 = x_2 + y_2$,
- alors il existe C tel que (x_1, y_1) et (x_2, y_2) sont sur γ_C ,
- alors $K(x_1, y_2) = K(x_2, y_2)$.

Donc il existe une fonction $\mathbb{R} \rightarrow \mathbb{R}$ telle que $K(x, y) = f(x + y)$.

Au final, la solution générale de l'équation est

$$u(x, y) = f(x + y) e^{-xy} \quad (34.25)$$

où f est une fonction à déterminer par les conditions initiales qui peuvent être données. Typiquement nous espérons que les conditions imposent une et une seule valeur de u sur chacune des courbes γ_C .

34.2.2 Un théorème d'existence et d'unicité

La méthode des caractéristiques donne essentiellement une preuve de l'unicité des solutions aux équations de transport, et une méthode pour construire cette solution. En effet, la procédure suivante permet de construire $u(x_0, y_0)$.

- Trouver la caractéristique passant par le point (x_0, y_0)
- Calculer en quel point elle passe par une condition initiale donnée.
- Attribuer à $u(x_0, y_0)$ la valeur trouvée sur la caractéristique là où elle passe par une condition initiale.

Rien ne permet a priori de savoir que cette procédure construit effectivement une solution. En particulier, comment calculer $\partial_x u$? Le quotient différentiel serait

$$\frac{\partial u}{\partial x}(x, y) = \lim_{\epsilon \rightarrow 0} \frac{u(x + \epsilon, y) - u(x, y)}{\epsilon}, \quad (34.26)$$

mais la caractéristique donnant la valeur de $u(x + \epsilon, y)$ est différente pour chaque ϵ . Rien a priori ne permet d'affirmer que le calcul soit simple, ni qu'il arrive à une solution du problème donné.

D'où la nécessité d'avoir un résultat un peu rigoureux donnant des conditions sous lesquelles les choses vont bien.

Proposition 34.4 (Équation de transport à coefficients variables[463, 454]).

Soit une fonction $c: \mathbb{R}^2 \rightarrow \mathbb{R}$ de classe C^2 en ses deux variables et uniformément Lipschitziennes en sa première variable¹ et $g \in C^1(\mathbb{R})$. Alors l'équation aux dérivées partielles de premier ordre

$$\begin{cases} \frac{\partial u}{\partial x}(x, t) + c(x, t) \frac{\partial u}{\partial t}(x, y) = 0 \\ u(x, 0) = h(y) \end{cases} \quad (34.27a)$$

$$(34.27b)$$

admet une unique solution de classe C^1 .

Cette solution est construite de la façon suivante². D'abord nous considérons la solution X au problème

$$\begin{cases} \frac{\partial X}{\partial s}(s; x, t) = c(X(s; x, t), s) \\ X(t; x, t) = x, \end{cases} \quad (34.28a)$$

$$(34.28b)$$

et ensuite le problème (34.27) a pour unique solution

$$u(x, t) = h(X(0; x, t)). \quad (34.29)$$

Démonstration. Nous commençons par étudier l'existence et l'unicité de la fonction X définie par le problème 34.28. La fonction c ici est dans les hypothèses de la fonction f du théorème de Cauchy-Lipschitz global 18.41. D'où l'existence et l'unicité de la fonction $s \mapsto X(s; x, t)$ sur \mathbb{R} pour chaque (x, y) donné³.

Le lemme 33.32 nous dit que X est de classe C^2 en (s, x, t) . Donc nous pourrions dériver et permuter les dérivées autant que nous voudrions (sans exagérer : ordre 2 au maximum).

Unicité Nous montrons que u doit être constante sur le chemin

$$\gamma_{(x,t)}(s) = \begin{pmatrix} X(s; x, t) \\ s \end{pmatrix}. \quad (34.30)$$

1. Dans [454], on ne demande que continue puis uniformément Lipschitz. Moi je crois que ce n'est pas assez pour assurer la dérivabilité de X par rapport à x , et encore moins pour permuter les dérivées dans $\partial_{ix}^2 Y$.

2. Le fait que la construction ait un sens fait partie des choses à prouver

3. Avec l'énoncé tel que donné dans [463], il faut utiliser la technique de 33.20 pour l'existence globale, parce que la fonction b là-dedans n'est pas dans les mêmes hypothèses.

En effet, en posant

$$\varphi(s) = u(\gamma(s)) = u(X(s; x, t), s), \quad (34.31)$$

et en dérivant nous obtenons

$$\varphi'(s) = \frac{\partial u}{\partial t}(X(s; x, t), s) \frac{\partial X}{\partial s}(s; x, t) + \frac{\partial u}{\partial t}(X(s; x, t)) = 0. \quad (34.32)$$

Par conséquent la valeur commune de tous les $u(\gamma_{(x,t)}(s))$ doit être celle en $\gamma_{(x,t)}(0) = h(X(0; x, t))$

Cela prouve l'unicité parce que la valeur de u est fixée en tout point. Nous devons encore vérifier que la fonction u ainsi construite est bien une solution du problème. C'est l'objet de la partie « existence » de la preuve.

Existence Même l'existence est divisée en plusieurs étapes.

Mise en place Nous prouvons que la fonction u donné par (34.29) est une solution du problème. Nous avons :

$$\frac{\partial u}{\partial t}(x, t) = h'(X(0; x, t)) \frac{\partial X}{\partial t}(0; x, t) \quad (34.33)$$

et

$$\frac{\partial u}{\partial x}(x, t) = h'(X(0; x, t)) \frac{\partial X}{\partial x}(0; x, t), \quad (34.34)$$

de sorte qu'en posant

$$g(s; x, t) = \frac{\partial X}{\partial t}(s; x, t) + c(x, t) \frac{\partial X}{\partial x}(s; x, t) \quad (34.35)$$

nous avons

$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = h'(X(0; x, t)) g(0; x, t). \quad (34.36)$$

Une équation différentielle pour g Nous allons prouver que

$$\frac{\partial g}{\partial s}(s; x, t) = \alpha_{(x,t)}(s) g(s; x, t) \quad (34.37)$$

avec ⁴

$$\alpha_{(x,t)}(s) = \frac{\partial c}{\partial x}(X(s; x, t), s). \quad (34.38)$$

D'abord nous avons

$$\frac{\partial g}{\partial s}(s; x, t) = \frac{\partial^2 X}{\partial s \partial t}(s; x, t) + c(x, t) \frac{\partial^2 X}{\partial s \partial x}(s; x, t). \quad (34.39)$$

Nous permutons les dérivées et nous tenons compte de (34.28) :

$$\frac{\partial g}{\partial s}(s; x, t) = \frac{\partial}{\partial t} \left(c(X(s; x, t), s) \right) + c(x, t) \frac{\partial}{\partial x} \left(c(X(s; x, t), s) \right). \quad (34.40)$$

Nous dérivons maintenant plus en profondeur. D'une part

$$\frac{\partial}{\partial t} \left(c(X(s; x, t), s) \right) = \frac{\partial c}{\partial x} \left(X(s; x, t), s \right) \frac{\partial X}{\partial t}(s; x, t) \quad (34.41)$$

et d'autre part,

$$\frac{\partial}{\partial x} \left(c(X(s; x, t), s) \right) = \frac{\partial c}{\partial x} \left(X(s; x, t), s \right) \frac{\partial X}{\partial x}(s; x, t), \quad (34.42)$$

de telle sorte que

$$\frac{\partial g}{\partial s}(s; x, t) = \frac{\partial c}{\partial x} \left(X(s; x, t), s \right) \left[\frac{\partial X}{\partial t}(s; x, t) + c(x, t) \frac{\partial X}{\partial x}(s; x, t) \right] \quad (34.43a)$$

$$= \alpha_{(x,t)}(s) g(s; x, t). \quad (34.43b)$$

4. La ligne suivante est une de celles qui me font penser qu'il manque des hypothèses dans [454]. Il faut bien pouvoir dériver c .

Une condition initiale pour g Nous montrons maintenant que $g(t; x, t) = 0$. La condition initiale pour X est $X(t; x, t) = x$ pour tout $t, x \in \mathbb{R}$. Nous dérivons cette dernière par rapport à x et à t :

$$\frac{\partial X}{\partial x}(t; x, t) = 1 \quad (34.44a)$$

$$\frac{\partial X}{\partial s}(t; x, t) + \frac{\partial X}{\partial t}(t; x, t) = 0. \quad (34.44b)$$

Mais $\partial_s X(t; x, t) = c(X(t; x, t), t)$, donc la relation (34.44b) donne

$$\frac{\partial X}{\partial t}(t; x, t) = -c(X(t; x, t), t) = -c(x, t) \quad (34.45)$$

où nous avons tenu compte du fait que $X(t; x, t) = x$.

Voyons à présent ce que (34.44a) et (34.45) donnent pour $g(t; x, t)$:

$$g(t; x, t) = \frac{\partial X}{\partial t}(t; x, t) + c(x, t) \frac{\partial X}{\partial x}(t; x, t) = -c(x, t) + c(x, t) = 0. \quad (34.46)$$

Conclusion pour g La fonction g vérifie l'équation différentielle

$$\begin{cases} \frac{\partial g}{\partial s}(s) = \alpha(s)g(s) \\ g(t) = 0. \end{cases} \quad (34.47a)$$

$$(34.47b)$$

Bien entendu, $g(s) = 0$ est une solution. Mais la solution à ce système est unique par Cauchy-Lipschitz 18.40. Ici nous utilisons le fait que

$$(s, y) \mapsto \alpha(s)y \quad (34.48)$$

est continue. C'est-à-dire entre autres que

$$s \mapsto \frac{\partial c}{\partial x}(X(s; x, t), s) \quad (34.49)$$

doit être continue. C'est le cas parce que c est de classe C^1 en sa première variable. □

Les hypothèses de la proposition 34.4 sont loin d'être optimales. Voici un exemple dans lequel c n'est même pas dérivable par rapport à t et qui se passe très bien quand même.

Exemple 34.5

Soit l'équation différentielle

$$\begin{cases} \frac{\partial u}{\partial t}(x, t) + |t-1| \frac{\partial u}{\partial x}(x, t) = 0 \\ u(x, 0) = h(x) \end{cases} \quad (34.50a)$$

$$(34.50b)$$

où h est une fonction bien régulière; mettons C^p . En suivant la méthode de la proposition nous devrions poser l'équation différentielle

$$\begin{cases} \frac{\partial X}{\partial s}(s; x, t) = |s-1| \\ X(t; x, t) = x. \end{cases} \quad (34.51a)$$

$$(34.51b)$$

Cela est la caractéristique passant par (x, t) . Cependant il sera plus simple de chercher les caractéristiques en demandant qu'elles passent par $(x_0, 0)$. Nous allons donc plutôt résoudre pour X_{x_0} l'équation différentielle

$$\begin{cases} \frac{\partial X}{\partial s}(s) = |s-1| \\ X_{x_0}(0) = x_0 \end{cases} \quad (34.52a)$$

$$(34.52b)$$

et les courbes caractéristiques seront les chemins

$$\gamma_{x_0}(s) = \begin{pmatrix} X_{x_0}(s) \\ s \end{pmatrix}. \quad (34.53)$$

La résolution donne d'abord

$$X_{x_0}(s) = \begin{cases} -\frac{s^2}{2} + s + K_1 & \text{si } s < 1 \\ \frac{s^2}{2} - s + K_2 & \text{si } s > 1. \end{cases} \quad (34.54)$$

Vu que la condition initiale est donnée pour $s = 0$, nous fixons K_1 pour la condition initiale et K_2 pour la continuité :

$$X_{x_0}(s) = \begin{cases} -\frac{s^2}{2} + s + x_1 & \text{si } s < 1 \\ \frac{1}{2} + x_0 & \text{si } s = 1 \\ \frac{s^2}{2} - s + 1 + x_0 & \text{si } s > 1. \end{cases} \quad (34.55)$$

Soit $(x, t) \in \mathbb{R}^2$. Quelle caractéristique passe par là? Nous allons déterminer la fonction $x_0(x, t)$ qui donne le x_0 tel que la caractéristique γ_{x_0} passe par (x, t) . Nous devons résoudre

$$\begin{pmatrix} X_{x_0}(s) \\ s \end{pmatrix} = \begin{pmatrix} x \\ t \end{pmatrix}. \quad (34.56)$$

Directement : $s = t$. Et ensuite $X_{x_0}(t) = x$. Nous avons

$$x_0(x, t) = \begin{cases} x + \frac{t^2}{2} - t & \text{si } t < 1 \\ x - \frac{1}{2} & \text{si } t = 1 \\ x - \frac{t^2}{2} + t - 1 & \text{si } t > 1. \end{cases} \quad (34.57)$$

Le truc presque étonnant est que x_0 est de classe C^1 . En effet le calcul de

$$\frac{\partial x_0}{\partial t}(1) = \lim_{\epsilon \rightarrow 0} \frac{x_0(x, 1 + \epsilon) - x_0(x, 1)}{\epsilon} \quad (34.58)$$

se fait en séparant les limites $\epsilon \rightarrow 0^+$ et $\epsilon \rightarrow 0^-$. Le résultat est que $\partial_t x_0(1) = 0$. Nous avons donc

$$\frac{\partial x_0}{\partial t}(t) = \begin{cases} t - 1 & \text{si } t < 1 \\ 0 & \text{si } t = 1 \\ -t + 1 & \text{si } t > 1. \end{cases} \quad (34.59)$$

Cela étant continu, la fonction x_0 est de classe C^1 en t , et la dérivée en x étant toujours 1, elle est de classe C^1 .

En ce qui concerne la solution de l'équation de départ,

$$u(x, t) = h(X_{x_0(x, t)}(0)) = h(x_0(x, t)). \quad (34.60)$$

Pourvu que h soit assez régulière, la fonction u est facilement de classe C^1 . \triangle

34.3 Méthode des caractéristiques pour l'ordre 2

34.3.1 Principe général

Soit l'opérateur différentiel agissant sur $C^2(\mathbb{R}^2)$:

$$D = a(x, y) \frac{\partial^2}{\partial x^2} + b(x, y) \frac{\partial^2}{\partial x \partial y} + c(x, y) \frac{\partial^2}{\partial y^2}. \quad (34.61)$$

Nous voulons résoudre des équations du type $Du = 0$ pour $u: \mathbb{R}^2 \rightarrow \mathbb{R}$.

Pour commencer [462], et c'est le point crucial, nous voyons D comme un polynôme en ∂_x et ∂_y et nous le factorisons : si

$$aX^2 + bXY + cY = (\alpha X + \beta Y)(\gamma X + \delta Y) \quad (34.62)$$

alors nous avons

$$D = \left(\alpha \frac{\partial}{\partial x} + \beta \frac{\partial}{\partial y} \right) \left(\gamma \frac{\partial}{\partial x} + \delta \frac{\partial}{\partial y} \right) + \text{termes d'ordre inférieurs.} \quad (34.63)$$

Les « termes d'ordre inférieurs » sont ceux de la forme $\alpha(x, y) \frac{\partial \delta}{\partial x} \frac{\partial}{\partial y}$.

L'astuce est de poser

$$v = (\gamma \partial_x + \delta \partial_y)u, \quad (34.64)$$

et de résoudre le système

$$\begin{cases} (\alpha \partial_x + \beta \partial_y)v = 0 & (34.65a) \\ (\gamma \partial_x + \delta \partial_y)u = v. & (34.65b) \end{cases}$$

Cela sont deux équations différentielles du premier ordre pour lesquelles nous avons déjà des techniques décrites en la section 34.2.

Afin que les fonctions α , β , γ et δ soient réelles, il faut que $b^2 - 4ac \geq 0$. Sachant que $a = \alpha\gamma$, $b = \alpha\delta + \beta\gamma$ et $c = \beta\delta$ cette condition sur a , b et c donne

$$(\alpha\delta + \beta\gamma)^2 - 4\alpha\gamma\beta\delta \geq 0. \quad (34.66)$$

Cela revient à

$$(\alpha\delta - \beta\gamma)^2 \geq 0. \quad (34.67)$$

Nous supposons à présent que l'inégalité soit stricte (cas hyperbolique). Nous avons en particulier que

$$\alpha\delta - \beta\gamma \neq 0. \quad (34.68)$$

Cette condition implique que les équations

$$\frac{dx}{dt} = \alpha(x, y) \quad \frac{dy}{dt} = \beta(x, y) \quad (34.69)$$

sont indépendantes des équations

$$\frac{dx}{dt} = \gamma(x, y) \quad \frac{dy}{dt} = \delta(x, y) \quad (34.70)$$

Ce sont les équations caractéristiques des équations (34.65).

34.3.2 Exemple : l'équation d'onde

Nous considérons l'équation aux dérivées partielles

$$\frac{\partial^2 u}{\partial t^2} - c^2 \frac{\partial^2 u}{\partial x^2} = 0 \quad (34.71)$$

où c est une constante réelle. Nous en cherchons des solutions de classe C^2 .

L'opérateur différentiel est donné par le polynôme $P(T, X) = T^2 - c^2 X^2$ qui se factorise en

$$P = (T + cX)(T - cX), \quad (34.72)$$

c'est-à-dire que nous pouvons récrire l'équation des ondes sous la forme

$$(\partial_t + c\partial_x)(\partial_t - c\partial_x)u = 0. \quad (34.73)$$

Nous posons donc $v = (\partial_t - c\partial_x)u$ et nous avons le système[453]

$$\begin{cases} (\partial_t + c\partial_x)v = 0 & (34.74a) \\ (\partial_t - c\partial_x)u = v. & (34.74b) \end{cases}$$

La méthode des caractéristiques est efficace pour résoudre la première, et pour trouver la solution générale de l'homogène associée à la seconde.

Nous nous lançons dans la résolution de (34.74a). Le flot est $v = \begin{pmatrix} 1 \\ c \end{pmatrix}$, et nous cherchons ses courbes intégrales sous la forme $\varphi(t) = (t, x(t))$. Immédiatement, $x'(t) = c$, ce qui donne

$$\gamma_C(t) = \begin{pmatrix} t \\ ct + C \end{pmatrix}. \quad (34.75)$$

Cela donne une caractéristique pour chaque valeur de C . En posant $\tilde{v}_C(t) = v(t, ct + C)$ nous avons

$$\tilde{v}'_C(t) = \frac{\partial v}{\partial t}(\gamma_C(t)) + c \frac{\partial v}{\partial x}(\gamma_C(t)) = 0. \quad (34.76)$$

Donc \tilde{v}_C est une fonction constante. Donc u est constant sur la courbe γ_C dont l'équation cartésienne est $x - ct = C$. Cela implique que

$$v(t, x) = f(x - ct) \quad (34.77)$$

où f est une fonction de classe C^1 . En effet si (t_1, x_1) et (t_2, x_2) vérifient $x_1 - ct_1 = x_2 - ct_2$ alors $v(t_1, x_1) = v(t_2, x_2)$. Le fait que f soit C^1 est une demande que u soit au final dans C^2 .

Nous devons maintenant résoudre l'équation (34.74b)

$$(\partial_t - c\partial_x)u = v. \quad (34.78)$$

Nous allons agir conformément à la stratégie expliquée par le lemme 33.6. Nous devons résoudre $Du = v$ avec

$$\begin{aligned} D: C^2(\mathbb{R}) &\rightarrow D(C^2(\mathbb{R})) \\ u &\mapsto (\partial_t - c\partial_x)u. \end{aligned} \quad (34.79)$$

Par la même méthode des caractéristiques que celle déjà menée plus haut nous trouvons $\ker(D)$ comme solution générale de $(\partial_t - c\partial_x)u_G = 0$. C'est-à-dire

$$u_G = g(x + ct) \quad (34.80)$$

où g est une fonction quelconque de classe C^2 .

Il nous faut maintenant une solution particulière de

$$(\partial_t - c\partial_x)u_P(t, x) = f(x - ct). \quad (34.81)$$

Si F est une primitive de f alors

$$u_P(t, x) = -\frac{1}{2c}F(x - ct) \quad (34.82)$$

fonctionne. Vu que f est quelconque dans $C^1(\mathbb{R})$, la fonction F est un élément quelconque de $C^2(\mathbb{R})$. Au final, la solution générale de l'équation des ondes est

$$u(t, x) = g_1(x + ct) + g_2(x - ct) \quad (34.83)$$

où g_1 et g_2 sont des éléments de $C^2(\mathbb{R})$.

34.4 Classification des équations du second ordre

Soit une équation générale d'ordre 2 sur \mathbb{R}^2 :

$$a \frac{\partial^2 u}{\partial x^2} + b \frac{\partial^2 u}{\partial x \partial y} + c \frac{\partial^2 u}{\partial y^2} + d \frac{\partial u}{\partial x} + e \frac{\partial u}{\partial y} + \beta u = f. \quad (34.84)$$

En ce qui concerne son symbole principal nous avons

$$\sigma(x, y, \xi_1, \xi_2) = a(x, y)\xi_1^2 + b(x, y)\xi_1\xi_2 + c(x, y)\xi_2^2, \quad (34.85)$$

ce qui donne l'équation des caractéristiques

$$a \left(\frac{\partial \phi}{\partial x} \right)^2 + b \frac{\partial \phi}{\partial x} \frac{\partial \phi}{\partial y} + c \left(\frac{\partial \phi}{\partial y} \right)^2 = 0. \quad (34.86)$$

Si nous nous posons sur un point (x_0, y_0) tel que $\nabla \phi(x_0, y_0) = 0$ et $\partial_x \phi(x_0, y_0) \neq 0$ alors, via le théorème de la fonction implicite⁵, la condition $\phi(x, y) = 0$ définit une fonction $y \mapsto x(y)$ vérifiant

$$\phi(x(y), y) = 0 \quad (34.87)$$

pour tout y dans un voisinage de y_0 .

Nous pouvons obtenir une équation différentielle ordinaire pour x de la façon suivante. D'abord nous posons $\varphi(y) = \phi(x(y), y)$ et ensuite nous calculons la dérivée de φ (qui est nulle par construction) :

$$0 = \varphi'(y) = \frac{\partial \phi}{\partial x} x' + \frac{\partial \phi}{\partial y}. \quad (34.88)$$

Nous pouvons donc remplacer $\partial_y \phi$ par $x' \partial_x \phi$ dans l'équation des caractéristiques (34.86) :

$$a \left(\frac{\partial \phi}{\partial x} \right)^2 + b x' \left(\frac{\partial \phi}{\partial x} \right)^2 + c \left(\frac{\partial \phi}{\partial x} \right)^2 (x')^2 = 0. \quad (34.89)$$

Vu que nous avons supposé $(\partial_x \phi) \neq 0$ sur un voisinage de (x_0, y_0) nous pouvons simplifier par $(\partial_x \phi)^2$ et avoir l'équation différentielle ordinaire

$$a(x(y), y) + b(x(y), y)x'(y) + c(x(y), y)x'(y)^2 = 0. \quad (34.90)$$

Notons que, conformément à ce que raconte le théorème des fonctions implicites, nous avons pris un voisinage de y_0 suffisamment petit pour que $x(y)$ reste dans un voisinage de x_0 . Ce voisinage étant, nous pouvons le restreindre pour nous assurer du signe de a , b et c . Cela est évidemment très théorique parce que le théorème de la fonction implicite parle de l'existence de voisinages, mais pas de façon de les construire.

Nous donnons la classification suivante.

Définition 34.6.

Si $b^2 - 4ac < 0$ alors l'équation est **elliptique**.

Si $b^2 - 4ac > 0$ alors l'équation est **hyperbolique**.

Si $b^2 - 4ac = 0$ alors l'équation est **parabolique**.

Exemple 34.7

Un exemple d'équation parabolique est l'équation de la chaleur

$$\frac{\partial u}{\partial t} - \alpha \frac{\partial^2 u}{\partial x^2} = 0 \quad (34.91)$$

où $\alpha > 0$ est une constante. Cette équation est avec $a = c = 0$, donc elle est parabolique. \triangle

5. Théorème 18.49.

Une équation aux dérivées partielles peut changer de nature selon le point.

Exemple 34.8

Soit l'équation

$$\frac{\partial^2 u}{\partial x^2} - (x^2 - y^2) \frac{\partial^2 u}{\partial y^2} = 0. \quad (34.92)$$

Son $b^2 - 4ac$ vaut $4(x^2 - y^2)$. Elle peut donc être hyperbolique, parabolique ou elliptique selon le point où l'on se trouve. \triangle

34.4.1 Problème au limite

Dans les définitions qui suivent nous considérons un ouvert $\Omega \subset \mathbb{R}^d$ assez régulier et possédant en particulier un vecteur normal extérieur $n(x)$ pour tout point $x \in \partial\Omega$.

Définition 34.9.

Un problème aux limites **de Dirichlet** est d'imposer la condition

$$u(x) = g(x) \quad (34.93)$$

pour tout $x \in \Gamma \subset \partial\Omega$. C'est-à-dire imposer la valeur de u sur une partie du bord du domaine.

Définition 34.10.

Un problème aux limites de **Von Neumann** est d'imposer

$$\frac{\partial u}{\partial n} \cdot x = g(x) \quad (34.94)$$

pour tout $x \in \Gamma \subset \partial\Omega$. C'est-à-dire imposer les valeurs de la dérivée normale de u sur une partie du bord.

Exemple 34.11

Lorsqu'on veut imposer un flux de chaleur aux bords d'un domaine pour l'équation de la chaleur, il s'agit de poser des conditions de type Von Neumann. \triangle

Définition 34.12.

Soit un domaine Ω de \mathbb{R}^d et un opérateur différentiel L sur une partie de $\text{Fun}(\Omega)$. Soit une fonction g sur $\partial\Omega$. Un problème **aux limites stationnaires** est un problème du type : trouver u définie sur Ω telle que

$$\begin{cases} L(u) = f & (34.95a) \\ u|_{\partial\Omega} = g. & (34.95b) \end{cases}$$

Définition 34.13.

Un problème **aux limites d'évolution** est du type : trouver $u \in \text{Fun}([0, \infty[\times \Omega)$ tel que

$$\begin{cases} \frac{\partial^m u}{\partial t^m} + L(u) = f & \text{sur }]0, \infty[\times \Omega \\ u(t, \cdot) = g(t, \cdot) & \text{sur }]0, \infty[\times \partial\Omega \\ u(0, \cdot) = u_0 & \text{sur } \Omega \end{cases} \quad (34.96)$$

où u_0 est une fonction sur Ω .

L'opérateur L ne doit pas opérer sur la partie « t » de u .

Définition 34.14 (Problème bien posé au sens de Hadamard).

Un problème aux limites est **bien posé au sens de Hadamard** si

(1) Il admet une unique solution.

(2) La solution dépend de façon continue en les données du problème.

La continuité est au sens des normes sur les fonctions sur $\partial\Omega$ et sur Ω en ce qui concerne les fonctions « données » du problème et des normes pour les fonctions sur Ω ou $\bar{\Omega}$ en ce qui concerne la solution.

Exemple 34.15 (Un problème de Dirichlet bien posé)

Trouver la fonction u définie sur $\Omega = [0, 1]^2$ telle que

$$\begin{cases} -\Delta u = 0 & \text{sur } \Omega \\ u = g & \text{sur } \partial\Omega. \end{cases} \quad (34.97)$$

△

34.5 Principe du maximum

Lemme 34.16 ([352]).

Soit un ouvert borné $\Omega \subset \mathbb{R}^n$. Soit une matrice symétrique strictement définie positive A telle que $A_{ij} \in C^0(\bar{\Omega})$ pour laquelle il existe $\lambda > 0$ minorant toutes les valeurs propres de toutes les matrices $A(x)$ pour $x \in \Omega$ ⁶.

Nous posons $L' = -\sum_{ij} A_{ij} \partial_{ij}$. Si $u \in C^2(\Omega)$ atteint un minimum local en $x_0 \in \Omega$, alors

$$(L'u)(x_0) \leq 0. \quad (34.98)$$

Démonstration. Nous allons bien entendu diagonaliser A . Si T est une matrice nous avons, en posant $u(x) = v(Tx)$:

$$\frac{\partial u}{\partial x_i}(x) = \sum_k \frac{\partial v}{\partial x_k}(Tx) T_{ki} \quad (34.99)$$

et

$$\frac{\partial^2 u}{\partial x_j \partial x_i}(x) = \sum_{kl} T_{ki} T_{lj} \frac{\partial^2 v}{\partial x_l \partial x_k}(Tx). \quad (34.100)$$

Si T est en particulier une matrice orthogonale diagonalisant A (théorème 11.189) nous avons $T^{-1} = T^t$ et

$$\sum_{ij} T_{ki} A_{ij} T_{jl}^{-1} = D_{kl} = \lambda_k \delta_{kl} \quad (34.101)$$

où les λ_k sont les valeurs propres de A . Notons que partout ici, tout est fonction de x sur Ω : tant A que T que les λ_k . Avec tous ces résultats nous calculons vite que

$$(Lu)(x_0) = -\sum_k \lambda_k (\partial_k^2 v)(Tx_0). \quad (34.102)$$

Si $x_0 \in \Omega$ est un minimum local de u , alors l'application v a un minimum local en $T^{-1}x_0$. Et donc

$$\frac{\partial^2 v}{\partial x_k^2}(Tx_0) \geq 0. \quad (34.103)$$

Du coup,

$$(Lu)(x_0) \leq 0. \quad (34.104)$$

□

6. Cela est plus que dire que toutes les $A(x)$ sont symétriques strictement définies positives.

Lemme 34.17 ([352]).

Soit un ouvert borné $\Omega \subset \mathbb{R}^n$. Soit une matrice symétrique strictement définie positive A telle que $A_{ij} \in C^0(\bar{\Omega})$ pour laquelle il existe $\lambda > 0$ minorant toutes les valeurs propres de toutes les matrices $A(x)$ pour $x \in \Omega$.

Nous posons

$$L = - \sum_{ij} A_{ij} \partial_{ij}^2 + \sum_i b_i \partial_i + c \tag{34.105}$$

où $b_i, c \in C^0(\bar{\Omega})$.

Soit $u \in C^0(\bar{\Omega}) \cap C^2(\Omega)$ telle que $Lu \geq 0$ sur Ω . Alors

(1) Si $c = 0$ alors

$$\min_{\bar{\Omega}}(u) = \min_{\partial\Omega}(u). \tag{34.106}$$

(2) Si $c \geq 0$ alors

$$\min_{\bar{\Omega}}(u) \geq \min_{\partial\Omega}(-u_-) \tag{34.107}$$

où u_- est défini par

$$u_-(x) = \begin{cases} 0 & \text{si } u(x) \geq 0 \\ -u(x) & \text{si } u(x) \leq 0. \end{cases} \tag{34.108}$$

Démonstration. D'abord, vu que Ω est borné, la fermeture $\bar{\Omega}$ est compacte et u y atteint son minimum. De plus $\partial\Omega$ est également compact (borné et le complémentaire est ouvert parce que Ω est ouvert). Donc u y atteint également son minimum. Cela pour dire que les minima écrits dans (34.106) et (34.107) ont un sens.

Pour (1).

$Lu \geq \eta > 0$ Nous supposons qu'il existe $\eta > 0$ tel que $(Lu)(x) \geq \eta$ pour tout $x \in \Omega$. Soit x_0 le point minimum sur $\bar{\Omega}$. Si $x_0 \in \Omega$ alors il est intérieur et $\partial_i u(x_0) = 0$ par la proposition 18.69.

Dans ce cas nous avons

$$(Lu)(x_0) = (L'u)(x_0) > 0, \tag{34.109}$$

ce qui contredit le lemme 34.16. Nous en déduisons que le minimum de u sur $\bar{\Omega}$ n'est pas atteint dans Ω , mais sur $\partial\Omega$. Pour la définition de la frontière, voir 8.52. Ici nous avons $\bar{\Omega} \setminus \Omega = \partial\Omega$.

Cela prouve (34.106) dans ce cas.

$Lu \geq 0$ sur Ω Nous prenons maintenant le cas général. Nous posons

$$u_{\gamma, \epsilon}(x) = u(x) - \epsilon e^{\gamma x_1}. \tag{34.110}$$

Nous avons⁷

$$L(e^{\gamma x_1}) = e^{\gamma x_1} (-A_{11}(x)\gamma^2 + b_1(x)\gamma). \tag{34.111}$$

Soit γ suffisamment grand pour que

$$\lambda\gamma^2 - \|b_1\|_{\bar{\Omega}}\gamma > 0. \tag{34.112}$$

Ici λ minore toutes les valeurs propres des $A(x)$ et nous notons que γ ne dépend pas de ϵ . En utilisant l'inégalité du lemme 11.19611.196, pour tout $\xi \in \mathbb{R}^n$ nous avons

$$\sum_{ij} A_{ij} \xi_i \xi_j \geq \lambda |\xi|^2, \tag{34.113}$$

ce qui donne avec $\xi = e_1$: $A_{11} \geq \lambda$, et même pour être plus précis : $A_{11}(x) \geq \lambda$ pour tout x . Ces inégalités donnent

$$-A_{11}(x)\gamma^2 + b_1(x)\gamma \leq -\lambda\gamma^2 + \|b_1\|_{\bar{\Omega}}\gamma < 0. \tag{34.114}$$

7. Nous abusons un peu de l'écriture parce que ce que nous calculons vraiment est $L(x \mapsto e^{\gamma x_1})$.

Nous avons donc

$$L(e^{\gamma x_1}) = e^{\gamma x_1} (-A_{11}(x)\gamma^2 + b_1(x)\gamma) \leq -\lambda\gamma^2 + \|b_1\|_{\bar{\Omega}}\gamma < 0. \quad (34.115)$$

Nous posons

$$\eta = \epsilon(\lambda\gamma^2 - \|b_1\|_{\bar{\Omega}}\gamma) \min_{\bar{\Omega}}(e^{\gamma x_1}), \quad (34.116)$$

où le minimum a un sens parce que Ω est borné. Nous avons $\eta > 0$.

C'est le moment de calculer ce que u_ϵ peut pour nous :

$$L(u_\epsilon) = Lu - \epsilon L(e^{\gamma x_1}) \quad (34.117a)$$

$$\geq -\epsilon L(e^{\gamma x_1}) \quad (34.117b)$$

$$= -\epsilon e^{\gamma x_1} (-A_{11}(x)\gamma^2 + b_1(x)\gamma) \quad (34.117c)$$

$$\geq -\epsilon e^{\gamma x_1} (\lambda\gamma^2 - \|b_1\|_{\bar{\Omega}}\gamma) \quad (34.117d)$$

$$\geq \eta > 0. \quad (34.117e)$$

Justification :

— (34.117b) parce que $Lu \geq 0$.

— (34.117d) par (34.115).

La fonction u_ϵ est donc dans le cas précédent et nous avons

$$\min_{x \in \bar{\Omega}} (u(x) - \epsilon e^{\gamma x_1}) = \min_{x \in \partial\Omega} (u(x) - \epsilon e^{\gamma x_1}). \quad (34.118)$$

Mentionnons le fait que le choix fait de γ ne dépend pas de ϵ . Nous pouvons donc encore faire varier ϵ sans toucher à γ et en maintenant toutes les inégalités prouvées jusqu'ici.

Vu que l'expression $e^{\gamma x_1}$ est majorable sur $\bar{\Omega}$ nous avons convergence uniforme

$$u_\epsilon \xrightarrow{\|\cdot\|_{\bar{\Omega}}} u \quad (34.119)$$

pour $\epsilon \rightarrow 0$.

Supposons qu'aucun point de $\partial\Omega$ ne réalise le minimum de u . Alors il existe $x_0 \in \Omega$ tel que $u(x_0) > u(x)$ pour tout $x \in \partial\Omega$. Mais comme $\partial\Omega$ est compact, il existe $\eta > 0$ tel que

$$u(x_0) < u(x) + \eta. \quad (34.120)$$

Soit ϵ tel que $\|u - u_\epsilon\|_{\bar{\Omega}} < \eta/2$. Nous avons

$$u(x_0) < u(x) - \eta \quad (34.121)$$

et donc aussi

$$u_\epsilon(x_0) < u(x_0) + \frac{\eta}{2} < u(x) - \frac{\eta}{2} < u_\epsilon(x), \quad (34.122)$$

ce qui signifierait que u_ϵ prend son minimum dans Ω . Or nous savons que ce n'est pas le cas.

Donc il existe un point de $\partial\Omega$ qui réalise le minimum de u .

Pour (2).

Supposons pour commencer que $u \geq 0$ sur Ω . Alors par continuité $u \geq 0$ sur $\bar{\Omega}$. Alors $u_- = 0$ et

$$\min_{\bar{\Omega}} u \geq 0, \quad (34.123)$$

ce qui fait que l'inégalité (34.107) est évidente.

Nous supposons donc que l'ensemble

$$\Omega_- = \{x \in \Omega \text{ tel que } u(x) < 0\} \quad (34.124)$$

est non vide. Notons que c'est également un ouvert.

Soit $\bar{L}u = Lu - cu$. Vu que $c \geq 0$ et que $Lu \geq 0$ nous avons $\bar{L}u \geq 0$ sur Ω_- . Par le point (1) nous avons

$$\min_{x \in \bar{\Omega}} u(x) = \min_{x \in \partial\Omega} u(x). \quad (34.125)$$

Mais le minimum de u est certainement atteint dans Ω_- , donc

$$\min_{x \in \bar{\Omega}_-} u(x) = \min_{x \in \bar{\Omega}} u(x). \quad (34.126)$$

Cela nous donne une première bonne égalité :

$$\min_{x \in \partial\Omega_-} u(x) = \min_{x \in \bar{\Omega}} u(x). \quad (34.127)$$

Nous pouvons la prolonger :

$$\min_{x \in \bar{\Omega}} u(x) = \min_{x \in \partial\Omega_-} u(x) \quad (34.128a)$$

$$= \min_{x \in \partial\Omega_-} (-u_-(x)) \quad (34.128b)$$

$$= \min_{x \in \partial\Omega_- \cap \partial\Omega} (-u_-(x)) \quad (34.128c)$$

$$= \min_{x \in \partial\Omega} (-u_-(x)) \quad (34.128d)$$

Justifications :

— Pour (34.128c) nous avons la décomposition

$$\partial\Omega_- = (\partial\Omega_- \cap \Omega) \cup (\partial\Omega_- \cap \partial\Omega). \quad (34.129)$$

Or sur $\partial\Omega_- \cap \Omega$ nous avons $u(x) = 0$ et donc pas le minimum.

— Pour (34.128d). Sur $\partial\Omega$, le minimum est atteint dans la partie $\partial\Omega_-$ parce que le reste ne contient que des valeurs positives de u .

Cela prouve ce que nous voulions. \square

Théorème 34.18.

(Principe du maximum fort) Soit un ouvert borné Ω de \mathbb{R}^n . Soit une matrice A dont

— A_{ij} est dans $C^0(\bar{\Omega})$

— $A(x)$ est symétrique strictement définie positive pour tout x .

— Il existe $\lambda > 0$ minimisant toutes les valeurs propres des toutes les matrices $A(x)$ sur $\bar{\Omega}$.

Soit $b_i, c \in C^0(\bar{\Omega})$ avec $c(x) \geq 0$ sur $\bar{\Omega}$.

Soit $u \in C^0(\bar{\Omega}) \cap C^2(\Omega)$ telle que

$$\begin{cases} -\sum_{ij} A_{ij} \partial_{ij} u(x) + \sum_i b_i \partial_i u(x) \geq 0 \\ u(x) \geq 0 \quad \forall x \in \partial\Omega. \end{cases} \quad (34.130a)$$

$$\quad (34.130b)$$

Alors $u \geq 0$ sur $\bar{\Omega}$.

Démonstration. Nous appliquons le lemme 34.17(2) :

$$\min_{\bar{\Omega}} u \geq \min_{\partial\Omega} (-u_-). \quad (34.131)$$

Vu que $u(x) \geq 0$ sur $\partial\Omega$, nous avons $u_- = 0$ sur $\partial\Omega$ et donc $\min_{\bar{\Omega}} u \geq 0$. \square

34.6 Quelques exemples

34.6.1 Un changement de variables

Soit l'équation différentielle

$$\frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2} = 0 \quad (34.132)$$

sur $\Omega =]0, \infty[\times \mathbb{R}$. Nous imposons la condition aux bords

$$u(0, x) = \begin{cases} 0 & \text{si } x < 0 \\ 1 & \text{si } x > 0. \end{cases} \quad (34.133)$$

Nous cherchons les solutions sous la forme

$$u(t, x) = f\left(\frac{x}{\sqrt{t}}\right). \quad (34.134)$$

Nous avons

$$\frac{\partial u}{\partial t} = -\frac{1}{2}xt^{-3/2}f'\left(\frac{x}{\sqrt{t}}\right) \quad (34.135)$$

ainsi que

$$\frac{\partial u}{\partial x} = \frac{1}{\sqrt{t}}f'\left(\frac{x}{\sqrt{t}}\right) \quad (34.136)$$

et

$$\frac{\partial^2 u}{\partial x^2} = \frac{1}{t}f''\left(\frac{x}{\sqrt{t}}\right). \quad (34.137)$$

En remettant le tout dans l'équation de départ et en simplifiant par $1/t$ (qui est permis parce que $t > 0$) :

$$-\frac{1}{2}f'\left(\frac{x}{\sqrt{t}}\right)\frac{x}{\sqrt{t}} - f''\left(\frac{x}{\sqrt{t}}\right) = 0. \quad (34.138)$$

Nous résolvons l'équation différentielle

$$\frac{z}{2}g'(z) + g''(z) \quad (34.139)$$

pour la fonction g de la variable réelle z . Cela fait, la réponse sera $f = g \circ z$ où z serait la fonction

$$z(x, t) = x/\sqrt{t}, \quad (34.140)$$

Pour résoudre (34.139) nous commençons par résoudre pour la dérivée $h = g'$, c'est-à-dire l'équation différentielle

$$\frac{z}{2}h(z) + h'(z) = 0 \quad (34.141)$$

qui donne

$$\frac{g'(z)}{g'z} = -z/2. \quad (34.142)$$

Une intégration fournit $\ln(h(z)) = -z^2/4 + K$ et donc

$$h(z) = Ke^{-z^2/4} \quad (34.143)$$

et

$$g(z) = C + K \int_0^z e^{-s^2/4} ds. \quad (34.144)$$

Nous ne pouvons pas aller plus loin parce que nous ne sommes pas capables de calculer la primitive demandée. Nous laissons donc g sous cette forme et nous posons

$$u(x, t) = g(x/\sqrt{t}), \quad (34.145)$$

en pleine confiance du fait que cela soit une solution de l'équation aux dérivées partielles (34.132).

Nous devons fixer K et C de telle façon à respecter les conditions aux bords (34.133). D'abord écrivons aussi explicitement que possible la fonction u :

$$u(t, x) = K + C \int_0^{x/\sqrt{t}} e^{-s^2/4} ds. \quad (34.146)$$

Il faut calculer $u(0, x)$ en termes de C et K . Pour cela nous calculons, pour un x fixé :

$$\lim_{t \rightarrow 0^+} u(t, x) = K + C \lim_{t \rightarrow 0^+} \int_0^{x/\sqrt{t}} e^{-s^2/4} ds. \quad (34.147)$$

En utilisant un petit changement de variables sur l'intégrale gaussienne de l'exemple 15.262, et en remarquant que la fonction est symétrique,

$$\lim_{t \rightarrow 0^+} u(t, x) = \begin{cases} K + C\sqrt{\pi} & \text{si } x > 0 \\ K - C\sqrt{\pi} & \text{si } x < 0. \end{cases} \quad (34.148)$$

À résoudre :

$$\begin{cases} K + C\sqrt{\pi} = 1 & (34.149a) \\ K - C\sqrt{\pi} = 0. & (34.149b) \end{cases}$$

Solution : $C = 1/2\sqrt{\pi}$ et $K = 1/2$. Au final,

$$u(t, x) = \frac{1}{2} + \frac{1}{2\sqrt{\pi}} \int_0^{x/\sqrt{t}} e^{-s^2/4} ds. \quad (34.150)$$

Notons que cela donne une valeur pour $u(t, 0)$:

$$u(t, 0) = \frac{1}{2}. \quad (34.151)$$

Chapitre 35

Numérique

D'autres lectures agréables dans [464].

35.1 Introduction

À quels types de problèmes peut-on s'attendre lorsqu'on se lance dans du calcul numérique, et en particulier dans la résolution numérique d'équations (algébrique, différentielles ou aux dérivées partielles, etc) ?

Quelques réflexions en vrac sur ce sujet.

- (1) Les erreurs de représentation de nombres : troncature et propagation de décimales (*drift*),
- (2) Erreur de compensation (*cancellation*),
- (3) Conditionnement, stabilité : les réponses peuvent fortement dépendre des paramètres,
- (4) Si on utilise une méthode itérative, comment savoir à quel moment on s'arrête ? Calculer la différence $|x_k - x_{k-1}|$ mène-t-il à une erreur de cancellation ?
- (5) Lors d'une implémentation, les matrices des systèmes à résoudre sont souvent très grandes et/ou très creuses. Cela pose la question de la manière de les enregistrer.
- (6) Pour la parallélisation, il faut faire attention au fait que parfois créer un nouveau processus demande plus de ressources que le mini-calcul qu'on voulait faire. Donc il ne faut pas toujours paralléliser tout ce qui est théoriquement parallélisable.
- (7) Le fait que certaines méthodes sont non-déterministes (Monté-Carlo) mène à des problèmes pour les tests unitaires des implémentations.

35.2 Représentations numériques

Dans cette section, les séquences de chiffres écrites entre crochet sont à comprendre comme des séquences de chiffres qui représentent une quantité suivant un codage donné.

35.2.1 Entier relatif en complément à deux (binaire)

Si nous avons m bits pour coder un entier relatif, une idée serait de prendre le premier bit pour le signe (0 pour positif et 1 pour négatif) et les autres pour la valeur absolue. Deux inconvénients :

- (1) Il y a deux codages pour le zéro, donc gaspillage.
- (2) L'algorithme pour faire la somme passe mal. Par exemple pour faire $1 + (-1)$, le 1 est codé comme [001] et le -1 par [101] et la somme se ferait naïvement comme

$$\begin{array}{r} 0 \ 0 \ 1 \\ 1 \ 0 \ 1 \\ \hline 1 \ 1 \ 0 \end{array}$$

Donc le résultat est [110] qui s'interprète comme -2 . Complètement faux.

Une solution est d'utiliser le **complément à deux**, qui est la façon usuelle de représenter des entiers signés.

Les **entiers positifs** se codent normalement, en laissant à zéro le premier bit (donc si nous disposons de m bits, nous codons sur $m - 1$ bits).

Les **entiers négatifs** se codent en trois étapes.

- coder la valeur absolue
- inverser tous les bits (d'où le nom de « complément à deux »)
- soustraire 1.

Exemple 35.1

Pour coder -1 nous faisons

- Nous codons 1 : [001]
- Nous inversons tous les bits : [110]
- Nous faisons -1 : [101].

△

Avec ce système, la somme passe bien : calculer $1 + (-1)$ donne

$$\begin{array}{r} 0 \ 0 \ 1 \\ 1 \ 0 \ 1 \\ \hline 1 \ 1 \ 0 \end{array}$$

La réponse est donc [110] qu'il faut interpréter via le complément à deux.

$$110 \xrightarrow{+1} 111 \xrightarrow{\text{complément}} 000. \quad (35.1)$$

Et ce dernier [000] s'interprète comme zéro.

Définition 35.2 (Entier signé en complément à deux[465]).

La suite de bits $[a_{m-1} \dots a_0]$ s'interprète via la formule

$$-a_{m-1}2^{m-1} + \sum_{i=0}^{m-2} a_i 2^i. \quad (35.2)$$

Le premier bit donne effectivement le signe du nombre, mais l'interprétation d'un nombre n'est pas aussi simple que ce que l'on pourrait croire de prime abord.

Exemple 35.3 (Entier signé en 8 bits)

Que pouvons nous faire avec 8 bits? Le plus grand nombre est codé par [01111111] qui vaut $\sum_{k=0}^6 2^k = 2^7 - 1 = 127$. (avez-vous utilisé la somme (12.142)?)

Le plus petit nombre codable en 8 bits n'est pas [11111111] mais bien [10000000] (cela est plus clair en regardant la formule (35.2) qu'en tentant de suivre la construction du complément à deux) qui signifie $-2^7 = -128$.

Nous pouvons donc coder tous les nombres de -128 à 127 . △

Plus généralement un système qui codes des entiers signés en N bits utilisant le complément à deux peut coder de $-(2^{N-1})$ à $2^{N-1} - 1$.

35.4 (Le dépassement).

Que se passe-t-il lorsque nous commettons un dépassement? Calculons sur 3 bits la somme [011] + [001] qui revient à ajouter 1 au nombre le plus grand :

$$\begin{array}{r} 0 \ 1 \ 1 \\ 0 \ 0 \ 1 \\ \hline 1 \ 0 \ 0 \end{array}$$

qui signifie $-2^2 = -4$. Lors d'un dépassement, nous retombons automatiquement sur le plus petit.

Ce phénomène est bien connu des personnes qui programment sans faire attention dans certains langages de programmation qui ne font pas attention à votre place.

Définition 35.5 (Représentation en virgule fixe).

Soit x un réel. On définit sa **représentation en virgule fixe** par

$$x = \{[x_n x_{n-1} \dots x_0, x_{-1} \dots x_{-m}], b, s\} \quad (35.3)$$

avec $b \in \mathbb{N}, b \geq 2, s \in \{0, 1\}$ et $x_j \in \mathbb{N}, x_j < b$ suivant la formule

$$x = (-1)^s \sum_{j=-m}^n x_j \cdot b^j. \quad (35.4)$$

35.2.2 Représentation en virgule flottante

Définition 35.6 (Représentation en virgule flottante).

La **représentation en virgule flottante normalisée** en base b d'un nombre est la donnée de

- (1) Un bit s pour le signe
- (2) Un entier non signé q de e chiffres pour l'exposant
- (3) Une suite de chiffres $[a_1 \dots a_m]$ pour la mantisse.

Ces données s'interprètent via la formule

$$\text{fl}(s, q, [a_1, \dots, a_m]) = (-1)^s \sum_{j=1}^m b^j a_j \times b^{q-d} \quad (35.5)$$

où $d = b^{e-1}$ est le **décalage**.

Une idée à retenir est que l'exposant est un entier non signé parce qu'il est plus simple d'introduire un décalage dans la formule (35.5) que de compliquer l'écriture de l'exposant.

35.2.3 Simple précision, IEEE-754

En écriture binaire, la représentation en virgule flottante est un peu différente parce qu'il y a une idée supplémentaire; la simple précision que nous allons voir maintenant n'est donc pas un cas particulier de 35.6 avec $b = 2$.

Nous commençons par une description informelle de la précision simple avant de donner la définition. La représentation en **précision simple** d'un nombre se fait sur 32 bits répartis comme suit :

- (1) 1 bit pour le signe,
- (2) 8 bits pour l'exposant interprété comme nombre entier non signé
- (3) 23 bits pour la mantisse

Soit le triple

$$(s, q, [a_1, \dots, a_{23}]) \quad (35.6)$$

Dans le cas générique, l'idée est de donner 24 bits pour la mantisse, mais en gardant en tête le fait que de toutes façons, le premier bit doit être 1, sinon il suffirait de décaler, c'est-à-dire changer l'exposant. Par conséquent la mantisse ne reçoit que 23 bits; il y a un « 1 » sous-entendu en première position. Donc la mantisse $[a_1, \dots, a_{23}]$ est à lire comme le nombre

$$1, a_1 \dots a_{23} = 1 + \sum_{j=1}^{23} a_j 2^{-j}. \quad (35.7)$$

Exemple 35.7

La mantisse $[011100\dots 0]$ signifie $1,0111 = 1 + 2^{-2} + 2^{-3} + 2^{-4} = 1 + \frac{1}{4} + \frac{1}{8} + \frac{1}{16}$. \triangle Cela pour justifier la formule

$$\text{sp}(s, q, [a_1, \dots, a_{23}]) = (-1)^s \left(1 + \sum_{j=1}^{23} a_j 2^{-j}\right) 2^{q-127}. \quad (35.8)$$

Notons :

- (1) Le « 1+ » dans la parenthèse correspond au 1 implicite en première position de la mantisse.
- (2) Il y a un décalage de 127 dans l'exposant, parce que q est un entier non signé.

Notons que cette règle du 1 implicite dans la mantisse empêche d'écrire le nombre 0, et ne permet pas d'écrire des nombres franchement petits parce que le 1 implicite est en *première* position dans la mantisse.

D'où l'idée de donner une règle particulière lorsque l'exposant vaut 0. Lorsque l'exposant est $q = 0$, alors nous ne considérons pas de 1 implicite dans la mantisse, et le décalage de l'exposant est -126 au lieu de -127 . D'où la formule

$$\text{sp}(s, q = 0, [a_1 \dots a_{23}]) = (-1)^s 2^{-216} \sum_{j=1}^{23} a_j 2^{-j}. \quad (35.9)$$

En particulier, si $q = 0$ et $a = [0\dots 0]$, nous avons le nombre zéro exact (il y a deux possibilités pour le code).

Enfin, nous avons des cas particuliers lorsque l'exposant est maximum, c'est-à-dire $q = [1111\ 1111] = 2^8 - 1 = 255$. Dans ce cas, le nombre codé est soit $+\infty$ soit NaN. Nous posons $\text{sp}(s, q = 255, a = 0) = +\infty$ et $\text{sp}(s, q = 255, a \neq 0) = NaN$. Il y a en réalité plusieurs valeurs différentes de NaN, mais nous n'entrons pas dans ces détails[466].

Définition 35.8 (Représentation en simple précision (binaire)).

La représentation en **précision simple** d'un nombre se fait sur 32 bits répartis comme suit :

- (1) 1 bit pour le signe,
- (2) 8 bits pour l'exposant interprété comme nombre entier non signé
- (3) 23 bits pour la mantisse

Un nombre est représenté par un triple

$$(s, q, [a_1, \dots, a_{23}]) \quad (35.10)$$

Selon que l'exposant $q - d$ soit égal à 0, $2^8 - 1 = 255$ ou autre chose, les règles d'interprétation sont différentes. Il y a donc trois cas.

Exposant q générique[467] Si $q \neq 0$ et $q \neq 255$ alors le nombre est **normalisé**. La règle de lecture est alors

$$\text{sp}(s, q, [a_1, \dots, a_{23}]) = (-1)^s \left(1 + \sum_{j=1}^{23} a_j 2^{-j}\right) 2^{q-127}. \quad (35.11)$$

Exposant q égal à 0 Le nombre est dit **dénormalisé** et la règle de lecture est

$$\text{sp}(s, q, [a_1, \dots, a_{23}]) = (-1)^s 2^{-126} \sum_{j=1}^{23} a_j 2^{-j}. \quad (35.12)$$

Exposant q égal à 255 La règle de lecture est alors au cas pas cas ou à peu près.

- (1) $\text{sp}(s, q = 255, a = 0) = +\infty$.
- (2) $\text{sp}(s, q = 255, a \neq 0) = NaN$.

Vous pouvez jouer avec la simple précision dans [468].

Exemple 35.9(Plus petit normalisé)

Pour faire un nombre normalisé, il faut au minimum $q = 1$. En prenant $a_j = 0$ nous obtenons le plus petit nombre normalisé possible en simple précision. La formule (35.11) donne

$$\text{sp}(1, q = 1, a = 0) = 2^{1-127} = 2^{-126} \simeq 1.17549435082229 \times 10^{-38}. \quad (35.13)$$

△

Exemple 35.10(Plus grand normalisé)

L'exposant q ne peut pas être maximum, sous peine de tomber dans les règles spéciales de $+\infty$ ou NaN. Donc $q = [11111110] = 2^8 - 2 = 254$. En ce qui concerne la mantisse, il faut la prendre maximale, c'est-à-dire $a_j = 1$ pour tout j . Nous avons alors le nombre

$$\text{sp}(1, q = 254, a = [1 \dots 1]) = \left(1 + \sum_{j=1}^{23} 2^{-j}\right) 2^{254-127} = \left(1 - \frac{1}{2^{24}}\right) 2^{128} \quad (35.14a)$$

$$= 3.40282346638528859811704183484516925440 \times 10^{38} \quad (35.14b)$$

où nous avons utilisé la somme (12.142) (et Sage pour le dernier calcul).

△

Notons ceci avec Sage :

```

1
2 SageMath Version 7.0, Release Date: 2016-01-19
3 Type "notebook()" for the browser-based notebook interface.
4 Type "help()" for help.
5
6 sage: A=(1- (1/2**24) )*2**(128)
7 sage: type(A)
8 <type 'sage.rings.rational.Rational'>
```

tex/sage/sageSnip003.sage

La précision du nombre donné en (35.14b) aurait été embarrassante si le type avait été un nombre en simple précision. Précision technique : en Python, le type `int` n'a pas de limite supérieure à part la mémoire.

Exemple 35.11(Plus petit non nul dénormalisé)

Pour être dénormalisé il faut $q = 0$ (ce qui est toutefois assez logique si nous voulons un petit nombre), et pour ne pas être nul, il faut une mantisse non nulle. Donc $a = [0 \dots 01]$. La formule (35.12) donne alors

$$\text{sp}(s = 0, q = 0, a = [0 \dots 01]) = 2^{-126} 2^{-23} = 2^{-149} \simeq 1.40129846432482 \times 10^{-45}. \quad (35.15)$$

△

Exemple 35.12(Plus grand dénormalisé)

Pour être dénormalisé il faut toujours $q = 0$, mais cette fois nous prenons la plus grande mantisse possible :

$$\text{sp}(s = 0, q = 0, a = [1 \dots 1]) = 2^{-126} \sum_{j=1}^{23} 2^{-j} = 2^{-216} (1 - 2^{-23}) = 1.17549421069244 \times 10^{-38} \quad (35.16)$$

△

Notons ceci avec Sage :

```

1 sage: B=2**(-126)*(1-2**(-23))
2 sage: A=2**(-126)
3 sage: n(A-B)
4 1.40129846432482e-45

```

tex/sage/sageSnip004.sage

Vu que $2^{-23} \simeq 1.2 \times 10^{-7}$, approximer la parenthèse par 1 donne une faute sur la septième décimale, ce qui est visible en simple précision.

35.3 Problèmes pour écrire des nombres

Définition 35.13.

L'*erreur relative* commise en remplaçant un nombre réel x par une valeur approchée \hat{x} est définie par

$$\epsilon_x := \left| \frac{x - \hat{x}}{x} \right|. \quad (35.17)$$

L'erreur relative n'est pas influencée par l'ordre de grandeur de x . En effet, l'ordre de grandeur de \hat{x} est certainement la même que celle de x , dans la majorité des cas sans problèmes. Du coup si $x' = 200x$ alors $\hat{x}' \simeq 200\hat{x}$ et le 200 se simplifie.

Le nombre de chiffres significatifs correct dans l'approximation est donné par $-\log_{10}(\epsilon_x)$. La partie entière de ce nombre est le nombre de chiffres tout à fait exacts et la partie décimale donne une idée sur le fait que le chiffre suivant est plus ou moins bien.

Remarque 35.14.

Si nous voulons donner $x \in \mathbb{R}$ à un ordinateur, nous sommes soumis à deux erreurs :

- (1) D'abord, vu que nous ne pouvons pas taper sur le clavier toutes les décimales de x , nous faisons une **erreur de troncature**.
- (2) L'ordinateur devant convertir cela en base deux, il commet une seconde erreur, dite **erreur d'assignation**.

35.3.1 Troncature : la base

Supposons que nous voulions écrire le nombre (écrit ici en base 10)

$$0.4567894251 \quad (35.18)$$

de façon plus facile à lire, on peut demander de ne laisser que t chiffres significatifs. Disons $t = 3$.

Technique de troncature On garde 3 chiffres significatifs : 0.456. Facile.

Technique d'arrondi Vu que le premier qu'on supprime est un 7, le dernier qu'on garde est majoré de 1 : on écrit 0.457.

Que faire si le premier chiffre rejeté est un 5 ? En première approximation, nous pouvons prendre la règle suivante : si le premier chiffre rejeté est un 5, il faut augmenter de 1 de dernier chiffre gardé parce qu'il y a presque certainement encore un chiffre non nul derrière.

Remarque 35.15.

Les ordinateurs travaillent tous en mode d'arrondi.

Exemple 35.16

Si on doit entrer le nombre 0.38358546 dans un ordinateur qui ne garde que 3 chiffres significatifs, il faut taper 0.384 au clavier (erreur classique dans les exercices). △

35.3.2 Troncature : le drift

Soit une machine ne pouvant retenir que 3 chiffres significatifs et effectuant les arrondis vers le haut lorsque le chiffre à éliminer est un 5. Nous notons \oplus et \ominus les opérations d'addition et soustraction avec arrondis[469]. Les égalités comprenant plus de trois chiffres significatifs sont des égalités au sens de la machine. Nous écrirons donc sans états d'âme :

$$1 \oplus 0.555 = 1.555 = 1.56. \quad (35.19)$$

Considérons la suite numérique

$$\begin{cases} x_0 = 1.00 & (35.20a) \\ x_n = (x_{n-1} \ominus y) \oplus y & (35.20b) \end{cases}$$

avec $y = -0.555$.

Nous avons

$$x_1 = (1 \oplus 0.555) \ominus 0.555 = 1.56 \ominus 0.555 = 1.005 = 1.01 \quad (35.21)$$

et ensuite

$$x_2 = (1.01 \oplus 0.555) \ominus 0.555 = 1.565 \ominus 0.555 = 1.57 \ominus 0.555 = 1.015 = 1.02. \quad (35.22)$$

Et ainsi de suite. La suite est donc croissante alors que la définition nous donnerait envie d'avoir $x_n = x_0$ pour tout n .

Remarque 35.17.

En réalité, cette suite se stabilise à $x_n = 10$ pour tout n à partir de $n = 845$. En effet,

$$(10 \oplus 0.555) \ominus 0.555 = 10.555 \ominus 0.555 = 10.6 \ominus 0.555 = 10.045 = 10. \quad (35.23)$$

Le fait est qu'à ce moment, l'erreur de troncature est assez loin dans les décimales pour que le premier chiffre négligé soit un "0" au lieu d'un "5".

Notons toutefois que cette stabilité n'est pas là pour nous rassurer parce qu'elle n'en est pas moins complètement fausse.

La règle de troncature adoptée dans Sage est d'arrondir au nombre pair le plus proche lorsque le premier nombre à négliger est un 5. Donc 12.5 s'arrondit à 12 plutôt que 13.

Exemple 35.18

Soient les expressions (algébriquement égales) :

$$(1) A = x(x + 1)$$

$$(2) B = x^2 + x$$

Nous savons que

$$x = \text{fl}(x) = 10^{-30} \quad (35.24)$$

et

$$1 = \text{fl}(1) \quad (35.25)$$

parce que pour 1 et 10^{-30} , il n'y a pas d'erreurs d'assignation.

En précision simple, $10^{-30} + 1 = 1$ parce qu'en précision simple, il n'y a que 7 ou 8 chiffres significatifs¹.

Nous avons $A = 10^{-30}$, mais x^2 donne un `underflow` parce que 10^{-60} ne peut pas être représenté en précision simple. En pratique, beaucoup de logiciels en font 0. Dans ce cas, en réalité B donne effectivement 10^{-30} après avoir fait $x^2 + x = 0 + x = 10^{-30}$. \triangle

1. Erreur de « relation normale ».

35.3.3 Quelques bonnes règles

- (1) Si on a plusieurs nombres à additionner ou soustraire, il vaut mieux commencer par sommer ou soustraire ceux dont on sait qu'ils ont le même ordre de grandeur. Il n'y a donc pas tout à fait « associativité » des erreurs.
- (2) Les opérations délicates sont l'addition et la soustraction. La multiplication et la division sont sans dangers, à part l'erreur de dépassement du maximum. Dans une multiplication, on perd au pire quelques chiffres significatifs, mais certainement les derniers, pas les premiers.

35.3.4 Erreur de “cancellation”

Lorsque deux nombres sont de même ordre de grandeur, avec plusieurs nombres significatifs identiques. La cancellation est le fait que, suite à la soustraction, tous les chiffres significatifs ou presque se sont simplifiés et qu'il ne reste plus que des chiffres non significatifs.

Exemple 35.19([470])

Sur une machine ne gardant que 4 chiffres significatifs, faire

$$0.5678 \times 10^6 - 0.5677 \times 10^6 = 0.0001 \times 10^6 = 0.1000 \times 10^3. \quad (35.26)$$

Le fait est que les trois derniers zéros ne sont pas significatifs, mais maintenant la machine nous fait croire qu'ils le sont.

Une autre façon de voir ce problème est d'imaginer qu'il faille calculer la différence

$$0.5678\,289798 \times 10^6 - 0.5677\,3136907 \quad (35.27)$$

sur cette machine. Certes la machine nous autorise à avoir 4 chiffres significatifs, donc au moment d'entrer les nombres nous perdons un beau paquet de chiffres. Mais au moment de faire la différence, nous perdons (presque) tout le reste. Donc là où nous pouvions espérer avoir 4 chiffres significatifs de la différence, nous n'en avons que 1. Les trois derniers zéros de la réponse (0.1000×10^3) sont faux. \triangle

Remarque 35.20.

L'erreur de cancellation provoque des chiffres significatifs faux, mais ne provoque pas de faute dans l'ordre de grandeur des réponses². Donc si nous voulons nous assurer que a et b sont égaux « à erreur numérique près », le test

$$|a - b| < \epsilon \quad (35.28)$$

est valide, malgré l'erreur de cancellation qui ne manquera pas de se produire dans le calcul de la différence.

Exemple 35.21

Soit à résoudre l'équation $ax^2 + bx + c = 0$ avec $a, b, c \neq 0$ et $b^2 - 4ac > 0$. Solution :

$$x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}. \quad (35.29)$$

Supposons que $|4ac| \ll b^2$ avec tout de même pas tellement petit qu'on se perd dans la précision. Bref, on suppose que seules quelques dernières décimales de $b^2 - 4ac$ sont différentes de zéro.

On a :

$$\sqrt{b^2 - 4ac} = \sqrt{\tilde{b}} = |\tilde{b}| \quad (35.30a)$$

$$x_1 = \frac{-b - \sqrt{b^2 - 4ac}}{2a} \quad (35.30b)$$

$$x_2 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \quad (35.30c)$$

2. Est-ce bien vrai, cela ?

Si $b > 0$, nous avons une erreur de cancellation dans x_2 parce qu'on fait la différence entre deux nombres presque égaux. Donc x_2 mal calculé. Par contre x_1 est bien calculé.

Si par contre $b < 0$, c'est le contraire.

Avec $a = 10^{-3}$, $b = 0.8$, $c = -1.2 \times 10^{-5}$. À la main nous obtenons : $x_1 = -800$, $x_2 = 1.5 \times 10^{-5}$, et un ordinateur se tromperait ...

```

1 SageMath Version 7.0, Release Date: 2016-01-19
2 Type "notebook()" for the browser-based notebook interface.
3 Type "help()" for help.
4
5
6 sage: f(x)=10**(-3)*x**2+0.8*x-1.2*10**(-5)
7 sage: solve(f(x)==0,x)
8 [x == -1/50*sqrt(400000030) - 400, x == 1/50*sqrt(400000030) - ←
9 400]
10 sage: numerical_approx(-1/50*sqrt(400000030))
11 -400.000015000000
12 sage: numerical_approx( 1/50*sqrt(400000030) - 400 )
13 0.0000149999996779115

```

tex/sage/sageSnip001.sage

Donc Sage ne tombe pas dans le piège. △

Comment résoudre ce problème ? Ou, autre façon de poser la question : comment Sage a fait pour résoudre le problème ?

Utilisons les relations coefficients-racines :

$$x_1 + x_2 = -b/a \quad (35.31a)$$

$$x_1 x_2 = c/a \quad (35.31b)$$

La première lie les deux racines par des opérations de addition et soustractions, et donc n'est pas intéressantes. La seconde est bien. Si nous connaissons x_1 , nous calculons

$$x_2 = \frac{c}{ax_1}. \quad (35.32)$$

Quitte à redéfinir x_1 et x_2 , la solution bien calculée est :

$$x_1 = \frac{-b - \operatorname{sgn}(b)\sqrt{b^2 - 4ac}}{2a}. \quad (35.33)$$

Exemple 35.22

Nous considérons :

$$f(x) = \cos(x + \delta) - \cos(x). \quad (35.34)$$

Cela a une erreur de cancellation lorsque $|\delta| \ll |x|$. On élimine l'erreur de cancellation par

$$f(x) = -2 \sin(\delta/2) \sin\left(x + \frac{\delta}{2}\right). \quad (35.35)$$

Problèmes et choses à faire

Pourquoi la condition pour avoir l'erreur est $\delta \ll x$ et non simplement $\delta \ll 1$?

△

Exemple 35.23

Pour

$$f(x) = \sqrt{x + \delta} - \sqrt{x}. \quad (35.36)$$

On fait la coup du binôme conjugué :

$$f(x) = \frac{\delta}{\sqrt{x + \delta} + \sqrt{x}}. \quad (35.37)$$

Plus d'erreur de cancellation, vu qu'au dénominateur nous avons une somme de deux positifs. \triangle

Les erreurs de cancellation ne se résolvent pas en augmentant la précision des nombres donnés.

Exemple 35.24(Dans la vie réelle)

La préparation de l'exemple 18.45 nous a porté à calculer la différence entre $\exp(x)$ et $f_{30}(x)$ où f_{30} est censée être une bonne approximation de l'exponentielle. Des erreurs de cancellation sont donc à craindre.

Et en effet, le code suivant produit un résultat non déterministe :

```

1 f=1/152444172305856930250752000000*x^28 + ←
  1/10888869450418352160768000000*x^27 + ←
  1/15511210043330985984000000*x^25 + ←
  1/310224200866619719680000*x^24 + 1/25852016738884976640000*x←
  ^23 + 1/51090942171709440000*x^21 + 1/1216451004088320000*x^20←
  + 1/121645100408832000*x^19 + 1/355687428096000*x^17 + ←
  1/10461394944000*x^16 + 1/1307674368000*x^15 + 1/6227020800*x←
  ^13 + 1/239500800*x^12 + 1/39916800*x^11 + 1/362880*x^9 + ←
  1/20160*x^8 + 1/5040*x^7 + 1/120*x^5 + 1/12*x^4 + 1/6*x^3 + x ←
  - cos(x) + 2 -exp(x)
2 a=numerical_approx(10)
3 print(f(a))

```

tex/sage/sageSnip016.sage

Voir la question ici :

<https://ask.sagemath.org/question/37946/undeterministic-numerical-approximation/>

\triangle

35.3.5 Calcul d'une dérivée

Pour calculer la dérivée de f en a , il est loisible d'utiliser la formule

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}. \quad (35.38)$$

Le numérateur est alors sujet à une erreur d'absorption dans le calcul de $a+h$ et ensuite une erreur de cancellation dans le calcul de la différence.

En utilisant la formule

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a-h)}{2h} \quad (35.39)$$

nous pouvons espérer avoir une erreur de cancellation plus petite.

35.3.6 Erreur d'absorption

L'addition d'un nombre avec un nombre très différent peut faire perdre de l'information sur le plus petit. Par exemple avec 4 chiffres significatifs,

$$0.5678 \oplus 0.0001237 = 0.5679 \quad (35.40)$$

où nous avons perdu presque toute l'information du petit nombre.

Une situation particulièrement ennuyeuse est celle où justement c'est le petit nombre qui nous intéresse parce que le grand est censé se simplifier :

$$(0.0001327 \oplus 0.5678) \ominus 0.5678 = 0.5679 \ominus 0.5678 = 0.0001 \quad (35.41)$$

qui ne possède qu'un seul chiffre significatif correct alors que voyant le calcul, la réponse aurait pu être trouvée.

Moralité : si certaines manipulations algébrique peuvent faire apparaître des simplifications avant de passer le calcul à la machine, il est bon de les effectuer.

35.4 Conditionnement et stabilité

Définition 35.25.

Soit F une fonction à valeurs réelles définie sur $X \times D$ où X et D sont des espaces vectoriels réels normés. Le problème de la recherche des solutions de

$$F(x, d) = 0 \quad (35.42)$$

est dit **stable** autour de $d_0 \in D$ si

- (1) la solution $x = x(d)$ existe et est unique pour tout d ;
- (2) Pour tout $\eta > 0$, et pour tout d_0 , il existe un nombre $K > 0$ tel que $\|d - d_0\| < \eta$ entraîne $\|x(d) - x(d_0)\| \leq K \|d - d_0\|$.

La seconde condition est le fait que x soit Lipschitz³ sur un voisinage de d_0 .

Exemple 35.26 (Stabilité de la différence)

Prenons le problème qui consiste à calculer la différence entre deux nombres : $x = a - b$. Cela se traduit par

$$\begin{aligned} F: \mathbb{R} \times \mathbb{R}^2 &\rightarrow \mathbb{R} \\ x &\mapsto x - a + b. \end{aligned} \quad (35.43)$$

Nous avons :

$$|x(a, b) - x(a', b')| = |a - b - a' + b'| \quad (35.44a)$$

$$\leq |a - a'| + |b - b'| \quad (35.44b)$$

$$= \|(a, b) - (a', b')\|_1 \quad (35.44c)$$

où nous avons utilisé la norme $\|\cdot\|_1$ sur \mathbb{R}^2 . Par la proposition 12.5 sur les équivalences de normes, le nombre $K = \sqrt{2}$ fonctionne pour toute valeurs de η .

La problème de la différence est donc un problème stable. △

Exemple 35.27 (Stabilité de la multiplication)

Si a est fixé, le problème de calculer ab (b est la donnée) est stable. En effet ce problème est donné par la fonction $F(x, b) = x - ab$, dont la solution est $x(b) = ab$. Nous avons donc

$$|x(b) - x(b')| = |ab - ab'| = |a||b - b'|. \quad (35.45)$$

La constante de Lipschitz de ce problème est donc $|a|$. △

3. Définition 13.254.

Définition 35.28.

Le nombre

$$K_{abs}(d_0, \eta) := \sup_{d \text{ tel que } |d_0 - d| < \eta} \frac{\|x(d) - x(d_0)\|_X}{\|d - d_0\|_D} \quad (35.46)$$

est appelé le **conditionnement absolu** du problème autour de d_0 .

Soit $F(x, d) = 0$ un problème stable de conditionnement absolu $K_{abs}(d, \eta)$. Le conditionnement relatif est défini par

$$K_{rel}(d, \eta) := K_{abs}(d, \eta) \frac{\|d\|_D}{\|x(d)\|_X}. \quad (35.47)$$

Le problème est dit **bien conditionné** près de d si $K_{rel}(d, \eta)$ est petit.

Exemple 35.29 (Mauvais conditionnement de la différence)

Reprenons le problème de la différence, mais en fixant a . Nous avons donc $x(b) = a - b$ et le conditionnement absolu est

$$\sup \frac{|x(b) - x(b_0)|}{|b - b_0|} = 1 \quad (35.48)$$

Le conditionnement relatif est :

$$K_{rel}(b_0, \eta) = \frac{|b|}{|a - b|}. \quad (35.49)$$

Et donc le problème est mal conditionné autour de a .

Autrement dit, si a' est un nombre proche de a , calculer la différence $a - a'$ est un problème mal conditionné. \triangle

Exemple 35.30 (Bon conditionnement de la multiplication)

Pour le problème $F(x, b) = x - ab$ nous avons

$$K_{abs} = \sup_{b'} \frac{|ab - ab'|}{|b - b'|} = |a|. \quad (35.50)$$

Et aussi

$$K_{rel} = a \frac{|b|}{|ab|} = 1. \quad (35.51)$$

Le conditionnement relatif du problème de la multiplication est donc toujours 1. Il est donc un toujours un problème bien conditionné. \triangle

Ne pas confondre :

Le conditionnement provient du problème lui-même.

La stabilité provient de l'algorithme de résolution.

Exemple 35.31 (Un problème mal conditionné)

Le système

$$\begin{cases} 2.1x + 3.5y = 8 & (35.52a) \\ 4.19x + 7.0y = 15 & (35.52b) \end{cases}$$

Solution : $x = 100$, $y = -57.714285\dots$ (périodique)

Perturbons : nous remplaçons 4.19 par 4.192. L'erreur relative est : 4.77×10^{-4} .

Solution : $\bar{x} = 125$, $\bar{y} = -72.714285\dots$, avec donc erreur relative de 0.26. Autrement dit : l'erreur relative sur la solution est grande même avec une petite erreur relative sur la donnée.

C'est un problème mal conditionné.

Le fait est que c'est une intersection de deux droites presque parallèles. Donc effectivement une petite perturbation d'une des deux droites donne une grande perturbation du point d'intersection.

Le fait est qu'un ordinateur effectue *toujours* une perturbation, au moins de l'ordre 10^{-16} pour ne fut-ce que représenter les nombres. C'est-à-dire une perturbation sur les six nombres définissant le système. Il n'y a donc pas d'espoir d'obtenir un algorithme donnant une bonne réponse. \triangle

Un résultat pratique pour étudier le conditionnement d'un problème est le suivant.

Corollaire 35.32.

Soit $x = x(d)$ un problème stable. Supposons \mathbb{D} de dimension finie, supposons que U est ouvert dans \mathbb{D} . Supposons encore $x: U \rightarrow \mathbb{R}$ différentiable en d_0 . Alors quand η est petit, on a

$$K_{abs}^\eta(d_0) \sim \|\nabla x(d_0)\|. \quad (35.53)$$

Lemme 35.33.

Tout problème de la forme $x = x(d)$ avec $d \in \mathbb{R}$ et $x \in C^1(\mathbb{R})$ est stable.

Démonstration. Il faut démontrer qu'une fonction C^1 sur \mathbb{R} vérifie automatiquement la condition (2) de la définition de la stabilité. Pour cela, remarquons qu'une fonction C^1 possède une dérivée continue, et donc bornée sur tout compact⁴

Prenons $\eta > 0$ et $d_0 \in \mathbb{R}$ et puis un d tel que $|d - d_0| < \eta$. Par le théorème des bornes atteintes, la fonction x' est bornée sur l'intervalle $[d_0 - \eta, d_0 + \eta]$. Appelons K un majorant de x' sur cet intervalle. La fonction

$$f(d) = x(d_0) + K|d - d_0| \quad (35.54)$$

majore $x(d)$, et donc on a

$$|x(d) - x(d_0)| \leq K|d - d_0|. \quad (35.55)$$

Attention : vérifier si ce raisonnement est correct avec $d_0 > d$, et adapter au besoin. \square

Exemple 35.34

Un exemple de problème stable de la forme $x = x(d)$ avec $d \in \mathbb{R}$ et $x \in C^0(\mathbb{R}) \setminus C^1(\mathbb{R})$.

La fonction

$$x(d) = \begin{cases} 0 & \text{si } x \geq 0 \\ x & \text{si } x < 0 \end{cases} \quad (35.56)$$

est continue, mais pas C^1 (non dérivable en $x = 0$). La dérivée est partout bornée par 1, et donc le problème est stable.

Un autre exemple très classique serait de prendre $x(d) = |d|$. Dans ce cas, on peut prendre n'importe que η et $K = 1$. Le calcul est que

$$|x(d) - x(d_0)| < K|d - d_0| \quad (35.57a)$$

$$||d| - |d_0|| < |d - d_0|. \quad (35.57b)$$

Cette dernière inéquation est correcte, comme on peut le voir en mettant au carré les deux membres.

\triangle

Exemple 35.35

Un exemple de problème instable de la forme $x = x(d)$ avec $d \in \mathbb{R}$ et $x \in C^0(\mathbb{R})$.

Un exemple assez classique de fonction dont la dérivée n'est pas bornée sans pour autant que la fonction aie un comportement immoral⁵ est $x \mapsto \sqrt{x}$. Afin d'avoir une fonction définie sur \mathbb{R} tout entier, nous regardons la fonction

$$x(d) = \sqrt{|d|}. \quad (35.58)$$

4. Un compact est un ensemble fermé et borné, typiquement un intervalle du type $[a, b]$.

5. Penser à $x \mapsto x \sin(1/x)$.

Si nous considérons maintenant $d_0 = 0$ et n'importe quel η , nous avons

$$\frac{|x(d) - x(d_0)|}{|d - d_0|} = \frac{\sqrt{d}}{d} = \frac{1}{\sqrt{d}}. \quad (35.59)$$

Il n'est pas possible de trouver un K qui majore ce rapport. Le problème est donc mal conditionné. Attention : dans ce calcul nous avons supposé $d > 0$. Pensez à adapter au cas $d < 0$. \triangle

Exemple 35.36 (Problème bien conditionné avec algorithme instable)

Soit à calculer

$$I_n = \frac{1}{e} \int_0^1 x^n e^x dx \quad (35.60)$$

avec $n \geq 0$. Par partie, nous obtenons :

$$I_n = 1 - nI_{n-1}. \quad (35.61)$$

D'autre part, $I_0 = \frac{e-1}{e}$, $I_1 = \frac{1}{e}$. Puis par récurrence, c'est tout en main.

Du côté de l'ordinateur, nous lui donnons forcément une approximation de I_1 , parce que nous lui donnons une approximation de e . Soit l'erreur ϵ_1 sur I_1 .

Sans démonstration :

Lemme 35.37.

Nous avons $\lim_{n \rightarrow \infty} I_n = 0$.

Mais numériquement, il n'est pas possible de rester longtemps sous ϵ_1 parce que nous n'espérons pas avoir une erreur plus petite que ça. Donc à partir du moment où $I_n < \epsilon_1$, les valeurs sont toutes complètement fausses. Cela est le mieux que l'on puisse espérer. Mais la réalité est pire.

En réalité, en lançant le calcul sur un ordinateur, les valeurs sont même croissantes avec n à partir d'un certain moment.

On peut étudier l'erreur et montrer que l'erreur est donnée par :

$$\epsilon_n = (-1)^{n-1} n! \epsilon_1. \quad (35.62)$$

Mais comme la factorielle est tellement forte que c'est sans espoir d'aller loin en essayant très fort de donner une petite erreur sur ϵ_1 .

\triangle

Il existe heureusement un algorithme stable pour cette intégrale. La formule est :

$$I_{n-1} = \frac{1}{n}(1 - I_n). \quad (35.63)$$

Si nous savons un I_N avec N grand, cette formule donne les I_i avec $i = N, N-1, \dots, 2$. Posons donc $I_N = a \in \mathbb{R}$ n'importe comment. Donc ϵ_N est grand. Mais il se trouve que l'erreur sur ϵ_1 est donnée par

$$\epsilon_1 = \frac{(-1)^{N-1}}{N!} \epsilon_N. \quad (35.64)$$

Donc même en prenant vraiment n'importe quoi pour I_N , nous obtenons de bonnes approximations pour I_i avec les petits i . Même avec $I_{20} = 1000$ (qui est complètement faux), nous trouvons énormément de chiffres significatifs corrects pour I_1 .

35.4.1 Comment choisir et penser le K ?

La formule (35.46) contient une formule qui ressemble étrangement à la dérivée. La stabilité d'un problème est très liée à la dérivée de F . La stabilité et la dérivée ne sont pas les mêmes choses,

mais il n'est pas mauvais de penser au K de la stabilité comme la dérivée. Ou plus précisément : le supremum de la dérivée.

Un fil conducteur du lemme 35.33 et des exemples 35.34, 35.35 est que l'on a un K qui fonctionne lorsque la dérivée est bornée sur l'intervalle $]d_0 - \eta, d_0 + \eta[$. Dans le cas où ce supremum existe, le prendre en guise de K fonctionne souvent.

Il faut cependant parfois faire acte d'imagination. La fonction $x \mapsto |x|$ n'est pas dérivable en 0. Il n'empêche que $K = 1$ fait fonctionner la définition de la stabilité. Remarquez que $K = 1$ est le supremum de la dérivée là où elle existe.

À partir du moment où c'est clair que le K est le supremum de la dérivée, on comprend pourquoi c'est le gradient qui arrive dans le corollaire 35.32. En effet, le gradient indique la direction de plus grande pente. C'est donc bien dans cette direction qu'il faut chercher la « plus grande dérivée ».

Proposition 35.38.

Pour le problème stable $x = x(d)$ avec $x \in C^1(\mathbb{R}^n, \mathbb{R})$, on a

$$K_{abs}(d) \sim \|dx_d\| \quad (35.65)$$

où dx_d désigne la différentielle de x en d et la norme est la norme opérateur.

35.5 Un peu de points fixes

35.5.1 Choix de la fonction à point fixe

Pour l'équation $f(x) = 0$, il existe une infinité de fonctions g pour lesquelles l'équation est équivalente à $x = g(x)$.

Exemple : $f(x) = x^2 - 2 - \ln(x)$, nous pouvons faire

(1) $x = x^2 - 2 - \ln(x) + x$

(2) Poser $x^2 = 2 + \ln(x)$ et donc

$$x = -\sqrt{2 + \ln(x)} \quad (35.66a)$$

$$x = \sqrt{2 + \ln(x)}. \quad (35.66b)$$

(3) Ou encore

$$x = \frac{2 + \ln(x)}{x} \quad (35.67)$$

où nous savons déjà que $x \neq 0$ parce que $x = 0$ n'est pas dans le domaine de f .

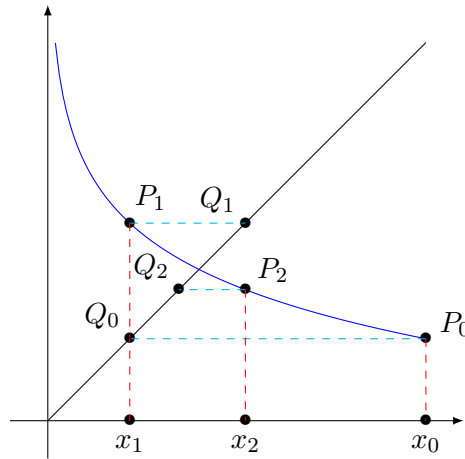
(4) Ou par l'exponentielle :

$$x = e^{x^2-2}. \quad (35.68)$$

Dans tous ces cas nous pouvons construire une suite (x_n) en posant un nombre arbitraire pour x_0 et ensuite la récurrence

$$x_{n+1} = g(x_n). \quad (35.69)$$

Graphiquement, la solution de l'équation est l'intersection entre les courbes $y = x$ et $y = g(x)$. Un petit dessin pour montrer la convergence :



Attention : cette méthode ne converge pas toujours. Parfois elle converge de façon monotone, et parfois pas. Le choix de la fonction g qui fait $x = g(x)$ peut énormément changer la vitesse de convergence.

Théorème 35.39 (Condition suffisante pour existence d'un point fixe).

Une fonction continue $f: [a, b] \rightarrow [a, b]$ admet au moins un point fixe dans $[a, b]$.

Théorème 35.40 (Condition suffisante pour l'unicité).

Soit f continue sur $[a, b]$ avec $g(x) \in [a, b]$ pour tout $x \in [a, b]$. Supposons qu'il existe $0 < k < 1$ tel que pour tout $x \in [a, b]$ nous ayons $|g'(x)| \leq k$ alors

- (1) La fonction g possède un unique point fixe dans $[a, b]$.
- (2) Pour tout $x_0 \in [a, b]$, tous les termes de la suite $x_{n+1} = g(x_n)$ sont dans $[a, b]$.
- (3) Ladite suite (x_n) converge vers le point fixe.

Théorème 35.41.

Soit f continue sur $[a, b]$ avec $g(x) \in [a, b]$ pour tout $x \in [a, b]$. Supposons

- (1) qu'il existe $0 < k < 1$ tel que pour tout $x \in [a, b]$ nous ayons $|g'(x)| \leq k$ et
- (2) g est p fois dérivable sur $[a, b]$.
- (3) $g'(\alpha) = g''(\alpha) = \dots = g^{(p-1)}(\alpha)$ et $g^{(p)}(\alpha) \neq 0$ où α est l'unique point fixe.

Alors la suite (x_n) converge avec un ordre p .

Exemple 35.42

Nous reprenons

$$f(x) = x^2 - 2 - \ln(x). \quad (35.70)$$

Et nous voulons résoudre $f(x) = 0$. Graphiquement c'est l'intersection entre $y = x^2 - 2$ et $y = \ln(x)$. Il est vite tracé de savoir qu'il y a deux solutions : $\alpha_1 \in [0, 1]$ et $\alpha_2 \in [\sqrt{2}, 2]$.

Déjà un petit problème : l'intervalle $[0, 1]$ ne va pas parce que f n'y est pas continue. Un petit raffinement d'analyse nous fournit $\alpha_1 \in [e^{-2}, 1]$.

Nous avons au moins les fonctions de points fixes suivantes :

$$g_1(x) = \sqrt{2 + \ln(x)} \quad (35.71a)$$

$$g_2(x) = e^{x^2 - 2}. \quad (35.71b)$$

Pour la première, il y avait un \pm qui a été négligé parce que nous savons que les deux solutions cherchées sont positives. Travaillons avec la première. D'abord

$$g_1'(x) = \frac{1}{2x\sqrt{2 + \ln(x)}}. \quad (35.72)$$

Nous avons $\lim_{x \rightarrow e^{-2}} g'_2(x) = +\infty$. Il ne sera donc pas possible de trouver $0 < k < 1$ tel que $|g'(x)| \leq k$. Tentons quand même la méthode :

$$x_0 = 0.5 \tag{35.73}$$

Il se fait que cela est plus proche de α_1 que de α_2 . Mais en réalité la suite converge vers α_2 .

Passons à la seconde méthode.

$$g'_2(x) = 2xe^{x^2-2}. \tag{35.74}$$

Sur l'intervalle $[e^{-2}, 1]$, g'_2 est croissante et prend toutes ses valeurs dans $[e^{-2}, 1]$. Nous pouvons prouver que

$$|g'_2(x)| \leq 2e^{-1} < 1. \tag{35.75}$$

Donc poser $k = 2e^{-1}$ fait fonctionner la proposition. Donc quel que soit le x_0 pris dans cet intervalle, nous aurons une suite convergente vers un point fixe à l'intérieur de l'intervalle. C'est-à-dire convergente vers α_1 .

Cela est un exemple de problème pour lequel changer de fonction g change réellement la vie.

△

35.5.2 Convergence quadratique

Définition 35.43.

Une suite (x_n) a une convergence **quadratique** vers α si elle converge vers α et s'il existe un C tel que pour tout n nous ayons

$$\|x_{n+1} - \alpha\| \leq C\|x_n - \alpha\|^2. \tag{35.76}$$

Il est bien entendu possible de parler de convergence quadratique si la relation (35.76) a lieu seulement à partir d'un certain indice.

Le lemme suivant donne l'importance du choix de point de départ lorsqu'on utilise une méthode itérative dont la convergence est quadratique.

Lemme 35.44.

Soit une suite $x_n \rightarrow \alpha$ de convergence quadratique. Si $\|x_0 - \alpha\| \leq r$ alors

$$\|x_n - \alpha\| \leq \frac{1}{C}(Cr)^{2^n} \tag{35.77}$$

Démonstration. Nous pourrions directement prouver la formule (35.77) par récurrence, mais nous allons la reconstruire un peu. Nous cherchons

$$\|x_n - \alpha\| \leq C^{k(n)}r^{2^n}. \tag{35.78}$$

Nous avons les inégalités

$$\|x_{n+1} - \alpha\| \leq C\|x_n - \alpha\|^2 \tag{35.79a}$$

$$\leq CC^{2k(n)}r^{2^{n+1}} \tag{35.79b}$$

$$= C^{2k(n)+1}r^{2^{n+1}} \tag{35.79c}$$

d'où nous voyons que la fonction k doit vérifier

$$\begin{cases} k(0) = 0 & (35.80a) \\ k(n+1) = 2k(n) + 1 & (35.80b) \end{cases}$$

La première équation est l'hypothèse $\|x_0 - \alpha\| \leq r$ comparée à la formule (35.78). Il est vite vérifié que $k(n) = 2^n - 1$. D'où le résultat. □

Si le point de départ est choisi de façon à avoir $Cr < 1$ alors nous avons là un très bon majorant parce qu'il s'agit d'un majorant convergeant très rapidement vers zéro. Si au contraire $Cr > 1$ alors ce majorant ne sert à rien.

35.45.

Le fait d'avoir une convergence quadratique signifie que le nombre décimales correctes double (environ) à chaque itération, dans n'importe quelle base. En effet supposons que x_n ait k décimales correctes ; cela signifie que $|x_n - \alpha| \sim 10^{-k}$. Donc

$$|x_{n+1} - \alpha| \lesssim M10^{-2k}. \quad (35.81)$$

Cela est le double de décimales correctes de $|x_n - \alpha|$, moins l'ordre de grandeur de M .

Pour la méthode de bisection, le nombre de décimales augmente de 1 à chaque itération, mais seulement en base 2. En base 10, de façon générique⁶ il faut entre 3 et 4 itérations pour avoir une décimale de plus.

35.46 (Condition d'arrêt[471]).

D'autre part, lorsqu'une méthode a une convergence quadratique, nous avons un test d'arrêt. Pour ce voir, nous avons la limite

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - \alpha|}{|x_n - \alpha|} \leq \lim_{n \rightarrow \infty} \frac{C|x_n - \alpha|^2}{|x_n - \alpha|} = 0. \quad (35.82)$$

Cette limite est alors également valable sans les valeurs absolues et si nous soustrayons $x_n - \alpha$ au numérateur, la limite devient -1 :

$$-1 = \lim_{n \rightarrow \infty} \frac{x_{n+1} - x_n}{x_n - \alpha}. \quad (35.83)$$

Ou encore

$$\lim_{n \rightarrow \infty} \frac{x_n - x_{n+1}}{x_n - \alpha} = 1. \quad (35.84)$$

Cela a pour conséquence que si n est grand,

- (1) x_{n+1} a le même ordre de grandeur que $x_n - \alpha$.
- (2) $x_n - x_{n+1}$ et $x_n - \alpha$ ont le même signe.

Donc si nous voulons une approximation de α avec une erreur ϵ , il suffit d'arrêter le calcul lorsque $|x_{n+1} - x_n| \leq \epsilon$. Et ce faisant nous savons de plus si l'approximation est par excès ou par défaut.

35.5.3 Convergence

Proposition 35.47 (Convergence d'une méthode de point fixe[471]).

Soit $g: \mathbb{R} \rightarrow \mathbb{R}$ de classe C^1 et α un point fixe attractif⁷ de g . Soit k tel que $|g'(\alpha)| < k < 1$ et δ tel que $\|g'\|_{B(\alpha, \delta)} < k$.

Alors

- (1) La fonction g est k -contractante⁸ sur $B(\alpha, \delta)$.
- (2) Nous avons $g(B(\alpha, \delta)) \subset B(\alpha, \delta)$.
- (3) Pour tout $x_0 \in B(\alpha, \delta)$ la suite $x_{n+1} = g(x_n)$ converge vers α et

$$|x_n - \alpha| \leq |x_0 - \alpha|k^n. \quad (35.85)$$

Si de plus $g'(\alpha) = 0$ et g est de classe C^2 alors nous avons convergence quadratique (définition 35.43).

Démonstration. Vu que α est un point fixe attractif de g nous pouvons considérer un k tel que $|g'(\alpha)| < k < 1$. Et comme g est de classe C^1 , la fonction g' est continue et donc bornée sur toute boule du type $\overline{B(\alpha, \delta)}$. Soit δ le plus grand nombre tel que $\|g'\|_{\overline{B(\alpha, \delta)}} \leq k$. Nous notons $I = \overline{B(\alpha, \delta)}$ pour cette valeur de δ .

6. C'est-à-dire sauf coup de malchance ou coup de chance.

7. Définition 18.30.

8. Définition 18.33

Pour tout $x \in I$ nous avons, en utilisant le théorème des accroissements finis 13.129(2) :

$$|g(x) - \alpha| = |g(x) - g(\alpha)| \quad (35.86a)$$

$$\leq \sup_{t \in I} |g'(t)| |x - \alpha| \quad (35.86b)$$

$$\leq k|x - \alpha| \quad (35.86c)$$

$$< \delta \quad (35.86d)$$

parce que $k < 1$ et $|x - \alpha| \leq \delta$. Par conséquent $g(x) \in B(\alpha, \delta)$. Cela prouve le point (2). Pour le point (1), soient $x, y \in B(\alpha, \delta)$ et

$$|g(x) - g(y)| \leq \sup_{a \in I} |g'(a)| |x - y| \leq k|x - y|. \quad (35.87)$$

Pout le point (3) nous avons $|g(x_n) - \alpha| \leq k|x_n - \alpha|$, c'est-à-dire

$$|x_{n+1} - \alpha| \leq k|x_n - \alpha|. \quad (35.88)$$

Le résultat annoncé s'obtient par récurrence sur n .

En ce qui concerne la convergence quadratique, c'est du Taylor (proposition 13.369). Développons $g(x_n)$ autour de $g(\alpha)$:

$$g(x_n) = g(\alpha) + g'(\alpha)(x_n - \alpha) + \frac{1}{2}(x_n - \alpha)^2 \epsilon(x_n - \alpha) \quad (35.89)$$

avec $\lim_{t \rightarrow 0} \epsilon(t) = 0$. En posant $C = \frac{1}{2} \sup_{t < \delta} |\epsilon(t)|$ nous avons $|g(x_n) - g(\alpha)| \leq C|x_n - \alpha|^2$, c'est-à-dire

$$|x_{n+1} - \alpha| \leq C|x_n - \alpha|^2. \quad (35.90)$$

□

Ce corollaire est une paraphrase de la proposition 35.47. Il en retient seulement les points intéressants en pratique.

Corollaire 35.48.

Soit α une solution de l'équation $x = g(x)$, avec g continue sur un voisinage de α et dérivable dans l'intérieur. Nous supposons que

$$|g'(\alpha)| < 1. \quad (35.91)$$

Alors il existe un rayon δ tel que si $x_0 \in B(\alpha, \delta)$, la suite (x_n) converge vers α .

Certes cette proposition demande moins d'hypothèses, mais en réalité, il ne donne pas de vrais moyens de choisir un point de départ x_0 . Avec les deux théorèmes précédents, nous pouvons prendre x_0 n'importe où dans $[a, b]$. Le fait est que pour choisir x_0 nous pouvons tracer et donner à la main un x_0 proche de ce qui semble être α . Si ça ne converge pas, il faut donner un x_0 plus proche. La proposition nous assure que si nous jouons bien à choisir x_0 très proche, la suite finira par converger.

Notons que le corollaire 35.48 a encore l'inconvénient de demander de calculer $g'(\alpha)$ alors que α est inconnu. La résolution de l'inéquation $|g'(x)| < 1$ nous donne un certain nombre d'intervalles dans \mathbb{R} .

Soient I_n les intervalles solutions de l'inéquation. Si $\alpha \in I_n$ alors la méthode converge. Sinon, c'est pas garanti. En tout cas nous ne devons pas savoir réellement α pour appliquer le théorème. Il suffit de savoir que α est dans un des I_n .

35.6 Méthode de Newton

L'objectif de la méthode de Newton est d'évaluer une racine α de l'équation $f(x) = 0$ lorsque nous avons déjà une approximation x_0 de la racine α .

C'est la méthode de Newton qui est à l'origine de la suite de la proposition 1.91 donnant une suite dans \mathbb{Q} qui converge vers \sqrt{A} .

Définition 35.49.

Le nombre α est une **racine simple** de l'équation $f(x) = 0$ si $f(\alpha) = 0$ et $f'(\alpha) \neq 0$. Le nombre α est une **racine multiple** d'ordre r de $f(x) = 0$ si

$$f(\alpha) = f'(\alpha) = \dots = f^{(r-1)}(\alpha) = 0 \quad (35.92)$$

et $f^{(r)}(\alpha) \neq 0$.

Exemple 35.50

La fonction $x \mapsto x^3$ en $x = 0$ est un racine d'ordre 3. △

35.6.1 « Justification » par la formule par Taylor

Soit une fonction f continue et dérivable sur $[a, b]$. Soit α une racine de f et x_n une de ses approximations. Nous notons l'erreur θ et nous avons $\alpha = x_n + \theta$. Du coup nous avons $f(x_n + \theta) = f(\alpha) = 0$.

Écrivons la série de Taylor du théorème 13.359 autour de x_n : il existe une fonction $\epsilon : \mathbb{R} \rightarrow \mathbb{R}$ telle que $\lim_{t \rightarrow 0} \epsilon(t) = 0$ telle que

$$f(\alpha) = f(x_n + \theta) = f(x_n) + \theta f'(x_n) + \frac{\theta^2}{2} \epsilon(\theta). \quad (35.93)$$

Nous isolons le θ du terme d'ordre 1 en nous souvenant que le membre de gauche est nul :

$$\theta = -\frac{f(x_n) - \theta^2 \epsilon(\theta)}{f'(x_n)} \quad (35.94)$$

Vu que $\alpha = x_n + \theta$, nous pouvons écrire

$$\alpha = x_n - \frac{f(x_n) + \theta^2 \epsilon(\theta)}{f'(x_n)}. \quad (35.95)$$

Il est donc raisonnable de poser

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \quad (35.96)$$

en espérant que cela soit une meilleur approximation de α que x_n .

En tout cas l'erreur sur x_{n+1} est

$$\alpha - x_{n+1} = x_n + \theta - x_n + \frac{f(x_n) + \theta^2 \epsilon(\theta)}{f'(x_n)} = \theta + \frac{f(x_n) + \theta^2 \epsilon(\theta)}{f'(x_n)}, \quad (35.97)$$

qui ne doit pas être fondamentalement plus grand que θ dès que θ est petit, surtout que si x_n est une approximation de α , nous pouvons espérer que $f(x_n)$ soit également petit. Là où les choses peuvent déraiper en grand, c'est si $f'(x_n)$ est petit.

Cette méthode de Newton ne converge pas toujours. Le pire est lorsque par malheur il y a une bosse pas loin de la racine. Alors il y a un risque de tomber sur $f'(x_{n+1}) = 0$ ou en tout cas très proche de zéro. Dans ce cas le point x_{n+2} est envoyé très loin.

35.6.2 « Justification » par points fixes

Nous savons que pour résoudre $f(x) = 0$ par une méthode de point fixe, il y a de nombreux choix possibles de fonctions g telles que $g(x) = x$ donne la même solution que $f(x) = 0$. Soit α une solution de $f(x) = 0$ et cherchons une fonction g de la forme

$$g(x) = x - kf(x). \quad (35.98)$$

Nous savons par la proposition 35.47 que la fonction g donne une convergence quadratique lorsque $g'(\alpha) = 0$. Pour la forme (35.98) nous avons $g'(\alpha) = 1 - kf'(\alpha)$, ce qui nous donne l'idée de poser $k = \frac{1}{f'(\alpha)}$.

Le fait est que $f'(\alpha)$ n'est pas connu, mais nous pouvons l'approximer par $f'(x)$ lorsque x est proche de α . D'où l'idée de considérer la fonction

$$g(x) = x - \frac{f(x)}{f'(x)}, \quad (35.99)$$

et donc la suite $x_{n+1} = g(x_n)$ c'est-à-dire

$$x_{n+1} = x - \frac{f(x_n)}{f'(x_n)}. \quad (35.100)$$

Dès que x_n est proche de α , sous l'hypothèse (raisonnable par continuité) que $f'(x_n)$ soit proche de $f'(\alpha)$, la méthode devrait donner une convergence quadratique.

Remarque 35.51.

Cette justification par points fixes n'est pas vraiment différente de celle par Taylor parce que Taylor est utilisé dans la preuve de la proposition 35.47.

Définition 35.52 (Méthode de Newton).

La *méthode de Newton* pour la fonction f est la suite définie par récurrence

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}. \quad (35.101)$$

Cette définition ne précise pas la valeur de x_0 , ni de condition d'arrêt.

35.6.3 Convergence de la méthode de Newton

Théorème 35.53 (Convergence quadratique de la méthode de Newton[471]).

Soit f une fonction continue vérifiant $f(\alpha) = 0$ et $f'(\alpha) \neq 0$. Nous considérons la fonction

$$g(x) = x - \frac{f(x)}{f'(x)} \quad (35.102)$$

que nous supposons être de classe C^2 .

Si C est une majoration de $\|g''\|$ sur un intervalle contenant α , alors en posant $\delta = 1/C$ nous avons

- (1) La boule $B(\alpha, \delta)$ est préservée par $g : g(B(\alpha, \delta)) \subset B(\alpha, \delta)$.
- (2) Pour tout $x_0 \in B(\alpha, \delta)$ nous avons convergence quadratique vers α de la suite définie par $x_{n+1} = g(x_n)$.
- (3) Nous avons l'estimation

$$|x_n - \alpha| \leq \frac{1}{C} (C|x_0 - \alpha|)^{2^n} \quad (35.103)$$

où C est la constante de la définition de convergence quadratique.

Démonstration. Nous commençons par calculer la dérivée de g :

$$g'(x) = -\frac{f(x)f''(x)}{f'(x)^2}, \quad (35.104)$$

d'où nous déduisons que $g'(\alpha) = 0$. Ensuite nous utilisons abondamment la formule des accroissements finis (théorème 13.252) en commençant par

$$|g(t) - g(\alpha)| \leq \|g'\|_{[t,\alpha]} |t - \alpha| \quad (35.105)$$

où par $\|f\|_A$ nous entendons la norme uniforme de f sur A , c'est-à-dire $\|f\|_A = \sup_{x \in A} \|f(x)\|$.
 Note : nous écrivons $[t, \alpha]$, mais ça pourrait être $[\alpha, t]$.

Si $x \in [t, \alpha]$ alors

$$|g'(x)| = |g'(x) - g'(\alpha)| \quad (35.106a)$$

$$\leq \|g''\|_{[x, \alpha]} |x - \alpha| \quad (35.106b)$$

$$\leq \|g''\|_{[x, \alpha]} |t - \alpha| \quad (35.106c)$$

$$\leq \|g''\|_{[t, \alpha]} |t - \alpha|. \quad (35.106d)$$

En particulier, $\|g'\|_{[t, \alpha]} \leq \|g''\|_{[t, \alpha]} |t - \alpha|$, et nous pouvons continuer les majorations (35.105) :

$$|g(t) - g(\alpha)| \leq \|g''\|_{[t, \alpha]} |t - \alpha|^2. \quad (35.107)$$

La fonction g étant de classe C^2 , la dérivée seconde g'' est bornée (nous supposons déjà travailler sur un compact contenant α). Soit C une borne. Nous sommes en mesure de prouver le point (1) avec $\delta = 1/C$. En effet si $t \in B(\alpha, 1/C)$ alors

$$|g(t) - \alpha| = |g(t) - g(\alpha)| \leq C |t - \alpha|^2 \leq C \frac{1}{C^2} = \frac{1}{C}, \quad (35.108)$$

ce qui prouve que $g(t) \in B(\alpha, 1/C)$.

Le point (2) se prouve de la même manière : si $x_n \in B(\alpha, 1/C)$ alors

$$|x_{n+1} - \alpha| = |g(x_n) - g(\alpha)| \leq C |x_n - \alpha|^2, \quad (35.109)$$

ce qui est bien la convergence quadratique.

La majoration du point (3) s'obtient par récurrence sur n . Pour $n = 0$, la relation (35.103) devient $|x_0 - \alpha| \leq |x_0 - \alpha|$ qui est vraie. Ensuite par la convergence quadratique et la récurrence,

$$|x_{n+1} - \alpha| \leq C |x_n - \alpha|^2 \leq C \left[\frac{1}{C} (C |x_0 - \alpha|)^{2^n} \right]^2 = \frac{1}{C} [M |x_0 - \alpha|]^{2^{n+1}}. \quad (35.110)$$

□

35.54.

Dans le cas pratiques, nous commençons souvent par résoudre l'équation $f(x) = 0$ par dichotomie. Au moment où nous sommes assez proche de la solution nous commençons Newton.

La raison est que la dichotomie fonctionne toujours : nous allons toujours nous approcher de la solution. Si par contre le point de départ est mal choisi, la méthode de Newton peut envoyer n'importe où, y compris très loin de la solution.

La proposition suivante nous indique que dans le cas d'une fonction convexe, le choix de point de départ de la méthode de Newton n'est pas tellement crucial parce que il sont tous bons. De plus la convergence se faisant de façon décroissante (si on part de la droite), nous savons que le résultat sera une approximation par excès de α .

Proposition 35.55 (Newton dans le cas convexe).

Soit f de classe C^2 et une racine α telle que $f'(\alpha) > 0$. Soit $b > \alpha$ tel que f soit convexe sur $[\alpha, b]$.

Alors pour tout $x_0 \in [\alpha, b]$ la suite de la méthode de Newton est

- (1) décroissante
- (2) reste dans $[\alpha, b]$
- (3) converge vers α .

Démonstration. Nous savons par la proposition 18.81(2) que la fonction f' est croissante, et par hypothèse $f'(\alpha) > 0$, donc sur $[\alpha, b]$ nous avons $f' > 0$. Par conséquent, nous avons aussi $f > 0$ sur $[\alpha, b]$.

Le graphe de f est au dessus de la tangente de f en $x = x_n$ (proposition 18.87). Si nous nommons t_x la fonction qui donne la tangente en x nous avons $t_{x_n}(\alpha) < 0$ parce que $f(\alpha) = 0$. Par conséquent

$$t_{x_n}(x) = 0 \quad (35.111)$$

pour $\alpha < x < x_n$. Cela prouve que $x_{n+1} \in [\alpha, b]$, et que (x_n) est une suite décroissante

Étant donné que (x_n) est une suite décroissante dans le compact $[\alpha, b]$, elle est convergente. Notons β sa limite. Nous avons la relation de récurrence

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}. \quad (35.112)$$

En passant à la limite $n \rightarrow \infty$ nous avons l'équation

$$\beta = \beta - \frac{f(\beta)}{f'(\beta)}. \quad (35.113)$$

Vu que $f(x) > 0$ sur $]\alpha, b]$ nous avons automatiquement $\beta = \alpha$. □

35.6.4 Formalisation de l'algorithme

La méthode de Newton consiste à exprimer la solution x de $f(x) = 0$ avec $f \in C^1(\mathbb{R})$ comme limite d'une suite $\{x_n\}_{n \in \mathbb{N}}$ définie par récurrence par la formule

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}. \quad (35.114)$$

où x_0 est arbitraire.

Si on veut exprimer cela en terme d'algorithme, nous disons que l'algorithme de Newton est donné par la suite de problèmes

$$F_n(x_{n+1}, x_n, f) = x_{n+1} - x_n + \frac{f(x_n)}{f'(x_n)}. \quad (35.115)$$

La donnée du problème est la fonction f , et rien que elle.

Plus précisément, une fois que la fonction f est donnée, il existe une infinité de problèmes : pour chaque $a \in \mathbb{R}$ nous avons le problème

$$G_a(x_n, f) = x - a + \frac{f(a)}{f'(a)}. \quad (35.116)$$

La méthode de Newton consiste à sélectionner une partie de ces problèmes de la façon suivante :

$$\left\{ \begin{array}{l} F_0 = G_{x_0} \\ F_n = G_{x_n}. \end{array} \right. \quad (35.117a)$$

$$(35.117b)$$

Le problème F_0 fournit un nombre x_1 qui nous permet de sélectionner le problème G_{x_1} qui va fournir le nombre x_2 , etc.

Au moment de calculer le conditionnement de F_n , nous ne devons pas voir x_{n-1} comme fonction de x_0 et de la donnée f . Il ne faut donc pas dériver à travers les x_n .

Proposition 35.56.

Si une racine est multiple, alors l'ordre de convergence de la méthode de Newton est 1.

Voici un algorithme possible :

```

1 def Newton(f, x0, toll, maxit):
2     fp=f.derivative()
3     n=0

```

```

4   x=x0
5   diff=toll+1
6   while abs(diff)>toll and n<maxit:
7       n=n+1
8       diff=-f(x)/fp(x)
9       x=x+diff
10  return x,n

```

tex/frido/codeSnip_2.py

Commentaires :

- (1) Notons que dans un langage vraiment numérique comme Matlab, il faut passer f' en argument.
- (2) Dans le `while` il faudrait mettre $x_{n+1} - x_n$ (en valeur absolue), mais cette différence est aussi utilisée pour calculer x_{n+1} donc on la calcule une seule fois.
- (3) Il faudrait faire une vérification sur $f(x_n) \neq 0$. Il n'y a pas tellement de choix que de changer le point initial.

35.6.5 Caractéristiques

L'algorithme de Newton a les caractéristiques suivantes :

- (1) Pour résoudre le problème numéro n , il faut avoir résolu le problème numéro $n - 1$.
- (2) Aucune des solutions x_n aux problèmes intermédiaires n'est une solution au problème de départ (à moins d'un coup de chance).
- (3) Étant donné que la donnée du problème F_n est la fonction f de départ, nous avons $d_m = d_n = d$ pour tout m et n .

Théorème 35.57.

Soit f continue sur un voisinage de α , racine simple. Alors il existe un voisinage de α de rayon σ tel que pour tout x_0 dans ce voisinage, la méthode converge vers α avec ordre de convergence $p = 2$.

Donc dès qu'on a continuité autour de la solution recherchée, il suffit de prendre x_0 assez proche pour que tout se passe bien. Cela se fait par localisation des racines, par exemples en traçant la fonction avec un bon niveau de zoom. Le fait est qu'on cherche disons 3 décimales à la main (travail sur ordinateur et graphique) et Newton donne les 20 décimales suivantes à la vitesse de la lumière.

35.6.6 Exemple de la racine carrée

Nous allons nous lancer dans un exemple : le cas de la racine carrée. Soit à calculer une approximation numérique de $\sqrt{2}$. Il s'agit d'une racine de la fonction $f(x) = x^2 - 2$. La fonction de la méthode de Newton associée est :

$$g(x) = x - \frac{f(x)}{f'(x)} = \frac{x^2 - 2}{2x}. \quad (35.118)$$

Cherchons un intervalle autour de $\sqrt{2}$ sur lequel nous avons convergence de la méthode de Newton. Cela s'obtient grâce à la proposition 35.47 qui nous informe qu'il suffit de trouver un intervalle autour de $\sqrt{2}$ sur lequel $|g'(x)| \leq 1$.

Nous avons

$$g'(x) = \frac{x^2 - 2}{2x^2}, \quad (35.119)$$

et nous cherchons à résoudre $|g'(x)| \leq 1$. D'abord $g'(x) = 1$ n'a aucune solutions alors que $g'(\sqrt{2}) = 0$. Donc nous avons $g'(x) \leq 1$ pour tout $x \in \mathbb{R}^+$. Par contre l'équation $g'(x) = -1$ a des solutions : $x = \pm\sqrt{2/3}$.

Nous avons donc convergence de la méthode de Newton pour x_0 dans un intervalle de la forme

$$[\sqrt{2/3}, \sqrt{2} + \dots] \quad (35.120)$$

où les trois points représentent l'expression qu'il faut pour que ce soit symétrique autour de $\sqrt{2}$. La valeur précise n'a pas tellement d'importance parce, vu que nous sommes en train de chercher $\sqrt{2}$, il est peu probable que nous ayons déjà en main une bonne approximation de nombres du type $\sqrt{2/3}$.

Proposition 35.58.

La méthode de Newton pour la fonction $f(x) = x^2 - 2$ converge vers $\sqrt{2}$ pour toute valeur de départ dans $]0, +\infty[$.

Démonstration. La fonction $f(x) = x^2 - 2$ est convexe et $f'(\sqrt{2}) = 2\sqrt{2} > 0$. Donc la méthode converge vers $\sqrt{2}$ pour tout $x_0 \geq \sqrt{2}$ par la proposition 35.55.

Si par contre $x_0 \in]0, \sqrt{2}[$ nous avons

$$x_1 = \frac{x_0^2 + 2}{2x_0}. \quad (35.121)$$

En posant $h(x) = (x^2 + 2)/2x$ et en résolvant $h'(x) = 0$ nous trouvons $x = \sqrt{2}$. Et là, $h(\sqrt{2}) = \sqrt{2}$. Donc $h(x)$ est toujours plus grand que $\sqrt{2}$ pour tout $x \in]0, \sqrt{2}[$.

En d'autres termes, si $x_0 \in]0, \sqrt{2}[$ alors $x_1 \geq \sqrt{2}$ et nous retombons dans le premier cas. \square

35.6.7 Si multiplicité

Supposons que α soit de multiplicité r (définition 35.49).

Cela se remarque en voyant que la méthode de Newton demande plutôt 20 itérations que 5. Le problème que cela pose est que chaque itération, les évaluations provoquent des erreurs. Donc moins d'itérations, c'est mieux.

Nous pouvons modifier la formule avec

$$x_{n+1} = x_n - r \frac{f(x_n)}{f'(x_n)}. \quad (35.122)$$

Il est possible de prouver que cette suite est à nouveau à convergence quadratique.

Ou alors on pose $F(x) = f^{(r-1)}(x)$ et α est une racine simple pour F . Donc faire Newton pour F est à nouveau quadratique, tout en donnant la même solution parce que $F(\alpha) = 0$ et $F'(\alpha) \neq 0$.

La seconde façon est bien parce que le théorème de localisation fonctionne 35.57

Et si r n'est pas connu ?

Il est toujours possible de faire $r = 2$ puis $r = 3$ et caetera jusqu'au moment où l'on remarque que le nombre d'itérations baisse un grand coup.

Mais ça demande beaucoup de calculs. Le mieux est de changer de méthode.

35.6.8 Et la dérivée ?

Un des problèmes de la méthode de Newton est que l'on doit pouvoir calculer la dérivée. Typiquement, il faut savoir f de façon analytique. Si cela n'est pas possible, nous pouvons changer de méthode et utiliser la méthode des sécantes décrite en 35.9.

35.6.9 Méthode de Newton : le cas général

Lemme 35.59.

Soient A et B deux matrices inversibles telles que la matrice $(A + \epsilon B)$ soit inversible pour tout ϵ assez petit. Alors il existe une matrice $X(\epsilon)$ telle que

$$(A + \epsilon B)^{-1} = (A^{-1} + \epsilon X) \quad (35.123)$$

et telle que $\lim_{\epsilon \rightarrow 0} X(\epsilon) = -A^{-1}BA^{-1}$.

Démonstration. Le candidat matrice X est relativement simple à trouver en écrivant

$$(A + \epsilon B)(A^{-1} + \epsilon X) = \mathbb{1} + \epsilon AX + \epsilon BA^{-1} + \epsilon^2 BX. \quad (35.124)$$

En imposant que cela soit $\mathbb{1}$, nous trouvons

$$X(\epsilon) = -(A + \epsilon B)^{-1}BA^{-1}. \quad (35.125)$$

La matrice $X(\epsilon)$ étant un inverse à droite de $(A + \epsilon B)$, son déterminant est non nul et X est inversible. Par conséquent elle est également inversible au sens usuel. Le calcul de la limite est direct :

$$\lim_{\epsilon \rightarrow 0} -(A + \epsilon B)^{-1}BA^{-1} = A^{-1}BA^{-1} \quad (35.126)$$

parce que l'inverse est une fonction continue sur $\mathbb{M}(n, \mathbb{R})$. \square

Remarque 35.60.

Un calcul naïf nous permet de trouver le même résultat de façon plus heuristique. En effet un développement usuel (dans \mathbb{R}) est

$$\frac{1}{a + \epsilon b} = \frac{1}{a} - \frac{\epsilon b}{a^2} + \dots \quad (35.127)$$

Si nous récrivons cela avec des matrices, nous écrivons (attention : passage heuristique!) :

$$(A + \epsilon B)^{-1} = A^{-1} - \epsilon A^{-1}BA^{-1} + \dots \quad (35.128)$$

Notons le choix de généraliser b/a^2 par $a^{-1}ba^{-1}$. Dans les réels les deux écritures sont équivalentes, mais pas dans les matrices.

Étudions si $A^{-1} - \epsilon A^{-1}BA^{-1}$ est bien un inverse à ϵ^2 près de $(A + \epsilon B)$:

$$(A + \epsilon B)(A^{-1} + \epsilon A^{-1}BA^{-1}) = 1 - \epsilon BA^{-1} + \epsilon BA^{-1} - \epsilon^2 BA^{-1}BA^{-1} = 1 - \epsilon^2 BA^{-1}BA^{-1}. \quad (35.129)$$

Par conséquent, à des termes en ϵ^2 près la matrice $A^{-1} - \epsilon A^{-1}BA^{-1}$ est bien un inverse de $A + \epsilon B$.

Théorème 35.61 (Méthode de Newton[472]).

Soit $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ une application de classe C^2 et un point $a \in \mathbb{R}^n$ tel que $f(a) = 0$. Nous supposons que df_a est inversible.

Alors il existe un voisinage V de a tel que pour tout $x_0 \in V$ la suite définie par récurrence

$$x_{n+1} = x_n - (df_a)^{-1}(f(x_n)) \quad (35.130)$$

converge vers a . De plus la vitesse est quadratique au sens où il existe $C > 1$ tel que

$$\|x_n - a\| \leq C^{-1-2^n}. \quad (35.131)$$

Démonstration. Étant donné que df_a est inversible et que df est continue, l'application df_x est continue⁹ pour tout x dans un voisinage de a . Nous prenons $r > 0$ tel que df_x est inversible pour tout $x \in B(a, r)$.

Nous considérons la fonction

$$F: B(a, r) \rightarrow \mathbb{R}^n \\ x \mapsto x - (df_x)^{-1}(f(x)). \quad (35.132)$$

Cela est une application C^1 . La clef est de montrer que l'application de F à un point $a + h$ rapproche de a pourvu que h soit assez petit. Nous avons la formule suivante :

$$F(a + h) - F(a) = h - (df_{a+h})^{-1}(f(a + h)). \quad (35.133)$$

9. Nous pouvons voir df comme l'application qui à x fait correspondre la matrice $df_x \in \mathbb{M}(n, \mathbb{R})$. Cette application étant continue et la non inversibilité d'une matrice étant donnée par l'annulation du déterminant, les matrices inversibles forment un ouvert dans l'ensemble des matrices.

Nous allons maintenant utiliser un développement de Taylor par rapport à h en suivant la formule (13.982). Nous avons

$$f(a+h) = f(a) + df_a(h) + \|h\|^2 \xi(h) \quad (35.134)$$

où $\xi: \mathbb{R}^n \rightarrow \mathbb{R}^n$ est une fonction qui tend vers une constante lorsque $h \rightarrow 0$. Nous avons aussi

$$df_{a+h} = df_a + \|h\| \tau(h) \quad (35.135)$$

où $\tau: \mathbb{R}^n \rightarrow \mathbb{M}(n, \mathbb{R})$ est une application qui tend vers une constante lorsque $h \rightarrow 0$. En ce qui concerne l'inverse nous utilisons le lemme¹⁰ 35.59 :

$$(df_a + \|h\| \tau(h))^{-1} = (df_a)^{-1} + \|h\| A(h) \quad (35.136)$$

où A est une autre matrice fonction de h qui tend vers une constante lorsque h tend vers zéro. En substituant le tout dans (35.133) nous trouvons

$$F(a+h) - F(a) = \|h\|^2 (df_a)^{-1} \xi(h) + \|h\| (A(h) \circ df_a)(h) + \|h\|^3 A(h) \xi(h). \quad (35.137)$$

En ce qui concerne la norme nous utilisons le fait que si T est un opérateur, $\|Tx\| \leq \|T\| \|x\|$. Nous trouvons

$$\|F(a+h) - F(a)\| \leq \|h\|^2 \|(df_a)^{-1}\| \|\xi(h)\| + \|h\|^2 \|A(h) \circ df_a\| + \|h\|^3 \|A(h)\| \|\xi(h)\| \quad (35.138a)$$

$$= \|h\|^2 \alpha(h) \quad (35.138b)$$

pour une certaine fonction $\alpha: \mathbb{R}^n \rightarrow \mathbb{R}$ qui tend vers une constante lorsque $h \rightarrow 0$.

En posant $C = \lim_{h \rightarrow 0} \alpha(h)$ nous avons la majoration

$$\|F(x) - a\| \leq C \|x - a\|^2. \quad (35.139)$$

Nous pouvons également supposer que $C > 1$. Afin de prouver la vitesse de convergence (35.131), nous allons encore redéfinir r en demandant $r < 1/C^2$. De cette manière nous avons

$$\|x_0 - a\| \leq \frac{1}{C^2} \quad (35.140)$$

et la récurrence sur n est :

$$\|x_{n+1} - a\| = \|F(x_n) - a\| \leq C \|x_n - a\|^2 \leq C (C^{-1-2^n})^2 = C^{-1-2^{n+1}}. \quad (35.141)$$

Note : ce dernier calcul est le lemme 35.44 appliqué à $r = (1/C^2)$. □

Remarque 35.62.

La valeur de la constante C a été fixée par l'équation (35.139). Certes nous pouvons toujours choisir C plus grand afin d'augmenter la vitesse de convergence, mais le point de départ x_0 devant être dans une boule de taille $1/C^2$ autour de a , demander C plus grand revient à demander un point de départ plus précis.

35.7 Estimation de l'ordre de convergence

Définition 35.63 ([473]).

Nous disons que la suite (x_n) de limite x est **convergente d'ordre q** pour $q > 1$ s'il existe $\mu > 0$ tel que

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x|}{|x_k - x|^q} = \mu. \quad (35.142)$$

En particulier :

¹⁰. Pour l'inversibilité de $\|h\| \tau(h)$, notons que df_a est inversible et que par hypothèse la somme $df_a + \|h\| \tau(h)$ est inversible.

- la convergence d'ordre 2 est dite quadratique,
- la convergence d'ordre 3 est dite cubique,
- la convergence d'ordre 4 est dite quartique.

Comment estimer numériquement l'ordre p de convergence de la méthode ? Soit une suite (x_n) convergente vers α . Considérons les 4 termes x_{n-3} , x_{n-2} , x_{n-1} , x_n . Alors nous pouvons écrire l'approximation

$$\frac{|x_n - x_{n-1}|}{|x_{n-1} - x_{n-2}|} \simeq \left(\frac{|x_{n-1} - x_{n-2}|}{|x_{n-1} - x_{n-3}|} \right)^p. \quad (35.143)$$

Cette approximation ne serait pas trop mauvaise tant que n est assez grand pour que la convergence soit bien engagée. Passons au logarithme :

$$\ln \frac{|x_n - x_{n-1}|}{|x_{n-1} - x_{n-2}|} \simeq p \ln \left(\frac{|x_{n-1} - x_{n-2}|}{|x_{n-1} - x_{n-3}|} \right). \quad (35.144)$$

et donc

$$p \simeq \frac{\ln \left(\frac{|x_n - x_{n-1}|}{|x_{n-1} - x_{n-2}|} \right)}{\ln \left(\frac{|x_{n-1} - x_{n-2}|}{|x_{n-1} - x_{n-3}|} \right)}. \quad (35.145)$$

Avec cette approximation, en réalité nous calculons une suite (p_i) qui sont les approximations de p à partir des termes i à $i + 3$ de la suite (x_n) . Il s'agit d'une suite d'estimations de p .

- (1) Dans le cas de la bisection, nous obtenons toujours $p_i = 1$.
- (2) Dans le cas de la méthode de Newton (35.6) nous avons $p = 2$. Mais les premières valeurs de p_i peuvent être aussi bien 0 que 7. Après quelques itérations pourtant les p_i se regroupent autour de 2.

En tout cas, le plus important est de savoir si $p > 1$ ou non. Rappel : nous voulons la superlinéarité parce que nous voulons utiliser le test d'arrêt de la différence entre deux termes, voir 35.46.

35.8 Autres méthodes

35.8.1 Méthode de Schröder

La formule est

$$x_{n+1} = x_n - \frac{f(x_n)f'(x_n)}{f'(x_n)^2 - f(x_n)f''(x_n)} \quad (35.146)$$

Cette méthode est d'ordre 2 pour toute racine et toute valeur de multiplicité. Le problème de cette méthode est qu'elle demande 3 évaluations de f . Son efficacité :

$$E = \sqrt[3]{2} \simeq 1.25 \quad (35.147)$$

Cela est donc moins efficace que Newton.

35.8.2 Halley

Il a $p = 3$ lorsque α est racine simple. Mais encore $p = 1$ pour les racines multiples. Plus efficace que Newton pour les racines simples, mais même problème pour les racines multiples.

$$x_{n+1} = x_n - \frac{2f(x_n)f'(x_n)}{2f'(x_n)^2 - f(x_n)f''(x_n)} \quad (35.148)$$

35.9 Méthode des sécantes variables

Si nous n'avons pas de formule analytique pour f , mais seulement la possibilité de calculer $f(x)$ pour tout x . Newton ne fonctionne pas, mais la bisection fonctionne.

Nous pouvons approximer

$$f'(x_n) = \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}. \quad (35.149)$$

En substituant dans la formule de Newton, nous obtenons

$$x_{n+1} = x_n - \frac{f(x_n)(x_n - x_{n-1})}{f(x_n) - f(x_{n-1})}. \quad (35.150)$$

Il s'agit de prendre la droite qui passe par $(x_{n-1}, f(x_{n-1}))$ et par $(x_n, f(x_n))$ et de prendre l'intersection de cette droite avec l'axe $y = 0$. Cela donne le x_{n+1} .

Pour cette méthode, il ne faut pas seulement x_0 mais également x_1 .

L'ordre de convergence est le nombre d'or

$$p = \frac{1 - \sqrt{5}}{2} \simeq 1.618. \quad (35.151)$$

Cela est donc superlinéaire.

La nombre d'évaluations est $s = 1$ (il y a deux apparitions de f dans la formule, mais l'une des deux est récupérée dans l'itération suivante). Donc l'efficacité est

$$E = p. \quad (35.152)$$

Donc bien efficace.

Proposition 35.64.

Si α est racine simple, il existe un voisinage de α tel que pour tout choix de x_0, x_1 dans ce voisinage, la méthode converge.

Psychologiquement, on est tenté de prendre x_0 et x_1 de part et d'autre de α (pensant à la bisection), mais en réalité ce n'est pas obligatoire du tout et n'a aucune influence. Il faut seulement les prendre très proches de α .

Remarque 35.65.

La méthode de la sécante est souvent écrite sous la forme

$$x_{n+1} = \frac{x_{n-1}f(x_n) - x_n f(x_{n-1})}{f(x_n) - f(x_{n-1})}. \quad (35.153)$$

C'est évidemment algébriquement équivalent.

Les formules (35.151) et (35.153) ont toutes deux des erreurs de cancellation. Laquelle est la plus grave ?

Dans la première, si la fraction est mal calculée, elle ne fait que modifier x_n . C'est-à-dire qu'on peut espérer qu'à la prochaine itération, ça aille mieux. En tout cas, dans ce cas si la fraction est mal calculée, ça ne détruit pas tout.

Dans la seconde, c'est la valeur elle-même qui risque d'être mal calculée. Et si la fraction est mal calculée, alors on casse complètement l'éventuel bonne approximation que nous avons déjà.

35.9.1 Aitken

La méthode du Δ^2 de Aitken est une méthode d'accélération de la convergence.

Soit (x_n) une suite qui converge. Nous voudrions une nouvelle suite (y_n) telle que

$$\lim_{n \rightarrow \infty} \frac{y_n - \alpha}{x_n - \alpha} \quad (35.154)$$

C'est la définition d'une convergence accélérée.

La façon de faire est :

$$y_n = \frac{x_{n+2}x_n - x_{n+1}^2}{x_{n+2} - 2x_{n+1} + x_n} = x_n - \frac{(x_{n+1} - x_n)}{x_{n+2} - 2x_{n+1} + x_n}. \quad (35.155)$$

La première expressions a deux cancellations (la seconde une seule) et de plus la première est y_n elle-même alors que la seconde est une correction.

Donc la seconde expression est numériquement meilleure.

L'opérateur Δ appliqué à une suite est :

$$(\Delta x)_n = x_{n+1} - x_n \quad (35.156)$$

Donc

$$(\Delta^2 x)_n = (\Delta x)_{n+1} - (\Delta x)_n = x_{n+2} - x_{n+1} - x_{n+1} + x_n = x_{n+2} - 2x_{n+1} + x_n. \quad (35.157)$$

L'accélération a alors la formule

$$y_n = \frac{(\Delta x)_n^2}{(\Delta^2 x)_n}. \quad (35.158)$$

Le problème est que ça accélère tellement que l'on arrive vite à des erreurs de cancellations, et donc à une précision en pics oscillants.

35.10 Équations algébrique

C'est une équation du type $P(x) = 0$ où P est un polynôme. Soit un polynôme de degré n . Nous en savons des choses.

- (1) L'équation a exactement n solutions dans \mathbb{C} en comptant les multiplicités.
- (2) Les racines complexes arrivent par paire complexes conjuguée. Elles sont donc toujours en nombre pair.

Si donc nous avons $n = 3$, nous ne pouvons pas avoir 2 racine réelles. Il y en a donc 1 ou 3 réelles. Pas zéro ni deux.

Quelques méthodes : Müller, matrice compagnon, Laguerre.

35.10.1 Résoudre un système linéaire

Pour résoudre un système linéaire d'équations, nous échelonons la matrice du système. Soit à résoudre le système $Ax = b$ où

$$A = \begin{pmatrix} 2 & 4 & -6 \\ 1 & 5 & 3 \\ 1 & 3 & 2 \end{pmatrix}, \quad \text{et } b = \begin{pmatrix} -4 \\ 10 \\ 5 \end{pmatrix}. \quad (35.159)$$

En termes de problèmes, on écrit $F(x, (A, b)) = Ax - b$. La donnée de ce problème est le couple (A, b) .

En ce qui concerne l'algorithme, on pose comme premier problème

$$F_1(x_1, (A_1, b_1)) = A_1 x_1 - b_1 = 0 \quad (35.160)$$

avec $A_1 = A$ et $b_1 = b$.

Ensuite, on commence à échelonner et le second problème est

$$F_2(x_2, (A_2, b_2)) = A_2 x_2 - b_2 = 0 \quad (35.161)$$

avec

$$A = \begin{pmatrix} 2 & 4 & -6 \\ 0 & 3 & 6 \\ 0 & 1 & 5 \end{pmatrix}, \quad \text{et} \quad b = \begin{pmatrix} -4 \\ 12 \\ 13 \end{pmatrix}. \quad (35.162)$$

Le troisième problème sera

$$F_3(x_3, (A_3, b_3)) = A_3 x_3 - b_3 = 0 \quad (35.163)$$

avec

$$A = \begin{pmatrix} 2 & 4 & -6 \\ 0 & 3 & 6 \\ 0 & 0 & 3 \end{pmatrix}, \quad \text{et} \quad b = \begin{pmatrix} -4 \\ 12 \\ 3 \end{pmatrix}. \quad (35.164)$$

Ce problème est facile à résoudre « à la main ». Nous nous arrêtons donc ici avec l'algorithme, et nous trouvons le x_3 qui résout le problème F_3 .

35.10.2 Caractéristiques

L'algorithme de résolution de systèmes linéaires d'équations a les propriétés suivantes, à mettre en contraste avec celles de Newton :

- (1) Pour résoudre le problème numéro n , il n'a pas fallu résoudre le problème numéro $n - 1$.
- (2) Toutes les solutions x_n des problèmes intermédiaires sont solutions du problème de départ. Nous avons $F_n(x, d_n) = 0$ pour tout n (ici, $d_n = (A_n, b_n)$).
- (3) D'un problème à l'autre, les données changent énormément : la matrice échelonnée peut être très différente de la matrice de départ.

35.10.3 Définitions

Nous allons maintenant formaliser en donnant quelques définitions pour nommer les propriétés que nous avons vues. D'abord, un algorithme est une suite de problèmes. Un **algorithme** pour résoudre un problème $F(x, d) = 0$ est une suite de problèmes $\{F_n(x_n, d_n) = 0\}_{n \in \mathbb{N}}$.

Définition 35.66.

Un tel algorithme est dit **fortement consistant** si pour toutes données admissibles d_n , on a

$$F_n(x, d_n) = 0 \quad \forall n, \quad (35.165)$$

où x est la solution de $F(x, d) = 0$.

L'algorithme des matrices est fortement consistant, mais pas l'algorithme de Newton.

Définition 35.67.

Un algorithme est **consistant** si $\lim_{n \rightarrow \infty} F_n(x, d_n) = 0$.

Dans le cas de l'algorithme de Newton, c'est plutôt une telle consistance qu'on attend.

L'algorithme est dit **stable** si pour tout n le problème correspondant est stable. Dans ce cas, on note K^{num} le **conditionnement relatif asymptotique** défini par

$$K^{\text{num}} = \limsup_n K_n \quad (35.166)$$

où K_n est le conditionnement relatif du problème $F_n(x_n, d_n) = 0$.

Définition 35.68.

Un algorithme est dit **convergeant** (en d) si pour tout $\epsilon > 0$, il existe $N = N(\epsilon)$ et $\delta = \delta(N, \epsilon)$ tels que pour $n \geq 0$ et $|d - d_n| < \delta$, on ait $|x(d) - x_n(d_n)| < \epsilon$.

Remarque 35.69.

Dans le cas de l'algorithme de Newton, nous avons vu que la donnée d_n du problème F_n était en fait la même que la donnée initiale d , donc nous avons $d_n = d$, et par conséquent nous avons toujours $|d - d_n| < \delta$. Dans ce cas, la définition de la convergence revient à demander que la suite numérique des x_n converge vers la solution x .

Remarque 35.70.

Dans le cas des matrices par contre, les données sont très différentes les unes des autres, nous avons donc en général que $|d - d_n| > \delta$. Mais en revanche nous savons que tous les problèmes intermédiaires F_n acceptent une solution unique¹¹ $x_n(d_n) = x(d)$. Par conséquent, $|x_n(d_n) - x(d)|$ est toujours plus petit que ϵ . L'algorithme des matrices est donc toujours un algorithme convergeant.

35.11 Équations non linéaire

Certains équations non linéaires sont résoluble explicitement, par exemples les polynômes de degré jusqu'à 4 ou des choses comme

$$\sin^2(x) + 3\sin(x) + 5 = 0. \quad (35.167)$$

Mais ces exemples sont très rares.

Nous allons étudier des équations du type $f(x) = 0$, dans \mathbb{R} .

- (1) Un problème écrit sous la forme $x = g(x)$ peut utiliser des théorèmes de points fixes.
- (2) Un problème sous la forme $f(x) = 0$ peut utiliser des méthodes de bisection, Newton ou autres.

Il y a évidemment beaucoup de façons de transformer un problème pour passer d'une forme à l'autre.

Exemple 35.71

Soit $f(x) = x^2 - a = 0$ avec $a > 0$. Nous pouvons l'écrire

$$x^2 + x - a = x \quad (35.168)$$

qui donne une forme $g(x) = x$ pour $g(x) = x^2 + x - a$.

Ou encore $x = \frac{a}{x}$ et donc $g(x) = a/x$ (si par ailleurs on sait que $x \neq 0$). Notons que $x \neq 0$ n'est pas une hypothèse très forte parce qu'on la vérifie directement sur a . \triangle

Exemple 35.72

Soit l'équation à résoudre

$$f(x) = x^2 - 2 - \ln(x) = 0 \quad (35.169)$$

Les solutions de cette équations peuvent être vues comme les intersections avec l'axe X du graphe $y = x^2 - 2 - \ln(x)$. Tracer peut donc aider. Par ailleurs, il faut noter que

$$\lim_{x \rightarrow \pm\infty} f(x) = \infty, \quad (35.170)$$

donc les solutions sont certainement contenues dans un compact de \mathbb{R} .

À part tracer nous pouvons écrire

$$x^2 - 2 = \ln(x). \quad (35.171)$$

Et là, ce sont deux fonctions dont nous pouvons tracer le graphe pour trouver graphiquement les points d'intersection. Une étude de fonction montre vite qu'il y a exactement deux solutions, qu'elles sont strictement positives. Pour trouver des bornes, il faut calculer par exemple pour $x = 2$ les valeurs de $\ln(x)$ et $x^2 - 2$ pour voir si le graphe de $x^2 - 2$ est déjà plus haut. \triangle

La majorité des méthodes numériques de résolution d'équation du type $f(x) = 0$ ou $x = g(x)$ seront sous la forme de suites. Avec questions à la clefs :

11. Nous n'envisageons que le cas où le déterminant est non nul.

- (1) Quel point de départ choisir ?
- (2) Convergence ?
- (3) Est-ce que la limite est bien une solution ?
- (4) Vu que la limite est unique, comment faire si l'équation a plusieurs solutions ? (souvent c'est le choix du point initial qui va jouer sur ce point)

35.73.

Si la fonction est très plate, il est possible d'avoir

$$|f(\tilde{\alpha})| \leq \epsilon \quad (35.172)$$

sans que $\tilde{\alpha}$ ne soit une bonne approximation.

Lorsqu'on fait tourner une méthode itérative résolvant $f(x) = 0$, il n'est pas suffisant de s'arrêter lorsque

$$f(x_n) \leq \epsilon_1. \quad (35.173)$$

Il faut aussi s'assurer que, si \bar{x} est la solution exacte, $|x_n - \bar{x}| \leq \epsilon_2$. Ici ϵ_1 et ϵ_2 sont deux « précisions » que nous nous fixons au départ.

Évidemment, vérifier la condition $|x_n - \bar{x}| \leq \epsilon_2$, il faudrait savoir \bar{x} . Et savoir \bar{x} c'est justement le problème. Nous sommes donc amenés à faire des estimations de $|x_n - \bar{x}|$.

35.74.

Lorsque nous effectuons une méthode itérative, il faut donc contrôler deux grandeurs :

$$|\bar{x} - x_n| \leq \epsilon_1 \quad (35.174a)$$

$$|x_{n+1} - x_n| \leq \epsilon_2. \quad (35.174b)$$

Proposition 35.75.

Soit p l'ordre de convergence de la suite (x_n) vers \bar{x} . Si $p > 1$ et $|x_{n+1} - x_n| \leq \epsilon_2$ alors $|\bar{x} - x_n| \leq \epsilon_2$.

35.11.1 Méthode de bisection

Il y a ce théorème des valeurs intermédiaires.

Théorème 35.76.

Soit f continue sur $[a, b]$ telle que $f(a)f(b) < 0$. Alors il existe au moins une solution à l'équation $f(x) = 0$ sur l'intervalle $]a, b[$.

Pour démarrer une bisection, il est toujours bon de prendre l'intervalle $[a, b]$ de façon à ne contenir qu'une seule solution.

Soit donc un premier intervalle $[a_0, b_0]$ tel que $f(a_0)f(b_0) < 0$ et ne contenant qu'une seule solution. À chaque itération nous considérons la moitié de l'intervalle précédent, mais la moitié contenant la solution.

Le test d'arrêt de la méthode de bisection se base uniquement sur la taille de l'intervalle qui reste. En effet si nous avons

$$|b_n - a_n| \leq \epsilon \quad (35.175)$$

nous avons certainement

$$|\bar{x} - x_n| \leq \frac{\epsilon}{2} \quad (35.176)$$

où x_n est le point du milieu de $[a_n, b_n]$.

35.77.

La fonction f n'intervient dans la méthode que via son signe, pas via ses valeurs exactes.

35.78.

Notons que le théorème des valeurs intermédiaires n'est pas très puissant pour choisir l'intervalle de départ ; penser à la fonction

$$f(x) = x^2 - 5 \quad (35.177)$$

sur l'intervalle $[-10, 10]$. Il y a bien deux solutions dans l'intervalle, mais elles sont invisibles du théorème des valeurs intermédiaires. La fonction $x \mapsto x^2$ a sa solution en $x = 0$, mais elle aussi n'est pas visible.

35.79.

Certes la méthode de bisection assure la convergence vers une solution, mais elle n'assure pas la convergence monotone. Il peut arriver que $|\bar{x} - x_n| < |\bar{x} - x_{n+1}|$. C'est le cas lorsque la solution est très proche du milieu de l'intervalle choisit. Le x_0 est alors proche de \bar{x} alors que x_1 sera à une distance de \bar{x} d'environ un quart de l'intervalle de départ.

Supposons déjà avoir trouvé un intervalle $[a, b]$ dans lequel se trouve une unique solution à $f(x) = 0$. Voici un algorithme possible.

```

1 from __future__ import division
2
3 def bisection(f, a, b, toll, mmax):
4     """
5     f : une fonction
6     a,b : les limites de l'intervalle
7     toll : la tolérance. C'est l'amplitude de l'intervalle à ←
8           partir duquel nous nous arrêtons.
9     nmax : le nombre maximum d'itérations.
10
11     Nous supposons que \(\ b>a\).
12
13     Retourne un tuple (x,n) où 'x' est la solution approchée et 'n←
14           ' est le nombre d'itérations effectuées
15     """
16     n = -1
17     amp = toll + 1 # Pour s'assurer que l'on entre dans le cycle
18     while amp > toll and n < mmax :
19         n = n + 1
20         amp=abs(b - a)
21         x = a + amp / 2
22         if f(a) * f(x) < 0:
23             b = x
24         elif f(a)*f(x) > 0:
25             a = x
26         else :          # Problème ZERO
27             amp = 0
28     return (x, n)

```

tex/frido/codeSnip_1.py

Plusieurs remarques :

- (1) Le fait de retourner le nombre d'itérations effectuées permet à l'utilisateur de savoir la précision et si le nombre maximum d'itérations est dépassé. Si ce n retourné est égal à n_{\max} , l'utilisateur sait que le x retourné n'est pas fiable.
- (2) La ligne `from __future__ import division` fait en sorte que l'opération `/` est bien la division usuelle. Sinon, le défaut en python 2 est que `/` soit la division *entière*, c'est-à-dire

que $1/2 = 0$ en python 2. En python 3, le symbole $/$ désigne bien la division usuelle, mais Sage utilise Python 2.

- (3) Même si l'intervalle $[a, b]$ contient plus d'une solution, la méthode fonctionne et donne une solution. Il est simplement éventuellement très compliqué de savoir laquelle.
- (4) Nous faisons `amp=tol+1` parce que nous voulons absolument lancer le cycle au moins une fois. Sinon, le `x` à retourner ne serait pas défini au moment de sortir du cycle (si le cycle n'est pas exécuté).
- (5) Calculer le point milieu d'un intervalle $[a, b]$ est par la formule $(a + b)/2$ sauf que cette opération est numériquement dangereuse parce qu'à cause de l'arithmétique en précision finie, il est possible que cela tombe *exactement* sur a ou b . D'où le fait de calculer le point milieu par

$$x = a + \frac{amp}{2}. \quad (35.178)$$

- (6) Dans le cas **Problème ZERO** nous déduisons $f(x) = 0$. Attention que c'est pas que $f(x) = 0$ mais simplement que en mettant x dans f , la *machine* retourne son zéro.
Il peut cependant avoir une fonction telle que $f(1) = 10^{-50}$ et $f(2) = 0$. L'algorithme de bisection risque de s'arrêter si $x_n = 1$. Parce que la machine risque de calculer $f(x_n) = 0$.
Quoi qu'il en soit, nous y mettons `amp=0` pour être sûr de sortir de la boucle dès la prochaine vérification.
- (7) Il y a moyen de sauver les valeurs de $f(a)$ et $f(x)$ pour ne pas les recalculer, et en particulier au moment de faire `b=x` nous pouvons poser `fa=fx`.

Si τ est la précision de la solution voulue, nous pouvons fixer a priori le nombre d'itérations à faire grâce à la formule

$$n \geq \left\lceil \log_2 \left(\frac{b-a}{\tau} \right) \right\rceil. \quad (35.179)$$

Il y a un "≥" et non une égalité parce qu'en arithmétique numérique, le nombre obtenu à droite pourrait ne pas être le bon à 1 près.

Ici pour $\nu \in \mathbb{R}$ le nombre $\lceil \nu \rceil$ est le plus petit entier à être plus grand ou égal à ν .

35.80.

Notons l'importance de la continuité de f . Par exemple que ferait la bisection sur la fonction $f(x) = 1/x$ pour l'intervalle $[-3, 1]$?

Il y a changement de signe sans avoir de racine.

Vu que 2^{10} est déjà 1024. Donc si on veut de la précision de l'ordre de $1/1000$, dix itérations suffisent. Si donc nous avons besoin de 200 itérations pour atteindre la précision voulue, c'est l'occasion de trouver un intervalle plus petit. Par exemple en traçant la fonction, en faisant un zoom et en trouvant des valeurs de a et b qui sont déjà proches.

35.81.

Dans le monde réel, il arrive souvent d'utiliser une méthode de bisection pour se donner un point de départ pour une autre méthode.

35.12 Efficacité

Définition 35.82.

L'*efficacité* est le nombre

$$E = \sqrt[s]{p} \quad (35.180)$$

où p est l'ordre de convergence de la méthode et s est le nombre de fois qu'il faut calculer une valeur de la fonction à chaque itération (nous ne comptons pas l'initialisation).

Que le nombre de d'évaluations de f intervienne est logique parce que chaque évaluation provoque une erreur possible.

Exemple 35.83(Bisection)

Pour la méthode de bisection, nous avons $s = 1$ parce que chercher x_{n+1} , il faut seulement calculer $f(x_n)$. \triangle

Exemple 35.84(Newton)

Pour l'algorithme de Newton nous avons $p = 2$ et il y a deux évaluations à chaque itération (une fois f et une fois f'), donc $s = 2$ et $E = \sqrt{2}$. \triangle

35.13 Exemples sous forme d'exercices

Exemple 35.85

Nous supposons une machine acceptant 5 chiffres significatifs. Elle retient les nombres sous la forme $\pm 0.x \cdots \times 10^{\cdots}$.

$$x = 1.403, y = 0.4112 \times 10^{-3}, z = -0.4111 \times 10^{-3}.$$

$$\text{Soient } a = (x \oplus y) \oplus z \text{ et } b = (y \oplus z) \oplus x.$$

- (1) D'abord la calculer à la main.
- (2) Quel est le calcul préférable ?
- (3) Donner l'erreur relative avec 3 chiffres significatifs.

Dans le cas du calcul à la main, il faut en faire un seul parce que, algébriquement, $a = b$.

Nous avons $x = 1.403$, $y = 0.0004112$ et $z = -0.0004111$. Et la somme donne :

$$a = b = 1.4030001 = 0.14030001 \times 10^1. \quad (35.181)$$

Faisons d'abord la normalisation de x , c'est-à-dire $\text{fl}(x)$.

$$\text{fl}(x) = 0.1403 \times 10^1. \quad (35.182)$$

et y est déjà normalisé :

$$\text{fl}(y) = 0.4112 \times 10^{-3}. \quad (35.183)$$

Il n'y a pas d'erreurs d'assignation pour ces deux nombres.

Pour faire la somme, il faudra déjà un peu casser les nombres pour les écrire de façon à pouvoir les sommer. En effet, il faut écrire les deux nombres avec le même exposant de 10 (le plus grand), pour pouvoir les mettre en colonne :

$$0.1404 \times 10^1 \rightarrow 0.1404 \times 10^1 \quad (35.184a)$$

$$0.4112 \times 10^{-3} \rightarrow 0.00004 \times 10^1. \quad (35.184b)$$

La somme donne 0.14034×10^1 . Et ça, c'est à nouveau arrondi. Le premier chiffre supprimé est un 4, donc

$$x \oplus y = 0.1403 \times 10^1. \quad (35.185)$$

Et là on remarque que nous avons la même chose que x . C'est un classique du calcul numérique.

Nous avons aussi

$$\text{fl}(z) = 0.4111 \times 10^{-3}. \quad (35.186)$$

Et pour faire la somme de cela avec $x \oplus y$ nous devons le remettre sous la forme d'un 10^1 :

$$\text{fl}(z) \rightarrow -0.00004 \times 10^1 \quad (35.187)$$

(erreur de conversion), et en sommant on trouve

$$(x \oplus y) \oplus z = 0.140216 \times 10^1, \quad (35.188)$$

qui est encore arrondi. Le premier chiffre supprimé est un 6, donc

$$\text{fl}(a) = 0.1403 \times 10^1, \quad (35.189)$$

Le nom de l'erreur qui consiste à avoir $x \oplus y = x$ est "relation anormale".

Calculons b .

Les nombres y et z ont même ordre de grandeur, donc pas d'erreur au moment de les mettre sous forme sommable.

$$\text{fl}(x) + \text{fl}(y) = 0.00010 \times 10^{-3}. \quad (35.190a)$$

Cela est renormalisé et arrondi : $\text{fl}(x) \oplus \text{fl}(y) = 0.1000 \times 10^{-6}$.

Notons que nous avons ici commis potentiellement une erreur de cancellation parce que entre y et z , il y a 3 chiffres sur 4 qui sont identiques. Seul le chiffre 1 est significatif en réalité.

Il faut maintenant ajouter x à cela. D'abord

$$\text{fl}(x) = 0.1403 \times 10^1. \quad (35.191)$$

Pour cette somme, il faudra remettre notre 0.1000×10^{-6} avec une puissance 10^1 . Et là, nous obtenons zéro parce que vraiment ce nombre est trop petit pour être écrit avec 10^1 . Résultat des courses :

$$\text{fl}(b) = 0.1403 \times 10^1. \quad (35.192)$$

Dans le premier calcul nous avons deux "relations anormales" et dans le second nous en avons une plus une cancellation.

Nous préférons avoir deux relations anormales, parce que l'erreur de cancellation est plus grave : elle consiste à une perte de chiffre significatifs. Le fait est que faisant la différence à l'ordinateur nous avons obtenu 0.1 qui est certes exact, mais qui est un coup de bol : la différence aurait aussi bien pu être 0.19 avec d'autres nombres, machinement égaux.

Note : avec les données ici, il n'y a en fait pas d'erreur de cancellation. Mais il y a une erreur potentielle de cancellation, potentiellement grave.

En ce qui concerne l'erreur relative. Dans la formule

$$\epsilon_r = \frac{|a - a^*|}{|a|}, \quad (35.193)$$

la différence ne peut pas être calculée à la calculatrice justement parce qu'elle est très potentiellement sujette à erreur de cancellation.

$$\epsilon_r = \frac{|0.1030001 \times 10^1 - 0.1403 \times 10^1|}{0.14030001 \times 10^1} = \frac{0.1 \times 10^{-6}}{0.14030001 \times 10^1} \simeq 0.712758 \times 10^{-7}. \quad (35.194)$$

En passant à 3 chiffres significatifs, 0.713×10^{-7} (le premier chiffre supprimé est un 7).

△

Exemple 35.86

Soient $x = 0.1 \times 10^{21}$ et $y = 0.5 \times 10^{20}$ et les expressions

- (1) $z_1 = \frac{x-y}{y} + \frac{x+y}{x}$
- (2) $z_2 = \frac{x^2+y^2}{xy}$.

Ces deux expressions sont algébriquement équivalentes.

- (1) Calculer les valeurs.

(2) On suppose une machine en précision simple. Laquelle des deux expressions est préférable ?

Pour z_1 , en arithmétique exacte :

$$z_1 = \frac{0.5 \times 10^{20}}{0.5 \times 10^{20}} + \frac{1.5 \times 10^{20}}{1 \times 10^{20}} = 0.25 \times 10^1. \quad (35.195)$$

Le calcul exact de z_2 donne la même chose.

Calcul de z_1 Les deux valeurs sont mémorisables et la différence $x - y$ se fait sans erreurs de cancellation. Idem pour la somme $x + y$. Idem pour les divisions.

Calcul de z_2 Pour faire x^2 , c'est pas possible parce que c'est de l'ordre de 10^{40} alors que nous sommes en précision simple. Idem pour le produit xy .

Morale : z_2 donne un **overflow** alors que z_1 fonctionne de façon exacte.

Remarque 35.87.

En réalité le z_1 n'est pas tout à fait calculable de façon exacte sur la machine parce qu'elle doit d'abord convertir en binaire, ce qui n'est pas toujours possible. Mais sur notre machine qui fonctionne en base 10, il n'y a pas de problèmes.

△

Exemple 35.88

Soit la fonction

$$f(x) = 2x^2 - 4x + 2 - e^{-x}. \quad (35.196)$$

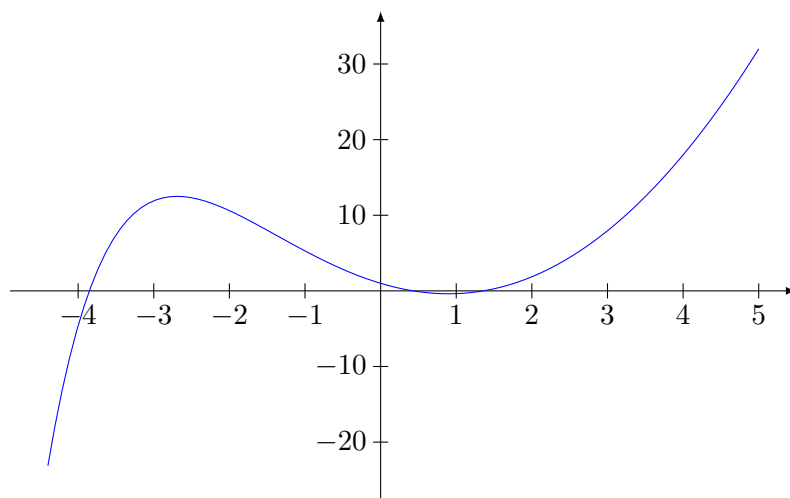
- (1) Identifier la plus grande des solutions réelles de $f(x) = 0$.
- (2) Effectuer une bisection pour la savoir.
- (3) Sachant que

$$\alpha \simeq 1.358500220734946, \quad (35.197)$$

quelle est l'erreur relative ?

Note : si c'est pour chercher à la main des approximations pour démarrer, il est évidemment préférable de dessiner $f_1(x) = 2x^2 - 4x + 2$ et $f_2(x) = -e^{-x}$ séparément.

Quoi qu'il en soit, voici un graphique :



Nous voyons trois racines : $\alpha_1 \in [0, 0.5,]$, $\alpha_2 \in [1, 1.5]$ et $\alpha_0 \in [-4, -3.5]$.

La plus grande solution est α_2 . Nous pouvons déjà remplir le tableau des précisions :

n	a_n	b_n	x_n	$g(x_n)$	$ b_n - a_n $
< ++ >	< ++ >	< ++ >	< ++ >	< ++ >	0.5
< ++ >	< ++ >	< ++ >	< ++ >	< ++ >	0.25
< ++ >	< ++ >	< ++ >	< ++ >	< ++ >	0.125

Et nous calculons les valeurs de f aux points d'extrémité de l'intervalle. Note que seul le signe nous importe :

$$f(1) \simeq -0.368 \quad (35.198a)$$

$$f(1.5) \simeq 0.278. \quad (35.198b)$$

Voici donc le tableau avec le signe de f indiqué :

n	a_n	b_n	x_n	$g(x_n)$	$ b_n - a_n $
0	1(+)	1.5(-)	< ++ >	-0.162×10^0	0.5
1	< ++ >	< ++ >	< ++ >	< ++ >	0.25
2	< ++ >	< ++ >	< ++ >	< ++ >	0.125

Puis :

n	a_n	b_n	x_n	$g(x_n)$	$ b_n - a_n $
0	1(+)	1.5(-)	1.25(-)	-0.162×10^0	0.5
1	< ++ >	< ++ >	< ++ >	< ++ >	0.25
2	< ++ >	< ++ >	< ++ >	< ++ >	0.125

Et enfin :

n	a_n	b_n	x_n	$g(x_n)$	$ b_n - a_n $
0	1(+)	1.5(-)	1.25(-)	-0.162×10^0	0.5
1	1.25(-)	1.5(+)	1.375(+)	$+0.284 \times 10^{-1}$	0.25
2	1.25(-)	1.375(+)	1.3125(-)	-0.738×10^{-1}	0.125

Note que les $f(x_n)$ restent toujours du même ordre de grandeur. Si un moment on voit un 1.6×10^7 , c'est qu'une erreur a été commise.

En ce qui concerne le calcul de l'erreur relative, la première chose à faire est de vérifier que le α proposé est dans l'intervalle qui nous reste. Sinon c'est qu'une erreur a été commise.

De plus notre approximation est $x_n = 1.3125$, dont déjà deux chiffres sont corrects. En deux itérations de bisection en partant de 0.5, nous ne pouvons pas nous attendre à mieux.

△

35.14 Approximations de fonctions

- (1) D'habitude on n'approxime pas une fonction sur tout son domaine, mais seulement sur une partie.
- (2) Il y a le problème du choix de la classe des fonctions qui vont approximer. Nous allons travailler avec des polynômes.
- (3) Il nous faut un critère disant si une approximation est bonne ou non.

35.14.1 Critère d'interpolation

À partir de $n + 1$ abscisses points distinctes x_i , nous calculons $y_i = f(x_i)$. Il y a ce théorème qui dit qu'il existe un unique polynôme de degré (au plus) $n + 1$ passant par les points (x_i, y_i) .

Proposition-définition 35.89 (Base de Lagrange).

Étant donnés $n + 1$ valeurs distinctes x_i , l'espace des polynômes de degré n admet la base

$$L_i^{(n)}(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j} \quad (35.199)$$

pour $i = 0, \dots, n$.

Soit par exemple les valeurs de f données dans

x_i	$y_i = f(x_i)$
$x_0 = 5$	1
$x_1 = -7$	-23
$x_2 = 6$	-54
$x_3 = 0$	-954

Les polynômes de Lagrange pour ces données dépendent seulement des x_i , pas des y_i . En particulier,

$$L_0^{(3)}(x) = \frac{(x - x_1)(x - x_2)(x - x_3)}{(x_0 - x_1)(x_0 - x_2)(x_0 - x_3)}, \quad (35.200)$$

etc.

La réponse est que

$$L_0^{(3)}(x) = \frac{x^3 + 13x^2 + 42x}{660} \quad (35.201a)$$

$$L_1^{(3)}(x) = \frac{-x^3 - x^2 + 30x}{84} \quad (35.201b)$$

$$\text{etc.} \quad (35.201c)$$

Ce qu'il y a de bien avec cette base est que en posant $a_i = f(x_i)$ alors le polynôme

$$\sum_{i=0}^n a_i L_i^{(n)}(x) \quad (35.202)$$

passé par les points $(x_i, f(x_i))$. Du coup il suffit d'écrire

$$P_3(x) = L_0^{(3)}(x) - 23L_1^{(3)}(x) - 54L_2^{(3)}(x) - 954L_3^{(3)}(x) = 4x^3 + 35x^2 - 84x - 954. \quad (35.203)$$

Un inconvénient de cette base est qu'elle est complètement dépendante des points choisis. Si on ajoute un point ou qu'on en prend un à peine différent, tous les coefficients changent. Mais en pratique, ajouter des points est quelque chose qui arrive souvent parce que souvent, après avoir vu le résultat d'un polynôme d'interpolation, on veut ajouter un point pour avoir un meilleur résultat.

35.90.

Une habitude : le premier et le dernier nœud se choisissent aux extrémités de l'intervalle sur lequel nous voulons une approximation.

Le but d'une approximation est d'avoir des approximations de $f(x^*)$ pour des valeurs de x^* qui ne soit pas une des abscisses données (parce que sur ces points, le polynôme et la fonction sont égaux). Nous considérons donc

$$f(x^*) \simeq P_n(x^*). \quad (35.204)$$

Si x^* est dans l'intervalle $I = [x_{min}, x_{max}]$ alors nous disons que nous calculons f par **interpolation**. Si au contraire x^* est en dehors de cet intervalle nous parlons d'**extrapolation**.

Si x^* est pris à l'extérieur de I , alors l'erreur risque d'être très grande, surtout parce que les polynômes tendent tous vers $\pm\infty$ lorsque $x \rightarrow \pm\infty$.

Autant l'interpolation via polynômes est le plus souvent valable, il faut garder à l'esprit que les extrapolations sont souvent mauvaises si x^* est trop loin des extrémités de I .

35.14.2 Base de Newton

Après la base canonique et la base de Lagrange, nous voyons la base de Newton. Soient encore $n + 1$ points donnés du graphe de f .

Définition 35.91.

La *base de Newton* pour les abscisses x_i est l'ensemble des polynômes suivants :

$$1, (x - x_0), (x - x_0)(x - x_1), \dots, (x - x_0)(x - x_1) \dots (x - x_{n-1}) \quad (35.205)$$

Notons que ces polynômes n'utilisent pas le dernier point des x_i . Le polynôme passant par les points est

$$P_n(x) = x_0 + \sum_{i=1}^n c_i \prod_{j=0}^{n-1} (x - x_j). \quad (35.206)$$

Le calcul des c_i n'est pas absolument évident. Mais si nous ajoutons un point d'interpolation, les polynômes déjà calculés sont encore bons ; en particulier

$$P_{n+1}(x) = P_n(x) + c_{n+1} \prod_{j=0}^n (x - x_j). \quad (35.207)$$

Et cela est bien, parce que ça donne une façon de les calculer par récurrence.

Il y a plusieurs façons de calculer les c_i .

Les différences divisées sont des façons d'approximer les dérivées.

Définition 35.92.

Soient $n + 1$ nœuds x_i pour la fonction f . La *différence divisée* sont :

Ordre 0

$$f[x_i] = f(x_i) \quad (35.208)$$

Ordre 1

$$f[x_i, x_j] = \frac{f[x_i] - f[x_j]}{x_i - x_j}. \quad (35.209)$$

Ordre 2

$$f[x_i, x_j, x_k] = \frac{f[x_i, x_j] - f[x_j, x_k]}{x_i - x_k}. \quad (35.210)$$

Ordre n

$$f[x_0, \dots, x_n] = \frac{f[x_0, \dots, x_{n-1}] - f[x_1, \dots, x_n]}{x_0 - x_n}. \quad (35.211)$$

Les ordres font référence à l'ordre de dérivation qui est approximé.

Nous avons alors

$$c_i = f[x_0, \dots, x_i]. \quad (35.212)$$

Cela donne effectivement une méthode de récurrence pour trouver les coefficients c_i .

Remarque 35.93.

Pour calculer c_0 , il faut seulement calculer $f[x_0] = f(x_0)$. Mais pour calculer c_1 il faut $f[x_0]$ et $f[x_1]$. Et pour c_2 il faut $f[x_0, x_1, x_2]$ qui demande $f[x_0, x_1]$ et $f[x_0, x_1]$, qui demande etc.

Il faut donc calculer en réalité tous les $f[x_i]$ pour terminer le calcul. Par contre, pour ajouter un point, il ne faut pas tout recalculer, et même pas tout conserver en mémoire. Il faut seulement garder en mémoire la dernière diagonale.

Exemple 35.94

Soit les nœuds

x	$f(x)$
3	1
1	-3
5	2
6	4

Trouver le polynôme d'interpolation via la base de Newton.

x_i	$f[x_i]$	$f[x_i, x_j]$	$f[x_i, x_j, x_k]$
3	1	$\frac{1+3}{2} = 2$	$\frac{f[x_0, x_1] - f[x_1, x_2]}{x_0 - x_2} = \frac{2 - \frac{5}{4}}{-2} = -\frac{3}{8}$
1	-3	$\frac{-3-2}{-4} = \frac{5}{4}$	$\frac{\frac{5}{4} - 2}{-5} = \frac{3}{20}$
5	2	$\frac{2-4}{-1} = 2$	$\frac{f[x_0, x_1, x_2] - f[x_1, x_2, x_3]}{x_0 - x_3} = \frac{4}{40}$
6	4	< ++ >	< ++ >

Le polynôme d'interpolation sera

$$P_3(x) = c_0 + c_1(x - x_0) + c_2(x - x_0)(x - x_1) + c_3(x - x_0)(x - x_1)(x - x_2). \quad (35.213)$$

△ < ++ >

Un exercice typique serait de donner tout pour 3 points puis de demander le polynôme qui aurait un quatrième point.

35.14.3 Méthode des minima quadratiques

Soient $m + 1$ points connus sur le graphe de la fonction f que nous devons approximer. Au lieu d'exiger que notre approximation ne passe par tous les points, nous allons chercher une approximation qui minimise la somme des carrés des erreurs sur ces points.

Soit \mathcal{F} une classe de fonctions dans laquelle nous allons chercher l'approximation. Nous cherchons $g \in \mathcal{F}$ qui minimise

$$E(g) = \sum_{i=0}^m (f(x_i) - g(x_i))^2 \omega_i \quad (35.214)$$

où $\omega_i > 0$ est une pondération. Souvent on prend $\omega_i = 1$, mais pas toujours. La fonction E sur \mathcal{F} est la **fonction d'erreur**.

35.95.

À part dans les exercices à la main, le nombre de points est grand, du type du milliard. Il est bien entendu pas envisageable de faire passer un polynôme *exactement* par un milliard de points, parce que cela demanderait un polynôme de degré un milliard.

Plus généralement, d'un point de vue scientifique, avoir n paramètres libres pour n données expérimentales, ça ne passe pas Popper.

Afin de faire de la science qui passe Popper nous nous restreignons à une classe de fonction \mathcal{F} dont la dimension n'est pas grande : $\dim(\mathcal{F}) \ll m$. Et nous notons $\dim(\mathcal{F}) = n + 1$.

Exemple 35.96

La qualité d'une expérience peut être influencée par des paramètres extérieurs comme l'humidité, le vent, etc. Donc il est normal d'avoir des expériences moins précises que d'autres. On le pèse moins. △

Exemple 35.97

Dans un questionnaire, il se met des questions volontairement contradictoires. Si quelqu'un répond

« oui » aux deux questions, il y a une indication que la personne a répondu un peu n'importe comment, et il faut moins peser ses réponses. \triangle

Soit $g \in \mathcal{F}$, et une base $\{g_i\}_{i=0,\dots,n}$ de \mathcal{F} . Nous écrivons

$$g = a_0g_0 + a_1g_1 + \dots + a_ng_n \quad (35.215)$$

La fonction donnée E donnée en (35.214), est, à partir du moment où \mathcal{F} et une base sont choisis, une fonction des paramètres a_i que nous nommons $F(a_0, \dots, a_n)$. Il faut minimiser F , c'est-à-dire poser

$$\frac{\partial F}{\partial a_j} = 0 \quad (35.216)$$

pour $j = 0, \dots, n$. Cela sont $n + 1$ équations pour $n + 1$ inconnues. Notons que ces équations sont linéaires parce que chacun des termes est du type

$$\left(f(x_i) - \sum_{j=0}^m a_j g_j(x_i) \right)^2, \quad (35.217)$$

et lors de la dérivation par rapport à a_j , nous obtenons du degré 1.

35.14.4 Notre espace de Hilbert

Nous allons maintenant formaliser un peu tout cela. Dans [474] il est expliqué que si ω est une fonction strictement positive, alors l'espace $L_\omega^2([a, b])$ dérivé de la norme

$$\|f\|_{L_\omega^2}^2 = \int_a^b |f(x)|^2 \omega(x) dx \quad (35.218)$$

est un espace de Hilbert. Nous allons tenter le coup avec $\omega = \sum_{i=0}^m \omega_i \delta_{x_i}$ où δ_a est la distribution de Dirac¹² centrée en a .

Sur l'ensemble des fonctions $\mathbb{R} \rightarrow \mathbb{R}$ nous considérons la relation d'équivalence $f \sim g$ si $f(x_i) = g(x_i)$ pour tout $i = 0, \dots, m$. Nous notons L_ω^2 cet ensemble.

Proposition 35.98.

La formule

$$\langle f, g \rangle = \sum_{i=0}^m f(x_i)g(x_i)\omega_i \quad (35.219)$$

définit un produit scalaire sur L_ω^2 . Ce dernier devient un espace de Hilbert.

Démonstration. Pour être un produit scalaire (définition 11.5), la forme considérée doit être symétrique et strictement définie positive. La symétrie de la formule (35.219) ne fait pas de doute. Le fait que ce soit semi-défini positif non plus. Pour le strict,

$$\langle f, f \rangle = \sum_{i=0}^m |f(x_i)|^2 \omega_i. \quad (35.220)$$

Étant donné que $\omega_i > 0$ pour tout i , l'annulation de $\langle f, f \rangle$ implique l'annulation de $f(x_i)$ pour tout i . Cela signifie que f est dans la classe de 0 et donc est nul dans L_ω^2 .

En ce qui concerne la complétude, la proposition 9.43 répond à notre place, étant donné que L_ω^2 est de dimension finie. Une base est donnée par exemple par $e_i(x) = \delta_{x, x_i}$. Ici le δ est celui de Kronecker, et non celui de Dirac. \square

Lemme 35.99.

Si la classe de fonctions \mathcal{F} est un sous-espace vectoriel de L_ω^2 et si $f \in L_\omega^2$ il existe un unique élément g de \mathcal{F} minimisant la distance à f .

12. Définition 31.34.

Démonstration. Le théorème de projection (au choix 13.233 ou 26.5) nous assure l'existence et l'unicité d'un élément de \mathcal{F} minimisant la distance à $f \in L_\omega^2$. \square

35.100.

Ce lemme est gentil, mais ne nous donne pas de méthodes pour trouver ce minimum. Nous allons donc écrire explicitement un système d'équations permettant de le trouver. Si $\{g_\alpha\}$ est une base (finie) de \mathcal{F} alors nous cherchons le minimisant sous la forme $f = \sum_\alpha a_\alpha g_\alpha$.

Nous devons minimiser

$$E(g) = \sum_{i=0}^m (f(x_i) - g(x_i))^2 \omega_i = \sum_{i=0}^m (f(x_i) - \sum_\alpha a_\alpha g_\alpha(x_i))^2 \omega_i. \quad (35.221)$$

Vu que cela est maintenant plutôt une fonction des coefficients a_α que de la fonction g nous la notons $F(a_0, \dots, a_n)$. Il s'agit d'étudier le système d'équations

$$\frac{\partial F}{\partial a_\alpha} = 0. \quad (35.222)$$

Un tout petit peu de calcul mène au système

$$\sum_i \sum_\beta a_\beta \omega_i g_\alpha(x_i) g_\beta(x_i) = \sum_i \omega_i g_\alpha(x_i) f(x_i). \quad (35.223)$$

À droite nous reconnaissons $\langle f, g_\alpha \rangle$. et à gauche, $\sum_\beta a_\beta \langle g_\alpha, g_\beta \rangle$. Donc le système s'écrit

$$\sum_\beta a_\beta \langle g_\alpha, g_\beta \rangle = \langle f, g_\alpha \rangle. \quad (35.224)$$

Il y a une équation pour chaque valeur de α .

La matrice $A \in \mathbb{M}(n+1, \mathbb{R})$ donnée par $\langle g_\alpha, g_\beta \rangle$ étant strictement définie positive (c'est un produit scalaire), le système a une unique solution. Et comme cette matrice est de plus symétrique, elle est diagonalisable par le théorème spectral 11.189. Toutes ses valeurs propres sont strictement positives.

Notons pour la curiosité que si l'on considère la matrice $B \in \mathbb{M}(m \times n)$ donnée par

$$B_{ij} = \sqrt{\omega_i} g_j(x_i), \quad (35.225)$$

alors nous avons $A = B^t B$.

35.14.5 Droite de régression

La droite de régression est le cas particulier $n = 1$, c'est-à-dire un système 2×2 . Nous cherchons $P = a_0 + a_1 x$. Et la base choisie est $g_0(x) = 1$, $g_1(x) = x$. Nous avons

$$\langle g_0, g_0 \rangle = \sum_i \omega_i g_0(x_i) g_0(x_i) = \sum_i \omega_i \quad (35.226a)$$

$$\langle g_0, g_1 \rangle = \sum_i \omega_i x_i \quad (35.226b)$$

$$\langle g_1, g_1 \rangle = \sum_i \omega_i x_i^2. \quad (35.226c)$$

Donc pour approximer une fonction f il faut résoudre le système

$$\begin{pmatrix} \sum_i \omega_i & \sum_i \omega_i x_i \\ \sum_i \omega_i x_i & \sum_i \omega_i x_i^2 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} \langle f, g_0 \rangle \\ \langle f, g_1 \rangle \end{pmatrix}. \quad (35.227)$$

Pour calculer les produits $\langle f, g_\alpha \rangle$ il suffit de savoir f sur les points x_i . Et encore heureux, parce que toute la méthode est basée sur le fait que nous ne connaissons pas f ailleurs. C'est pour cela que nous avons défini L_ω^2 comme un ensemble quotient.

Exemple 35.101

Faisons la droite de régression pour les données avec tous les poids $\omega_i = 1$.

x_i	$f(x_i)$
-5	18
-3	7
1	0
3	7
4	16
6	50
8	67

Nous avons

$$\langle g, g_0 \rangle = \sum_i f(x_i) = 165 \quad (35.228a)$$

$$\langle f, f_1 \rangle = \sum_i f(x_i)x_i = 810. \quad (35.228b)$$

et donc le système

$$\begin{pmatrix} 7 & 14 \\ 14 & 160 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} 165 \\ 810 \end{pmatrix} \quad (35.229)$$

dont résolution donne la droite de régression. △

Proposition 35.102.

Si tous les poids sont identiques, alors la droite de régression passe par le barycentre des points donnés :

$$\begin{cases} x_M = \frac{1}{m+1} \sum_{i=0}^m x_i \\ y_M = \frac{1}{m+1} \sum_{i=0}^m y_i. \end{cases} \quad (35.230a)$$

$$\quad (35.230b)$$

Cela donne une vérification possible de la réponse trouvée.

Définition 35.103.

L'**erreur quadratique** est la fonction $F(a_0, \dots, a_n)$ dont il est question plus haut. Et si une solution est connue, son erreur quadratique est la valeur de F pour cette solution.

35.15 Conditionnement d'une matrice

Soit le système d'équations linéaires $Au = b$ avec la matrice inversible A ainsi que le système perturbé $(A + \Delta A)u' = (b + \Delta b)$. Nous notons $\Delta u = u' - u$ et nous voudrions pouvoir dire des choses de l'erreur relative $\frac{\|\Delta u\|}{\|u\|}$.

Exemple 35.104 ([156])

Soit la matrice

$$A = \begin{pmatrix} 10 & 7 \\ 7 & 5 \end{pmatrix} \quad (35.231)$$

et $b = \begin{pmatrix} 32 \\ 23 \end{pmatrix}$. La solution de $Au = b$ est $u = \begin{pmatrix} -1 \\ 6 \end{pmatrix}$. Si nous conservons la même matrice mais nous considérons $b = \begin{pmatrix} 32.1 \\ 22.9 \end{pmatrix}$. La solution devient $u' = \begin{pmatrix} 0.2 \\ 4.3 \end{pmatrix}$

En norme $\|\cdot\|_\infty$ nous avons ¹³

$$\frac{\|\Delta b\|}{\|b\|} = \frac{0.1}{32} = 0.003125 \quad (35.232)$$

et

$$\frac{\|\Delta u\|}{\|u\|} = \frac{1.7}{6} = 0.28. \quad (35.233)$$

Cela montre environ amplification d'un facteur 100 entre l'erreur sur b et l'erreur sur la solution.

△

Définition 35.105.

Le **conditionnement** de la matrice inversible $A \in \text{GL}(n, \mathbb{C})$ est le nombre positif

$$\text{Cond}(A) = \|A\| \|A^{-1}\|. \quad (35.234)$$

Cette dénomination sera justifié par le corollaire 35.110 parce qu'il est évident que le conditionnement d'une matrice est lié au conditionnement du problème de résolution d'un système linéaire.

Remarque 35.106.

Le conditionnement dépend de la norme choisie, mais cette dépendance est contrôlée par la proposition 12.5 qui nous indique que si le conditionnement d'une matrice est grand dans une norme, il sera grand dans une autre norme.

D'autre part, lorsque nous écrivons $\|A\|$ nous supposons toujours que $\|\cdot\|$ est une norme d'algèbre ¹⁴ et donc que nous avons toujours

$$\|AB\| \leq \|A\| \|B\|. \quad (35.235)$$

De plus nous supposons toujours avoir une norme subordonnée à une norme sur l'espace \mathbb{C}^n , de telle sorte à avoir

$$\|Au\| \leq \|A\| \|u\| \quad (35.236)$$

pour tout $u \in \mathbb{C}^n$. Voir aussi le lemme 12.17.

Proposition 35.107 ([156]).

Si A est une matrice inversible et si $\alpha \in \mathbb{C}$ nous avons :

- (1) $\text{Cond}(A) \geq 1$
- (2) $\text{Cond}(A) = \text{Cond}(A^{-1})$
- (3) $\text{Cond}(\alpha A) = \text{Cond}(A)$.

Si $Q \in \text{O}(n)$ alors

- (1) Nous avons $\text{Cond}_2(Q) = 1$ où Cond_2 est le conditionnement pour la norme $\|\cdot\|_2$.
- (2) Nous avons aussi

$$\text{Cond}_2(A) = \text{Cond}_2(AQ) = \text{Cond}_2(QA). \quad (35.237)$$

Démonstration. Nous savons que $\text{Cond}(\mathbb{1}) = 1$ et donc

$$1 = \|\mathbb{1}\| \leq \|A\| \|A^{-1}\| \quad (35.238)$$

parce que la norme utilisée est une norme matricielle.

Les deux autres formules sont évidentes à partir du fait que la définition du conditionnement de A est symétrique entre A et A^{-1} .

13. La proposition 12.5(3) montre que si nous voulions des estimations en norme $\|\cdot\|_2$, il y aurait au maximum un facteur $\sqrt{2}$ par-ci par là.

14. Définition 12.14.

En ce qui concerne les formules relatives à la matrice orthogonale Q nous savons par la proposition 11.83(3) qu'une matrice orthogonale est une bijection de l'ensemble $\{x \in \mathbb{R}^n \text{ tel que } \|x\| = 1\}$. Par conséquent

$$\|AQ\| = \sup_{x \text{ tel que } \|x\|=1} \|AQx\| = \sup_{Q^{-1}x \text{ tel que } \|x\|=1} \|AQQ^{-1}x\| = \|A\|. \quad (35.239)$$

Donc $\|AQ\| = \|A\|$. Les assertions s'ensuivent immédiatement en remarquant que Q^{-1} est également orthogonale. \square

Soit une matrice inversible $A \in \text{GL}(n, \mathbb{C})$. La matrice A^*A est hermitienne¹⁵ et le théorème 11.183 nous assure que ses valeurs propres sont réelles. Par la remarque 11.184, ses valeurs propres sont même positives.

Proposition 35.108 ([156]).

Soit une matrice inversible $A \in \text{GL}(n, \mathbb{C})$, et $\mu_1 \leq \dots \leq \mu_n$ les valeurs propres de A^*A . Alors nous avons la formule

$$\text{Cond}_2(A) = \sqrt{\frac{\mu_n}{\mu_1}}. \quad (35.240)$$

Démonstration. Par le théorème 12.27, la norme de A est liée au au rayon spectral de A^*A par

$$\|A\|_2 = \sqrt{\rho(A^*A)} = \sqrt{\mu_n}. \quad (35.241)$$

Vu que le spectre de AA^* est le même que celui de A^*A (lemme 11.187) nous avons aussi

$$\|A^{-1}\|_2 = \sqrt{\rho((A^{-1})^*A^{-1})} = \sqrt{\rho((A^*A)^{-1})} = \frac{1}{\sqrt{\mu_1}} \quad (35.242)$$

parce que la plus grande valeur propre de $(A^*A)^{-1}$ est l'inverse de la plus petite de A^*A .

Ces deux calculs étant,

$$\text{Cond}_2(A) = \|A\|_2 \|A^{-1}\|_2 = \sqrt{\frac{\mu_n}{\mu_1}}. \quad (35.243)$$

\square

Problèmes et choses à faire

À mon avis ce qui est dans la proposition 18.109 est le conditionnement de la matrice ou sa racine carrée ou un truc du genre. Il faut voir le lien entre les valeurs propres de A et celles de AA^* .

35.15.1 Perturbation du vecteur

Proposition 35.109 (Système linéaire : perturbation du vecteur[156]).

Soit une matrice inversible A et les systèmes d'équations linéaires

$$Au = b \quad (35.244a)$$

$$Au' = b'. \quad (35.244b)$$

En notant $\Delta u = u' - u$ et $\Delta b = b' - b$ nous avons

$$\frac{\|\Delta u\|}{\|u\|} \leq \text{Cond}(A) \frac{\|\Delta b\|}{\|b\|}. \quad (35.245)$$

Démonstration. En soustrayant les équations (35.244) nous avons $\Delta b = A\Delta u$, et donc $\Delta u = A^{-1}\Delta b$. D'une part nous avons alors

$$\|\Delta u\| \leq \|A^{-1}\| \|\Delta b\|. \quad (35.246)$$

15. Définition 11.76.

Et d'autre part, $\|b\| \leq \|A\| \|u\|$, ce qui donne

$$\frac{\|b\|}{\|A\|} \leq \|u\|. \quad (35.247)$$

En mettant les deux ensemble,

$$\frac{\|\Delta u\|}{\|u\|} \leq \frac{\|A^{-1}\| \|\Delta b\|}{\|b\|} \|A\| = \text{Cond}(A) \frac{\|\Delta b\|}{\|b\|}. \quad (35.248)$$

□

Le corollaire suivant justifie le nom « conditionnement » au conditionnement d'une matrice.

Corollaire 35.110.

Soit $A \in \text{GL}(n, \mathbb{C})$ fixée et le problème de résoudre $Au = b$, c'est-à-dire la fonction

$$F(u, b) = Au - b. \quad (35.249)$$

(1) Ce problème est stable pour toute valeur de b .

(2) Nous avons une majoration pour le conditionnement relatif¹⁶ :

$$K_{rel}(\eta, b_0) \leq \text{Cond}(A). \quad (35.250)$$

Démonstration. **Stabilité** Vu que A est inversible, il existe une solution unique à tout système de la forme $Au = b'$. De plus $u(b) = A^{-1}b$, donc

$$\|u(b) - u(b_0)\| = \|A^{-1}(b - b_0)\| \leq \|A^{-1}\| \|b - b_0\|, \quad (35.251)$$

de telle sorte que la condition 35.25(2) fonctionne avec $K = \|A^{-1}\|$.

Conditionnement En partant de la définition 35.47, et en utilisant la majoration de la proposition 35.109 sous la forme

$$\|u(b) - u(b_0)\| \leq \text{Cond}(A) \|u(b_0)\| \frac{\|\Delta b\|}{\|b_0\|}, \quad (35.252)$$

nous obtenons :

$$K_{rel}(b_0, \eta) = K_{abs}(b_0, \eta) \frac{\|b_0\|}{\|u(b_0)\|} \quad (35.253a)$$

$$= \sup_{\|b-b_0\| \leq \eta} \frac{\|u(b) - u(b_0)\|}{\|b - b_0\|} \frac{\|b_0\|}{\|u(b_0)\|} \quad (35.253b)$$

$$\leq \sup_b \text{Cond}(A) \frac{\|u(b_0)\|}{\|b_0\|} \|\Delta b\| \frac{1}{\|b - b_0\|} \frac{\|b_0\|}{\|u(b_0)\|} \quad (35.253c)$$

$$= \text{Cond}(A). \quad (35.253d)$$

□

Remarque 35.111.

La notion de conditionnement relatif dépend aussi de la norme choisie. Dans la formule (35.250) il faut prendre le conditionnement $\text{Cond}(A)$ pour la norme dans laquelle le K_{rel} est écrit. Encore une fois, toutes les normes étant équivalentes, cette majoration est à constante près bonne pour toutes les normes. Si la dimension est très grande, cette constante peut par contre être grande.

16. Si vous doutez de la norme à prendre, lisez la remarque 35.111

35.15.2 Perturbation de la matrice

Proposition 35.112 (Système linéaire : perturbation de la matrice[156]).

Soient les systèmes linéaires

$$Au = b \quad (35.254a)$$

$$A'u' = b \quad (35.254b)$$

avec A et A' inversibles. Nous notons $\Delta A = A' - A$. Alors

(1)

$$\frac{\|\Delta u\|}{\|u'\|} \leq \text{Cond}(A) \frac{\|\Delta A\|}{\|A\|} \quad (35.255)$$

(2)

$$\frac{\|\Delta u\|}{\|u\|} \leq \text{Cond}(A) \frac{\|\Delta A\|}{\|A\|} (1 + \alpha(\|\Delta A\|)) \quad (35.256)$$

où $\lim_{x \rightarrow 0} \alpha(x) = 0$.

Démonstration. D'abord nous avons

$$0 = Au' - Au \quad (35.257a)$$

$$= (A' - A)u' - Au' - Au \quad (35.257b)$$

$$= \Delta Au' + A\Delta u. \quad (35.257c)$$

Par conséquent, $\Delta u = -A^{-1}(\Delta A)u'$ et

$$\|\Delta u\| \leq \|A^{-1}\| \|\Delta A\| \|u'\|. \quad (35.258)$$

Donc

$$\frac{\|\Delta u\|}{\|u'\|} \leq \|A^{-1}\| \|A\| \frac{\|\Delta A\|}{\|A\|} = \text{Cond}(A) \frac{\|\Delta A\|}{\|A\|}. \quad (35.259)$$

Cela est (1).

Pour l'autre inégalité, nous avons $A' = A + \Delta A$ et donc

$$\|A'^{-1}\| = \|(A + \Delta A)^{-1}\| \quad (35.260)$$

Nous repartons alors de (35.258) en changeant le rôle de A et A' (et donc aussi de u et u'). Ce changement étant, $\|\Delta u\|$ et $\|\Delta A\|$ ne changent pas. Nous avons :

$$\frac{\|\Delta u\|}{\|u\|} \leq \|A'^{-1}\| \|\Delta A\| \quad (35.261a)$$

$$= \|(A + \Delta A)^{-1}\| \|\Delta A\| \frac{\text{Cond}(A)}{\|A\| \|A^{-1}\|} \quad (35.261b)$$

$$= \frac{\|(A + \Delta A)^{-1}\|}{\|A^{-1}\|} \frac{\|\Delta A\|}{\|A\|} \text{Cond}(A). \quad (35.261c)$$

Il reste à voir que

$$\lim_{\|\Delta A\| \rightarrow 0} \frac{\|(A + \Delta A)^{-1}\|}{\|A^{-1}\|} = 1, \quad (35.262)$$

ou autrement dit que

$$\lim_{A \rightarrow A'} \frac{\|A'^{-1}\|}{\|A^{-1}\|} = 1 \quad (35.263)$$

où la limite est celle dans $\text{GL}(n, \mathbb{C})$. Par définition de la topologie, la norme est continue (quelle qu'elle soit par l'équivalence de norme 12.6). Par le théorème 12.159, l'application $A \mapsto A^{-1}$ est également continue et commute donc avec la limite. Nous avons donc

$$\lim_{A' \rightarrow A} \|A'^{-1}\| = \|(\lim_{A' \rightarrow A} A')^{-1}\| = \|A^{-1}\|. \quad (35.264)$$

Donc la limite du quotient (35.263) est bien 1. \square

35.16 Système linéaires (généralités)

Soit un système d'équations linéaires $Ax = b$ avec $A \in \mathbb{M}(n, \mathbb{R})$. Le problème est évidemment de savoir s'il existe une unique solution x et de la déterminer. Nous supposons l'existence et l'unicité. C'est-à-dire que les conditions équivalentes¹⁷ sont vérifiées :

- (1) A est inversible, c'est-à-dire qu'il existe une matrice notée A^{-1} telle que $AA^{-1} = A^{-1}A = \mathbb{1}$.
- (2) $\det(A) \neq 0$.

Note : si nous avons un système pas carré du type $Bx = v$ avec $B \in \mathbb{M}(n \times m)$ alors nous pouvons nous ramener à un système carré en écrivant

$$B^t Bx = B^t v. \quad (35.265)$$

Mais attention : bien que $B^t B$ soit symétrique et semi-définie positive, certaines valeurs propres peuvent être nulles.

35.113.

Deux choses générales en calcul numérique :

- (1) On ne calcule pas l'inverse d'une matrice.
- (2) On ne calcule même pas son déterminant.

Par conséquent nous ne faisons pas $x = A^{-1}v$.

Il faut garder en tête le fait que dans la pratique, la matrice A possède des millions de lignes et colonnes, si pas pire. Pour une matrice de taille de l'ordre du million, il y a 1000 milliards d'entrées. Si on compte 32 bits par nombre (précision simple, définition 35.8), c'est-à-dire 4 octets, il faut 4000 giga-octets pour enregistrer la matrice. Même pour la mémoire actuellement disponible, ce n'est pas rien. Surtout que souvent, la précision simple n'est pas utilisée, mais la précision double, ce qui donne 8000 giga pour enregistrer la matrice.

Heureusement, dans la majorité des cas pratiques, les matrices géantes qui apparaissent sont pleines de zéros.

Définition 35.114.

Une matrice est **creuse** si elle possède beaucoup de zéros. Une matrice non creuse est dite **dense**.

Notons que lorsqu'on parle de matrice comprenant beaucoup de « zéros », nous pensons à des éléments très petits, et non de vrai zéros.

Les matrices creuses ne sont pas mémorisées entièrement, mais plutôt comme un dictionnaire (i, j, v) qui donne la valeur v de A_{ij} .

Définition 35.115.

Une matrice est de « grande dimension » si elle ne peut pas être mise en mémoire sur un ordinateur donné. Sur certains ordinateurs, ça commence à 5000 inconnues. Mais sur des plus forts, on peut aller jusqu'au million ou le milliard.

Si la matrice est de petite dimension, il est possible d'utiliser des méthodes dites « directes ». Sinon, il faudra utiliser des méthodes itératives.

35.16.1 Les méthodes directes

Une méthode directe consiste à successivement transformer un système $A^{(0)}x = b^{(0)}$ en de nouveaux systèmes $A^{(i)}x = b^{(i)}$ dont la solution est identique jusqu'à obtenir un système $A^{(n-1)}x = b^{(n-1)}$ qui est à résolution immédiate.

L'avantage d'une méthode directe est qu'elle fournit une réponse exacte, pour autant que les calculs intermédiaires soient bien faits (ce qui n'est pas le cas sur un ordinateur).

Une méthode directe fonctionne en général avec un nombre de pas fixés par la taille du système. Par exemple pour un système $n \times n$, la méthode de Gauss demande exactement n pas, et il n'y a

17. L'équivalence est la proposition 11.52(2).

pas moyen de faire mieux. Or chaque pas demande de recalculer tous les éléments de la matrice. Encore une fois, si la matrice a une taille de l'ordre du milliard, cela fait 10^{18} éléments à recalculer un milliard de fois (sans compter les éléments du vecteur b). Infaisable.

Souvent une méthode directe passe par une factorisation $A = BC$ avec $B, C \in \mathbb{M}(n \times n)$.

Quelques types de matrices dont la résolution est immédiate :

- Matrice diagonale.
- Matrice orthogonale parce que si A est orthogonale alors $Ax = v$ se résout par $x = A^t v$ qui n'est pas particulièrement lourd à faire numériquement.
- Matrice triangulaire.

Remarque 35.116.

Pour une matrice diagonale, le déterminant et l'inverse sont faciles. Mais également pour la triangulaire. Pour une matrice triangulaire, le déterminant est le produit des éléments diagonaux, et il se fait qu'il y a une algorithmes facile pour calculer l'inverse.

Donc en fait les matrices à résolution immédiates sont des matrices pour lesquelles l'inverse et le déterminant sont facile à calculer.

35.16.2 Méthodes itératives

Si la matrice est trop grande, il n'est pas possible de faire des manipulations de matrices à chaque itération.

En général, les méthodes itératives ne convergent pas toujours. Mais lorsqu'une méthode converge, c'est une propriété de la matrice, et donc la convergence aura lieu pour tout vecteur de départ x_0 . Cela est très différent du cas des équations non linéaires type Newton pour lesquelles la convergence peut fortement dépendre du point de départ.

35.17 Système linéaires (méthodes directes)

Les matrices que nous sommes autorisés à inverser sont les matrices

- orthogonales : l'inverse est la transposée
- diagonales : l'inverse est diagonale avec les inverses sur la diagonale
- triangulaires : nous en parlons maintenant.

35.17.1 Inversion de matrice triangulaire

Si T est une matrice triangulaire (mettons supérieure pour fixer les idées), il est possible d'en calculer l'inverse sans trop d'efforts. Notons B la matrice inverse que nous allons construire ligne par ligne. Vu que $BT = \mathbb{1}$ nous avons

$$\delta_{1j} = \sum_{k=1}^n B_{1k} T_{kj} = \sum_{k=1}^j B_{1k} T_{kj} \quad (35.266)$$

parce que $T_{kj} = 0$ pour $k > j$. Donc nous pouvons calculer les éléments B_{1j} un par un parce que chacun ne dépend que des précédents. Le même procédé fonctionne pour les autres lignes :

$$\delta_{ij} = \sum_{k=1}^j B_{ik} T_{kj}. \quad (35.267)$$

Et tu notes que le calcul peut être parallélisé : le calcul de la ligne numéro j ne dépend pas du résultat des autres lignes.

35.17.2 Transformation gaussienne

Définition 35.117 (Transformation gaussienne[475]).

Soit $x \in \mathbb{R}^n$ avec $x_k \neq 0$. La k^{e} **transformation gaussienne** pour x est la matrice

$$M_k(x) = \mathbb{1} - T_k(x) \quad (35.268)$$

où $T_k(x)$ est la matrice unité à qui on a ajouté le vecteur

$$\tau_k(x) = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ x_{k+1}/x_k \\ \vdots \\ x_n/x_k \end{pmatrix} \quad (35.269)$$

à la k^{e} colonne.

Autrement dit, la matrice $M_k(x)$ est la matrice

$$M_k(x) = \begin{pmatrix} 1 & & & & & \\ 0 & \ddots & & & & \\ \vdots & 0 & & 1 & & \\ \vdots & \vdots & -x_{k+1}/x_k & 1 & & \\ \vdots & \vdots & \vdots & 0 & \ddots & \\ 0 & 0 & -x_n/x_k & 0 & & 1 \end{pmatrix} \quad (35.270)$$

En coordonnées nous avons

$$(M_k(x))_{ij} = \delta_{ij} - \tau_k(x)_i \delta_{kj}. \quad (35.271)$$

35.118.

Les matrices de transformation gaussienne sont des matrices triangulaires de diagonale unitaire (c'est-à-dire avec des 1 sur la diagonale).

Lemme 35.119.

Si $x \in \mathbb{R}^n$ alors nous avons

$$M_k(x)x = \begin{pmatrix} x_1 \\ \vdots \\ x_k \\ 0 \\ \vdots \\ 0 \end{pmatrix}. \quad (35.272)$$

Démonstration. Nous avons

$$(M_k(x)x)_i = \sum_l M_k(x)_{il} x_l = \sum_l (\delta_{il} - \tau_k(x)_i \delta_{kl}) x_l = x_i - \tau_k(x)_i x_k. \quad (35.273)$$

Si $i \leq k$ nous avons $\tau_k(x)_i = 0$ et donc $(M_k(x)x)_i = x_i$. Si par contre $i \geq k+1$ alors $\tau_k(x)_i = \frac{x_i}{x_k}$ et alors $(M_k(x)x)_i = 0$. \square

Lemme 35.120.

Si $y \in \mathbb{R}^n$ vérifie $y_i = 0$ pour $i > k$ alors $M_{k+1}(x)y = y$.

Démonstration. C'est une simple vérification :

$$(M_{k+1}(x)y)_i = \sum_l (\delta_{il} - \tau_{k+1}(x)_i \delta_{k+1,l}) y_l = y_i - \tau_{k+1}(x)_i y_{k+1}. \quad (35.274)$$

Mais comme $y_{k+1} = 0$ il nous reste automatiquement y_i . \square

Le sens de ce lemme est si un vecteur est déjà « gaussiannisé » au niveau k , alors en lui appliquant une transformation gaussienne de niveau plus élevé que k , il ne change pas. Ce fait est important parce qu'il assure que lorsque l'on avance dans le processus de Gauss, chaque étape ne détruit pas les précédentes.

Le lemme suivant nous indique que l'inverse d'une matrice de transformation gaussienne est facile à calculer ¹⁸.

Lemme 35.121.

L'inverse de la transformation gaussienne

$$M_k(x)_{ij} = \delta_{ij} - \tau_k(x)_i \delta_{kj}. \quad (35.275)$$

est la matrice donnée par

$$M_k(x)_{ij}^{-1} = \delta_{ij} + \tau_k(x)_i \delta_{kj}. \quad (35.276)$$

Autrement dit, il suffit de changer le signe de la partie non diagonale.

Démonstration. Il s'agit d'une simple vérification, utilisant le produit matriciel explicite, et en remarquant que $\tau_k(x)_k = 0$ pour tout k . \square

35.17.3 Méthode de Gauss pour résoudre des systèmes d'équations linéaires

Pour résoudre un système d'équations linéaires, on procède comme suit :

- (1) Écrire le système sous forme matricielle.

$$\text{p.ex. } \begin{cases} 2x + 3y = 5 \\ x + 2y = 4 \end{cases} \Leftrightarrow \left(\begin{array}{cc|c} 2 & 3 & 5 \\ 1 & 2 & 4 \end{array} \right)$$

- (2) Se ramener à une matrice avec un maximum de 0 dans la partie de gauche en utilisant les transformations admissibles :

- (2a) Remplacer une ligne par elle-même + un multiple d'une autre ;

$$\text{p.ex. } \left(\begin{array}{cc|c} 2 & 3 & 5 \\ 1 & 2 & 4 \end{array} \right) \xrightarrow{L_1 - 2L_2 \rightarrow L'_1} \left(\begin{array}{cc|c} 0 & -1 & -3 \\ 1 & 2 & 4 \end{array} \right)$$

- (2b) Remplacer une ligne par un multiple d'elle-même ;

$$\text{p.ex. } \left(\begin{array}{cc|c} 0 & -1 & -3 \\ 1 & 2 & 4 \end{array} \right) \xrightarrow{-L_1 \rightarrow L'_1} \left(\begin{array}{cc|c} 0 & 1 & 3 \\ 1 & 2 & 4 \end{array} \right)$$

- (2c) Permuter des lignes.

$$\text{p.ex. } \left(\begin{array}{cc|c} 0 & 1 & 3 \\ 1 & 0 & -2 \end{array} \right) \xrightarrow{L_1 \leftrightarrow L'_2 \text{ et } L_2 \leftrightarrow L'_1} \left(\begin{array}{cc|c} 1 & 0 & -2 \\ 0 & 1 & 3 \end{array} \right)$$

- (3) Retransformer la matrice obtenue en système d'équations.

$$\text{p.ex. } \left(\begin{array}{cc|c} 1 & 0 & -2 \\ 0 & 1 & 3 \end{array} \right) \Leftrightarrow \begin{cases} x = -2 \\ y = 3 \end{cases}$$

18. Elle rentre d'ailleurs dans la catégorie des matrices triangulaires dont nous avons déjà discuté l'inverse.

Remarques :

— Si on obtient une ligne de zéros, on peut l'enlever :

$$\text{p.ex. } \left(\begin{array}{ccc|c} 3 & 4 & -2 & 2 \\ 4 & -1 & 3 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right) \Leftrightarrow \left(\begin{array}{ccc|c} 3 & 4 & -2 & 2 \\ 4 & -1 & 3 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right)$$

— Si on obtient une ligne de zéros suivie d'un nombre non-nul, le système d'équations n'a pas de solution :

$$\text{p.ex. } \left(\begin{array}{ccc|c} 3 & 4 & -2 & 2 \\ 4 & -1 & 3 & 0 \\ 0 & 0 & 0 & 7 \end{array} \right) \Leftrightarrow \begin{cases} \dots \\ \dots \\ 0x + 0y + 0z = 7 \end{cases} \Rightarrow \text{Impossible}$$

— Si on a moins d'équations que d'inconnues, alors il y a une infinité de solutions qui dépendent d'un ou plusieurs paramètres :

$$\text{p.ex. } \left(\begin{array}{ccc|c} 1 & 0 & -2 & 2 \\ 0 & 1 & 3 & 0 \end{array} \right) \Leftrightarrow \begin{cases} x - 2z = 2 \\ y + 3z = 0 \end{cases} \Leftrightarrow \begin{cases} x = 2 + 2\lambda \\ y = -3\lambda \\ z = \lambda \end{cases}$$

35.17.4 Méthode de Gauss sans pivot (décomposition LU)

La méthode de Gauss est encore utilisée aujourd'hui dans les vrais problèmes.

La méthode de Gauss est souvent aussi appelée méthode « LU » qui va décomposer $A = LU$ où L est triangulaire inférieure et U est triangulaire supérieure. La décomposition est même plus précise que cela : on demande que L ait seulement des 1 sur la diagonale.

Si A est une matrice nous notons

$$\Delta_k(A) = (A_{ij})_{1 \leq i, j \leq k} \quad (35.277)$$

la matrice tronquée dont nous ne gardons que le carré $k \times k$ en haut à gauche.

Lemme 35.122.

Soit S une matrice triangulaire inférieure. Soient également A et B telles que $B = SA$. Alors

$$\Delta_k(B) = \Delta_k(S)\Delta_k(A). \quad (35.278)$$

Démonstration. En effet nous avons

$$\Delta_k(B)_{ij} = \sum_{l=1}^n S_{il}A_{lj}. \quad (35.279)$$

Dans la somme sur l il ne reste que les termes $l \leq i$. Mais dans le calcul des éléments de matrice $\Delta_k(B)_{ij}$, nous avons évidemment $i, j \leq k$. Donc $l \leq i \leq k$. Les seuls éléments de matrice de A qui sont utilisés dans la somme (35.279) sont les éléments A_{lj} avec $l, j \leq k$.

Nous pouvons donc limiter la somme à $l = k$ au lieu de $l = n$ et écrire $\Delta_k(A)_{lj}$ au lieu de A_{lj} .

Même chose en ce qui concerne S . À partir du moment où l est limité à k , les éléments S_{il} et $\Delta_k(S)_{il}$ sont les mêmes. \square

Théorème 35.123 (Décomposition LU[476, 1]).

Soit une matrice A inversible telles que $\det(\Delta_k(A)) \neq 0$ pour tout k . Alors il existe un unique couple de matrices (L, U) telles que

- U soit triangulaire supérieure
- L soit triangulaire inférieure, de diagonale unité

— $A = LU$.

De plus pour tout $k \leq n$ nous avons

$$\Delta_k(A) = \Delta_k(L)\Delta_k(U). \quad (35.280)$$

Démonstration. Nous allons prouver par récurrence le fait suivant : pour tout $1 \leq k \leq n - 1$ il existe des matrices E_i ($i = 1, \dots, k$) telles que en posant

$$A_k = E_k \dots E_1 A, \quad (35.281)$$

- E_j est une transformation gaussienne pour la j^{e} colonne,
- pour tout $j \leq k$, $A_{ij} = 0$ dès que $i > j$. Autrement dit la matrice A_k est triangulaire supérieure jusqu'à y compris la $(k + 1)^{\text{e}}$ colonne (laquelle est quelconque). Exemple pour fixer les idées : pour une matrice $A \in \mathbb{M}(4 \times 4)$, la matrice A_2 doit avoir la forme

$$A_2 = E_2 E_1 A = \begin{pmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & 0 & \textcircled{*} & * \\ 0 & 0 & * & * \end{pmatrix} \quad (35.282)$$

où les éléments notés * sont a priori non nuls,

- l'élément de matrice $(A_k)_{k+1, k+1}$ est non nul (celui entouré dans l'exemple).

La chose un peu triste dans cette démonstration est que l'initialisation va être très ressemblante au pas de récurrence.

Initialisation : $k = 1$ Vu que $\Delta_1(A)$ est inversible, l'élément A_{11} est non nul. Il existe donc une transformation gaussienne E_1 telle que la première colonne de la matrice $A_1 = E_1 A$ soit nul sauf la première composante. En particulier $(A_1)_{21} = 0$.

Par le lemme 35.122, nous avons $\Delta_2(A_1) = \Delta_2(E_1)\Delta_2(A)$, donc¹⁹

$$\det(\Delta_2(A_1)) = \det(\Delta_2(E_1)) \det(\Delta_2(A)). \quad (35.283)$$

Étant donnée la forme (35.270), toutes les matrices du type $\Delta_k(E_i)$ ont un déterminant unité, et par hypothèse $\Delta_2(A)$ est inversible, donc de déterminant non nul. Par conséquent $\det(\Delta_2(A_1)) \neq 0$. Mais comme ce déterminant est le produit des éléments diagonaux (c'est une matrice triangulaire), ces derniers ne sont pas nuls. Finalement, $(A_1)_{22} \neq 0$.

Le pas de récurrence Nous supposons avoir $A_k = E_k \dots E_1 A$ avec $(A_k)_{k+1, k+1} \neq 0$. Alors il existe une transformation gaussienne E_{k+1} de la $(k + 1)^{\text{e}}$ colonne telle que $A_{k+1} = E_{k+1} A_k$ soit une matrice dont la $(k + 1)^{\text{e}}$ colonne n'ait que des zéros en dessous de la $(k + 1)^{\text{e}}$ position. Vu le lemme 35.120, cette transformation n'affecte pas les colonnes précédentes.

La matrice A_{k+1} est donc triangulaire supérieure jusqu'à la $(k + 2)^{\text{e}}$ colonne.

Vu que le produit $E_k \dots E_1$ est une matrice triangulaire inférieure, le lemme 35.122 fonctionne encore et nous avons

$$\Delta_{k+1}(A_k) = \Delta_{k+1}(E_k \dots E_1)\Delta_{k+1}(A). \quad (35.284)$$

En ce qui concerne les déterminants, par hypothèse, nous avons $\det(\Delta_{k+1}(A)) \neq 0$ ainsi que $\det(\Delta_{k+1}(E_k \dots E_1)) = 1$. Donc

$$\det(\Delta_{k+1}(A_k)) \neq 0. \quad (35.285)$$

Cette matrice étant triangulaire, ses éléments diagonaux sont non nuls et nous avons $(A_k)_{k+1, k+1} \neq 0$.

19. Le déterminant est multiplicatif, proposition 11.52(1).

En poussant la récurrence jusqu'au bout, la matrice

$$A_{n-1} = E_{n-1} \dots E_n A \quad (35.286)$$

est triangulaire supérieure.

Nous posons alors $L = (E_{n-1} \dots E_n)^{-1}$ et $U = A_{n-1}$. Cela prouve l'existence parce que

$$A = (E_{n-1} \dots E_1)^{-1} A_{n_1}. \quad (35.287)$$

Encore une fois, le lemme 35.122 nous donne

$$\Delta_k(A) = \Delta_k\left((E_{n_1} \dots E_1)^{-1}\right) \Delta_k(A_{n-1}), \quad (35.288)$$

ou encore $\Delta_k(A) = \Delta_k(L) \Delta_k(U)$.

En ce qui concerne l'unicité, si $A = L_1 U_1 = L_2 U_2$ alors $L_2^{-1} L_1 = U_2 U_1^{-1}$. Vu qu'à gauche nous avons une matrice triangulaire inférieure et que à droite nous avons une triangulaire inférieure, nous savons que les deux membres représentent une matrice diagonale. Mais à gauche, la diagonale est unitaire. Donc les deux membres représentent la matrice unité. \square

35.124.

En pratique, pour résoudre $Ax = b$, il faut seulement appliquer les transformations gaussiennes à la matrice élargie $(A|b)$ pour finir sur un système du type

$$Ux = b' \quad (35.289)$$

qui est immédiatement soluble. Autrement dit, en effectuant les annulations de colonnes, la matrice U est « gratuite ».

Il n'est pas indispensable de calculer la matrice L qui, elle, demande à chaque étape de se souvenir de la matrice E_i utilisée. S'il faut résoudre plusieurs systèmes $Ax_i = b_i$, nous pouvons encore travailler avec la matrice encore plus élargie $(A|b_1 \dots b_m)$.

Si par contre nous ne connaissons pas à l'avance l'ensemble des vecteurs b avec lesquels il faudra résoudre le système, il est bon de calculer la décomposition $A = LU$ in extenso, c'est-à-dire de garder une trace des matrices L et U séparément. Dans ce cas, résoudre $Ax = b$ revient à résoudre $Ly = b$, et ensuite $Ux = y$. Ce sont deux systèmes de résolution directe parce que les matrices sont triangulaires.

35.125.

Le fait que

$$\Delta_k(A) = \Delta_k(L) \Delta_k(U) \quad (35.290)$$

nous dit que si après avoir calculer L et U nous remarquons que le système est un peu plus petit ou un peu plus grand que prévu, tout le travail n'est pas perdu. En particulier si le système est plus petit que prévu, l'adaptation de L et U est immédiate.

Notons que U et L sont inversibles, et que $\det(L) = 1$. Donc $\det(U) = \det(A)$.

Exemple 35.126

Pour travailler la méthode de Gauss pour le système $Ax = b$, nous introduisons la matrice un peu augmentée $(A|b)$. Nous faisons un exemple. Soit à résoudre

$$\begin{pmatrix} 2 & 1 & 3 \\ 4 & 3 & 10 \\ -2 & 1 & 73 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 11 \\ 28 \\ 3 \end{pmatrix}. \quad (35.291)$$

Nous introduisons la matrice augmentée

$$(A|b)^{(0)} = \begin{pmatrix} 2 & 1 & 3 & 11 \\ 4 & 3 & 10 & 28 \\ -2 & 1 & 7 & 3 \end{pmatrix}. \quad (35.292)$$

Le premier pas consiste à annuler tous les éléments sous la diagonale de la première colonne. Autrement dit, nous prenons le 2 comme pivot. Nous introduisons les multiplicateurs $l_{ij} = \frac{A_{ij}}{A_{i1}}$. La nouvelle matrice est :

$$(A|b)^{(1)} = \begin{pmatrix} 2 & 1 & 3 & 11 \\ 0 & 1 & 4 & 6 \\ 0 & 2 & 10 & 14 \end{pmatrix} \quad (35.293)$$

où nous avons utilisé les multiplicateurs $l_{21} = 2$, $l_{31} = -1$.

Et la matrice suivante est :

$$(A|b)^{(2)} = \begin{pmatrix} 2 & 1 & 3 & 11 \\ 0 & 1 & 4 & 6 \\ 0 & 0 & 2 & 2 \end{pmatrix} \quad (35.294)$$

où nous avons utilisé le multiplicateur $l_{32} = 2$.

Cela est un système de résolution immédiate :

$$\begin{cases} 2x + y + 3z = 11 & (35.295a) \\ y + 4z = 6 & (35.295b) \\ 2z = 2. & (35.295c) \end{cases}$$

La troisième donne $z = 1$. Ensuite $y + 4 = 6$, donc $y = 2$. Et la première donne : $2x + 2 + 3 = 11$, c'est-à-dire $2x = 6$, enfin : $x = 3$.

Solution : $(x, y, z) = (3, 2, 1)$.

Nous notons surtout que dans $(A|b)^{(2)}$ nous avons une matrice triangulaire supérieure. Où est la matrice triangulaire inférieure ? En réalité la matrice L est la matrice des multiplicateurs :

$$L = \begin{pmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & 2 & 1 \end{pmatrix}. \quad (35.296)$$

△

Le problème de cette méthode est que faisant ainsi nous risquons d'avoir un zéro sur un des pivots. Par exemple tomber sur

$$(A|b) = \begin{pmatrix} 2 & 1 & 3 & 11 \\ 0 & 0 & 4 & 6 \\ 0 & 2 & 10 & 14 \end{pmatrix}. \quad (35.297)$$

Le zéro sur la deuxième ligne nous ennuie si nous voulons tout faire dans l'ordre. Mais notons qu'en échangeant les deux dernières lignes, tout va bien : le système donné par

$$(A|b) = \begin{pmatrix} 2 & 1 & 3 & 11 \\ 0 & 2 & 10 & 14 \\ 0 & 0 & 4 & 6 \end{pmatrix} \quad (35.298)$$

fonctionne très bien. Et même tellement bien qu'il est de résolution immédiate, dans ce cas.

Un autre problème est que si un des pivots est 10^{-14} , le multiplicateur sera de l'ordre 10^{14} , qui est mal représenté en mémoire. Il est donc bon de prendre les pivots le plus grand possible. Si le pivot est le plus grand nombre en valeur absolue d'une colonne, alors les nombres x_{k+i}/x_k qui entrent dans la matrice de transformation gaussienne sont des nombres dans $[-1, 1]$ qui sont bien représentés en mémoire.

Tout cela nous incite à développer une méthode de Gauss qui permet de tenir une trace des permutations.

35.17.5 Matrice de permutation élémentaires

Définition 35.127.

Une *matrice de permutation élémentaire* est une matrice obtenue en permutant deux lignes de la matrice identité. Nous notons P_{ij} la matrice obtenue en inversant les lignes i et j de la matrice identité.

Exemple 35.128

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \rightarrow \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} = P_{12}. \quad (35.299)$$

△

Lemme 35.129.

La matrice $P_{ij}A$ est la matrice A avec ses lignes i et j inversées.

Démonstration. Il suffit d'écrire

$$(P_{ij}A)_{kl} = \sum_m (P_{ij})_{km} A_{ml} \quad (35.300)$$

et de faire trois cas selon que $k = i$, $k = j$ ou k différent de i et j . Si $k = i$ alors $(P_{ij})_{im} = \delta_{mj}$ et si k est différent de i et j alors $(P_{ij})_{mk} = \delta_{km}$ (troisième cas similaire au premier). □

Et la matrice AP_{12} est la A avec ses deux premières colonnes échangées.

Avec ces notations, notre matrice $(A|b)^{0'}$ est

$$P_{12}(A|b)^{(0)}. \quad (35.301)$$

Puis la matrice $(A|b)^{(1')}$ est

$$P_{23}(A|b)^{(1)}. \quad (35.302)$$

Et la matrice P qui arrive dans $PA = LU$ est la matrice $P = P_{23}P_{21}$, qui est une matrice de permutation non élémentaire. Elle vaut :

$$P = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}. \quad (35.303)$$

Lemme 35.130 ([477]).

Si $i, j > k$ alors les matrices de permutation élémentaires ont la relation de « commutation » suivante avec les transformations gaussiennes :

$$M_k(x)P_{ij} = P_{ij}M_k(P_{ij}(x)). \quad (35.304)$$

Démonstration. Il suffit de calculer les éléments de matrice :

$$(P_{ij}M_k(x))_{st} = (P_{ij})_{st} - \sum_m (P_{ij})_{sm} \tau_k(x)_m \delta_{kt}, \quad (35.305)$$

mais

$$\sum_m (P_{ij})_{sm} \tau_k(x)_m = (P_{ij}\tau_k(x))_s = \tau_k(P_{ij}(x))_s \quad (35.306)$$

parce que $i, j > k$ implique que dans $P_{ij}\tau_k(x)$ nous inversons deux élément non nuls de $\tau_k(x)$, tout en laissant le k^{e} élément. Le dénominateur ne change pas et il s'agit réellement d'une inversion de ligne. Donc

$$(P_{ij}M_k(x))_{st} = (P_{ij})_{st} - \tau_k(P_{ij}(x))_s \delta_{kt}. \quad (35.307)$$

De l'autre côté,

$$(M_k(y)P_{ij})_{st} = (P_{ij})_{st} - \tau_k(y)_s(P_{ij})_{kt}. \quad (35.308)$$

Mais comme $i, j > k$ la k^{e} ligne de P_{ij} est la même que celle de la matrice unité, donc $(P_{ij})_{kt} = \delta_{kt}$.

$$(M_k(y)P_{ij})_{st} = (P_{ij})_{st} - \tau_k(y)_s\delta_{kt}. \quad (35.309)$$

Cela correspond bien à (35.307). \square

35.18 Méthode de Gauss avec pivot partiel (décomposition PLU)

35.18.1 L'idée

À chaque pas, nous faisons une permutation de ligne. Nous permutons à chaque pas la première ligne avec celle qui a le pivot le plus grand (en valeur absolue). Donc :

$$(A|b)^{(0)} = \begin{pmatrix} 2 & 1 & 3 & 11 \\ 4 & 3 & 10 & 28 \\ -2 & 1 & 7 & 3 \end{pmatrix} \quad (35.310)$$

Nous commençons par déplacer des lignes :

$$(A|b)^{(0')} = \begin{pmatrix} 4 & 3 & 10 & 28 \\ 2 & 1 & 3 & 11 \\ -2 & 1 & 7 & 3 \end{pmatrix}. \quad (35.311)$$

Les multiplicateurs sont $l_{21} = 1/2$ et $l_{31} = -1/2$. Le fait est que les multiplicateurs ont toujours le plus grand dénominateur possible et nous avons alors toujours $0 \leq |l_{ij}| \leq 1$, qui sont des nombres relativement petits, et bien représentés en mémoire.

Nous avons la nouvelle matrice

$$(A|b)^{(1)} = \begin{pmatrix} 4 & 3 & 10 & 28 \\ 0 & -1/2 & -2 & -3 \\ 0 & 5/2 & 12 & 17 \end{pmatrix}. \quad (35.312)$$

Le pivot serait $-1/2$. Nous cherchons un pivot plus grand en dessous de ce $-1/2$ (et pas au dessus, sinon on casserait les zéros déjà trouvés). Nous trouvons le $5/2$ qui est plus grand. Nous permutons donc les deux dernières lignes :

$$(A|b)^{(1')} = \begin{pmatrix} 4 & 3 & 10 & 28 \\ 0 & 5/2 & 12 & 17 \\ 0 & -1/2 & -2 & -3 \end{pmatrix} \quad (35.313)$$

où le pivot est maintenant $l_{32} = -1/5$. La matrice suivante :

$$(A|b)^{(1)} = \begin{pmatrix} 4 & 3 & 10 & 28 \\ 0 & 5/2 & 12 & 17 \\ 0 & 0 & 2/5 & 2/5 \end{pmatrix} \quad (35.314)$$

Dans ce cas, la matrice L n'est pas aussi simple à construire parce que nous avons permuté des choses. Dans ce cas, la matrice L est encore de la forme

$$L = \begin{pmatrix} 1 & 0 & 0 \\ . & 1 & 0 \\ . & . & 1 \end{pmatrix}. \quad (35.315)$$

Mais vu que nous avons permuté les lignes 2 et 3 au deuxième pas, nous devons permuter l_{21} et l_{31} avant de remplir la matrice L avec les multiplicateurs :

$$L = \begin{pmatrix} 1 & 0 & 0 \\ -1/2 & 1 & 0 \\ 1/2 & -1/5 & 1 \end{pmatrix}. \quad (35.316)$$

Notons que ces L et U ne sont pas les mêmes que le LU obtenu sans pivot. Où est l'unicité? Elle est que en fait maintenant nous n'avons pas $A = LU$, mais

$$PA = LU \quad (35.317)$$

où P est une matrice de permutation.

35.18.2 Le théorème

Proposition 35.131 (Méthode de Gauss avec pivot partiel[477]).

Soit une matrice inversible $A \in \mathbb{M}(n, \mathbb{C})$. Il existe

- une matrice de permutation P
- une matrice triangulaire inférieure de diagonale unitaire L ,
- une matrice triangulaire supérieure inversible U

telles que

$$PA = LU. \quad (35.318)$$

Notons que cette proposition ne demande que l'hypothèse d'inversibilité pour A . Il n'y a pas d'hypothèses sur tous les mineurs comme c'était le cas avec Gauss sans pivot.

Démonstration. Nous prouvons par récurrence qu'il existe des matrices Q_k, E_1, \dots, E_k et A_k telles que

$$Q_k A = E_1 \dots E_k A_k \quad (35.319)$$

avec

- (1) Q_k est une matrice de permutation
- (2) E_i est une transformation gaussienne sur la i^{e} colonne
- (3) A_k est triangulaire supérieure jusqu'à la k^{e} colonne.

Sachant que $\det(Q_k) = \pm 1$, et que $\det(E_i) = 1$, le passage au déterminant dans (35.319) nous donne $\det(A_k) \neq 0$ et si nous notons $\Omega_k(A)$ la matrice tronquée de A , ne gardant que les entrées plus grandes que k , nous avons

$$\det(A_k) = \prod_{i=1}^k (A_k)_{ii} \det(\Omega_{k+1}(A_k)). \quad (35.320)$$

Donc : $(A_k)_{ii} \neq 0$ pour $i \leq k$ et $\det(\Omega_{k+1}(A_k)) \neq 0$.

Pour fixer les idées, voici une image de $k = 2$:

$$\begin{array}{c} \Delta_k(A_2) \\ \left(\begin{array}{ccccc} * & * & * & * & * \\ 0 & * & * & * & * \\ 0 & 0 & * & * & * \\ 0 & 0 & * & * & * \\ 0 & 0 & * & * & * \end{array} \right) \\ \Omega_{k+1}(A_2) \end{array} \quad (35.321)$$

Étant donné que $\det(\Omega_{k+1}(A_k)) \neq 0$, parmi les nombres $(A_k)_{i,k+1}$ ($i \geq k+1$), au moins un est non nul et nous posons r_{k+1} tel que $|(A_k)_{r_{k+1},k+1}|$ soit maximum parmi ces éléments.

Le nombre r_{k+1} est enregistré parce qu'il servira à écrire la matrice P plus tard. Les matrices E_i ne sont pas enregistrées, parce que nous verrons qu'elles vont encore changer. Seule la dernière sera enregistrée.

La composante $(k+1, k+1)$ de la matrice

$$P_{r_{k+1},k+1} A_k \quad (35.322)$$

est non nulle et peut donc servir de pivot. Soit M_{k+1} la transformation gaussienne pour la $(k+1)$ ^e colonne de la matrice $P_{r_{k+1},k+1}A_k$. La matrice

$$A_{k+1} = M_{k+1}P_{r_{k+1},k+1}A_k \quad (35.323)$$

est alors une matrice triangulaire supérieure jusqu'à la $(k+1)$ ^e colonne. En posant $E_{k+1} = M_{k+1}^{-1}$ nous avons

$$P_{r_{k+1},k+1}E_{k+1}A_{k+1} = A_k, \quad (35.324)$$

et nous nous sentons en droit de récrire l'équation de départ (35.319) :

$$Q_k A = E_1 \dots E_k A_k = E_1 \dots E_k P_{r_{k+1},k+1} E_{k+1} A_{k+1}. \quad (35.325)$$

Le lemme 35.130 nous permet de ramener la matrice $P_{r_{k+1},k+1}$ en première position, quitte à modifier un peu (pas beaucoup) chacune des matrices E_i ($i = 1, \dots, k$). C'est pour cela que nous n'enregistrons pas les matrices E_i . Nous avons donc

$$P_{r_{k+1},k+1} Q_k A = E'_1 \dots E'_k E_{k+1} A_{k+1} \quad (35.326)$$

où

— Le produit $P_{r_{k+1},k+1} Q_k$ est encore une matrice de permutation, et mieux : elle vaut

$$\prod_{i=1}^{k+1} P_{r_i, i}. \quad (35.327)$$

Cela montre qu'il est suffisant d'enregistrer les nombres r_i pour reconstituer cette partie.

— La matrice E'_i est une transformation gaussienne pour la i ^e colonne.

— La matrice A_{k+1} est triangulaire supérieure jusqu'à la $k+1$ ^e colonne.

La récurrence est maintenant finie et nous pouvons écrire avec $k = n$:

$$Q_n A = E_1 \dots E_n A_n \quad (35.328)$$

où le produit $E_1 \dots E_n$ est triangulaire inférieure et A_n est triangulaire supérieur.

Maintenant nous enregistrons la matrice $U = A_n$, le produit $L = \prod_{i=1}^n E_n$ et les nombres r_i qui permettent de retrouver P . \square

Note : dans l'équation (35.328) nous avons bien entendu massivement renommé les E'_i en E_i . En réalité la matrice E_1 vient avec n primes sur la tête.

Dans les exemples 35.134, 35.135 et 35.136, nous allons résoudre le système

$$\begin{pmatrix} 10^{-9} & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad (35.329)$$

d'abord de façon exacte, et ensuite en supposant une machine ne tenant que 8 chiffres significatifs en utilisant la méthode de Gauss avec ou sans pivot.

Commençons par voir comment se passe en pratique la décomposition $PA = LU$ de Gauss avec pivot partiel.

Exemple 35.132

Décomposons la matrice

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 5 & 0 \\ 3 & 8 & 0 \end{pmatrix}. \quad (35.330)$$

Sur la première colonne, le plus grand nombre est 3. Nous commençons par permuter la première et la troisième ligne en utilisant la matrice de permutation $P_1 = P_{3,1}$ et nous enregistrons $r_1 = 3$. Nous avons alors la matrice

$$A'_0 = \begin{pmatrix} 3 & 8 & 0 \\ 2 & 5 & 0 \\ 1 & 3 & 3 \end{pmatrix}. \quad (35.331)$$

Pour trouver la matrice A_1 nous suivons l'équation (35.323). Bien que le résultat net soit des combinaisons de lignes : $L_2 \rightarrow L_2 - 2L_1/3$ et $L_3 \rightarrow L_3 - L_1/3$ (que nous pourrions savoir dès à présent), il est important de passer par la matrice gaussienne pour obtenir la matrice L_1 .

La matrice de transformation gaussienne pour la première colonne de (35.331) est :

$$M_1 = \begin{pmatrix} 1 & 0 & 0 \\ -2/3 & 1 & 0 \\ -1/3 & 0 & 1 \end{pmatrix} \quad (35.332)$$

et $L_1 = M_1^{-1}$. Le lemme 35.121 nous dit comment calculer facilement cet inverse :

$$L_1 = \begin{pmatrix} 1 & 0 & 0 \\ 2/3 & 1 & 0 \\ 1/3 & 0 & 1 \end{pmatrix} \quad (35.333)$$

En suivant l'équation (35.323) nous posons $A_1 = M_1 A'_0$:

$$A_1 = \begin{pmatrix} 1 & 0 & 0 \\ -2/3 & 1 & 0 \\ -1/3 & 0 & 1 \end{pmatrix} \begin{pmatrix} 3 & 8 & 0 \\ 2 & 5 & 0 \\ 1 & 3 & 3 \end{pmatrix} = \begin{pmatrix} 3 & 8 & 0 \\ 0 & -1/3 & 0 \\ 0 & -2/3 & 3 \end{pmatrix} \quad (35.334)$$

et nous avons

$$Q_1 A = L_1 A_1 \quad (35.335)$$

où L_1 , A_1 et $r_1 = 3$ sont enregistrés. La matrice Q_1 peut être retrouvée en sachant r_1 parce que P est la matrice de permutation $P_{r_1,1}$.

Nous travaillons maintenant sur la deuxième colonne de A_1 . Le plus grand élément en valeur absolue (sur ou sous la diagonale) est $-2/3$. Nous posons $r_2 = 3$ et

$$A'_1 = \begin{pmatrix} 3 & 8 & 0 \\ 0 & -2/3 & 3 \\ 0 & -1/3 & 0 \end{pmatrix} \quad (35.336)$$

et la matrice gaussienne pour la deuxième colonne est

$$M_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -1/2 & 1 \end{pmatrix} \quad (35.337)$$

Le $-1/2$ provient du calcul $-((-1/3)/(-2/3))$. L'inverse de cette matrice est facile :

$$L_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1/2 & 1 \end{pmatrix} \quad (35.338)$$

et la matrice suivante à enregistrer est

$$A_2 = M_2 P_{3,2} A_1 = M_2 A'_1 = \begin{pmatrix} 3 & 8 & 0 \\ 0 & -2/3 & 3 \\ 0 & 0 & -3/2 \end{pmatrix}. \quad (35.339)$$

Notons toutefois que pour calculer cette matrice, seul le dernier élément demande un calcul. La première colonne ne change pas (par construction), la seconde gagne un zéro en dernière ligne (la matrice M_2 sert à ça) et sur la dernière colonne, seule la dernière ligne est sujette à changement.

Avec la matrice A_2 , la trigonalisation supérieure est faite. La décomposition n'est cependant pas terminée. Nous devons encore trouver la partie triangulaire inférieure. Nous en sommes à

$$Q_1 A = L_1 A_1 = L_1 P_{3,2} L_2 A \quad (35.340)$$

où Q_1 est la première matrice de permutation.

Utilisant le lemme 35.130, il est facile de permuter L_1 avec $P_{3,2}$:

$$L_1 P_{3,2} = P_{3,2} \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 1/3 & 1 & 0 \\ 2/3 & 0 & 1 \end{pmatrix}}_{L'_1} \quad (35.341)$$

Nous avons donc

$$P_{3,2} P_{3,1} A = L'_1 L_2 A \quad (35.342)$$

Deux multiplications matricielles plus tard nous terminons :

$$PA = LU \quad (35.343)$$

avec

$$P = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}, \quad L = \begin{pmatrix} 1 & 0 & 0 \\ 1/3 & 1 & 0 \\ 2/3 & 1/2 & 1 \end{pmatrix}, \quad U = A_2 = \begin{pmatrix} 3 & 8 & 0 \\ 0 & -2/3 & 3 \\ 0 & 0 & -3/2 \end{pmatrix}. \quad (35.344)$$

△

Notons que Sage utilise la méthode de Gauss avec pivots :

```

1 sage: A=matrix([ [1,2,3],[2,5,0],[3,8,0] ])
2 sage: A
3 [1 2 3]
4 [2 5 0]
5 [3 8 0]
6 sage: A.LU()
7 (
8 [0 1 0] [ 1 0 0] [ 3 8 0]
9 [0 0 1] [1/3 1 0] [ 0 -2/3 3]
10 [1 0 0], [2/3 1/2 1], [ 0 0 -3/2]
11 )

```

tex/sage/sageSnip006.sage

Mais attention : Sage crée une décomposition $A = PLU$ et non $PA = LU$. D'où le fait que la matrice de permutation de Sage est l'inverse de celle donnée ici.

35.18.3 D'un point de vue algorithmique

Problèmes et choses à faire

Je ne suis pas certain de l'optimalité de ce que je raconte ici. Je décris simplement ce que j'ai fait pour écrire mon programme `finitediff`.

Si vous êtes expert en calcul numérique, n'hésitez pas à donner votre avis.

Nous décrivons à présent la décomposition $A = PLU$ (du théorème 35.131, avec le P à droite). En suivant l'exemple 35.132 nous voyons assez bien comment créer les matrices U et P au fur et à mesure. La construction de L est peut-être moins évidente.

Écrivons un exemple très explicite pour

$$A = \begin{pmatrix} 2 & 1 & 3 \\ 4 & 3 & 10 \\ -2 & 1 & 7 \end{pmatrix}. \quad (35.345)$$

Nous commençons par permuter des lignes pour avoir un grand pivot :

$$P_{12}A = \begin{pmatrix} 4 & 3 & 10 \\ 2 & 1 & 3 \\ -2 & 1 & 7 \end{pmatrix}. \quad (35.346)$$

Et nous effectuons l'élimination avec la matrice

$$M_1 = \begin{pmatrix} 1 & 0 & 0 \\ -1/2 & 1 & 0 \\ 1/2 & 0 & 1 \end{pmatrix}. \quad (35.347)$$

Cela donne le premier résultat :

$$M_1 P_{12} A = \begin{pmatrix} 4 & 3 & 10 \\ 0 & -1/2 & -2 \\ 0 & 5/2 & 12 \end{pmatrix} \quad (35.348)$$

Nous continuons avec P_{23} pour avoir un nouveau grand pivot :

$$P_{23} M_1 P_{12} A = \begin{pmatrix} 4 & 3 & 10 \\ 0 & 5/2 & 12 \\ 0 & -1/2 & -2 \end{pmatrix}. \quad (35.349)$$

Nous utilisons la matrice

$$M_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1/5 & 1 \end{pmatrix} \quad (35.350)$$

et au final :

$$M_2 P_{23} M_1 P_{12} A = \begin{pmatrix} 4 & 3 & 10 \\ 0 & 5/2 & 12 \\ 0 & 0 & 2/5 \end{pmatrix} = U. \quad (35.351)$$

L'égalité obtenue est

$$M_2 P_{23} M_1 P_{12} A = U. \quad (35.352)$$

Pour avoir la décomposition PLU il faut écrire

$$A = P_{12} M_1^{-1} P_{23} M_2^{-1} U, \quad (35.353)$$

et permuter P_{23} avec M_1^{-1} , ce qui est facile par le lemme 35.130.

Remarque 35.133.

Nous ne devons permuter la matrice M_k avec une matrice de permutation qu'à partir de la deuxième étape. En effet l'équation (35.348) revient à

$$A = M_1^{-1} P_{12} m_U \quad (35.354)$$

qui est dans le bon ordre. Ce n'est qu'à partir de la seconde étape que des matrices de permutation apparaissent à droite des matrices gaussiennes.

Cependant dans un cas 4×4 , cette méthode deviendrait fastidieuse parce que nous aurions encore des étapes à faire. En repartant de (35.353), mais avec m_U (la matrice pas encore tout à fait triangularisée) au lieu de U , nous aurons, pour un certain $k > 3$:

$$M_3 P_{3k} M_2 P_{23} M_1 P_{12} A = U, \quad (35.355)$$

ce qui fait :

$$A = P_{12} M_1^{-1} P_{23} M_2^{-1} P_{3k} M_3^{-1} U. \quad (35.356)$$

Tous les P_{ij} peuvent être mis à gauche parce que leurs indices sont toujours strictement supérieurs à ceux des M_l placés devant eux. Mais c'est fastidieux.

Nous allons donc permuter à chaque étape pour ne retenir que l'important. Si à une certaine étape nous avons

$$A = P_{1,r_1} \dots P_{k,r_k} M_1^{-1} \dots M_k^{-1} m_U \quad (35.357)$$

avec

$$m_U = P_{k+1, r_{k+1}} M_{k+1}^{-1} U \quad (35.358)$$

alors nous allons directement permuter $P_{k+1, r_{k+1}}$ avec tous les M_i^{-1} . Si nous notons P_k la permutation (pas élémentaire) à l'étape k et L_k la matrice triangulaire inférieure à de l'étape k ,

$$A = P_k L_k m_U = P_k L_k P_{k+1, r_{k+1}} M_{k+1}^{-1} m'_U. \quad (35.359)$$

Nous enregistrons alors $P_{k+1} = P_k P_{k+1, r_{k+1}}$ et pour L_{k+1} nous partons de L_k et nous faisons deux opérations suivantes :

- nous permutons, sur ses colonnes non triviales, les indices $k+1$ et r_{k+1} ,
- nous multiplions par M_{k+1}^{-1} , ce qui revient à simplement lui ajouter une colonne non triviale.

Notons que $r_{k+1} \geq k+1$, de telle sorte que sur les colonnes non triviales (qui sont jusqu'au numéro k), la permutation des lignes $k+1$ et r_{k+1} ne change pas l'aspect de la matrice : elle reste multi-gaussienne de dernière colonne k .

35.18.4 Exemples

Nous nous lançons dans la résolution du système

$$\begin{pmatrix} 10^{-9} & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}. \quad (35.360)$$

Exemple 35.134

Nous commençons de façon exacte, par la méthode de Gauss sans pivot. La première transformation gaussienne est

$$E_1 = \begin{pmatrix} 1 & 0 \\ -10^9 & 1 \end{pmatrix} \quad (35.361)$$

et nous calculons

$$E_1 A = \begin{pmatrix} 1 & 0 \\ -10^9 & 1 \end{pmatrix} \begin{pmatrix} 10^{-9} & 1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 10^{-9} & 1 \\ 0 & 1 - 10^{-9} \end{pmatrix}. \quad (35.362)$$

Vu que cette dernière est triangulaire supérieure, nous avons fini la méthode de Gauss et $U = E_1 A$. En ce qui concerne la matrice L , elle est donnée par $L = E_1^{-1}$, c'est-à-dire

$$L = \begin{pmatrix} 1 & 0 \\ -10^9 & 1 \end{pmatrix}^{-1} = \begin{pmatrix} 1 & 0 \\ 10^9 & 1 \end{pmatrix}. \quad (35.363)$$

Au final nous avons la décomposition $A = LU$ exacte suivante :

$$L = \begin{pmatrix} 1 & 0 \\ 10^9 & 1 \end{pmatrix} \quad U = \begin{pmatrix} 10^{-9} & 1 \\ 0 & 1 - 10^9 \end{pmatrix}. \quad (35.364)$$

Résoudre le système $Ax = b$ revient à résoudre $LUx = b$ et donc résoudre successivement les systèmes

$$\begin{cases} Ly = b \\ Ux = y. \end{cases} \quad (35.365a)$$

$$(35.365b)$$

D'abord le système

$$\begin{pmatrix} 1 & 0 \\ 10^9 & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad (35.366)$$

donne $y_1 = 1$ et $y_2 = 2 - 10^9$.

Ensuite nous résolvons

$$\begin{pmatrix} 10^{-9} & 1 \\ 0 & 1 - 10^9 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 - 10^9 \end{pmatrix}. \quad (35.367)$$

Cela donne

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} -\frac{10^9}{1-10^9} \\ \frac{2-10^9}{1-10^9} \end{pmatrix} \simeq \begin{pmatrix} 1 \\ 1 \end{pmatrix}. \quad (35.368)$$

C'est également le résultat que trouve Sage :

```

1 sage: var('y')
2 y
3 sage: solve([10**(-9)*x+y==1, x+y==2], [x, y])
4 [[x == (1000000000/999999999), y == (999999998/999999999)]]

```

tex/sage/sageSnip007.sage

△

Exemple 35.135([477])

Nous recommençons tout le calcul avec une précision limitée à 8 chiffres significatifs, sans pivot.

Nous avons à nouveau la transformation gaussienne

$$E_1 = \begin{pmatrix} 1 & 0 \\ -10^9 & 1 \end{pmatrix}, \quad (35.369)$$

mais pour calculer U nous effectuons le produit matriciel

$$U = E_1 A = \begin{pmatrix} 1 & 0 \\ -10^9 & 1 \end{pmatrix} \begin{pmatrix} 10^{-9} & 1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 10^{-9} & 1 \\ 0 & * \end{pmatrix}. \quad (35.370)$$

Nous détaillons à présent le calcul de l'élément noté *. Le calcul de $10^9 \ominus 1$ donne

$$999999999 = 9.99999999 \times 10^8, \quad (35.371)$$

mais la précision étant limitée à 8 chiffres, un arrondi arrive. Étant donné que le premier chiffres supprimé est un 9 nous retombons sur 10^9 , et donc notre machine à précision limitée donnera

$$U = \begin{pmatrix} 10^{-9} & 1 \\ 0 & -10^9 \end{pmatrix}. \quad (35.372)$$

Ensuite le calcul de $L = E_1^{-1}$ ne cause pas de problèmes :

$$L = \begin{pmatrix} 1 & 0 \\ -10^9 & 1 \end{pmatrix}. \quad (35.373)$$

Maintenant il s'agit de résoudre les systèmes $Ly = b$ et $Ux = y$. Du système

$$\begin{pmatrix} 1 & 0 \\ 10^9 & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad (35.374)$$

nous tirons tout de suite $y_1 = 1$ et ensuite $10^9 + y_2 = 2$, c'est-à-dire $y_2 = 2 - 10^9$, qui en précision limitée donne encore $y_2 = -10^9$. À résoudre maintenant :

$$\begin{pmatrix} 10^{-9} & 1 \\ 0 & -10^9 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ -10^9 \end{pmatrix}. \quad (35.375)$$

Cela donne immédiatement $x_2 = 1$ et ensuite

$$10^{-9}x_1 + 1 = 1, \quad (35.376)$$

donc $x_1 = 0$. La solution trouvée est

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad (35.377)$$

qui est complètement faux au niveau de la première variable. \triangle

Exemple 35.136([477])

Nous résolvons encore le même système en précision limitée, mais en utilisant cette fois la méthode de Gauss avec pivot partiel.

Le plus grand élément de la première colonne est 1 ; nous utilisons donc la permutation $P_{1,2}$:

$$P_{1,2}A = \begin{pmatrix} 1 & 1 \\ 10^{-9} & 1 \end{pmatrix}. \quad (35.378)$$

La matrice de transformation gaussienne pour la première colonne de cette matrice est

$$M_1 = \begin{pmatrix} 1 & 0 \\ -10^{-9} & 1 \end{pmatrix} \quad (35.379)$$

et nous posons

$$A_1 = M_1 P_{1,2} A = \begin{pmatrix} 1 & 0 \\ -10^{-9} & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 10^{-9} & 1 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 0 & -10^{-9} + 1 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} = U \quad (35.380)$$

où un arrondi a eu lieu pour $-10^{-9} + 1 = 1$. En inversant M_1 nous avons

$$L_1 = M_1^{-1} = \begin{pmatrix} 1 & 0 \\ 10^{-9} & 1 \end{pmatrix}. \quad (35.381)$$

La décomposition est

$$A = \underbrace{\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}}_P \underbrace{\begin{pmatrix} 1 & 0 \\ 10^{-9} & 1 \end{pmatrix}}_L \underbrace{\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}}_U \quad (35.382)$$

Le moment de résoudre est venu. Vu que $PLUx = b$ nous devons résoudre les systèmes

$$\begin{cases} Pz = b & (35.383a) \end{cases}$$

$$\begin{cases} Ly = z & (35.383b) \end{cases}$$

$$\begin{cases} Ux = y. & (35.383c) \end{cases}$$

Pour z c'est facile :

$$z = \begin{pmatrix} 2 \\ 1 \end{pmatrix}. \quad (35.384)$$

Pour y il y a un arrondi :

$$\begin{pmatrix} 1 & 0 \\ 10^{-9} & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \end{pmatrix}. \quad (35.385)$$

Tout de suite : $y_1 = 2$ et ensuite $2 \times 10^{-9} + y_2 = 1$, ce qui donne $y_2 = 1 \ominus 2 \times 10^{-9} = 1$. Donc

$$y = \begin{pmatrix} 2 \\ 1 \end{pmatrix}. \quad (35.386)$$

Et enfin pour x c'est le système

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \end{pmatrix}. \quad (35.387)$$

Nous avons $x_2 = 1$ et ensuite $x_1 + 1 = 2$ c'est-à-dire $x_1 = 1$. Au final la solution trouvée est

$$x = \begin{pmatrix} 1 \\ 1 \end{pmatrix}. \quad (35.388)$$

Cette solution est considérablement meilleure que (35.377). \triangle

35.137.

L'utilisation du pivot non seulement assure le fait que la trigonalisation va bien se passer (on évite les zéros en pivot), mais aussi et surtout, en choisissant de prendre le plus grand pivot possible, nous obtenons une meilleure stabilité numérique.

35.19 Résolution de systèmes linéaires (suite)**35.19.1 Déterminant**

Pour calculer un déterminant lorsque nous avons la décomposition $A = LU$ nous pouvons faire

$$\det(A) = \det(LU) = \det(L) \det(U) = \det(U) \quad (35.389)$$

parce que L est triangulaire avec des 1 sur la diagonale.

Si par contre nous avons fait des pivots, nous avons $PA = LU$. Il nous faut le déterminant de P , qui n'est autre que ± 1 . Nous avons

$$\det(P) = (-1)^s \quad (35.390)$$

où s est le nombre de permutations effectives effectuées. Nous précisons « effectives » parce qu'il ne faut pas compter le pas où nous n'avons pas permuté (les cas où le bon pivot était présent du premier coup). Nous avons alors

$$\det(A) = (-1)^s \det(U). \quad (35.391)$$

35.19.2 Plusieurs termes indépendants

Mettons un système $Ax = b$ qu'il faut résoudre pour plusieurs b différents. C'est typiquement le cas où l'on voudrait calculer l'inverse de A . Mais on va directement se calmer. Soient donc à résoudre $Ax_1 = b_1, \dots, Ax_n = b_n$.

Les opérations (avec ou sans pivot) que nous faisons ne dépendent que de la matrice A , mais aucune décisions concernant les pivots ou la matrice des multiplicateurs ne dépend de b . Autre façon de dire : si le système $(A|b_1)$ devient $(U|y_1)$, le système $(A|b_i)$ devient $(U|y_i)$ avec le même U .

Nous ne sommes donc pas obligés de faire tout le travail autant de fois qu'il n'y a de systèmes à résoudre. Donc si on a plusieurs systèmes à résoudre avec la même matrice, on fait mieux de retenir une fois pour toute la décomposition LU (avec ou sans pivots), avant de vraiment résoudre.

Ou alors on peut aussi faire que, au lieu de faire $(A|b_i)$ plein de fois, faire une seule fois

$$(A|b_1 \dots b_n). \quad (35.392)$$

Et on fait tout le travail sur tous les vecteurs d'un en même temps.

Soit e_i la base canonique. Si nous notons x_n les solutions des problèmes $Ax_i = e_i$, tous les problèmes $Ax_i = e_i$ s'écrivent d'un seul coup

$$AX = Y \quad (35.393)$$

où X est la matrice des x_i en colonnes, et Y est celle des e_i en colonnes. Oh, mais $Y = \mathbb{1}$ évidemment. Donc

$$AX = \mathbb{1}. \quad (35.394)$$

Si nous supposons que A est inversible, alors ce X est l'inverse.

Donc pour calculer l'inverse d'une matrice de dimension non trop grande, il suffit d'utiliser la méthode de Gauss sur les vecteurs de la base canonique. Cette idée est la base du calcul de l'inverse par matrice companion. En effet, si nous partons du problème

$$(A|\mathbb{1}) \quad (35.395)$$

et nous appliquons la méthode de Gauss avec pivot, nous arrivons à

$$(U|L^{-1}P). \quad (35.396)$$

Attention : le produit $L^{-1}P$ est une permutation des *colonnes* de L^{-1} . Vu que L est triangulaire inférieure avec des 1 sur la diagonale, L^{-1} est triangulaire inférieure avec des 1 sur la diagonale. Donc si la matrice n'est pas trop grande, on peut assez facilement remettre les colonnes de $L^{-1}P$ dans l'ordre pour recomposer une matrice triangulaire inférieure avec 1 sur la diagonale.

Une autre façon de calculer l'inverse, si $A = LU$ est connue, il suffit de faire

$$A^{-1} = U^{-1}L^{-1}. \quad (35.397)$$

Et il existe un algorithme facile pour l'inverse d'une matrice triangulaire.

35.19.3 Cholesky

Le commandant Cholesky travaillait sur le tir de canon (chose éminemment liée à de nombreuses mathématiques ingénieuses). La méthode de Cholesky est encore utilisée aujourd'hui dans les vrais problèmes.

La méthode de Gauss s'applique sans hypothèses sur la matrice A , à part qu'elle doit être de petite dimension, comme pour toute méthode directe. Souvent nous savons des choses sur la matrice. Ici nous allons supposer que A est symétrique et définie positive.

Comment numériquement vérifier ces hypothèses ? En ce qui concerne la symétrie, il suffit de faire le test complet :

$$A^t = A. \quad (35.398)$$

La vérification de cela coûte au maximum n^2 comparaisons (et en fait la moitié de ça moins la diagonale).

Le fait que A soit définie positive est facile à vérifier pour utiliser Cholesky parce que il suffit de le faire, et s'il n'y a pas de nombres complexes qui arrivent, c'est que la matrice était définie positive.

Un lemme très simple à mettre en oeuvre numériquement nous permet de traiter certains cas.

Lemme 35.138.

Une matrice symétrique possédant un élément négatif sur la diagonale n'est pas définie positive.

Démonstration. Un simple calcul ou effort d'imagination montre que $\langle Me_k, e_k \rangle = M_{kk}$. Donc si M doit être définie positive, M_{kk} doit être positive par le lemme 11.196. \square

Ce lemme est un moyen déjà de faire quelques vérifications. Et si les éléments diagonaux de A sont tous négatifs, on peut prendre $-A$.

Lemme 35.139.

Si A est une matrice symétrique strictement définie positive, alors pour tout k , la matrice tronquée $\Delta_k(A)$ l'est également.

Démonstration. Le fait que $\Delta_k(A)$ soit symétrique est évidemment. Le fait qu'elle soit définie positive l'est moins. Soit $y \in \mathbb{R}^k$ et le vecteur $\tau y \in \mathbb{R}^n$, qui est « complété » avec des zéros.

Nous avons $\langle \Delta_k(A)y, y \rangle_k = \langle A\tau y, \tau \rangle_n$. En effet

$$\langle \Delta_k(A)y, y \rangle = \sum_{i=1}^k \sum_{l=1}^k A_{il}y_l y_i. \quad (35.399)$$

Et à droite :

$$\langle A\tau y, \tau y \rangle = \sum_{i=1}^n (A\tau y)_i (\tau y)_i = \sum_{i=1}^k (A\tau y)_i y_i = \sum_{i=1}^k \sum_{l=1}^n A_{il}(\tau y)_l y_i = \sum_{i=1}^k \sum_{l=1}^k A_{il}y_l y_i \quad (35.400)$$

où nous avons utilisé le fait que $(\tau y)_i = 0$ dès que $i > k$ et que $(\tau y)_i = y_i$ sinon.

En conséquence de quoi $\langle \Delta_k(A)y, u \rangle > 0$ pour tout $y \in \mathbb{R}^k$ et la matrice $\Delta_k(A)$ est strictement définie positive. \square

Lemme 35.140.

Si T est une matrice triangulaire, alors $(T_{ii})^{-1} = (T^{-1})_{ii}$.

Démonstration. Il suffit de se rendre compte que le coefficient ii de l'égalité $\mathbb{1} = TT^{-1}$ donne

$$1 = \sum_l T_{il}(T^{-1})_{li}. \quad (35.401)$$

Dans la somme il ne reste que le terme $l = i$. \square

Nous allons chercher une décomposition de type LU sous la forme $A = LL^t$, c'est-à-dire $U = L^t$. Attention : maintenant nous n'avons plus des 1 sur la diagonale. Ce n'est donc pas exactement la décomposition LU dont nous parlions plus haut. C'est pour cela que nous n'allons pas la noter LL^t mais BB^t .

Théorème 35.141 (Cholesky[476]).

Soit une matrice réelle symétrique strictement définie positive. Il existe une unique matrice réelle B telle que

- B est triangulaire inférieure,
- la diagonale de B est positive,
- $A = BB^t$.

Démonstration. Par la décomposition LU du théorème 35.123 nous avons des matrices L et U telles que $A = LU$. Soit D la matrice diagonale donnée par

$$D_{ii} = \sqrt{U_{ii}}. \quad (35.402)$$

Cette définition fonctionne parce que $U_{ii} > 0$. En effet nous savons que $\Delta_k(A) = \Delta_k(L)\Delta_k(U)$, et en passant au déterminant,

$$\det(\Delta_k(A)) = \det(\Delta_k(U)). \quad (35.403)$$

Vu que $\Delta_k(A)$ est strictement définie positive par le lemme 35.139, son déterminant est strictement positif²⁰ et nous avons

$$\det(\Delta_k(U)) > 0. \quad (35.404)$$

En appliquant cela à $k = 1$ nous avons $U_{11} > 0$ puis de proche en proche, $U_{ii} > 0$ pour tout i .

Nous posons :

$$B = LD \quad \text{qui est triangulaire inférieure} \quad (35.405a)$$

$$C = D^{-1}U \quad \text{qui est triangulaire supérieure.} \quad (35.405b)$$

Nous avons bien entendu $A = BC$ et nous allons prouver que $C = B^t$. Vu que $A = A^t$ nous pouvons identifier BC et C^tB^t :

$$BC = C^tB^t. \quad (35.406)$$

En mettant les matrices triangulaires supérieures à gauche et inférieures à droite :

$$C(B^t)^{-1} = B^{-1}C^t, \quad (35.407)$$

qui sont donc deux matrices diagonales. Nous montrons que cette diagonale est en réalité l'identité.

D'abord

$$B_{ii} = \sum_{l=1}^n L_{il}D_{li} = L_{ii}\sqrt{U_{ii}} = \sqrt{U_{ii}} \quad (35.408)$$

²⁰. Le théorème 11.189 donne une diagonalisation par des matrices de déterminant 1. Vu que les valeurs propres forment sur la diagonale, et qu'elles sont toutes positives, el déterminant est positif.

parce que $L_{ii} = 1$. Notons en passant que la diagonale de B est positive. Ensuite

$$C_{ii} = \sum_{l=1}^n (D^{-1})_{il} U_{li} = (D^{-1})_{ii} U_{ii} = \frac{1}{\sqrt{U_{ii}}} U_{ii} = \sqrt{U_{ii}}. \quad (35.409)$$

Donc B et C ont des diagonales égales. Calculons alors la diagonale de $B^{-1}C^t$:

$$(B^{-1}C^t)_{ii} = \sum_l (B^{-1})_{il} (C^t)_{li} = (B^{-1})_{ii} C_{ii} \quad (35.410)$$

parce que encore une fois, de la somme il ne reste que le terme $l = i$.

Mais B est une matrice triangulaire qui tombe sous le coup du lemme 35.140. Donc $(B^{-1})_{ii} = (B_{ii})^{-1} = (C_{ii})^{-1}$. Nous avons alors

$$(B^{-1}C^t)_{ii} = 1. \quad (35.411)$$

Cela conclut l'existence de la décomposition de Cholesky.

En ce qui concerne l'unicité, soient $A = BB^t = CC^t$. Nous regroupons les supérieures et les inférieures :

$$B^t(C^t)^{-1} = B^{-1}C. \quad (35.412)$$

Ces deux matrices sont donc diagonales et nous posons $D = B^{-1}C$, c'est-à-dire $C = BD$. Nous remplaçons donc C par BD dans (35.412) :

$$A = BB^t = BD(BD)^t = BDD^tB^t. \quad (35.413)$$

Donc $DD^t = \mathbb{1}$, ce qui signifie que les éléments diagonaux de D sont ± 1 . Nous montrons qu'ils sont positifs : à partir de $C = BD$ nous déballons

$$C_{ii} = \sum_l B_{il} D_{li}, \quad (35.414)$$

et donc

$$B_{ii} D_{ii} = C_{ii}. \quad (35.415)$$

En sachant que les conditions de la décomposition de Cholesky demandent les éléments diagonaux positifs nous en déduisons que D_{ii} est positif et donc égal à 1. Finalement $D = \mathbb{1}$ et $B = C$. \square

Prenons la matrice

$$A = \begin{pmatrix} 4 & 2 & -2 \\ 2 & 10 & -7 \\ -2 & -7 & 9 \end{pmatrix} \quad (35.416)$$

Elle est symétrique et définie positive. Nous posons

$$\{ l_{11} = \sqrt{a_{11}} l_{i1} = a_{i1} / l_{11} \quad (35.417a)$$

pour $i = 2, \dots, n$. Et aussi

$$\left\{ \begin{array}{l} l_{jj} = (a_{jj} - \sum_{k=1}^{j-1} l_{jk}^2)^{1/2} \end{array} \right. \quad (35.418a)$$

$$\left\{ \begin{array}{l} l_{ij} = (a_{ij} - \sum_{k=1}^{j-1} l_{ik} l_{jk}) / a_{jj} \end{array} \right. \quad (35.418b)$$

pour $i = j + 1, \dots, n$.

Les formules (35.417) nous disent comment remplir la première colonne. Cela donne la matrice

$$L = \begin{pmatrix} \sqrt{4} = 2 & . & . \\ 2/2 = 1 & . & . \\ -2/2 = -1 & . & . \end{pmatrix} \quad (35.419)$$

Les formules (35.418) donnent les autres colonnes en fonction des précédentes.

Dans Sage :

```

1 sage: A=matrix( [ [4,2,-2],[2,10,-7],[-2,-7,9] ] )
2 sage: A
3 [ 4  2 -2]
4 [ 2 10 -7]
5 [-2 -7  9]
6 sage: A.cholesky()
7 [ 2  0  0]
8 [ 1  3  0]
9 [-1 -2  2]

```

tex/sage/sageSnip005.sage

35.20 Système linéaire (méthodes itératives)

Nous trouvons des méthodes itératives lorsque les matrices sont grandes, ce qui arrive lorsque l'on discrétise une équation différentielle.

Nous allons chercher des méthodes de la forme $x_{n+1} = Bx_n + q$; ce sont des méthodes stationnaires. La convergence d'une méthode est toujours liée à la matrice B et en général, la convergence ne dépend pas du choix du vecteur initial. Nous faisons donc souvent $x_0 = 0$ et donc $x_1 = q$. Voilà donc une itération de faite gratuitement.

Nous notons e_k le **vecteur d'erreur** qui est défini par $e_k = x - x_k$. Et le **vecteur résidu** $r_k = b - Ax_k$. Attention : ici k n'est pas un indice mais un numéro de vecteur.

Notons que si x est solution, alors $b - Ax = 0$, ce qui motive le vecteur résidu.

Les conditions d'arrêt d'un algorithme seraient

$$\left\{ \begin{array}{l} \|e_k\|_\infty \ll \epsilon_1 \\ \|r_k\|_\infty < \epsilon_2 \end{array} \right. \quad (35.420a)$$

$$\left\{ \begin{array}{l} \|e_k\|_\infty \ll \epsilon_1 \\ \|r_k\|_\infty < \epsilon_2 \end{array} \right. \quad (35.420b)$$

où ϵ_1 et ϵ_2 sont des précisions décidées à l'avance par l'utilisateur.

Proposition 35.142.

Si A est une matrice inversible, alors

$$\lim_{k \rightarrow \infty} e_k = \lim_{k \rightarrow \infty} r_k = 0. \quad (35.421)$$

Vu que $r_k = Ae_k$, si la matrice A est mal conditionnée, il peut arriver que r_k reste grand alors que e_k est déjà petit.

Remarque 35.143.

Dans les méthode stationnaires, nous avons $x_{n+1} = Bx_n + q$ avec B et q fixés au départ de l'algorithme. Il existe des méthodes non stationnaires pour lesquelles l'itération prend la forme $x_{n+1} = B_n x_n + q_n$ avec B_n et q_n qui changent avec les étapes.

Proposition 35.144.

Pour la méthode $x_{n+1} = Bx_n + q$ nous avons équivalence de

- (1) La méthode converge pour tout x_0
- (2) B est une matrice convergente²¹
- (3) $\rho(B) < 1$ (rayon spectral).

De plus si $\|B\| < 1$ alors la méthode converge (quelle que soit la norme algébrique).

La norme d'une matrice (en tout cas, certaines normes) est quelque chose de facile à calculer à l'ordinateur. Typiquement $\|\cdot\|_\infty$ est un simple maximum. Cependant si après avoir calculé $\|B\|_i$ pour des dizaines de normes i différentes, nous avons toujours $\|B\|_i \geq 1$, alors nous ne pouvons rien conclure.

21. C'est-à-dire $\lim_{k \rightarrow \infty} B^k = 0$.

35.20.1 La méthode générale

Nous décomposons la matrice A sous la forme $A = M - N$ avec M inversible. Le système $Ax = b$ devient

$$Mx - Nx = b \quad (35.422)$$

puis $Mx = Nx + b$ et finalement

$$x = M^{-1}Nx + M^{-1}b, \quad (35.423)$$

et voilà une méthode stationnaire avec $B = M^{-1}N$ et $q = M^{-1}b$.

Mais ici nous voyons que M doit être non seulement inversible, mais en plus doit être facilement calculable. En sachant que nous travaillons avec des grandes matrices, il n'est pas question d'inverser M avec une méthode de Gauss.

En bref, il faut choisir M triangulaire parce que c'est en gros la seule que nous pouvons inverser facilement ²².

Remarque 35.145.

La matrice B ne doit pas spécialement être inversible. Si elle ne l'est pas, ce n'est pas un problème.

35.20.2 Jacobi

Nous décomposons

$$A = D - E - F \quad (35.424)$$

où D est la diagonale de A , $-F$ est la partie triangulaire supérieure (sans la diagonale) et $-E$ la triangulaire inférieure (sans la diagonale). Donc D , E et F sont simplement des extractions de parties de la matrice A (et quelques changements de signes).

La méthode de Jacobi prend $M = D$ et $N = (E + F)$. L'inverse de M est facile à calculer parce que M est diagonale. Nous notons B_J la matrice B de la méthode de Jacobi.

Remarque 35.146.

Il se peut que la matrice A ait des zéros sur la diagonale, même si elle est inversible. Et cela est un problème parce qu'alors la matrice D ici construite n'est pas inversible. Dans ce cas, avant de nous lancer dans la méthode de Jacobi, il faut permuter deux lignes de A et donc de b .

Attention cependant que l'on pourrait vouloir effectuer ces permutations en mettant sur la diagonale des nombres les plus grands possibles (parce qu'ensuite, ce qui rentre dans les calculs, c'est D^{-1} qui aura alors des petits nombres). Mais il faut toutefois faire en sorte que le rayon spectral de la matrice B résultante reste plus petit que 1.

Chaque changement dans A induit des changements dans B et donc sur la convergence de la méthode.

35.20.3 Gauss-Seidel

Nous partons de la même décomposition $A = D - E - F$ que dans (35.424). La méthode de Gauss-Seidel prend $M = (D - E)$ et $N = F$.

35.20.4 Autres

Voir la méthode des gradients, et des gradients conjugués.

35.21 Indices connectés, matrice irréductible

Définition 35.147.

Soit $A \in \mathbb{M}(n, \mathbb{R})$ et $i, j \in \{1, \dots, n\}$. Nous disons que les indices i et j sont **directement connectés** si $A_{ij} \neq 0$ ou $A_{ji} \neq 0$.

²². Les matrices orthogonales sont aussi facilement inversibles, mais ne se prêtent pas bien à une décomposition de type somme.

Définition 35.148.

Soit $A \in \mathbb{M}(n, \mathbb{R})$ et $i, j \in \{1, \dots, n\}$. Nous disons que les indices i et j sont **connectés** si il existe un ensemble d'indices $i_0 = i, i_1, \dots, i_{r-1}, i_r = j$ tels que $A_{i_k, i_{k+1}} \neq 0$ pour tout $0 \leq k \leq r$.

Par exemple pour que les indices 1 et 4 soient connectés, on peut avoir les éléments A_{13}, A_{32}, A_{24} non nuls.

Définition 35.149 ([478]).

Une matrice carrée A est **réductible** si il existe une permutation σ telle que

$$\sigma^t A \sigma = \begin{pmatrix} K & L \\ 0 & M \end{pmatrix} \quad (35.425)$$

où K et M sont carrées.

Notons que par définition de la matrice d'une application linéaire,

$$B_{ij} = \langle e_i, B e_j \rangle = \langle e_i, \sigma^t A \sigma e_j \rangle = \langle \sigma e_i, A \sigma e_j \rangle = A_{\sigma(i), \sigma(j)}. \quad (35.426)$$

Proposition 35.150 ([478]).

Soit une matrice carrée A . Les faits suivants sont équivalents :

- (1) A est réductible.
- (2) Il existe une partition non triviale I, J de $\{1, \dots, n\}$ telle que $I \cup J = \{1, \dots, n\}$, $I \cap J = \emptyset$ et pour tout $i \in I$, et pour tout $j \in J$, $A_{ij} = 0$.
- (3) La matrice A admet des indices non connectés (définition 35.148).

Démonstration. Dans plusieurs sens...

(1) implique (2) Nous notons j^* la taille de la matrice K dans (35.425). Nous avons $B_{ij} = 0$ si

$$\begin{cases} J^* + 1 \leq i \leq n \\ 1 \leq j \leq j^* \end{cases} \quad (35.427a)$$

$$\begin{cases} J^* + 1 \leq i \leq n \\ 1 \leq j \leq j^* \end{cases} \quad (35.427b)$$

Donc en posant $I = \sigma\{j^* + 1, \dots, n\}$ et $J = \sigma\{1, \dots, j^*\}$ nous avons une partition non triviale de $\{1, \dots, n\}$ telle que si $i \in I$ et $j \in J$ alors $i = \sigma(i_0)$, $j = \sigma(j_0)$ et

$$A_{ij} = A_{\sigma(i_0), \sigma(j_0)} = B_{i_0, j_0} = 0. \quad (35.428)$$

(2) implique (1) Soit une partition I, J comme indiquée dans l'hypothèse. Soit j^* le nombre d'éléments dans J . Soit σ une permutation de $\{1, \dots, n\}$ telle que $\sigma\{j^* + 1, \dots, n\} = I$ et $\sigma\{1, \dots, j^*\} = J$. Nous posons ensuite $B = \sigma^t A \sigma$. Par construction si $i \in I$ et $j \in J$ alors $A_{ij} = 0$.

Mais si

$$\begin{cases} J^* + 1 \leq i \leq n \\ 1 \leq j \leq j^* \end{cases} \quad (35.429a)$$

$$\begin{cases} J^* + 1 \leq i \leq n \\ 1 \leq j \leq j^* \end{cases} \quad (35.429b)$$

alors $B_{ij} = A_{\sigma(i)\sigma(j)} = 0$. Donc B a la bonne forme.

(3) implique (2) Soient i et j deux indices non connectés : il n'existe pas de chaînes partant de i et arrivant à j . Nous notons I l'ensemble des indices connectés à i , et J les autres. Par hypothèses ces ensembles sont non vides.

Si $k \in i$ et $l \in J$ alors $A_{kl} = 0$ parce que sinon on aurait une chaîne de i à k puis de k à l et donc de i à l , ce qui signifierait que l est connecté à i .

(2) implique (3) Soit une partition I, J comme dans l'hypothèse. Si $j \in J$ est connecté à $i \in I$ alors il existe une chaîne

$$i = i_0, i_1, \dots, i_r = j. \quad (35.430)$$

Si i_s est le premier dans J alors $i_{s-1} \in I$ et $A_{i_{s-1}, i_s} = 0$, ce qui empêche la chaîne de connecter j à i .

□

35.22 Localisation des valeurs propres

Sur l'ensemble $\mathbb{M}(n, \mathbb{R})$ des matrices $n \times n$ à coefficients réels nous introduisons l'ordre partiel²³ donné par $A \geq B$ lorsque $A_{ij} \geq B_{ij}$ pour tout i et j . Nous définissons de façon similaire les relations $A \leq B$, $A < B$ et $A > B$.

Si $x \in \mathbb{R}^n$ nous notons $|x| = (|x_1|, \dots, |x_n|)$ et $x \leq y$ lorsque $x_i \leq y_i$ pour tout i .

Proposition 35.151.

Soit $A \in \mathbb{M}(n, \mathbb{R})$ et $x, y \in \mathbb{R}^n$.

(1) Si $A \geq 0$ et si $x \leq y$ alors $Ax \leq Ay$.

(2) Si $A \geq 0$ alors $Ax \leq |Ax| \leq A|x|$.

Démonstration. Pour la première inégalité, pour tout i et k nous avons $A_{ik}x_k \leq A_{ik}y_k$ et donc

$$(Ax)_i = \sum_k A_{ik}x_k \leq \sum_k A_{ik}y_k = (Ay)_i. \quad (35.431)$$

Pour la seconde, d'abord l'inégalité $Ax \leq |Ax|$ est évidente. Ensuite vu que $A_{ik} \geq 0$ nous avons

$$|Ax|_i = \left| \sum_k A_{ik}x_k \right| \leq \sum_k A_{ik}|x_k| = (A|x|)_i. \quad (35.432)$$

□

Soit une matrice $A \in \mathbb{M}(n, \mathbb{R})$. Nous notons

$$r_i = \sum_{j \neq i} |A_{ij}|. \quad (35.433)$$

Notons la somme sur la ligne i , pas sur la colonne : la somme est horizontale.

Définition 35.152.

Les ensembles

$$D_i = \{z \in \mathbb{C} \text{ tel que } |z - A_{ii}| \leq r_i\} \quad (35.434)$$

sont les **disques de Gershgorin**. Nous allons également noter $B_i = \text{Int}(D_i)$ les boules ouvertes correspondantes.

Théorème 35.153 (Gershgorin).

Soit $A \in \mathbb{M}(n, \mathbb{R})$. Si $\lambda \in \mathbb{C}$ est valeur propre de A alors $\lambda \in D_i$ pour un certain i .

Démonstration. Soit une valeur propre λ et un de ses vecteurs propres $u \in \mathbb{R}^n$: $Au = \lambda u$ avec $u \neq 0$. Soit i un indice réalisant le maximum $|u_i| = \max\{|u_k|\}_k$. Nous écrivons la i ème ligne de $Au = \lambda u$:

$$\sum_k A_{ik}u_k = \lambda u_i, \quad (35.435)$$

c'est-à-dire $A_{ii}u_i + \sum_{k \neq i} A_{ik}u_k = \lambda u_i$, ou encore

$$A_{ii} + \sum_{k \neq i} A_{ik} \frac{u_k}{u_i} = \lambda, \quad (35.436)$$

qui donne

$$|A_{ii} - \lambda| \leq \sum_{k \neq i} |A_{ik}| \frac{|u_k|}{|u_i|} \leq \sum_{k \neq i} |A_{ik}| \quad (35.437)$$

pare que $|u_i| \geq |u_k|$. Notons que sur la ligne précédente, $|\cdot|$ est le module dans \mathbb{C} , pas la valeur absolue dans \mathbb{R} . □

23. Définition 1.8.

Théorème 35.154 (Gershgorin 2[478]).

Soit une matrice irréductible $A \in \mathbb{M}(n, \mathbb{R})$ et une valeur propre λ de A . Si elle est sur la frontière de l'union des disques de Gershgorin, alors elle est sur le bord de tous les disques.

Démonstration. Soit une valeur propre λ de A telle que $\lambda \in \partial(\bigcup_i D_i)$. Alors λ n'est dans l'intérieur d'aucune boule et nous avons $|\lambda - A_{ii}| \geq r_i$ pour tout i .

Soit un vecteur propre u de A tel que $\|u\|_\infty = 1$. Nous posons $I = \{1 \leq i \leq n \text{ tel que } |u_i| = 1\}$ et $J = \{1 \leq j \leq n \text{ tel que } |u_j| < 1\}$. Par hypothèse I n'est pas vide, et de plus $I \cap J = \emptyset$ et $I \cup J = \{1, \dots, n\}$ parce qu'aucune composante de u n'a un module²⁴ plus grand que 1.

La i^{e} composante de la relation $Au = \lambda u$ peut s'écrire

$$(A_{ii} - \lambda)u_i + \sum_{k \neq i} A_{ik}u_k = 0. \quad (35.438)$$

Forts de cela nous écrivons les inégalités suivantes :

$$r_i \leq |\lambda - A_{ii}| = |(\lambda - A_{ii})u_i| = \left| \sum_{k \neq i} A_{ik}u_k \right| \leq \sum_{k \neq i} |A_{ik}| |u_k| \leq \sum_{k \neq i} |A_{ik}| = r_i. \quad (35.439)$$

Donc les inégalités sont des égalités :

$$r_i = |\lambda - A_{ii}| = |(\lambda - A_{ii})u_i| = \left| \sum_{k \neq i} A_{ik}u_k \right| = \sum_{k \neq i} |A_{ik}| |u_k| = \sum_{k \neq i} |A_{ik}|. \quad (35.440)$$

En particulier l'égalité $\sum_{k \neq i} |A_{ik}| |u_k| = \sum_{k \neq i} |A_{ik}|$ donne

$$\sum_{k \neq i} |A_{ik}| (|u_k| - 1) = 0. \quad (35.441)$$

Donc pour tout $k \in J$ nous avons $A_{ik} = 0$. Vu que A est irréductible, cela donnerait une partition impossible $\{1, \dots, n\} = I \cup J$. Nous en déduisons que J est vide et donc que $|u_j| = 1$ pour tout j . En repartant de (35.440) nous avons alors

$$r_i = |(\lambda - A_{ii})u_i| = |\lambda - A_{ii}| |u_i| = |\lambda - A_{ii}|. \quad (35.442)$$

Cela prouve que $\lambda \in \partial D_i$ pour tout i . □

Exemple 35.155

Soit la matrice

$$B = \begin{pmatrix} 2 & 0 & 1 \\ 0 & 1 & -1/2 \\ -1 & 0 & 3 \end{pmatrix}. \quad (35.443)$$

D'abord nous rappelons que si vous voulez entrer cette matrice dans Sage (ou plus généralement dans Python2²⁵), vous devez faire attention au 1/2 qui, tel quel, est évalué à 0. Nous vous rappelons donc que tous vos codes Sage doivent commencer par ceci :

```

1 # -*- coding: utf8 -*-
2 from __future__ import unicode_literals
3 from __future__ import division

```

tex/sage/sageSnip015.sage

24. Les composantes de u sont a priori dans \mathbb{C} , et non spécialement dans \mathbb{R} , même si A est une matrice réelle.

25. Que vous n'avez aucune raison d'utiliser autre que Sage.

Les éléments non nuls hors diagonale sont B_{13} , B_{31} et B_{23} . Elle n'est donc pas irréductible; nous avons par exemple la partition $I = \{1, 3\}$, $J = \{2\}$ pour le critère de la proposition 35.150(2).

Les disques de Gershgorin sont

$$D_1 = \{z \in \mathbb{C} \text{ tel que } |z - 2| \leq 1\} \quad (35.444a)$$

$$D_2 = \{z \in \mathbb{C} \text{ tel que } |z - 1| \leq 1/2\} \quad (35.444b)$$

$$D_3 = \{z \in \mathbb{C} \text{ tel que } |z - 3| \leq 1\} \quad (35.444c)$$

Les valeurs propres de la matrice sont sur des bords de disques de Gershgorin, sans être sur tous les bords, comme ça aurait été le cas par le théorème 35.154 si la matrice avait été irréductible. Elles sont sur la figure 35.1; notez en particulier les valeurs propres λ_2 et λ_3 qui sont sur le bord de deux disques mais pas sur le bord des trois disques en même temps.

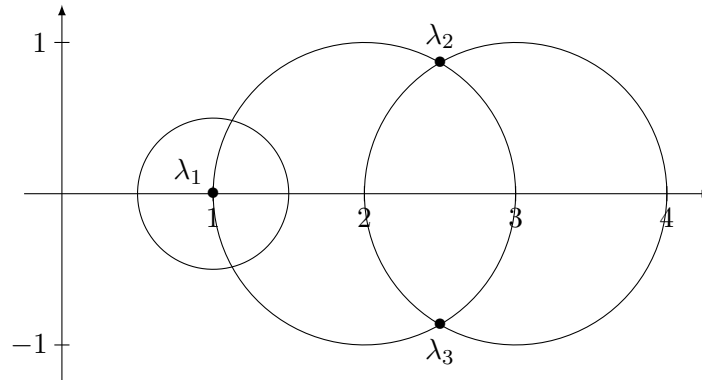


FIGURE 35.1 – Les disques de Gershgorin et les valeurs propres pour l'exemple 35.155.

△

Exemple 35.156

Soit la matrice

$$A = \begin{pmatrix} 0 & -1 & 0 \\ 0 & 1 & 2 \\ 3 & 0 & 2 \end{pmatrix}. \quad (35.445)$$

Nous avons

$$D_1 = \{z \in \mathbb{C} \text{ tel que } |z| \leq 1\}, \quad (35.446a)$$

$$D_2 = \{z \in \mathbb{C} \text{ tel que } |z - 1| \leq 2\}, \quad (35.446b)$$

$$D_3 = \{z \in \mathbb{C} \text{ tel que } |z - 2| \leq 3\}, \quad (35.446c)$$

Le polynôme caractéristique est

$$\chi(\lambda) = -\lambda^3 + 3\lambda^2 - 2\lambda - 6. \quad (35.447)$$

Une fois remarqué que $\lambda_1 = -1$ est une racine, les autres sont faciles à trouver (division euclidienne de $\chi(\lambda)$ par $\lambda + 1$): $\lambda_2 = 2 + i\sqrt{2}$ et $\lambda_3 = 2 - i\sqrt{2}$.

La matrice A est irréductible. En effet les éléments non diagonaux non nuls sont A_{12} , A_{23} et A_{31} . Ils peuvent former une chaîne reliant tous les indices entre eux.

Les contraintes sur la localisation des valeurs propres est donc qu'elles doivent être dans ou sur les disques de Gershgorin, mais que celles qui sont sur le bord d'un disque doivent être sur le bord de tous les disques en même temps. C'est cela que nous observons sur la figure 35.2. Notez en particulier la position de la valeur propre λ_1 .

△

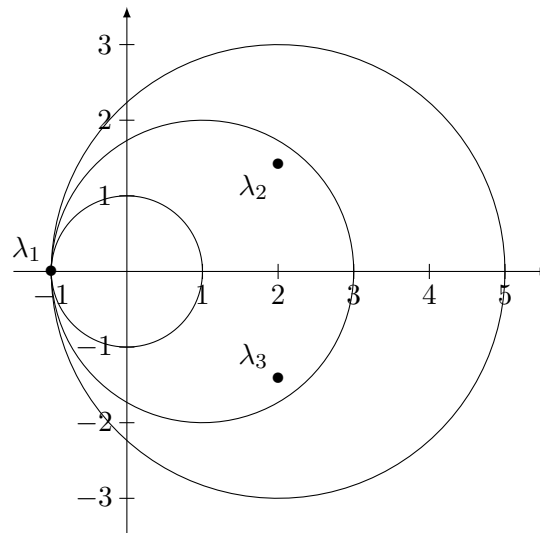


FIGURE 35.2 – Les disques de Gershgorin et les valeurs propres pour l'exemple 35.156.

35.22.1 Matrices à diagonale dominante

Définition 35.157 ([479, 478]).

La matrice $A \in \mathbb{M}(n, \mathbb{K})$ est à **diagonale dominante** si pour tout i ,

$$\sum_{j \neq i} |A_{ij}| \leq |A_{ii}| \quad (35.448)$$

où $|\cdot|$ est la module dans \mathbb{C} ou la valeur absolue dans \mathbb{R} .

Elle est à **diagonale fortement dominante** si elle est à diagonale dominante et si il existe un i tel que

$$\sum_{j \neq i} |A_{ij}| < |A_{ii}|. \quad (35.449)$$

Elle est à **diagonale strictement dominante** si

$$\sum_{j \neq i} |A_{ij}| < |A_{ii}| \quad (35.450)$$

pour tout i (entier entre 1 et n).

Nous avons les inclusions suivantes :

$$\text{strictement dominante} \subset \text{fortement dominante} \subset \text{dominante}. \quad (35.451)$$

Lemme 35.158.

Si A est dans un des deux cas suivant :

- diagonale strictement dominante,
- diagonale dominante et irréductible²⁶

alors $A_{ii} \neq 0$ pour tout i .

Démonstration. Si A est à diagonale strictement dominante, alors l'inégalité stricte (35.450) n'est pas possible.

Si A est à diagonale dominante, alors si $A_{ii} = 0$, toute la ligne est nulle. Dans ce cas, la matrice ne peut pas être irréductible. \square

Proposition 35.159 ([479]).

Une matrice à diagonale strictement dominante est inversible.

26. Définition 35.149.

Démonstration. Soit une matrice A à diagonale strictement dominante. Soit x tel que $Ax = 0$. Le but est de montrer que $x = 0$. Soit un indice i_0 réalisant la norme maximum :

$$|x_{i_0}| = \|x\|_\infty. \quad (35.452)$$

Nous écrivons la composante i_0 de l'égalité $Ax = 0$:

$$\sum_k A_{i_0 k} x_k = 0, \quad (35.453)$$

et nous séparons le terme $k = i_0$ des autres :

$$\sum_{k \neq i_0} A_{i_0 k} x_k + A_{i_0 i_0} x_{i_0} = 0. \quad (35.454)$$

Nous prenons le module et majorons les sommes :

$$|A_{i_0 i_0}| |x_{i_0}| \leq \sum_{k \neq i_0} |A_{i_0 k}| |x_k| \leq \sum_{k \neq i_0} |A_{i_0 k}| |x_{i_0}|. \quad (35.455)$$

Si $|x_{i_0}|$ est non nul nous pouvons simplifier :

$$|A_{i_0 i_0}| \leq \sum_{k \neq i_0} |A_{i_0 k}|. \quad (35.456)$$

Hélas, l'hypothèse de diagonale strictement dominante implique l'inégalité stricte dans le sens inverse. Impossible. Nous en déduisons que $|x_{i_0}| = 0$. Donc $\|x\|_\infty = 0$, ce qui signifie que $x = 0$.

Le fait que le noyau de A se réduise à $\{0\}$ implique l'inversibilité de A . \square

Proposition 35.160.

Soit $A \in \mathbb{M}(n, \mathbb{C})$ une matrice qui est dans un des deux cas suivants :

- à diagonale strictement dominante
- à diagonale dominante et irréductible

Si $A = D - M$ où D est la diagonale de A (et M est « le reste ») alors D est inversible et

$$\rho(D^{-1}M) < 1 \quad (35.457)$$

où ρ est le rayon spectral (thème 41).

Démonstration. Le lemme 35.158 nous dit que les éléments diagonaux de A sont non nuls. Cela donne déjà le fait que la matrice D est inversible et que la produit $D^{-1}M$ ait un sens. Nous posons $T = D^{-1}M$. Nous avons alors

$$T_{ii} = \sum_k (D^{-1})_{ik} M_{ki}. \quad (35.458)$$

Si $k = i$ alors $M_{ki} = 0$ et si $k \neq i$ alors $D_{ik} = 0$. Donc $T_{ii} = 0$ pour tout i .

En ce qui concerne les autres éléments de T ,

$$T_{ij} = \sum_k (D^{-1})_{ik} M_{kj} = \sum_k \frac{1}{A_{ik}} \delta_{ik} M_{kj} = -\frac{A_{ij}}{A_{ii}}. \quad (35.459)$$

Notes :

- Les hypothèses sur A jouent pour dire que $A_{ii} \neq 0$.
- Le signe moins est dû au fait que $M_{ij} = -A_{ij}$ lorsque $i \neq j$.

En faisant la somme des modules :

$$\sum_{j \neq i} |T_{ij}| = \sum_{j \neq i} \frac{|A_{ij}|}{|A_{ii}|} = \frac{1}{|A_{ii}|} \sum_{j \neq i} |A_{ij}| \leq 1. \quad (35.460)$$

La dernière inégalité est le fait que A soit à diagonale dominante.

Si A est à diagonale strictement dominante Alors nous avons l'inégalité stricte

$$\sum_{j \neq i} |T_{ij}| < 1. \quad (35.461)$$

Et le théorème de Gershgorin 35.153 dit que le spectre de T est contenu dans l'union des disques

$$D_i = \{z \in \mathbb{C} \text{ tel que } |z - T_{ii}| \leq r_i\} \quad (35.462)$$

où

$$r_i = \sum_{j \neq i} |T_{ij}|. \quad (35.463)$$

Mais nous avons prouvé que pour tout i , $T_{ii} = 0$ et $\sum_{j \neq i} |T_{ij}| < 1$. Donc toutes ces boules sont contenues dans $B(0, 1)$. Cela prouve que $\rho(T) < 1$.

Diagonale dominante, irréductible La matrice T est alors également irréductible parce que les éléments non nuls de A et de T sont les mêmes : $T_{ij} = -A_{ij}/A_{ii}$. Nous utilisons alors le second théorème de Gershgorin 35.154. Si λ est une valeur propre de T , alors soit

$$\lambda \in \bigcup_i B(0, r_i) \quad (35.464)$$

soit

$$\lambda \in \bigcap_i \partial B(0, r_i). \quad (35.465)$$

Vu que $r_i \leq 1$ pour tout i , dans le premier cas λ est dans l'union des boules *ouvertes* de rayon 1. Le nombre λ est donc une la boule ouverte de rayon 1. Bref, $|\lambda| < 1$.

Dans le second cas, l'intersection de deux cercles de même centre sont soit vide soit tout le cercle (auquel cas les rayons sont égaux). Dans le second cas, ledit rayon est certainement strictement plus petit que 1 parce que

$$r_i = \sum_{j \neq i} |T_{ij}| = \sum_{j \neq i} \frac{|A_{ij}|}{|A_{ii}|} < 1. \quad (35.466)$$

□

35.22.2 M-matrice

Définition 35.161.

Une matrice $A \in \mathbb{M}(n, \mathbb{R})$ est une **M-matrice** si

- (1) $A_{ii} > 0$ pour tout i ,
- (2) $A_{ij} \leq 0$ si $i \neq j$
- (3) A est inversible et $A^{-1} \geq 0$.

Proposition 35.162.

Soit $A \in \mathbb{M}(n, \mathbb{R})$ telle que $A_{ii} > 0$ pour tout i et $A_{ij} \leq 0$ pour tout $i \neq j$. Nous posons $A = D - M$ où D est la diagonale de A .

La matrice A est une M-matrice si et seulement si $\rho(D^{-1}M) < 1$.

Démonstration. En deux morceaux.

Si $\rho(D^{-1}M) < 1$ Nous posons encore $T = D^{-1}M$. Par le théorème 16.113, la matrice $\mathbb{1} - T$ est inversible et

$$(\mathbb{1} - T)^{-1} = \sum_{k=0}^{\infty} T^k. \quad (35.467)$$

D'autre part, via des calculs déjà faits, et les hypothèses sur les signes des éléments de A ,

$$T_{ij} = -\frac{A_{ij}}{A_{ii}} \geq 0. \quad (35.468)$$

Donc tous les éléments de T sont positifs (ou nuls). Par conséquent $T^k \geq 0$ pour tout k et $(\mathbb{1} - T)^{-1}$ est positive.

Mais $A = D - M = D(\mathbb{1} - D^{-1}M) = D(\mathbb{1} - T)$. Vu que D et $\mathbb{1} - T$ sont inversibles, nous savons que A est inversible et

$$A^{-1} = (\mathbb{1} - A)^{-1}D^{-1}, \quad (35.469)$$

qui est un produit de matrices positives. Donc $A^{-1} \geq 0$.

Au final, A est une M-matrice.

Si A est une M-matrice Soit une valeur propre λ de $T = D^{-1}M$ est un vecteur propre u : $Tu = \lambda u$. Vu que $T \geq 0$ nous avons d'une part $|\lambda u| = |\lambda|u|$ et d'autre part $|\lambda u| = |Tu| \leq T|u|$, ce qui donne

$$|\lambda||u| \leq T|u|. \quad (35.470)$$

Dans cette inégalité nous substituons T par $\mathbb{1} - (\mathbb{1} - T)$ pour avoir

$$|\lambda||u| \leq |u| - (\mathbb{1} - T)|u| \quad (35.471)$$

ou encore

$$(\mathbb{1} - T)|u| \leq (1 - |\lambda|)|u|. \quad (35.472)$$

Mais $(\mathbb{1} - T)^{-1} = A^{-1}D \geq 0$ parce que A et D sont positives. Donc en appliquant $(\mathbb{1} - T)^{-1}$ à l'inégalité (35.472), elle est conservée (proposition 35.151(2)) :

$$|u| \leq (\mathbb{1} - T)^{-1}(1 - |\lambda|)|u|. \quad (35.473)$$

Si $|\lambda| \geq 1$ alors toutes les composantes de $(1 - |\lambda|)|u|$ sont négatives et l'inégalité n'est possible qu'avec $|u| = 0$. Dans ce cas, λ n'est pas une valeur propre (le vecteur propre soit être non nul).

Nous en déduisons que $|\lambda| < 1$ et donc que $\rho(T) = \rho(D^{-1}M) < 1$.

□

Le théorème suivant résume ce que nous avons vu en donnant une condition suffisante facile à vérifier pour être une M-matrice.

Théorème 35.163.

Soit $A \in \mathbb{M}(n, \mathbb{R})$ telle que

- (1) $A_{ii} > 0$
- (2) $A_{ij} \leq 0$ pour $i \neq j$
- (3) vérifiant une des deux conditions suivantes :
 - à diagonale strictement dominante
 - à diagonale dominante et irréductible.

Alors A est une M-matrice.

Démonstration. Au vu de la proposition 35.162, il suffira de montrer que $\rho(D^{-1}M) < 1$ où D et M sont la décomposition $A = D - M$ habituelle. C'est le cas grâce à la proposition 35.160. □

Proposition 35.164.

Soit $A \in \mathbb{M}(n, \mathbb{R})$, une M-matrice irréductible. Alors $A^{-1} > 0$.

Démonstration. Nous posons $T = D^{-1}M$. En comparant la définition 35.161 de M-matrice et la caractérisation de la proposition 35.162, nous avons $\rho(D^{-1}M) < 1$. Par conséquent

$$(\mathbb{1} - T)^{-1} = \sum_{k=0}^{\infty} T^k \quad (35.474)$$

par la proposition 35.162. D'autre part, $A^{-1} = (\mathbb{1} - A)^{-1}D$ où les éléments D sont strictement positifs. Donc nous devons encore prouver que $(\mathbb{1} - T)^{-1} > 0$. Nous savons que $T \geq 0$, et vu que

$$\left(\sum_k T^k\right)_{ij} = \sum_k (T^k)_{ij} \quad (35.475)$$

il nous suffit de prouver que pour chaque (ij) , un des $(T^k)_{ij}$ est strictement positif. Soient donc deux indices i et j . Vu que A est irréductible, ils sont connectés par une suite d'indice $i = i_0, i_1, \dots, i_r = j$ tels que

$$T_{i_k, i_{k+1}} = -\frac{A_{i_k, i_{k+1}}}{A_{i_k, i_k}} > 0. \quad (35.476)$$

Or les indices i_k sont choisis de telle sorte que les numérateurs soient non nuls et donc strictement négatifs. Nous avons, en général :

$$(T^k)_{ij} = \sum_{l_1, \dots, l_{r-1}} T_{i, l_1} T_{l_1, l_2} \dots T_{l_{r-1}, j}. \quad (35.477)$$

Chacun des termes est positif ou nul, mais pour $k = r$, il y a entre autres le terme

$$T_{i, i_1} T_{i_1, i_2} \dots T_{i_r, j} \neq 0. \quad (35.478)$$

Donc $(T^r)_{ij} > 0$ et $\sum_{k=0}^{\infty} (T^k)_{ij} > 0$. Et par conséquent

$$A^{-1} = (\mathbb{1} - T)^{-1}D > 0. \quad (35.479)$$

□

Théorème 35.165.

Soit une M-matrice $A \in \mathbb{M}(n, \mathbb{R})$ et $g \in \mathbb{R}^n$ tel que $(Ag)_i \geq 1$ pour tout i . Alors $\|A^{-1}\|_{\infty} \leq \|g\|_{\infty}$.

Démonstration. Nous posons $u = (1, \dots, 1)$ et considérons $x \in \mathbb{R}^n$. Vu que A est une M-matrice, nous avons $A^{-1} \geq 0$, donc

$$|A^{-1}x| \leq A^{-1}|x| \leq \|x\|_{\infty} A^{-1}u \leq \|x\|_{\infty} g. \quad (35.480)$$

Justifications :

- La première inégalité est la proposition 35.151(2).
- La seconde provient de

$$(B|x|)_i = \sum_k B_{ik} |x_k| \leq \sum_k B_{ik} \|x\|_{\infty} = \|x\|_{\infty} \sum_k B_{ik} u_k = \|x\|_{\infty} Bu. \quad (35.481)$$

- Étant donné que $A^{-1} \geq 0$ nous conservons l'inégalité et $Ag \geq u$ implique $g \geq A^{-1}u$ (c'est la proposition 35.151(1)).

En ce qui concerne la norme de A^{-1} nous avons donc

$$\|A^{-1}\|_{\infty} = \sup_{|x|_{\infty}=1} \|A^{-1}x\|_{\infty} \leq \sup_{\|x\|_{\infty}=1} \|x\|_{\infty} \|g\|_{\infty} = \|g\|_{\infty}. \quad (35.482)$$

□

Proposition 35.166.

Une matrice de $\mathbb{M}(n, \mathbb{R})$ qui

- (1) est symétrique,
 (2) Vérifie une des deux conditions suivantes
 — est irréductible à diagonale fortement dominante
 — est à diagonale strictement dominante,
 (3) vérifie $A_{ii} > 0$ pour tout i

est strictement définie positive.

Démonstration. D'après le théorème de Gershgorin 35.153, chaque valeur propre de A est dans un des disques fermés

$$D_i = \{z \in \mathbb{C} \text{ tel que } |z - A_{ii}| \leq r_i\}. \quad (35.483)$$

Par hypothèse, les centres de ces disques sont réels et strictement positifs. Mais le fait que A soit à diagonale dominante donne que le rayon de ces cercles sont plus petits que A_{ii} . Donc D_i n'intersecte pas $]-\infty, 0[$. Mais le fait que A soit symétrique implique que les valeurs propres soient réelles (théorème 11.189(1)). Cela montre que les valeurs propres de A sont toutes dans $[0, \infty[$.

Si la matrice A est à diagonale strictement dominante, alors les inégalités sont strictes et le théorème est prouvé.

Sinon nous sommes dans le cas irréductible à diagonale fortement dominante et nous avons le théorème de Gershgorin numéro 2 35.154. Soit une valeur propre λ . Soit elle est dans un des disques ouvert (qui est inclus dans $]0, \infty[$), soit elle est dans l'intersection des bords des disques. Mais au moins un des disques n'intersecte pas 0 (parce que la diagonale est strictement dominante). Dans ce cas non plus λ ne peut pas être nul.

Nous en déduisons que dans tous les cas, les valeurs propres sont toutes réelles strictement positives. \square

Chapitre 36

Méthode des différences finies

36.1 Problèmes de dimension un

Soit $u: \mathbb{R} \rightarrow \mathbb{R}$ une fonction et $h > 0$. Nous définissons les opérations suivantes (qui sont supposées approximer la dérivée $u'(x)$ lorsqu'elle existe).

Définition 36.1.

La *différence progressive* est

$$(D_h^+ u)(x) = \frac{u(x+h) - u(x)}{h}, \quad (36.1)$$

la *différence régressive* est

$$(D_h^- u)(x) = \frac{u(x) - u(x-h)}{h}, \quad (36.2)$$

la *différence centrée* est

$$(D_h^0 u)(x) = \frac{u(x+h) - u(x-h)}{2h}. \quad (36.3)$$

Nous ne noterons pas toujours la dépendance en h , c'est-à-dire que nous noterons $D^+ u$ au lieu de $D_h^+ u$ lorsque cela ne pose pas de problèmes.

Notons que u'' peut être approximé par $D^+ D^+ u$, $D^0 D^+$, $D^+ D^-$, et encore de nombreuses autres possibilités.

Voici un lemme qui dit que tout cela n'est pas si mal, pourvu que u soit assez régulière.

Lemme 36.2.

Soit un ouvert connexe Ω de \mathbb{R} , soit $x \in \Omega$ et $h > 0$ tel que $\overline{B(x, h)} \subset \Omega$.

(1) Si $u \in C^2(\Omega)$ alors

$$|u'(x) - D^+ u(x)| \leq \frac{h}{2} \|u''\|_{\bar{\Omega}} \quad (36.4)$$

et

$$|u'(x) - D^- u(x)| \leq \frac{h}{2} \|u''\|_{\bar{\Omega}}. \quad (36.5)$$

(2) Si $u \in C^3(\bar{\Omega})$ alors

$$|u'(x) - D^0(x)| \leq \frac{h^2}{2} \|u^{(3)}\|_{\bar{\Omega}} \quad (36.6)$$

(3) Si $u \in C^4(\bar{\Omega})$ alors

$$|u''(x) - D^- D^+ u(x)| \leq \frac{h^2}{12} \|u^{(4)}\|_{\bar{\Omega}}. \quad (36.7)$$

Démonstration. Nous prouvons le point (3). D'abord nous regardons de quoi nous avons besoin :

$$D^- D^+ u(x) = \frac{(D^+ u)(x) - (D^+ u)(x-h)}{h} = \frac{u(x+h) - 2u(x) + u(x-h)}{h^2} \quad (36.8)$$

Nous allons y mettre les approximations de $u(x+h)$ et $u(x-h)$ par Taylor, proposition 13.370 :

$$u(x+h) = u(x) + hu'(x) + \frac{h^2}{2}u''(x) + \frac{h^3}{6}u'''(x) + \frac{h^4}{24}u^{(4)}(x+\theta_1h) \quad (36.9)$$

avec $\theta_1 \in [0, 1]$. De même,

$$u(x-h) = u(x) - hu'(x) + \frac{h^2}{2}u''(x) - \frac{h^3}{6}u'''(x) + \frac{h^4}{24}u^{(4)}(x-\theta_2h) \quad (36.10)$$

avec $\theta_2 \in [0, 1]$.

Donc

$$u(x+h) + u(x-h) - 2u(x) = h^2u''(x) + \frac{h^4}{4!} \left(u^{(4)}(x+\theta_1h) + u^{(4)}(x-\theta_2h) \right), \quad (36.11)$$

ce qui donne

$$(D^-D^+u)(x) = u''(x) + \frac{h^2}{4!} \left(u^{(4)}(x+\theta_1h) + u^{(4)}(x-\theta_2h) \right). \quad (36.12)$$

Chacun des deux termes dans la parenthèse peut être majoré par $\|u^{(4)}\|_{\bar{\Omega}}$. Notons que c'est une majoration très sauvage parce que $x+\theta_1h$ ne prend ses valeurs que dans $[x, x+h]$ avec h supposé petit. Quoiqu'il en soit nous ne pouvons pas dire mieux que

$$|u''(x) - D^-D^+u(x)| \leq \frac{h^2}{12} \|u^{(4)}\|_{\bar{\Omega}}. \quad (36.13)$$

□

Remarque 36.3 ([1]).

Lorsque nous écrivons

$$|u'(x) - D^+u(x)| \leq \delta \quad (36.14)$$

pour tout x , nous ne pouvons pas écrire

$$\|u' - D^+u\|_{\infty} \leq \delta \quad (36.15)$$

parce que l'inégalité (36.14) n'est valable que pour les x tels que $[x-h, x+h] \subset \Omega$, de telle sorte que l'inégalité n'est pas spécialement correcte sur $\bar{\Omega}$.

36.1.1 Un schéma à cinq points

36.1.1.1 Poser le système

Soit $\Omega =]0, 1[$ et l'équation différentielle

$$\begin{cases} -u''(x) + c(x)u(x) = f & \text{sur } \Omega \\ u(0) = \alpha \\ u(1) = \beta \end{cases} \quad (36.16)$$

où c est une fonction positive et $\alpha, \beta \in \mathbb{R}$. Nous considérons $h > 0$ assez petit pour que le reste ait un sens. Si nous cherchons des solutions dans $C^4(\bar{\Omega})$, le lemme 36.2 nous dit que

$$|u''(x) - D^-D^+u(x)| = \eta(h^2) \quad (36.17)$$

où η est une fonction telle que $\lim_{t \rightarrow 0} \eta(t) = 0$. Nous pouvons récrire l'équation différentielle sous la forme

$$-D^-D^+u(x) + c(x)u(x) = f(x) + \eta(h^2). \quad (36.18)$$

Si nous négligeons le terme $\eta(h^2)$ qui est supposé être petit nous pouvons tenter de résoudre pour la fonction u_h

$$-D^-D^+u_h(x) + c(x)u_h(x) = f(x). \quad (36.19)$$

Notons ici l'importance de la notion de problème bien posé parce que en remplaçant le paramètre (fonctionnel) f par $f + \eta(h^2)$, nous modifions les solutions. Dans la mesure où le problème est bien posé, cette petite modification ne modifiera pas trop la solution et nous pouvons espérer que $\|u - u_h\|$ soit petit pour une norme ou une autre.

Utilisant l'expression (36.8) pour D^-D^+ nous avons l'équation suivante pour u_h :

$$\frac{1}{h^2} \left(2u_h(x) - u_h(x+h) - u_h(x-h) \right) + c(x)u_h(x) = f(x). \quad (36.20)$$

Avons-nous gagné quelque chose? Pas encore. L'idée de la discrétisation est de ne considérer u_h qu'en certains points, et de prendre ces points à intervalles réguliers de taille h . Soient donc N un nombre entier et $h = 1/(N+1)$. Nous posons

$$x_k = kh \quad (36.21)$$

pour $i = 0, \dots, N+1$. Avec cela nous avons

$$\bar{\Omega} = \bigcup_{k=0}^{N-1} [x_k, x_{k+1}] \quad (36.22a)$$

$$x_0 = 0 \quad (36.22b)$$

$$x_{N+1} = 1. \quad (36.22c)$$

Nous posons surtout

$$\Omega_h = \{x_i\}_{i=1, \dots, N} \quad (36.23)$$

et

$$\bar{\Omega}_h = \{x_i\}_{i=0, \dots, N+1}. \quad (36.24)$$

Enfin, nous ne considérons plus u_h que comme une fonction $u_h: \bar{\Omega}_h \rightarrow \mathbb{R}$. C'est-à-dire que u_h est un vecteur à $N+2$ composantes.

L'équation (36.20) devient

$$\frac{1}{h^2} \left(2u_h(x_i) - u_h(x_{i+1}) - u_h(x_{i-1}) \right) + c(x_i)u_h(x_i) = f(x_i) \quad (36.25)$$

pour $i = 1, \dots, N$. Sur les bords, cette équation n'est pas possible parce que x_{i-1} ou x_{i+1} n'existerait pas. Au contraire, sur les bords nous avons les conditions aux bords

$$u_h(x_0) = \alpha \quad (36.26)$$

et

$$u_h(x_{N+1}) = \beta. \quad (36.27)$$

En posant $c_i = c(x_i)$ et $u_i = u_h(x_i)$, les inconnues du problème sont les nombres u_i ($i = 1, \dots, N+1$). Elles sont immédiatement résolues pour u_0 et u_{N+1} . Pour les autres, il faut écrire et résoudre un système d'équation linéaire.

L'écriture du système linéaire à résoudre consiste essentiellement à écrire (36.25) en séparant les cas $i = 0$ et $i = N$ parce que nous savons déjà les valeurs de u_0 et u_N . Le système que nous avons est :

$$\begin{cases} \left(\frac{2}{h^2} + c_1 \right) u_1 - \frac{1}{h^2} u_2 = f_1 + \frac{\alpha}{h^2} & i = 1 \\ \left(\frac{2}{h^2} + c_N \right) u_N - \frac{1}{h^2} u_{N-1} = f_N + \frac{\beta}{h^2} & i = N \\ \left(\frac{2}{h^2} + c_i \right) u_i - \frac{1}{h^2} u_{i+1} - \frac{1}{h^2} u_{i-1} = f_i. & \text{autres} \end{cases} \quad (36.28a)$$

$$\left(\frac{2}{h^2} + c_N \right) u_N - \frac{1}{h^2} u_{N-1} = f_N + \frac{\beta}{h^2} \quad i = N \quad (36.28b)$$

$$\left(\frac{2}{h^2} + c_i \right) u_i - \frac{1}{h^2} u_{i+1} - \frac{1}{h^2} u_{i-1} = f_i. \quad \text{autres} \quad (36.28c)$$

Cela se met sous la forme matricielle

$$L_h U_h = F_h \quad (36.29)$$

pour

$$F_h = \left(f_1 + \frac{\alpha}{h^2}, f_2, \dots, f_{N-1}, f_N + \frac{\beta}{h^2} \right) \quad (36.30)$$

et les éléments non nuls de L_h sont :

$$(L_h)_{i,i-1} = -\frac{1}{h^2} \quad \text{pour } i = 2, \dots, N \quad (36.31a)$$

$$(L_h)_{i,i+1} = -\frac{1}{h^2} \quad \text{pour } i = 1, \dots, N-1 \quad (36.31b)$$

$$(L_h)_{i,i} = \frac{2}{h^2} + c_i \quad \text{pour } i = 1, \dots, N. \quad (36.31c)$$

Cette matrice est pleine de zéros, à part les trois diagonales centrales, et il existe des méthodes efficaces pour résoudre le système d'équation correspondant.

36.1.1.2 Propriétés du système

La matrice est la suivante :

$$L_h = \begin{pmatrix} \frac{2}{h^2} + c_1 & -1/h^2 & 0 & 0 & \dots & 0 \\ -1/h^2 & \frac{2}{h^2} + c_2 & -1/h^2 & 0 & \dots & 0 \\ 0 & -1/h^2 & \frac{2}{h^2} + c_3 & -1/h^2 & \ddots & \vdots \\ 0 & 0 & -1/h^2 & \ddots & \ddots & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & -1/h^2 \\ 0 & 0 & \dots & 0 & -1/h^2 & \frac{2}{h^2} + c_N \end{pmatrix} \quad (36.32)$$

où nous avons déjà posé l'hypothèse $c_i \geq 0$ pour tout i .

Lemme 36.4.

La matrice L_h est irréductible¹ à diagonale fortement dominante².

Démonstration. Nous décomposons la preuve en plusieurs parties.

La première ligne Sur la première ligne, seuls deux éléments sont non nuls et nous avons

$$L_{11} = \frac{2}{h^2} + c_1 \geq \frac{2}{h^2} > \frac{1}{h^2} = L_{12}. \quad (36.33)$$

La dernière ligne Elle est semblable à la première.

Les autres lignes Sur les autres lignes nous avons trois éléments non nuls et

$$\sum_{j \neq i} |A_{ij}| = \frac{2}{h^2} \leq \frac{2}{h^2} + c_i = L_{ii}. \quad (36.34)$$

Diagonale fortement dominante Nous avons prouvé jusqu'à présent que L_h était une matrice à diagonale fortement dominante.

Irréductible Nous allons utiliser la caractérisation de la proposition 35.150(3). Pour cela nous considérons le chaîne d'éléments non nuls

$$A_{12}, A_{23}, \dots, A_{N-1,N} = -\frac{1}{h^2}. \quad (36.35)$$

Soient deux indices i et j avec $i < j$. Cette suite d'indices (ou un sous-suite) rend i et j connectés.

Si par contre $i > j$, il faut considérer la suite inversée grâce au fait que L_h est symétrique :

$$A_{N,N-1}, A_{N-1,N-2}, \dots, A_{32}, A_{21} = -\frac{1}{h^2}. \quad (36.36)$$

1. Caractérisation 35.150.

2. Définition 35.157.

□

Proposition 36.5.

Soit le problème

$$\begin{cases} -u''(x) + c(x)u(x) = f & \text{sur } \Omega \\ u(0) = \alpha \\ u(1) = \beta \end{cases} \quad (36.37)$$

où c est une fonction positive et $\alpha, \beta \in \mathbb{R}$. Nous considérons $h > 0$ assez petit pour que le reste ait un sens. Et nous approximations u'' par D^-D^+u .

La matrice L_h des différences finies associée à ce problème est

- (1) une M-matrice,
- (2) strictement définie positive,
- (3) d'inverse $L_h^{-1} > 0$.

Démonstration. Le théorème 35.163 dit que L_h est une M-matrice. La Proposition 35.166 nous donne aussi que L_h est strictement définie positive.

Le lemme 36.4 dit que L_h est irréductible, ce qui permet à la proposition 35.164 de conclure que $L_h^{-1} > 0$. □

Cela étant rappelé, nous pouvons continuer.

Lemme 36.6.

Soit $\Omega =]0, 1[$, soit $N \in \mathbb{N}$ et $h = 1/(N + 1)$. La solution $w_h: \Omega_h \rightarrow \mathbb{R}$ du problème discrétisé

$$\begin{cases} -(D^-D^+w_h)(x_i) = 1 & (36.38a) \\ w_h(0) = 0 & (36.38b) \\ w_h(1) = 0 & (36.38c) \end{cases}$$

pour tout $x_i = ih$ ($i = 1, \dots, N$) donne les valeurs exactes des $w(x_i)$ lorsque w est la solution de

$$\begin{cases} -w''(x) = 1 & (36.39a) \\ w(0) = 0 & (36.39b) \\ w(1) = 0. & (36.39c) \end{cases}$$

Démonstration. Un enseignement de la proposition 36.5 est que le système (36.38) peut être écrit sous la forme d'un système linéaire $L_h^0 w_h = F_h$ où L_h^0 est inversible. Il y a donc unicité de la solution.

D'autre part, la solution du système (36.39) est $w(x) = \frac{1}{2}(x - x^2)$, qui est de classe C^∞ . Le lemme 36.2(3) dit que $D^-D^+w = w''$. Donc les valeurs $w(x_i)$ résolvent aussi le système (36.38). □

Lemme 36.7 (Quelques estimations).

La matrice L_h du problème sus-mentionné en (36.37) vérifie³ :

- (1) $\|L_h\|_\infty \leq \frac{4}{h^2} + \|c\|_\infty$
- (2) $\|L_h^{-1}\|_\infty \leq \frac{1}{8}$.

Démonstration. Nous nous souvenons de la formule (12.24) :

$$\|A\|_\infty = \max_{i=1, \dots, n} \sum_{j=1}^n |A_{ij}|. \quad (36.40)$$

La première ligne a pour somme : $\frac{3}{h^2} + c_1$, la dernière a pour somme $\frac{3}{h^2} + c_n$ et les autres sont pour somme $\frac{4}{h^2} + c_i$. Elles sont donc toutes majorées par $\frac{4}{h^2} + \|c\|_\infty$.

3. Dans le CTES d'analyse numérique de Marseille, l'estimation donnée est $\|L_h^{-1}\|_\infty \leq \frac{1}{4}$.

Pour l'estimation de $\|L_h^{-1}\|_\infty$ nous allons nous appuyer sur le théorème 35.165.

Commençons par considérer le problème

$$\begin{cases} -w'' = 1 & (36.41a) \\ w(0) = w(1) = 0. & (36.41b) \end{cases}$$

La première équation dit que w est un polynôme de degré 2. En écrivant $w(x) = ax^2 + bx + c$ et en imposant toutes les contraintes, nous trouvons l'unique solution

$$w(x) = -\frac{1}{2}(x^2 - x). \quad (36.42)$$

Le lemme 36.6 nous dit que la fonction w prise aux points $x_i = ih$ donne les valeurs de w_h .

La matrice L_h^0 est une M-matrice et le vecteur w_h vérifie $L_h^0 w_h = \mathbb{1}$. Donc le théorème 35.165 s'applique et

$$\|(L_h^0)^{-1}\| \leq \|w_h\|_\infty = \frac{1}{8}. \quad (36.43)$$

L'obtention de $1/8$ n'est rien d'autre que la recherche du maximum (en valeur absolue) de la parabole $x \mapsto (x - x^2)/2$ pour $x \in [0, 1]$. Le maximum est atteint pour $x = 1/2$; calcul de dérivée et tout ça ...

Nous retournons maintenant à notre matrice originale L_h . Nous avons

$$L_h - L_h^0 = \text{diag}(c_1, \dots, c_n) \geq 0, \quad (36.44)$$

et aussi

$$L_h^{-1} - (L_h^0)^{-1} = \underbrace{L_h^{-1}}_{\geq 0} \underbrace{(L_h^0 - L_h)}_{\leq 0} \underbrace{(L_h^0)^{-1}}_{\geq 0} \quad (36.45)$$

parce que L_h est une M-matrice. Donc tous les coefficients de $L_h^{-1} - (L_h^0)^{-1}$ sont négatifs. Cela implique

$$L_h^{-1} \leq (L_h^0)^{-1}. \quad (36.46)$$

Mais nous savons que les coefficients de L_h^{-1} sont positifs, donc le maximum de ses coefficients en valeur absolue est plus petit que ceux de $(L_h^0)^{-1}$, c'est-à-dire

$$\|L_h^{-1}\|_\infty \leq \|(L_h^0)^{-1}\|_\infty \leq \frac{1}{8}. \quad (36.47)$$

□

36.1.2 Exemple

Soit $\Omega =]0, 1[$ et une fonction $u: \bar{\Omega} \rightarrow \mathbb{R}$ de classe C^4 vérifiant

$$\begin{cases} -u''(x) + u(x) = \sin(x) & (36.48a) \end{cases}$$

$$\begin{cases} u(0) = 0 & (36.48b) \end{cases}$$

$$\begin{cases} u(1) = 0. & (36.48c) \end{cases}$$

Nous allons écrire la méthode des différences finies pour $h = 1/4$. Nous posons donc les points

$$\begin{cases} x_0 = 0 & (36.49a) \end{cases}$$

$$\begin{cases} x_1 = 1/4 & (36.49b) \end{cases}$$

$$\begin{cases} x_2 = 1/2 & (36.49c) \end{cases}$$

$$\begin{cases} x_3 = 3/4 & (36.49d) \end{cases}$$

$$\begin{cases} x_4 = 1. & (36.49e) \end{cases}$$

Vu que nous avons supposé u de classe C^4 , le lemme 36.2(3) nous donne⁴

$$u''(x) = (D^- D^+ u)(x) + \alpha(h) \quad (36.50)$$

avec $\lim_{h \rightarrow 0} \alpha(h)/h = 0$. L'équation discrétisée serait alors

$$\begin{cases} -(D^- D^+ u)(x) + u(x) = \sin(x) & (36.51a) \\ u(0) = u(1) = 0. & (36.51b) \end{cases}$$

où nous n'avons pas précisé l'indice h au bas des opérateurs D^+ et D^- . Les équations (36.51) ne doivent être posées que pour x_1 , x_2 et x_3 parce que les valeurs en x_0 et x_4 sont déjà connues.

Pour x_1

$$\frac{u_2 - 2u_1 + u_0}{h^2} + u_1 = \sin(x_1) \quad (36.52)$$

Pour x_2

$$\frac{u_3 - 2u_2 + u_1}{h^2} + u_2 = \sin(x_2) \quad (36.53)$$

Pour x_3

$$\frac{u_4 - 2u_3 + u_2}{h^2} + u_3 = \sin(x_3). \quad (36.54)$$

Nous tenons compte du fait que $u_0 = u_4 = 0$ et que $h = 1/4$ pour écrire le système

$$\begin{pmatrix} -31 & 16 & 0 \\ 16 & -31 & 16 \\ 0 & 16 & -31 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} = \begin{pmatrix} s_1 \\ s_2 \\ s_3 \end{pmatrix} \quad (36.55)$$

A où les s_i sont des nombres parfaitement connus : par exemple $s_1 = \sin(x_1) = \sin(1/4) \simeq 0.247403959254523$.

36.2 Problèmes de dimension deux

Nous allons considérer le système

$$\begin{cases} -\Delta u = f & \text{sur } \Omega \\ u = g & \text{sur } \partial\Omega \end{cases} \quad (36.56)$$

où $\Omega =]0, a[\times]0, b[$.

Remarque 36.8.

Pourquoi un signe moins devant le laplacien ? Pour avoir la proposition 36.5 qui dira que la matrice correspondant aux différences finies appliquées à ce système est une M-matrice. Sinon, c'est la matrice $-L_h$ qui en serait une.

36.2.1 Discrétisation en croix

Nous allons maintenant déduire une discrétisation du laplacien en discrétisant les opérations ∂_x^2 et ∂_y^2 . Nous discrétisons Ω en mailles carrés de côté h : $x_k = kkh$ et $y_k = kh$. L'opération de dérivée partielle ∂_x est discrétisée par

$$(D_x^+ u)(x, y) = \frac{u(x+h, y) - u(x, y)}{h} \quad (36.57)$$

ou

$$(D_x^- u)(x, y) = \frac{u(x, y) - u(x-h, y)}{h} \quad (36.58)$$

4. Nous ferions n'importe quoi pour ne pas écrire $u''(x) = (D^- D^+ u)(x) + o(h^2)$. Notez que vous faites ce que vous voulez : écrivez avec la notation « petit o » si cela vous chante.

ou

$$(D_x^0 u)(x, y) = \frac{u(x+h, y) - u(x-h, y)}{2h} \quad (36.59)$$

où le h est sous-entendu dans les opérateurs D^0 , D^+ et D^- .

La dérivée partielle seconde $\partial_x^2 u$ peut être approximée par toutes les combinaisons imaginable, par exemple

$$(D_x^- D_x^+ u)(x, y) = \frac{u(x+h, y) - 2u(x, y) + u(x-h, y)}{h^2}. \quad (36.60)$$

Pour évaluer la différence entre $(\partial_x^2 u)(x, y)$ et $(D^- D^+ u)(x, y)$, il est possible de faire du Taylor en deux dimensions, mais nous pouvons également recycler ce qui a été fait. Nous posons $u_y(x) = u(x, y)$ et alors $(\partial_x^2 u)(x, y) = u_y''(x)$ et le lemme 36.2(3) donne, si u_y est de classe C^4 ,

$$|u_y''(x) - D^- D^+ u_y(x)| \leq \frac{1}{12} h^2 \|u_y^{(4)}\|_\infty. \quad (36.61)$$

Là, les opérateurs D^+ et D^- sont ceux à une dimension. Mais nous avons $(D^- D^+ u_y)(x) = (D^- D^+ u)(x, y)$ (à droite ce sont les opérateurs à deux dimensions), donc

$$|(\partial_x^2 u)(x, y) - (D^- D^+ u)(x, y)| \leq \frac{1}{12} h^2 \|\partial_x^4 u\|_\infty \quad (36.62)$$

et nous pouvons écrire

$$(\partial_x^2 u)(x, y) = (D^- D^+ u)(x, y) + h^2 R(x, y, h) \quad (36.63)$$

où R est une fonction qui dépend de x , y et h , mais aussi de u . Le point important est que R soit majoré par une quantité indépendante de h , de telle sorte que nous ayons quelques garanties que négliger ce terme soit une bonne approximation lorsque $h \rightarrow 0$.

Au niveau de la discrétisation, nous considérons x_i avec $i = 0, \dots, N_x$ et y_j avec $j = 0, \dots, N_y$. La discrétisation de $-(\Delta u)(x, y) = f(x, y)$ donne, pour $i = 1, \dots, N_x - 1$ et $j = 1, \dots, N_y - 1$,

$$\frac{1}{h^2} (-u_{i+1,j} + 4u_{i,j} - u_{i-1,j} - u_{i,j+1} + u_{i,j-1}) = f_{i,j}. \quad (36.64)$$

Les équations avec i ou j valant 0 ou N_x , N_y sont les valeurs au bords.

36.9.

Nous notons pour référence ultérieure la discrétisation suivante du laplacien :

$$(\Delta_h u)(x_i, y_j) = \frac{1}{h^2} (-4u_{i,j} + u_{i+1,j} + u_{i-1,j} + u_{i,j-1} + u_{i,j+1}). \quad (36.65)$$

Elle vérifie

$$\Delta_h f = \Delta f + h^2 \alpha(h). \quad (36.66)$$

Cette discrétisation est dite « en croix » parce que les points exploités forment une croix.

36.2.2 Discrétisation en carré

L'opérateur laplacien est défini par $\Delta = \partial_x^2 + \partial_y^2$, mais il existe de nombreuses autres façons de l'écrire.

Lemme 36.10.

Le laplacien est invariant par changement de coordonnées orthogonales. Plus précisément, si A est une matrice orthogonale, en posant $u_i = \sum_k A_{ik} e_k$ nous avons

$$\sum_i \frac{\partial^2 f}{\partial u_i^2} = \Delta f. \quad (36.67)$$

Démonstration. Nous avons :

$$\frac{\partial f}{\partial u_i} = \sum_k A_{ik} \frac{\partial f}{\partial x_k}, \quad (36.68)$$

et donc

$$\sum_i \frac{\partial^2 f}{\partial u_i^2} = \sum_i \frac{\partial}{\partial u_i} \left(\sum_k A_{ik} \frac{\partial f}{\partial x_k} \right) = \sum_{ijk} (A^t)_{ji} A_{ik} \partial^2 f = \sum_{jk} (A^t A)_{jk} \partial_{jk}^2 f. \quad (36.69)$$

En particulier si A est une matrice orthogonale, $(A^t A)_{jk} = \delta_{jk}$ et le résultat est prouvé. \square

Les plus convaincus diront que $\Delta = \nabla \cdot \nabla$ et que le produit scalaire est invariant sous changement de coordonnées orthogonales.

Nous avons déjà déduit la discrétisation (36.65) du laplacien :

$$(\Delta_h u)(x_i, y_j) = \frac{1}{h^2} (-4u_{i,j} + u_{i+1,j} + u_{i-1,j} + u_{i,j-1} + u_{i,j+1}). \quad (36.70)$$

Nous allons maintenant en déduire une par l'idée de décomposer le laplacien dans la base $u = e_1 + e_2$, $v = e_1 - e_2$. Pour cela nous introduisons les opérations (le nombre h dont dépendent ces opérateurs est sous-entendu)

$$(D_u^+ f)(x, y) = \frac{f(x+h, y+h) - f(x, y)}{h} \quad (36.71a)$$

$$(D_v^+ f)(x, y) = \frac{f(x+h, y-h) - f(x, y)}{h} \quad (36.71b)$$

$$(D_u^- f)(x, y) = \frac{f(x, y) - f(x-h, y-h)}{h} \quad (36.71c)$$

$$(D_v^- f)(x, y) = \frac{f(x, y) - f(x-h, y+h)}{h}. \quad (36.71d)$$

Vu que

$$\partial_u^2 + \partial_v^2 = 2\Delta, \quad (36.72)$$

nous discrétisons le laplacien par

$$\Delta'_h = \frac{1}{2} (D_u^- D_u^+ + D_v^- D_v^+). \quad (36.73)$$

Un peu de calcul donne :

$$(\Delta'_h f)(x, y) = \frac{1}{2h^2} \left(-4f(x, y) + f(x+h, y+h) + f(x-h, y+h) \right. \quad (36.74a)$$

$$\left. + f(x+h, y-h) + f(x-h, y-h) \right). \quad (36.74b)$$

36.2.3 Résolution de la discrétisation en croix

Les équations (36.64) forment un système d'équations linéaires à résoudre. Certaines peuvent être simplifiées parce qu'elles « touchent » le bord. Nous verrons cela un peu plus tard.

Nous allons d'abord numéroter correctement les équations de façon à ne pas avoir deux mais un seul indice. Notre fonction de numérotation sera

$$\varphi(i, j) = (j-1)(N_x-1) + i \quad (36.75)$$

avec $i = 1, \dots, N_x-1$ et $j = 1, \dots, N_y-1$. Cela correspond à numéroter les points de l'intérieur du quadrillage ligne par ligne en bas en haut et de gauche à droite. Avec cela les équations (36.64) vont être numérotées par un seul indice I allant de $\varphi(1, 1) = 1$ à $\varphi(N_x-1, N_y-1) = (N_x-1)(N_y-1)$.

Si $I = \varphi(i, j)$ alors nous avons vite

$$\varphi(i + 1, j) = I + 1 \tag{36.76a}$$

$$\varphi(i, j + 1) = I + N_x - 1 \tag{36.76b}$$

$$\varphi(i - 1, j) = I - 1 \tag{36.76c}$$

$$\varphi(i, j - 1) = I - N_x + 1. \tag{36.76d}$$

Nous posons $U_I = u_{\varphi^{-1}(I)}$, et l'équation (36.64) devient

$$\frac{1}{h^2}(-U_{I+1} + 4U_I - U_{I-1} - U_{I+N_x-1} - U_{I-N_x+1}) = f_I. \tag{36.77}$$

Pour savoir la matrice représentant ce système, nous devons simplifier les équations qui doivent l'être. Par exemple avec $I = 1$, le terme $U_{I-1} = U_0$ vaut $u_{0,1} = f_0$. Ce n'est donc pas réellement une inconnue de notre problème.

Nous voulons mettre les équations sous la forme du système

$$L_h U = F. \tag{36.78}$$

Sur la ligne numéro I de L_h , les éléments non nuls sont :

$$L_{I,I} = 4 \tag{36.79a}$$

$$L_{I,I+1} = -1 \tag{36.79b}$$

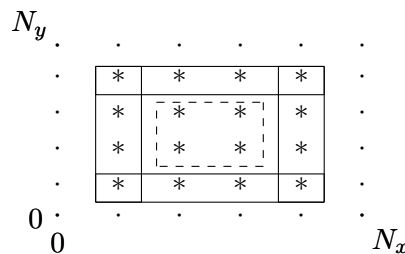
$$L_{I,I-1} = -1 \tag{36.79c}$$

$$L_{I,I+N_x-1} = -1 \tag{36.79d}$$

$$L_{I,I-N_x+1} = -1 \tag{36.79e}$$

pour peu qu'ils existent. Par exemple pour $I = 1$, il n'y a pas d'éléments $L_{I,I-1}$. Les indices I et J de $L_{I,J}$ vont de 1 à $\varphi(N_x - 1, N_y - 1) = (N_y - 1)(N_x - 1)$.

Voici un dessin de notre situation :



À chaque élément du quadrillage correspond une équation.

- Aux points simples sur le bord, correspondent des équations triviales parce que la fonction u y est directement donnée par les conditions aux bords.
- Aux points étoilés entourés en traits continus correspondent des équations « incomplètes » parce que certains termes de l'équation (36.64) sont donnés par les conditions aux bords. Elle correspondent aussi aux lignes incomplète de la matrice L_h où certains éléments donnés en (36.79) n'existent pas.
Le membre de droite de ces équations est par contre enrichi de ce qui à gauche est « donné ».
- Au points étoilés du centre entourés en traits discontinus correspondent des équations complètes.

Notons que f_{00} ne joue aucun rôle dans notre histoire parce que dans les équations (36.64), chaque point (i, j) du maillage n'est liée qu'aux quatre points situés « à côté ».

Proposition 36.11.

La matrice L_h est

- (1) irréductible et à diagonale fortement dominante⁵,
- (2) une M-matrice,
- (3) inversible avec $L_h > 0$,
- (4) symétrique,
- (5) strictement définie positive.

Démonstration. On divise la preuve.

Irréductible Une matrice $n \times n$ dont les deux premières diagonales sont entièrement composées d'éléments non nuls est toujours irréductible. En effet, la première lie l'élément $(1, 2)$ à l'élément $(n - 1, n)$ et donc permet de dire que tous les $i < j$ sont connectés.

La seconde diagonale lie l'élément $(n, n - 1)$ à l'élément $(2, 1)$.

Diagonale fortement dominante En ce qui concerne la dominance de la diagonale, il faut sommer sur les lignes. Or chaque ligne contient (en valeur absolue) un 4 sur la diagonale et au plus quatre éléments qui valent 1. D'où

$$|L_{II}| \geq \sum_{J \neq I} |L_{IJ}|. \quad (36.80)$$

La première ligne n'est jamais complète : elle contient un 4 sur l'élément $(1, 1)$ au maximum deux 1 plus à droite. Donc la matrice L_h est à diagonale fortement dominante.

M-matrice D'après ce que nous venons de voir (proposition 36.11), le théorème 35.163 fonctionne et L_h est une M-matrice⁶.

Inverse strictement positif La proposition 35.164 nous assure qu'une M-matrice irréductible est d'inverse strictement positif. Donc $L_h^{-1} > 0$.

Symétrique La ligne numéro I est

$$\left(\dots, \underbrace{-1}_{I-N_x+1}, \dots, -1, 4, -1, \dots, \underbrace{-1}_{I+N_x-1}, \dots \right) \quad (36.81)$$

Prenons par exemple l'élément $(I, I - N_x + 1)$ qui vaut -1 . Son symétrique est l'élément $(I - N_x + 1, I)$ qui se trouve sur la ligne $I - N_x + 1$. Sur cette dernière ligne nous avons un -1 sur la colonne $I - N_x + 1 + N_x - 1 = I$. Donc l'élément $(I - N_x + 1, I)$ vaut bien -1 et la matrice est symétrique.

Strictement définie positive Vu que la matrice L_h est symétrique, irréductible à diagonale fortement dominante (proposition 36.11), vu que ses éléments diagonaux sont strictement positifs (ils valent 4), la proposition 35.166 nous dit que L_h est strictement définie positive. \square

36.3 Consistance, convergence

36.3.1 Définitions, mise en place

Soit un ouvert $\Omega \subset \mathbb{R}^n$ et un opérateur différentiel L sur Ω . Nous considérons le problème qui consiste à trouver une fonction u sur Ω telle que

$$Lu = f \quad (36.82)$$

pour une fonction f donnée.

Problèmes et choses à faire

La définition suivante est une invention personnelle, n'est pas précise et mérite des commentaires de la part du lecteur.

5. Définition 35.157. Le cas 1×1 est discutablement à diagonale fortement dominante, il faut avouer.

6. Notons que c'est ici que nous sommes content d'avoir posé $-\Delta u = f$ dans le système (36.56), avec un signe négatif devant le laplacien. Sinon tous les signes auraient changé et la matrice $-L_h$ aurait été une M-matrice au lieu de L_h .

Définition 36.12 ([1]).

Un **schéma numérique** de pas h pour $Lu = f$ est la donnée de

- (1) un nombre $h > 0$ supposé petit,
- (2) une quantité N de points x_i dans Ω formant l'ensemble discret Ω_h ,
- (3) une matrice L_h de taille $N \times N$,
- (4) une solution $u_h: \Omega_h \rightarrow \mathbb{R}$ de l'équation $(L_h u_h)(x_i) = f_i$ où nous avons posé $f_i = f(x_i)$.

36.13.

Évidemment pour qu'un schéma mérite le nom de schéma de pas h pour l'équation $Lu = f$, il faut que le nombre h soit lié au choix des points x_i , et que la matrice L_h soit liée à l'opérateur L . La définition n'impose pas formellement de tels liens, parce qu'il y a de nombreuses façons d'approximer une équation différentielle en un système linéaire, sans compter que même l'équation $(L_h u_h)_i = f_i$ peut se résoudre de beaucoup de façons, exacte ou approchées.

Cela pour dire que le lien entre la solution exacte u et la solution approchée n'a rien d'évident, et va dépendre des choix faits lors de la discrétisation et lors de la résolution du système linéaire. Nous allons supposer dans un premier temps que l'équation $L_h u_h = f$ est résolue exactement (nous avons un peu parlé de ces problèmes dans les sections 35.15 et suivantes).

Définition 36.14.

L'**erreur de consistance** d'un schéma numérique est la fonction $\tau_h: \Omega_h \rightarrow \mathbb{R}$ définie par

$$\tau_h(x_i) = (L_h u)_i - (Lu)(x_i). \quad (36.83)$$

Il y a un jeu de notation pas tout à fait évident dans la définition (36.83). En effet, L_h est une matrice, et ne s'applique donc a priori pas immédiatement à une fonction. Ce que signifie la notation $(L_h u)_i$ est que l'on applique la matrice L_h au vecteur $j \mapsto u(x_j)$ et que l'on prend la composante i du résultat.

Définition 36.15.

Nous disons que le schéma est **consistant** avec l'opérateur différentiel L lorsque

$$\lim_{h \rightarrow 0^+} \|\tau_h\| = 0 \quad (36.84)$$

où la norme $\|\cdot\|$ est souvent la norme uniforme, c'est-à-dire $\|\tau_h\| = \max_i \tau_h(x_i)$.

Notons que le lien entre h et le choix des x_i fait partie de la définition des schéma. Sur un segment de longueur L , lorsque h n'est pas un diviseur de L , le schéma devrait expliquer ce que l'on fait pour que la limite (36.84) ait un sens.

Définition 36.16.

Le schéma (Ω_h, L_h) est **consistant à l'ordre p** avec l'opérateur différentiel L pour la norme $\|\cdot\|$ si il existe une constante C indépendante de h telle que

$$\|\tau_h\| \leq Ch^p. \quad (36.85)$$

Définition 36.17.

L'**erreur de discrétisation** entre la solution u du problème $Lu = f$ et la solution approchée u_h sur Ω_h est la fonction

$$e_h: \Omega_h \rightarrow \mathbb{R} \\ x_i \mapsto u(x_i) - u_i. \quad (36.86)$$

où $u_i = u_h(x_i)$ est la solution approchée.

Le schéma discret $(L_h u_h)(x_i) = f_i$ est **convergeant** si $\lim_{h \rightarrow 0} \|e_h\| = 0$. Si de plus il existe une constante C et $p > 0$ tels que

$$\|e_h\| \leq Ch^p, \quad (36.87)$$

alors nous disons que le schéma est convergent à l'ordre p .

Si l'erreur de consistance est petite, le *problème* est bien approximé par la système linéaire. Cela n'implique cependant pas que la solution trouvée soit bien approximée.

Exemple 36.18(Deux opérateurs différentiels proches dont les solutions sont loin)
Soit la partie $\Omega =]0, \infty[$, et les problèmes

$$\begin{cases} L_1 u = u' = 0 \\ u(0) = 1 \end{cases} \quad (36.88a)$$

$$\begin{cases} L_2 v = v' - \epsilon v = 0 \\ v(0) = 1. \end{cases} \quad (36.88b)$$

et

$$\begin{cases} L_2 v = v' - \epsilon v = 0 \\ v(0) = 1. \end{cases} \quad (36.89a)$$

$$\begin{cases} L_2 v = v' - \epsilon v = 0 \\ v(0) = 1. \end{cases} \quad (36.89b)$$

Les solutions exactes sont $u(x) = 1$ et $v(x) = e^{\epsilon x}$.

En ce qui concerne les opérateurs, quelle que soit la norme utilisée nous avons

$$\|L_1 - L_2\| = \sup_{\|f\|=1} \|L_1(f) - L_2(f)\| \quad (36.90a)$$

$$= \sup_{\|f\|=1} \|\epsilon f\| \quad (36.90b)$$

$$= \epsilon. \quad (36.90c)$$

Donc lorsque ϵ est petit, l'opérateur L_2 approxime bien l'opérateur L_1 . Pour toutes les normes. Mais

$$|u(x) - v(x)| = |1 - e^{\epsilon x}|, \quad (36.91)$$

donc quel que soit ϵ nous avons $\|u - v\|_\infty = \infty$. Et d'ailleurs, quelle que soit la norme raisonnable que nous mettons sur l'espace des fonctions, avoir $\|u - v\| = \infty$ semble inévitable.

Donc deux opérateurs différentiels proches peuvent avoir des solutions lointaines. \triangle

36.3.2 Exemple

Soit l'opérateur différentiel L donné par

$$Lu = -u'' + cu \quad (36.92)$$

où c est une fonction. Nous considérons sur $\Omega =]0, 1[$ l'équation différentielle

$$Lu = 0. \quad (36.93)$$

En ce qui concerne la discrétisation, nous définissons le maillage $\Omega_h = \{x_i = ih\}$ avec $i = 0, \dots, N+1$. La solution approchée discrètement sera le vecteur v qui peut être vu comme fonction $v: \Omega_h \rightarrow \mathbb{R}$. Les nombres v_0 et v_{N+1} sont a priori donnés par les conditions aux bords. Pour les autres v_i nous avons les équations

$$(L_h v)_i = -\frac{v_{i+1} - 2v_i + v_{i-1}}{h^2} + c(x_i)v_i. \quad (36.94)$$

Cela est la définition de l'opérateur L_h , et le vecteur v solution de $L_h v = 0$ est la solution du problème au sens de la méthode des différences finies (pour peu qu'il existe, soit unique et tout ça).

Pour calculer l'erreur de consistance, nous considérons une fonction u et nous posons $u_i = u(x_i)$. Le vecteur (u_i) ainsi construit est approximé par v (on espère). Nous avons :

$$\tau_h(x_i) = -\frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} - c(x_i)u_i - (Lu)(x_i). \quad (36.95)$$

Pour étudier cela nous développons $u_{i+1} = u(x_i + h)$ et $u_{i-1} = u(x_i - h)$ à l'ordre 4 : il existe $\alpha_i \in [x_i, x_i + h]$ et $\beta_i \in [x_i - h, x_i]$ tels que

$$u_{i+1} = u(x_i) + hu'(x_i) + \frac{h^2}{2}u''(x_i) + \frac{h^3}{3!}u^{(3)}(x_i) + \frac{h^4}{4!}u^{(4)}(\alpha_i) \quad (36.96)$$

et

$$u_{i-1} = u(x_i) - hu'(x_i) + \frac{h^2}{2}u''(x_i) - \frac{h^3}{3!}u^{(3)}(x_i) + \frac{h^4}{4!}u^{(4)}(\beta_i). \quad (36.97)$$

Après simplification de plusieurs termes,

$$\tau_h(x_i) = -\frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} - c_i u_i + u''(x_i) + c_i u_i = \frac{h^2}{4!}(u^{(4)}(\alpha_i) + u^{(4)}(\beta_i)). \quad (36.98)$$

Parler de la consistance du schéma demande d'étudier $\lim_{h \rightarrow 0^+} \|\tau_h\|$, et pour cela, il faut préciser la norme avec laquelle nous voulons travailler. L'ordre de consistance va dépendre de la norme utilisée.

Pour la norme $\|\cdot\|_\infty$, les nombres $u^{(4)}(\alpha_i)$ et $u^{(4)}(\beta_i)$ se majorent par $\|u^{(4)}\|_\infty$ et nous avons

$$\|\tau_h\|_\infty \leq \frac{h^2}{12}\|u^{(4)}\|_\infty. \quad (36.99)$$

Nous avons consistance d'ordre 2.

Remarque 36.19.

La valeur de $\|\tau_h\|_\infty$ dépend de la fonction u sur laquelle nous la calculons. Cependant nous avons convergence $\|\tau_h\|_\infty \rightarrow 0$ pour toute fonction (de classe disons C^4).

La constante C pour laquelle nous avons $\|\tau_h\| \leq Ch^2$ et donc qui nous vaut de pouvoir dire que la consistance est d'ordre 2 ne dépend pas de h , ni des valeurs ponctuelles de u ou de ses dérivées, mais dépend des normes de u et de ses dérivées (en l'occurrence seulement de la norme de $u^{(4)}$.)

Étudions la consistance pour la norme L_1 :

$$\|\tau_h\|_1 = \sum_i |\tau_h(x_i)| \leq \frac{h^2}{12} \sum_i \|u^{(4)}\|_\infty \quad (36.100)$$

où nous avons majoré chacun des $u^{(4)}(\alpha_i)$ par $\|u^{(4)}\|_\infty$. Combien de termes dans la somme? Nous avons $h = 1/(N-1)$ et donc $N = (1+h)/h$, ce qui donne

$$\|\tau_h\|_1 \leq N \frac{h^2}{12} \|u^{(4)}\|_\infty = (1+h)Ch. \quad (36.101)$$

La constante $1+h$ se majore par n'importe quelle constante strictement plus grande que 1. Nous pouvons donc la rentrer dans C et écrire

$$\|\tau_h\|_1 \leq Ch \quad (36.102)$$

et donc avoir la consistance à l'ordre 1.

36.3.3 Consistance, stabilité et convergence

Lemme 36.20.

Soit un opérateur différentiel L , soit u la solution de $Lu = f$ et un schéma numérique (L_h, Ω_h) pour cette équation. Nous notons u_h la solution de $L_h u_h = f$. Alors nous avons

$$L_h e_h = \tau_h \quad (36.103)$$

Et si de plus L_h est inversible,

$$\|e_h\| \leq \|L_h^{-1}\| \|\tau_h\|. \quad (36.104)$$

Démonstration. Par définition u_h est solution de $L_h u_h = f$ en tant que fonctions sur Ω_h . Nous avons donc

$$L_h e_h = L_h u_h - L_h u \quad (36.105)$$

où u doit être compris comme la restriction de u à Ω . En appliquant au point x_i ,

$$(L_h e_h)(x_i) = \underbrace{(L_h u_h)(x_i)}_{=f_i} - (L_h u)(x_i), \quad (36.106)$$

mais $f_i = (Lu)(x_i)$ parce que u est solution de $Lu = f$. Donc

$$(L_h e_h)(x_i) = (Lu)(x_i) - (L_h u)(x_i) = \tau_h(x_i). \quad (36.107)$$

Si la matrice L_h est inversible nous avons $e_h = L_h^{-1} \tau_h$ et donc

$$\|e_h\| \leq \|L_h^{-1}\| \|\tau_h\| \quad (36.108)$$

par le lemme 12.17. □

Bien entendu, en tant qu'opérateur linéaire sur un espace de dimension finie, l'opérateur L_h^{-1} est borné pour chaque h . Mais si il n'y a pas une borne uniforme en h , alors le lemme 36.20 dit qu'il n'y a pas d'espoir de majorer $\|e_h\|$ de façon à passer à la limite $\lim_{h \rightarrow 0} \|L_h^{-1}\|$.

Définition 36.21.

Un schéma numérique est **stable** si il existe une constante $C > 0$ indépendante de h telle que $\|L_h^{-1}\| \leq C$.

Théorème 36.22.

En deux parties.

- (1) Si un schéma discret est consistant et stable, alors il est convergent.
- (2) Si de plus il est consistant à l'ordre p , alors il est convergent à l'ordre p .

Démonstration. Nous savons du lemme 36.20 (qui s'applique parce que l'inversibilité de L_h est dans la définition de la stabilité) que $\|e_h\| \leq \|L_h^{-1}\| \|\tau_h\|$ et que $\|L_h^{-1}\| \leq C$. En passant à la limite ⁷,

$$\lim_{h \rightarrow 0} \|e_h\| \leq C \lim_{h \rightarrow 0} \|\tau_h\| = 0. \quad (36.109)$$

La dernière limite est le fait que le schéma soit consistant. Le schéma est donc convergent.

Si de plus il est consistant à l'ordre p , alors

$$\|e_h\| \leq C \|\tau_h\| \leq C' h^p \quad (36.110)$$

Il est donc également convergent à l'ordre p . □

36.3.4 Exemple : schéma à cinq points, laplacien en croix

Nous avons développé le schéma dont l'opérateur sur Ω_h est (voir (36.64))

$$(L_h u_h)(x_i, y_j) = \frac{1}{h^2} (-u_{i+1,j} - u_{i,j+1} + 4u_{i,j} - u_{i-1,j} - u_{i,j-1}). \quad (36.111)$$

Proposition 36.23.

Le schéma est :

- (1) consistant à l'ordre 2,
- (2) stable pour la norme uniforme et

$$\|L_h^{-1}\|_\infty \leq \frac{1}{8}, \quad (36.112)$$

- (3) convergent à l'ordre 2 pour la norme $\|\cdot\|_\infty$.

7. Toutes les limites $h \rightarrow 0$ sont en réalité des limites $h \rightarrow 0^+$, mais nous allégeons cette notation.

Démonstration. Cet opérateur avait été construit de telle sorte à avoir (voir (36.63))

$$(\Delta u)(x_i, y_j) = (L_h u)(x_i, y_i) + h^2 R(x, y, h) \quad (36.113)$$

où R peut être majoré indépendamment de h . En tant que fonctions sur Ω_h nous avons

$$\tau_h = \Delta u - L_h u = h^2 R(x, y, h), \quad (36.114)$$

et donc $\|\tau_h\|_\infty \leq Ch^2$, parce que le lemme 36.2(3) donne aussi

$$\|R\|_\infty \leq C \max\left\{\left\|\frac{\partial^4 u}{\partial x^4}\right\|_\infty, \left\|\frac{\partial^4 u}{\partial y^4}\right\|_\infty\right\}. \quad (36.115)$$

En ce qui concerne la stabilité nous allons utiliser le théorème 35.165. Nous considérons la fonction

$$g(x, y) = -\frac{1}{4}(x^2 + y^2), \quad (36.116)$$

qui vérifie $-\Delta g = 1$ sur le carré $[0, 1]^2$. Nous considérons le vecteur g_h d'indices $(i, j) \mapsto g_{ij} = g(x_i, y_j)$ sur lequel nous calculons L_h :

$$(L_h g_h)_{ij} = \frac{1}{h^2}(-g_{i+1,j} - g_{i,j+1} + 4g_{ij} - g_{i-1,j} - g_{i,j-1}); \quad (36.117)$$

en remplaçant les g par leurs valeurs en termes de $x_i, x_{i-1}, x_{i+1}, y_j, y_{j-1}$ et y_{j+1} , et ne tenant compte du fait que $x_k = kh$ et $y_l = lh$ nous avons :

$$(L_h g_h)_{ij} = \frac{1}{4}((i+1)^2 - j^2 + i^2 + (j+1)^2 - 4i^2 - 4j^2 + (i-1)^2 + j^2 + i^2 + (j-1)^2) = 1. \quad (36.118)$$

Donc $L_h g_h = 1$.

Vu que L_h est une M-matrice (proposition 36.11(2)), le théorème 35.165 nous dit alors que L_h^{-1} vérifie

$$\|L_h^{-1}\|_\infty \leq \|g_h\|_\infty. \quad (36.119)$$

Mais

$$\|g_h\|_\infty \leq \|g\|_\infty = g\left(\frac{1}{2}, \frac{1}{2}\right) = \frac{1}{8}. \quad (36.120)$$

Notons que cela est bien une inégalité et non une égalité parce que rien n'assure que le point $(1/2, 1/2)$ soit sur le maillage ; donc rien n'assure que la valeur $g(1/2, 1/2)$ ne soit parmi les valeurs du vecteur discrétisé g_h .

Notre schéma numérique est stable et consistant à l'ordre 2 pour la norme $\|\cdot\|_\infty$. Le théorème 36.22 dit alors que le schéma est convergent à l'ordre 2 pour la même norme. \square

36.4 Autres laplaciens

Nous avons vu le laplacien en croix (36.65)

$$(\Delta_h f)(x, y) = \frac{1}{h^2}(-4f(x, y) + f(x+h, y) + f(x-h, y) \quad (36.121a)$$

$$+ f(x, y+h) + f(x, y-h)) \quad (36.121b)$$

qui vérifie

$$\Delta_h f = \Delta f + Kh^2, \quad (36.122)$$

ainsi que le laplacien en carré (36.74)

$$(\Delta'_h f)(x, y) = \frac{1}{2h^2} \left(-4f(x, y) + f(x+h, y+h) + f(x-h, y+h) \quad (36.123a)$$

$$+ f(x+h, y-h) + f(x-h, y-h) \right) \quad (36.123b)$$

qui vérifie également $\Delta'_h f = \Delta f + Kh^2$.

A priori toutes combinaisons de la forme

$$a\Delta_h + b\Delta'_h \quad (36.124)$$

avec $a + b = 1$ est valable comme tentative de discrétiser le laplacien. Ce sont des schémas à 9 points. Évidemment la matrice L_h correspondante va être moins creuse, mais nous pouvons espérer ajuster a et b de telle sorte à obtenir une consistance d'un ordre supérieur à 2.

Nous allons développer les $\Delta_h f$ et $\Delta'_h f$ à l'ordre 4 (reste à l'ordre 6). Quelques remarques avant de commencer.

- (1) Allez relire la proposition 13.274 et les notations qui vont avec pour comprendre les différentielles.
- (2) Écrivez les formules du type (13.689) pour $d^2 f$ et $d^4 f$.
- (3) Allez relire le développement de Taylor du théorème 13.363.
- (4) À l'ordre zéro, il n'y a rien parce que le terme $-4f(x, y)$ compense les quatre termes d'ordre zéro des autres termes.
- (5) Aux ordres impairs, il n'y a rien. En effet, prenons un nombre impair l et la formule

$$(d^l f)_x(h, \dots, h) = \sum_{i_1, \dots, i_l} h_{i_1} \dots h_{i_l} \frac{\partial^l f}{\partial x_{i_1} \dots \partial x_{i_l}}(x). \quad (36.125)$$

Nous avons

$$(d^l f)_x(h, \dots, h) + (d^l f)_x(-h, \dots, -h) = 0. \quad (36.126)$$

Or dans les expressions (36.123) et (36.121), les termes arrivent par paires opposées.

Commençons par calculer $h^2(\Delta_h f)(x, y)$.

Ordre 4 . Le premier terme est :

$$(d^2 f)_{(x,y)}((h, 0), (h, 0)) = h^2(d^2 f)_{(x,y)}((1, 0), (1, 0)). \quad (36.127)$$

La formule (13.689) à peine adaptée permet de calculer ça explicitement.

Il y a encore les termes du même type avec $(0, 1)$, $(-1, 0)$ et $(0, -1)$.

Ordre 4 Cette fois, ce sont 4 termes du type

$$h^4(d^4 f)_{(x,y)}((1, 0), (1, 0), (1, 0), (1, 0)) \quad (36.128)$$

à calculer.

Cela fait beaucoup de termes à calculer. Je vous laisse vous persuader que le programme suivant en Sage nous donne les coefficients.

```

1  #! /usr/bin/env python3
2  # -*- coding: utf8 -*-
3
4
5  x,y=var("x,y")
6
7  C=[]
8  C.append(x)
9  C.append(y)
10
11 def coef4(h):
12     S=0
13     for i in [0,1]:

```

```

14     for j in [0,1]:
15         for k in [0,1]:
16             for l in [0,1]:
17                 S=S+h[i]*h[j]*h[k]*h[l]*C[i]*C[j]*C[k]*C[l]
18     return S
19
20 def coef2(h):
21     S=0
22     for i in [0,1]:
23         for j in [0,1]:
24             S=S+h[i]*h[j]*C[i]*C[j]
25     return S
26
27 cross=[]
28 square=[]
29
30 cross.append( [1,0] )
31 cross.append( [-1,0] )
32 cross.append( [0,1] )
33 cross.append( [0,-1] )
34
35 square=[]
36 square.append( [-1,-1] )
37 square.append( [1,-1] )
38 square.append( [-1,1] )
39 square.append( [1,1] )
40
41
42 print("Cross scheme :")
43
44 K=sum( coef2(v) for v in cross )
45 L=sum( coef4(v) for v in cross )
46 print(K)
47 print(L)
48
49 print("square scheme :")
50
51 K=sum( coef2(v) for v in square )
52 L=sum( coef4(v) for v in square )
53 print(K)
54 print(L)

```

tex/sage/coefs.sage

Le résultat est que, en utilisant la formule

$$\partial_x^4 f + 2\partial_{xxyy}^4 f + \partial_y^4 f = \Delta\Delta f, \quad (36.129)$$

nous avons

$$(\Delta_h f)(x, y) = \frac{1}{2}(2\partial_x^2 f + 2\partial_y^2 f)(x, y) + \frac{1}{4!}2h^2(\partial_x^4 f + \partial_y^4 f)(x, y) + Kh^4 \quad (36.130a)$$

$$= (\Delta f)(x, y) + \frac{1}{12}h^2(\partial_x^4 f + \partial_y^4 f)(x, y) + Kh^4 \quad (36.130b)$$

$$= \Delta f + \frac{h^2}{12}\Delta\Delta f - \frac{h^2}{12}2\partial_{xxyy}^4 f + Kh^4 \quad (36.130c)$$

$$= \Delta f + \frac{h^2}{12}\Delta\Delta f - \frac{h^2}{6}\partial_{xxyy}^4 f + Kh^4 \quad (36.130d)$$

où K est une constante qui peut être majorée en terme des dérivées quatrièmes de f . En particulier la plus grande des normes supremum de ces dérivées.

Le même genre de calculs donnent

$$(\Delta'_h f)(x, y) = \frac{1}{2}\left[\frac{1}{2}4\Delta f + \frac{h^2}{4!}(4\partial_x^4 f + 24\partial_x^2\partial_y^2 f + 4\partial_y^4 f)\right] + Kh^4. \quad (36.131)$$

Ça donne :

$$(\Delta'_h f) = \Delta f + \frac{h^2}{12}\Delta\Delta f + \frac{h^2}{3}\partial_{xxyy}^4 f + Kh^4 \quad (36.132)$$

avec redéfinition du K ; nous ne le précisons plus à chaque fois.

Nous avons donc le résultat proposé dans [480] :

$$a\Delta_h f + b\Delta'_h f = (a+b)\Delta f + (a+b)\frac{h^2}{12}\Delta^2 f + h^2\frac{1}{6}(a-2b)\partial_{xxyy}^4 f + Kh^4. \quad (36.133)$$

L'idée est d'appliquer ça à une fonction u qui vérifie l'équation différentielle $-\Delta u = f$ (attention au clash de notation pour f). Le mieux est de supprimer le terme en $\partial_{xxyy}^4 f$ en demandant $a-2b=0$. Nous avons donc à résoudre le système

$$\begin{cases} a+b=1 \\ a-2b=0. \end{cases} \quad (36.134a)$$

$$(36.134b)$$

Qui propose une décomposition PLU pour résoudre ce système linéaire ? Quelle que soit la manière, la solution est

$$a = \frac{2}{3}, \quad b = \frac{1}{3}. \quad (36.135)$$

Nous allons donc étudier la discrétisation à neuf points

$$L_h = \frac{2}{3}\Delta_h + \frac{1}{3}\Delta'_h. \quad (36.136)$$

En faisant quelques additions nous trouvons que l'opération

$$(L_h u)(x_i, y_j) = \frac{1}{6h^2}\left(-20u_{ij} + 4(u_{i+1,j} + u_{i-1,j} + u_{i,j+1} + u_{i,j-1}) + u_{i+1,j+1} + u_{i-1,j+1} + u_{i+1,j-1} + u_{i-1,j-1}\right) \quad (36.137)$$

vérifie

$$L_h f = \Delta f + \frac{h^2}{12}\Delta^2 f + Kh^4. \quad (36.138)$$

36.4.1 Travail avec le laplacien à 9 points

Nous allons écrire un schéma numérique pour l'équation différentielle $-\Delta u = f$ utilisant la discrétisation à 9 points du laplacien. Nous recopions ses propriétés fondamentales (36.136), (36.137), (36.138) :

$$L_h = \frac{2}{3}\Delta_h + \frac{1}{3}\Delta'_h, \quad (36.139)$$

et

$$(L_h u)(x_i, y_j) = \frac{1}{6h^2} \left(-20u_{ij} + 4(u_{i+1,j} + u_{i-1,j} + u_{i,j+1} + u_{i,j-1}) \right. \\ \left. + u_{i+1,j+1} + u_{i-1,j+1} + u_{i+1,j-1} + u_{i-1,j-1} \right), \quad (36.140)$$

et

$$L_h f = \Delta f + \frac{h^2}{12} \Delta^2 f + Kh^4. \quad (36.141)$$

Nous appliquons (36.141) à u et nous isolons Δu :

$$\Delta u = L_h u - \frac{h^2}{12} \Delta^2 u + Kh^4 = \frac{1}{6h^2} T_h u + \frac{h^2}{12} \Delta f + Kh^4 \quad (36.142)$$

où nous avons utilisé $\Delta^2 u = -\Delta f$ et avons noté

$$T_h u = -20u_{ij} + 4(u_{i+1,j} + u_{i-1,j} + u_{i,j+1} + u_{i,j-1}) + u_{i+1,j+1} + u_{i-1,j+1} + u_{i+1,j-1} + u_{i-1,j-1}. \quad (36.143)$$

Nous imposons maintenant $\Delta u = -f$ en écrivant

$$\frac{1}{6h^2} T_h u = -f - \frac{h^2}{12} \Delta f + \alpha(h)h^4. \quad (36.144)$$

Une idée est de remplacer Δf par son approximation en croix (36.122) :

$$T_h u = -6h^2 f - \frac{h^4}{2} (\Delta_h f + Kh^2) + \alpha(h)h^6 \quad (36.145)$$

Avec quelques calculs nous trouvons le schéma numérique suivant :

$$20u_{ij} - 4(u_{i+1,j} + u_{i-1,j} + u_{i,j+1} - u_{i,j-1}) - u_{i+1,j+1} - u_{i-1,j+1} - u_{i+1,j-1} - u_{i-1,j-1} \quad (36.146a)$$

$$= \frac{h^2}{2} (8f_{ij} + f_{i+1,j} + f_{i-1,j} + f_{i,j+1} + f_{i,j-1}) + \alpha(h)h^6. \quad (36.146b)$$

En oubliant le terme en $\alpha(h)h^6$, nous obtenons un système d'équations linéaires.

Problèmes et choses à faire

Il me semble que ce schéma donne une convergence d'ordre 6. C'est correct ?

Chapitre 37

Variables aléatoires et théorie des probabilités

37.1 Espace de probabilité

Définition 37.1.

Une **mesure de probabilité** sur un espace mesurable¹ (Ω, \mathcal{A}) est une mesure positive telle que $P(\Omega) = 1$. Dans ce cas, le triple (Ω, \mathcal{A}, P) est un **espace de probabilité**.

Un point $\omega \in \Omega$ est une **observation**, une partie mesurable $A \in \mathcal{A}$ est un **événement**. L'ensemble $A \cup B$ représente l'événement A ou B tandis que l'ensemble $A \cap B$ représente l'événement A et B .

Si les A_n sont des événements, nous avons défini en 8.27 limite supérieure et la limite inférieure de la suite A_n par

$$\limsup_{n \rightarrow \infty} A_n = \bigcap_{n \geq 1} \bigcup_{k \geq n} A_k \quad (37.1)$$

et

$$\liminf_{n \rightarrow \infty} A_n = \bigcup_{n \geq 1} \bigcap_{k \geq n} A_k \quad (37.2)$$

Si $\omega \in \liminf A_n$, alors ω réalise tous les A_n sauf un nombre fini.

Nous avons

$$\limsup A_n = \{\omega \in \Omega \text{ tel que } \omega \in A_n \text{ pour une infinité de } n\}. \quad (37.3)$$

Théorème 37.2 (Borel-Cantelli).

Si

$$\sum_{n=1}^{\infty} P(A_n) < \infty \quad (37.4)$$

alors $P(\limsup A_n) = 0$.

Démonstration. La condition $\sum_{n \geq 1} P(A_n) < \infty$ signifie que la fonction

$$\varphi = \sum_{n \geq 1} \mathbb{1}_{A_n} \quad (37.5)$$

est P -intégrable. Par conséquent, elle est finie presque partout (au sens de P), c'est-à-dire

$$P(\varphi = \infty) = 0. \quad (37.6)$$

Les points ω sur lesquels $\varphi(\omega) = \infty$ sont ceux tels que

$$\sum_{n \geq 1} \mathbb{1}_{A_n}(\omega) = \infty, \quad (37.7)$$

1. Espace mesurable : 15.1, mesure positive : 15.20.

c'est-à-dire les ω qui appartiennent à une infinité d'ensembles A_n , ou encore les $\omega \in \limsup A_n$. Nous avons donc montré que

$$\{\omega \text{ tel que } \varphi(\omega) = \infty\} = \{\omega \in \Omega \text{ tel que } \omega \in A_n \text{ pour une infinité de } n\} = \limsup A_n. \quad (37.8)$$

Or l'hypothèse signifie que la probabilité du membre de gauche est nulle. \square

Corollaire 37.3.

Si $\sum_{n=1}^{\infty} P(\complement A_n) < \infty$, alors presque sûrement tous les B_n sont réalisés à l'exception d'un nombre fini.

37.2 Variables aléatoires

Définition 37.4.

Une **variable aléatoire** est une application mesurable

$$X: (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d)). \quad (37.9)$$

Nous convenons que $\mathbb{R}^1 = \bar{\mathbb{R}}$, c'est-à-dire que dans le cas où la variable aléatoire X est réelle, nous acceptons les valeurs $\pm\infty$.

Définition 37.5.

Une variable aléatoire réelle X est **absolument continue** s'il existe une fonction positive et intégrable $f: \mathbb{R} \rightarrow \mathbb{R}$ telle que pour tout intervalle $I \subset \mathbb{R}$,

$$P(X \in I) = \int_I f(t) dt. \quad (37.10)$$

Nous disons alors que f est la **densité** de X .

Cela ne devrait pas être sans rappeler la définition 18.20.

37.2.1 Indépendance

La définition suivante vient de l'instructive motivation de [481]. La définition d'indépendance de deux événements se généralise à n événements de la façon suivante.

Définition 37.6.

Nous disons que les événements A_1, \dots, A_n sont **indépendants** si pour tout choix $\{i_1, \dots, i_k\} \subset \{1, \dots, n\}$ nous avons

$$P(A_{i_1} \cap \dots \cap A_{i_k}) = P(A_{i_1}) \dots P(A_{i_k}). \quad (37.11)$$

Les sous tribus $\mathcal{A}_1, \dots, \mathcal{A}_n$ sont **indépendantes** si pour tout choix $A_i \in \mathcal{A}_i$, les événements A_i sont indépendants.

Exemple 37.7

Soit $\Omega = [0, 1] \times [0, 1]$ muni de la mesure de Lebesgue. Soient $A = [0, a] \times [0, 1]$ et $B = [0, 1] \times [0, b]$. Nous avons $P(A) = a$ et $P(B) = b$ ainsi que $P(A \cap B) = ab$. \triangle

Lemme 37.8.

Les tribus $\mathcal{A}_1, \dots, \mathcal{A}_n$ sont indépendantes si et seulement si

$$P(A_1 \cap \dots \cap A_n) = P(A_1) \dots P(A_n) \quad (37.12)$$

pour tout $A_i \in \mathcal{A}_i$.

Démonstration. L'implication dans le sens direct découle immédiatement des définitions.

Nous supposons avoir un choix $(A_i)_{i=1,\dots,n}$ avec $A_i \in \mathcal{A}_i$ et nous devons montrer que ces événements sont indépendants, c'est-à-dire que si $J \subset \{1, \dots, n\}$ alors les événements $(A_j)_{j \in J}$ sont indépendants. Sans perte de généralité, nous pouvons supposer que si $i \notin J$, $A_i = \Omega$. Alors nous avons

$$P\left(\bigcap_{j \in J} A_j\right) = P\left(\bigcap_{i=1}^n A_i\right) = \prod_{i=1}^n P(A_i) = \prod_{j \in J} P(A_j) \quad (37.13)$$

parce que $P(A_i) = P(\Omega) = 1$ lorsque i n'est pas dans J . \square

Si A est un événement, la **tribu engendrée** par A est

$$\sigma(A) = \{\emptyset, A, \complement A, \Omega\}. \quad (37.14)$$

Soit $X: \Omega \rightarrow \mathbb{R}^d$ une variable aléatoire. Conformément à la définition 15.46, la **tribu engendrée** est

$$\mathcal{A}_X = \{X^{-1}(B) \text{ tel que } B \in \mathcal{B}(\mathbb{R}^d)\}. \quad (37.15)$$

Cela est la plus petite tribu sous tribu de \mathcal{A} pour laquelle X est mesurable. Elle sera aussi (le plus souvent notée $\sigma(X)$).

Définition 37.9.

Nous disons que les variables aléatoires $X_k: \Omega \rightarrow \mathbb{R}^d$ sont **indépendantes** lorsque les tribus engendrées $\mathcal{A}_{X_1}, \dots, \mathcal{A}_{X_n}$ le sont.

Remarque 37.10.

Il n'a de sens de dire que X_1 et X_2 sont indépendants que si X_1 et X_2 sont des applications dont l'espace de départ est identique.

Si nous voulons modéliser le jet de deux pièce indépendantes, le mauvais choix est de faire $\Omega = \{0, 1\}$, y mettre la mesure d'équiprobabilité, et de considérer les deux variables aléatoires

$$X_i(\omega) = \begin{cases} f & \text{si } \omega = 0 \\ p & \text{si } \omega = 1. \end{cases} \quad (37.16)$$

Ces deux variables sont évidemment pas indépendantes. Il faut poser $\Omega = \{0, 1\} \times \{0, 1\}$, y mettre la mesure d'équiprobabilité et poser

$$X_1(x, y) = \begin{cases} f & \text{si } x = 0 \\ p & \text{si } x = 1 \end{cases}, \quad (37.17)$$

$$X_2(x, y) = \begin{cases} f & \text{si } y = 0 \\ p & \text{si } y = 1 \end{cases}, \quad (37.18)$$

Ces variables aléatoires sont indépendantes. Par exemple

$$X_1^{-1}\{p\} = \{(1, 0), (1, 1)\} \quad (37.19a)$$

$$X_2^{-1}\{p\} = \{(0, 1), (1, 1)\} \quad (37.19b)$$

et on a bien

$$P(X_1^{-1}\{p\} \cap X_2^{-1}\{p\}) = P\{(1, 1)\} = \frac{1}{4} \quad (37.20)$$

ainsi que

$$P\{X_p^{-1}(p)\} = \frac{1}{2} \quad (37.21a)$$

pour $i = 1$ et $i = 2$.

Proposition 37.11.

Soient $(X_k : \Omega \rightarrow \mathbb{R}^{d_k})$ des variables aléatoires indépendantes.

(1) Si $B_k \in \mathcal{B}or(\mathbb{R}^{d_k})$. Alors

$$P(X_k \in B_k \forall k \leq n) = P(X_1 \in B_1) \dots P(X_n \in B_n). \quad (37.22)$$

(2) Les événements $\{X_i \in B_i\}$ sont indépendants.

(3) Les tribus engendrées par des X_i et d'autres sont indépendantes. Plus précisément, si I et J sont deux ensembles disjoints de \mathbb{N} alors les tribus

$$\sigma(\{X_i, i \in I\}) \quad (37.23)$$

et

$$\sigma(\{X_i, i \in J\}) \quad (37.24)$$

sont indépendantes.

Démonstration. Lorsque nous écrivons $X_i \in B_i$, nous parlons de l'événement

$$(X_i \in B_i) = \{\omega \in \Omega \text{ tel que } X_i(\omega) \in B_i\} = X_i^{-1}(B_i) \in \mathcal{A}_{X_i}. \quad (37.25)$$

Vu que par hypothèse les tribus (\mathcal{A}_i) sont indépendantes, le lemme 37.8 nous montre que

$$P\left(\bigcap_{i=1}^n X_i \in B_i\right) = \prod_i P(X_i \in B_i). \quad (37.26)$$

Il reste à voir que l'ensemble $X_i^{-1}(B_i)$ fait partie de la tribu \mathcal{A} de départ. Cela est la définition du fait que l'application X_i soit une variable aléatoire : elle doit être mesurable en tant qu'application

$$X_i : (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}^d, \mathcal{B}or(\mathbb{R}^d)). \quad (37.27)$$

Les affirmations (2) et (3) ne sont que des façons alternatives d'exprimer la même chose. \square

Lemme 37.12.

Les événements $(A_i)_{i=0, \dots, n}$ sont indépendants si et seulement si les événements que nous obtenons en remplaçant certains des A_i par $\mathcal{C}A_i$ le sont.

Démonstration. Sans perte de généralité, nous pouvons nous contenter de prouver que les événements $\mathcal{C}A_0, A_1, \dots, A_n$ sont indépendants sous l'hypothèse que les événements A_0, A_1, \dots, A_n sont indépendants. Soit I un sous-ensemble de $\{1, \dots, n\}$. Nous avons

$$P(\mathcal{C}A_0 \bigcap_{i \in I} A_i) = P\left(\bigcap_{i \in I} A_i \setminus \bigcap_{i \in I} A_i \cap A_0\right) \quad (37.28a)$$

$$= P\left(\bigcap_{i \in I} A_i\right) - P\left(\bigcap_{i \in I} A_i \cap A_0\right) \quad (37.28b)$$

$$= P\left(\bigcap_{i \in I} A_i\right) (1 - P(\mathcal{C}A_0)) \quad (37.28c)$$

$$= P\left(\bigcap_{i \in I} A_i\right) P(\mathcal{C}A_0). \quad (37.28d)$$

\square

Proposition 37.13.

Les événements $(A_i)_{i=1, \dots, n}$ sont indépendants si et seulement si les variables aléatoires associées $\mathbb{1}_{A_1}, \dots, \mathbb{1}_{A_n}$ le sont.

Démonstration. La tribu engendrée par la variable aléatoire $\mathbb{1}_{A_k}$ est

$$\mathcal{A}_{\mathbb{1}_{A_k}} = \{\emptyset, A_k, \complement A_k, \Omega\}. \quad (37.29)$$

En effet si $1 \in B$, alors $A_i \subset \mathbb{1}_{A_i}^{-1}(B)$, et si $0 \in B$, alors $\complement A_i \subset \mathbb{1}_{A_i}^{-1}(B)$. Les éléments 0 et 1 sont tous deux soit dans B , soit hors de B . Cela donne les 4 possibilités énumérées dans (37.29).

Supposons que les événements (A_i) sont indépendants. Nous devons vérifier que les tribus le soient, c'est-à-dire que les événements A_i et $\complement A_j$ sont indépendants. Cela est une conséquence du lemme 37.12. \square

Théorème 37.14 (Doob[236]).

Soit $X: \Omega \rightarrow \mathbb{R}^d$ une variable aléatoire. Une fonction $Y: \Omega \rightarrow \mathbb{R}^p$ est une variable aléatoire \mathcal{A}_X -mesurable si et seulement s'il existe une fonction borélienne $f: \mathbb{R}^d \rightarrow \mathbb{R}^p$ telle que $Y = f(X)$.

Proposition 37.15.

Soient des variables aléatoires $X_k: \Omega \rightarrow \mathbb{R}^{d_k}$ des variables aléatoires indépendantes et des fonctions boréliennes $f_k: \mathbb{R}^{d_k} \rightarrow \mathbb{R}^{p_k}$. Alors les variables aléatoires $f_k(X_k)$ sont indépendantes.

Démonstration. Le théorème 37.14 assure que les applications

$$f_k \circ X_k: \Omega \rightarrow \mathbb{R}^{p_k} \quad (37.30)$$

sont \mathcal{A}_{X_k} -mesurables. En particulier pour tout borélien $B \subset \mathbb{R}^{p_k}$, nous avons $X_k^{-1} \circ f_k^{-1}(B) \in \mathcal{A}_{X_k}$. Nous avons donc

$$\sigma(f_k \circ X_k) \subset \sigma(X_k), \quad (37.31)$$

et par conséquent les tribus $\sigma(f_k \circ X_k)$ sont indépendantes étant donné que les tribus $\sigma(X_k)$ le sont. \square

Lemme 37.16 (Lemme de regroupement).

Soit (Ω, \mathcal{A}, P) un espace de probabilité et $(\mathcal{A})_{i \in I}$ une famille de tribus indépendantes dans \mathcal{A} . Si $(M_j)_{j \in J}$ est une partition de I , alors les tribus

$$\mathcal{B}_j = \sigma\left(\bigcup_{i \in M_j} \mathcal{A}_i\right) \quad (37.32)$$

sont indépendantes.

Si les variables aléatoires $\{X_1, X_2, X_3, X_4, X_5\}$ sont indépendantes, et si f et g sont des fonctions mesurables, alors les variables aléatoires $f(X_2, X_3, X_5)$ et $g(X_1, X_4)$ sont indépendantes.

Une preuve a l'air d'être donnée dans [482].

37.2.2 Lois conjointes et indépendance

Définition 37.17.

Deux événements A et B sont dits **indépendants** si

$$P(A \cap B) = P(A)P(B). \quad (37.33)$$

Si nous considérons n variables aléatoires réelles $X_1, \dots, X_n: \Omega \rightarrow \mathbb{R}$, la loi du n -uplet $X = (X_1, \dots, X_n)$ est une variable aléatoire $X: \Omega \rightarrow \mathbb{R}^n$ appelée la **loi conjointe** des lois X_i . Dans ce cas, les variables aléatoires X_i elles-mêmes sont dites lois **marginales** de X .

Proposition 37.18.

Les variables aléatoires $\{X_i\}$ sont indépendantes si et seulement si

$$P_{(X_1, \dots, X_n)} = P_{X_1} \otimes \dots \otimes P_{X_n}. \quad (37.34)$$

Définition 37.19.

Soient $\{X_i\}_{1 \leq i \leq n}$ des variables aléatoires réelles (pas spécialement indépendantes). La **densité conjointe** de X_1, \dots, X_n est la fonction $f: \mathbb{R}^n \rightarrow \mathbb{R}$ qui satisfait

- (1) $f(x_1, \dots, x_n) \geq 0$ pour tout $(x_1, \dots, x_n) \in \mathbb{R}^n$,
- (2) $\int_{\mathbb{R}^n} f = 1$,
- (3) pour tout $A_i \subset \mathbb{R}$ nous avons

$$P\left(\bigcap_{i=1}^n X_i \in A_i\right) = \int_{\prod_i A_i} f(x_1, \dots, x_n) dx_1 \dots dx_n. \quad (37.35)$$

Proposition 37.20.

Si les variables aléatoires X_1, \dots, X_n sont indépendantes et ont des densités f_{X_1}, \dots, f_{X_n} , alors la variable aléatoire conjointe $X = (X_1, \dots, X_n)$ a pour densité conjointe la fonction

$$f_X(x_1, \dots, x_n) = f_{X_1}(x_1) \dots f_{X_n}(x_n). \quad (37.36)$$

Démonstration. En partant de la définition de l'indépendance et de la fonction de densité conjointe, ainsi qu'en utilisant le théorème de Fubini,

$$\begin{aligned} \int_{A_1 \times \dots \times A_n} f_X(x_1, \dots, x_n) dx_1 \dots dx_n &= P(X_1 \in A_1, \dots, X_n \in A_n) \\ &= P(X_1 \in A_1) \dots P(X_n \in A_n) \\ &= \left(\int_{A_1} f_{X_1}(x_1) dx_1 \right) \dots \left(\int_{A_n} f_{X_n}(x_n) dx_n \right) \\ &= \int_{A_1 \times \dots \times A_n} f_{X_1}(x_1) \dots f_{X_n}(x_n) dx_1 \dots dx_n. \end{aligned} \quad (37.37)$$

La fonction $(x_1, \dots, x_n) \mapsto f_{X_1}(x_1) \dots f_{X_n}(x_n)$ vérifie donc la condition (3) de la définition 37.19. La vérification des autres conditions est immédiate. \square

La proposition suivante provient du fait que la mesure d'une loi conjointe est le produit des mesures lorsque les variables aléatoires sont indépendantes (proposition 37.18).

Proposition 37.21 ([236]).

Si les variables aléatoires réelles X_1, \dots, X_n sont intégrables et indépendantes, alors leur produit est intégrable et l'espérance du produit est égal au produit des espérances :

$$E(X_1 \dots X_n) = E(X_1) \dots E(X_n). \quad (37.38)$$

37.2.3 Somme et produit de variables aléatoires indépendantes

Soient X et Y , deux variables aléatoires réelles indépendantes. Nous voudrions étudier la loi de la variable aléatoire $S = X + Y$. Nous commençons par calculer la fonction de répartition en utilisant le résultat de la proposition 37.20 :

$$F_{X+Y}(z) = P(X + Y \leq z) = \int_{x+y \leq z} f_{X,Y}(x, y) dx dy \quad (37.39a)$$

$$= \int_{-\infty}^{\infty} dx \int_{-\infty}^{z-x} dy f_X(x) f_Y(y) \quad (37.39b)$$

$$= \int_{\mathbb{R}} \left(\int_{-\infty}^{z-x} f_Y(y) dy \right) f_X(x) dx \quad (37.39c)$$

$$= \int_{\mathbb{R}} F_Y(z-x) f_X(x) dx. \quad (37.39d)$$

Pour calculer la fonction de densité de S , nous dérivons la fonction de répartition :

$$f_{X+Y}(z) = \frac{dF_{X+Y}}{dz}(z) \quad (37.40a)$$

$$= \int_{\mathbb{R}} f_Y(z-x)f_X(x)dx, \quad (37.40b)$$

ce qui nous amène à dire que la densité de la somme est le produit de convolution² des densités :

$$f_{X+Y}(x) = \int_{\mathbb{R}} f_Y(x-t)f_X(t)dt, \quad (37.41)$$

ou encore $f_{X+Y} = f_X * f_Y$.

Notez que nous avons passé sous le silence la difficulté d'inverser la dérivée et l'intégrale. Un exemple sera donné au point 37.5.8.

Lemme 37.22.

Soient X et Y , deux variables aléatoires indépendantes. Alors

$$E(XY) = E(X)E(Y). \quad (37.42)$$

Démonstration. Par indépendance et par proposition 37.20, la fonction de densité conjointe de X et Y vaut $f_{X,Y} = f_X f_Y$. Par conséquent l'utilisation de Fubini sous la forme (15.811) entraîne

$$E(XY) = \int_{\mathbb{R} \times \mathbb{R}} xy f_{X,Y}(x,y) dx dy = E(X)E(Y). \quad (37.43)$$

□

Nous dirons tout un tas de chose sur l'indépendance et la variance en 37.5.13, mais pour l'instant nous allons mentionner et démontrer déjà ceci :

Lemme 37.23.

Soient X et Y deux variables aléatoires indépendantes et identiquement distribuées. Alors

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y). \quad (37.44)$$

Démonstration. Par définition, $\text{Var}(X + Y) = E([X + Y - E(X) - E(Y)]^2)$. En développant le carré et en utilisant le lemme 37.22,

$$\text{Var}(X + Y) = E(X^2) - E(X)^2 + E(Y^2) - E(Y)^2 = \text{Var}(X) + \text{Var}(Y). \quad (37.45)$$

□

Exemple 37.24

Deux variables aléatoires non indépendantes dont la covariance est nulle. Nous considérons la variable aléatoire

$$Z: \Omega \rightarrow \{(1, 0), (-1, 0), (0, 1), (0, -1)\} \quad (37.46)$$

de loi uniforme. C'est-à-dire que $P(Z = z) = \frac{1}{4}$ pour tout z . Ensuite nous considérons les variables aléatoires $X = \text{proj}_1 \circ Z$ et $Y = \text{proj}_2 \circ Z$. Toute personne étant capable de compter jusqu'à 4 voit que

$$P(X = 1) = P(X = -1) = \frac{1}{4} \quad (37.47a)$$

$$P(X = 0) = \frac{1}{2}, \quad (37.47b)$$

2. Définition 28.52.

et les mêmes probabilités pour Y . De même $E(X) = E(Y) = 0$. Par conséquent

$$\text{Cov}(X, Y) = E(XY) = 0 \quad (37.48)$$

parce que pour tout $\omega \in \Omega$ nous avons soit $X(\omega) = 0$ soit $Y(\omega) = 0$. Ces variables aléatoires X et Y ne sont donc pas corrélées.

Mais elles ne sont pas indépendantes pour autant, comme nous allons le voir pas plus tard qu'immédiatement. Nous avons

$$P(X = 0|Y = 0) = \frac{P(X = 0, Y = 0)}{P(Y = 0)} = 0 \quad (37.49)$$

parce que X et Y ne peuvent pas être simultanément nulles, tandis que

$$P(X = 0)P(Y = 0) = \frac{1}{4}. \quad (37.50)$$

△

37.2.4 Espérance

Nous dirons que la variable aléatoire X a un **moment d'ordre** p si $X \in L^p(\Omega, \mathcal{A}, P)$ ($1 \leq p < \infty$). Si X est **intégrable** (c'est-à-dire si $X \in L^1$), alors nous définissons l'**espérance** de X par

$$E(X) = \int_{\Omega} X dP \in \mathbb{R}^d. \quad (37.51)$$

Si $E(X) = 0$ nous disons que la variable aléatoire est **centrée**. La variable aléatoire $X - E(X)$ est la variable aléatoire centrée associée à X .

Le **moment** d'ordre p de la variable aléatoire X est l'espérance

$$m_n(X) = E(X^n). \quad (37.52)$$

Proposition 37.25.

Si X et Y sont deux variables aléatoires (pas spécialement indépendantes), nous avons

$$E(X + Y) = E(X) + E(Y). \quad (37.53)$$

Nous donnons la preuve dans le cas de variables aléatoires indépendantes. Le cas plus général de variable aléatoires non indépendantes peut être trouvé dans [483].

Démonstration. Nous avons le calcul suivant :

$$E(X + Y) = \int_{\mathbb{R}} x f_{X+Y}(x) dx \quad (37.54a)$$

$$= \int_{\mathbb{R}} x \int_{\mathbb{R}} f_Y(x - t) f_X(t) dt dx \quad (37.54b)$$

$$= \int_{\mathbb{R}} f_X(t) \underbrace{\int_{\mathbb{R}} x f_Y(x - t) dx}_{=E(Y)+t} dt \quad (37.54c)$$

$$= \int_{\mathbb{R}} f_X(t) (E(Y) + t) dt \quad (37.54d)$$

$$= E(Y) + \int_{\mathbb{R}} t f_X(t) dt \quad (37.54e)$$

$$= E(Y) + E(X) \quad (37.54f)$$

où nous avons utilisé la proposition 37.41 et le fait que l'intégrale sur \mathbb{R} d'une densité vaut 1. □

Une application de l'inégalité de Hölder (proposition 28.33) est la suivante. Si X et Y sont des variables aléatoires intégrables alors

$$E(XY) \leq E(X^2)^{1/2} E(Y^2)^{1/2}. \quad (37.55)$$

En effet

$$E(XY) \leq \|XY\|_{L^1(\Omega)} \leq \|X\|_{L^2(\Omega)} \|Y\|_{L^2(\Omega)}. \quad (37.56)$$

37.2.5 Variance

Si $X \in L^2(\Omega, \mathcal{A}, P)$ alors nous définissons la **variance** de X par

$$\text{Var}(X) = E([X - E(X)]^2). \quad (37.57)$$

Proposition 37.26.

La variance de la variable aléatoire X peut être exprimée par la formule

$$\text{Var}(X) = E(X^2) - [E(X)]^2 \quad (37.58)$$

où $X^2 = X \cdot X$ et $E(X)^2 = E(X) \cdot E(X)$ sont des produits scalaires dans \mathbb{R}^d .

Démonstration. De façon explicite, nous avons

$$E([X - E(X)]^2) = \int_{\Omega} (X(\omega) - E(X)) \cdot (X(\omega) - E(X)) dP(\omega) \quad (37.59)$$

où $E(X) \in \mathbb{R}^d$ est une constante. En développant le produit scalaire nous avons

$$E([X - E(X)]^2) = E(X^2 - 2X \cdot E(X) + E(X)^2) \quad (37.60a)$$

$$= E(X^2) - 2E(X)^2 + E(X)^2 \quad (37.60b)$$

$$= E(X^2) - E(X)^2. \quad (37.60c)$$

□

Nous définissons l'**écart-type** de X par

$$\sigma_X = \sqrt{\text{Var}(X)}. \quad (37.61)$$

En d'autres termes,

$$\sigma_X = \|X - E(X)\|_{L^2}. \quad (37.62)$$

On définit encore la **moyenne quadratique** de X par

$$\|X\|_{L^2} = [E(X^2)]^{1/2}. \quad (37.63)$$

La variable aléatoire

$$\bar{V}_n = \frac{1}{n} \sum_i (X_i - \bar{X}_n)^2 \quad (37.64)$$

est la **variance empirique** de l'échantillon (X_i) .

Lemme 37.27.

Si X est une variable aléatoire,

(1) $\text{Var}(ax) = a^2 \text{Var}(X)$ pour tout $a \in \mathbb{R}$;

(2) Si de plus Y est une variable aléatoire indépendante de X , alors $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$.

Démonstration. Nous avons

$$\text{Var}(X + Y) = E(X^2 + Y^2 + 2XY) - (E(X) + E(Y))^2 \quad (37.65a)$$

$$= E(X^2) + E(Y^2) + 2E(XY) - E(X)^2 - E(Y)^2 + 2E(X)E(Y). \quad (37.65b)$$

Étant donné que X et Y sont indépendantes nous avons $E(XY) = E(X)E(Y)$ par le lemme 37.22. \square

Si les X_1, \dots, X_n sont des variables aléatoires on considère la **moyenne empirique**

$$\bar{X}_n = \frac{X_1 + \dots + X_n}{n}. \quad (37.66)$$

37.2.6 Covariance

Soient X et Y , deux variables aléatoires réelles. Leur **covariance** est définie par

$$\text{Cov}(X, Y) = E\left[(X - E(X))(Y - E(Y))\right] \quad (37.67)$$

L'idée est que la covariance devient grande si X et Y s'écartent de leurs moyennes dans le même sens. Il existe une formule alternative :

$$\text{Cov}(X, Y) = E(XY) - E(X)E(Y) \quad (37.68)$$

En ce qui concerne les dimensions plus hautes, si $X : \Omega \rightarrow \mathbb{R}^d$ est un vecteur aléatoire de carré intégrable, nous définissons

$$\text{Cov}(X) = E\left[(X - E(X)) \otimes (X - E(X))\right] \quad (37.69)$$

où par $a \otimes b$ nous entendons la matrice $(a \otimes b)_{ij} = a_i b_j$. Cela peut aussi être noté $a^t b$ si l'on fait bien attention à qui est un vecteur colonne et qui est un vecteur ligne.

Proposition 37.28.

Si X et Y sont deux variables aléatoires non spécialement indépendantes, nous avons

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y). \quad (37.70)$$

Démonstration. Il s'agit d'un calcul en partant de

$$\begin{aligned} \text{Var}(X + Y) &= E((X + Y)^2) - E(X + Y)^2 \\ &= E(X^2) + E(Y^2) + 2E(XY) \\ &\quad + (E(X) + E(Y))^2 - 2E(X)^2 - 2E(X)E(Y) \\ &\quad - 2E(Y)E(X) - 2E(Y)^2. \end{aligned} \quad (37.71)$$

À partir d'ici il s'agit de recombinaison tous les termes pour former la formule annoncée. \square

Plus généralement nous avons la formule

$$\text{Var}\left(\sum_i X_i\right) = \sum_i \text{Var}(X_i) + 2 \sum_{1 \leq i < j \leq n} \text{Cov}(X_i, X_j). \quad (37.72)$$

37.2.7 Probabilité conditionnelle : événements

Proposition-définition 37.29.

Soit (Ω, \mathcal{A}, P) un espace de probabilité et $B \in \mathcal{A}$ avec $P(B) > 0$. Alors avec la formule

$$P_B(A) = \frac{P(A \cap B)}{P(B)}, \quad (37.73)$$

l'espace $(\Omega, \mathcal{A}, P_B)$ est un espace probabilisé. Nous notons $P_B(A)$ le nombre $P(A|B)$ et nous le nommons **probabilité conditionnelle** de A sachant B .

Démonstration. On vérifie que $(\Omega, \mathcal{A}, P_B)$ est un espace de probabilité parce que $P_B(\Omega) = 1$ et

$$P_B\left(\bigcup_i A_i\right) = \sum_i P_B(A_i) \quad (37.74)$$

si les A_i sont deux à deux disjoints. □

Une conséquence immédiate de (37.73) est que si A et B sont des événements indépendants alors

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = P(A). \quad (37.75)$$

La probabilité conditionnelle à B est quelque chose qui ne tient compte que de ce qui se passe dans B . Si K est un événement tel que $A \cap B = K \cap B$, alors

$$P(A|B) = P(K|B). \quad (37.76)$$

Théorème 37.30.

Soient $(B_n)_{n \geq 1}$ une partition finie de Ω telle que $P(B_i) > 0$. Soit $A \in \mathcal{A}$ tel que $P(A) > 0$.

(1) Si A, B et C sont des événements, alors

$$P(A \cap B|C) = P(A|B \cap C)P(B|C). \quad (37.77)$$

(2) Si $P(B) > 0$, alors $P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$.

(3) On a la **formule des probabilités totales** :

$$P(A) = \sum_{i=1}^n P(A|B_i)P(B_i) = \sum_i P(A \cap B_i). \quad (37.78)$$

(4) On a la **formule de Bayes** :

$$P(B_k|A) = \frac{P(A|B_k)P(B_k)}{\sum_i P(A|B_i)P(B_i)}. \quad (37.79)$$

Démonstration. (1) En développant le membre de droite,

$$\begin{aligned} P(A \cap B|C) &= \frac{P(A \cap B \cap C)}{P(B \cap C)} \frac{P(B \cap C)}{P(C)} \\ &= P(A \cap B|C). \end{aligned} \quad (37.80)$$

(2) C'est la définition de $P(A|B)$ et $P(B|A)$.

(3) Vu que les B_i forment une partition, nous avons

$$P(A) = \sum_i P(A \cap B_i) = \sum_i P(A|B_i)P(B_i). \quad (37.81)$$

(4) En utilisant les deux premiers points, nous trouvons

$$\begin{aligned} P(A|B_k)P(B_k) &= P(A \cap B_k) \\ &= P(B_k|A)P(A) \\ &= P(B_k|A) \sum_i P(A|B_i)P(B_i). \end{aligned} \quad (37.82)$$

□

37.2.8 Espérance conditionnelle

Théorème-définition 37.31 (Définition de l'espérance conditionnelle[385]).

Soit un espace de probabilité (Ω, \mathcal{A}, P) et une variable aléatoire intégrable $X : \Omega \rightarrow \mathbb{R}$. Pour chaque sous tribu \mathcal{F} de \mathcal{A} , il existe une (presque partout) unique variable aléatoire $Y : \Omega \rightarrow \mathbb{R}$ telle que

- (1) Y est \mathcal{F} -mesurable
- (2) Y est P -intégrable
- (3) pour tout $B \in \mathcal{F}$,

$$\int_B X dP = \int_B Y dP. \quad (37.83)$$

Cette variable aléatoire sera notée $E(X|\mathcal{F})$ pour des raisons qui apparaîtront plus tard.

Démonstration. Remarquons que prendre $Y = X$ ne fonctionne pas parce qu'en général si \mathcal{O} est mesurable dans \mathbb{R} , alors $X^{-1}(\mathcal{O})$ est dans la tribu \mathcal{A} , mais n'est pas automatiquement dans la tribu \mathcal{F} . Il faudra donc un peu plus travailler.

Unicité Si Y_1 et Y_2 vérifient tous les deux les conditions, l'ensemble $\{Y_1 < Y_2\}$ est un élément de \mathcal{F} et nous avons

$$\int_{\{Y_1 < Y_2\}} X = \int_{Y_1 < Y_2} Y_1 = \int_{Y_1 < Y_2} Y_2. \quad (37.84)$$

En particulier nous avons $\int_{\{Y_1 < Y_2\}} (Y_1 - Y_2) = 0$ et donc

$$(Y_1 - Y_2)\mathbb{1}_{Y_1 < Y_2} = 0 \quad (37.85)$$

presque partout. Le corollaire 15.180 montre alors que $Y_1 - Y_2 \geq 0$ presque partout. De la même manière, l'ensemble $\{Y_2 < Y_1\}$ est dans \mathcal{F} et nous trouvons que $Y_2 - Y_1 \geq 0$ presque partout. Par conséquent $Y_1 = Y_2$ presque partout.

Existence dans le cas de carré intégrable Nous supposons que $X \in L^2(\Omega, \mathcal{A}, P)$ et nous considérons K , le sous-ensemble de $L^2(\Omega, \mathcal{A}, P)$ des fonctions \mathcal{F} -mesurables. Le théorème des projections 26.5 nous indique que

$$L^2(\Omega, \mathcal{A}, P) = K \oplus K^\perp \quad (37.86)$$

par la décomposition $X = \text{proj}_K X + (X - \text{proj}_K X)$. La variable aléatoire $Y = \text{proj}_K X$ a les propriétés d'être \mathcal{F} -mesurable et $\langle Y - X, Z \rangle = 0$ pour tout $Z \in K$. Soit $A \in \mathcal{F}$, si nous considérons $Z = \mathbb{1}_A$, la dernière condition signifie que

$$\int_\Omega X \mathbb{1}_A = \int_\Omega Y \mathbb{1}_A, \quad (37.87)$$

ou encore

$$\int_A Y = \int_A X. \quad (37.88)$$

La variable aléatoire $Y = \text{proj}_K(X)$ répond donc à la question lorsque $X \in L^2(\Omega, \mathcal{F}, P)$.

Existence en général Nous considérons maintenant que $X \in L^1(\Omega, \mathcal{A}, P)$. Quitte à décomposer X en deux fonctions positives X_+ et X_- telles que $X = X_+ + X_-$, nous pouvons supposer que X est positive. Par hypothèse $X \in L^1(\Omega, \mathcal{A}, P)$; pour chaque $n \in \mathbb{N}$ nous posons

$$X_n(\omega) = \min\{X(\omega), n\}. \quad (37.89)$$

Étant donné que la mesure P est une mesure de probabilité, les constantes sont intégrables et $X_n \in L^2(\Omega, \mathcal{A}, P)$. De plus la suite (X_n) est croissante et

$$\lim_{n \rightarrow \infty} X_n(\omega) = X(\omega). \quad (37.90)$$

Si nous notons encore K l'ensemble des variables aléatoires dans $L^2(\Omega, \mathcal{A}, P)$ qui sont \mathcal{F} -mesurables, pour chaque n nous avons donc la variable aléatoire

$$Y_n = \text{proj}_K X_n = E(X_n | \mathcal{F}) \quad (37.91)$$

qui est \mathcal{F} -mesurable et telle que

$$\int_A X_n = \int_A Y_n \quad (37.92)$$

pour tout $A \in \mathcal{F}$. Nous voudrions prouver que la variable aléatoire $Y = \lim_n Y_n$ existe et est la solution au problème, c'est-à-dire est $E(X | \mathcal{F})$.

Commençons par prouver que $Y_n \geq 0$ presque partout. Pour cela nous remarquons que l'ensemble $\{Y_n < 0\}$ est mesurable et

$$0 \geq \int_{Y_n < 0} Y_n = \int_{Y_n < 0} X_n \geq 0. \quad (37.93)$$

La première inégalité est évidente et la dernière est due au fait que X_n est positive. Par conséquent

$$\int_{Y_n < 0} Y_n = 0 \quad (37.94)$$

et le lemme 15.180 conclut que $P(Y_n < 0) = 0$.

Soit $Z: \Omega \rightarrow \mathbb{R}$ une variable aléatoire positive dans $L^2(\Omega, \mathcal{A}, P)$. Montrons que $\text{proj}_K Z$ est encore positive. Pour cela nous considérons l'ensemble $A = \{\text{proj}_K Z < 0\}$ et les inégalités

$$0 \leq \int_A Z = \int_A \text{proj}_K Z \leq 0, \quad (37.95)$$

ce qui montre que $\int_A \text{proj}_K Z = 0$ et par conséquent que $P\{\text{proj}_K(Z) < 0\} = 0$. Cela nous montre que la projection depuis L^2 conserve la positivité.

Étant donné que $X_{n-1} - X_n \geq 0$ nous avons aussi

$$Y_{n-1} - Y_n \geq 0 \quad (37.96)$$

La suite de fonctions

$$n \mapsto Y_n = E(X_n | \mathcal{F}) \quad (37.97)$$

est croissante et vérifie le théorème de la convergence monotone :

$$\int_A X = \lim_{n \rightarrow \infty} \int_A X_n = \lim_{n \rightarrow \infty} \int_A E(X_n | \mathcal{F}) = \int_A \lim_{n \rightarrow \infty} E(X_n | \mathcal{F}) = \int_A Y. \quad (37.98)$$

Par conséquent $E(X | \mathcal{F})$ existe et

$$Y = \lim_{n \rightarrow \infty} E(X_n | \mathcal{F}) = E(X | \mathcal{F}). \quad (37.99)$$

□

37.32.

Vu la définition 37.31 nous pourrions croire que la variable aléatoire $E(X | \mathcal{F}) = X$ fait l'affaire. Il n'en est rien parce que la variable aléatoire X n'est pas spécialement \mathcal{F} -mesurable alors qu'il est requis que $E(X | \mathcal{F})$ le soit. Avec la tribu $\mathcal{F} = \{\emptyset, \Omega\}$, nous n'avons en général pas que $X^{-1}(B) \in \mathcal{F}$ pour tout borélien B .

Par contre si $\sigma(X)$ est la tribu engendrée par la variable aléatoire X , alors $E(X | \sigma(X)) = X$.

Définition 37.33.

Soit Z une variable aléatoire. L'*espérance conditionnelle* « X sachant Z » est la variable aléatoire

$$E(X | Z) = E(X | \sigma(Z)) \quad (37.100)$$

où $\sigma(Z)$ est la tribu engendrée par Z . Le membre de droite est une variable aléatoire définie en 37.31.

Définition 37.34.

Soient $A \in \mathcal{A}$ un événement et \mathcal{F} une sous-tribu de \mathcal{A} . Nous définissons $P(A|\mathcal{F})$ par

$$P(A|\mathcal{F}) = E(\mathbb{1}_A|\mathcal{F}). \quad (37.101)$$

Notons que cela est une variable aléatoire et non un réel. Le membre de droite est l'espérance conditionnelle de la variable aléatoire $\mathbb{1}_A$ par rapport à \mathcal{F} définie en 37.31.

Et l'espérance conditionnelle d'un événement par rapport à une variable aléatoire est :

$$E(A|X) = E(A|\sigma(A)). \quad (37.102)$$

Proposition 37.35.

Soit une espace probabilisé (Ω, \mathcal{A}, P) ainsi qu'une variable aléatoire X à valeurs dans \mathbb{R}^d , et un événement A . Alors

$$E(P(A|X)) = P(A). \quad (37.103)$$

Démonstration. Tout le point de la preuve est de remarquer que $E(\mathbb{1}_A) = E(\mathbb{1}_A|X)$.

La formule $E(\mathbb{1}_A) = E(\mathbb{1}_A|X)$ La notation $E(\mathbb{1}_A|X)$ est un raccourcis pour écrire la variable aléatoire $E(\mathbb{1}_A|\sigma(X))$. Cette dernière est l'application $\Omega \rightarrow \mathbb{R}^d$ telle que

$$\int_B E(\mathbb{1}_A|\sigma(X)) = \int_B \mathbb{1}_A \quad (37.104)$$

pour tout borélien B de \mathbb{R}^d tout en étant $\sigma(X)$ -mesurable. Comme expliqué en 37.32, il est tentant de dire $E(\mathbb{1}_A|\sigma(X)) = \mathbb{1}_A$, mais ce n'est pas le cas parce qu'il n'y a aucune raisons que $\mathbb{1}_A$ soit une application $\sigma(X)$ -mesurable. Au niveau des espérances, par contre, l'égalité tient :

$$E(E(\mathbb{1}_A|X)) = \int_{\Omega} E(\mathbb{1}_A|X) = \int_{\Omega} \mathbb{1}_A = E(\mathbb{1}_A) \quad (37.105)$$

où nous avons utilisé le fait que Ω lui-même soit $\sigma(X)$ -mesurable.

La preuve Nous avons alors

$$P(A) = E(\mathbb{1}_A) = E(E(\mathbb{1}_A|X)), \quad (37.106)$$

alors que $E(\mathbb{1}_A|X) = P(A|X)$. En mettant l'un dans l'autre :

$$P(A) = E(P(A|X)). \quad (37.107)$$

□

Proposition 37.36 (Transitivité de l'espérance conditionnelle).

Si $\mathcal{B}_2 \subseteq \mathcal{B}_1 \subset \mathcal{A}$ alors

$$E(E(X|\mathcal{B}_1)|\mathcal{B}_2) = E(X|\mathcal{B}_2). \quad (37.108)$$

Démonstration. Si $B \in \mathcal{B}_2$, nous avons

$$\int_B E(E(X|\mathcal{B}_1)|\mathcal{B}_2) dP = \int_B E(X|\mathcal{B}_1) dP = \int_B dP. \quad (37.109)$$

La première égalité est la définition de l'espérance conditionnelle par rapport à \mathcal{B}_2 . La seconde égalité est celle de l'espérance conditionnelle par rapport à \mathcal{B}_1 et le fait que $B \in \mathcal{B}_2 \subset \mathcal{B}_1$. Ce que nous avons prouvé est que

$$E(E(X|\mathcal{B}_1)|\mathcal{B}_2) \quad (37.110)$$

est une variable aléatoire \mathcal{B}_2 -mesurable vérifiant la condition

$$\int_B E(E(X|\mathcal{B}_1)|\mathcal{B}_2) = \int_B E(X|\mathcal{B}_2) \quad (37.111)$$

pour tout $B \in \mathcal{B}_2$. C'est donc $E(X|\mathcal{B}_2)$ par la partie unicité du théorème 37.31. □

Proposition 37.37.

Soit (Ω, \mathcal{F}, P) un espace de probabilité, soit \mathcal{A} une sous tribu de \mathcal{F} et X , une variable aléatoire \mathcal{F} -mesurable et intégrable. Alors la variable aléatoire $E(X|\mathcal{A})$ du théorème 37.31 est l'unique (presque partout) variable aléatoire à être \mathcal{A} -mesurable telle que nous ayons

$$E(E(X|\mathcal{A})Y) = E(XY). \quad (37.112)$$

pour toute variable aléatoire Y \mathcal{A} -mesurable.

Démonstration. Supposons pour commencer que Y soit une fonction simple positive, alors $Y = \sum_{i=1}^n a_i \mathbb{1}_{E_i}$ et nous avons

$$\int_{\Omega} E(X|Y) = \sum_i a_i \int_{E_i} E(X|\mathcal{A}) \quad (37.113a)$$

$$= \sum_i a_i \int_{E_i} X \quad (37.113b)$$

$$= \int_{\Omega} XY. \quad (37.113c)$$

Maintenant si Y est mesurable et bornée, elle est limite croissante de fonctions étagées bornées (proposition 15.105) et le résultat tient par la convergence monotone, théorème 15.160.

Si Y n'est pas positive, nous séparons $Y = Y_+ - Y_-$.

Pour l'unicité, soit Z et Z' deux variables aléatoires telles que pour toute variable aléatoire Y ,

$$\int_{\Omega} ZY = \int_{\Omega} XY = \int_{\Omega} Z'Y. \quad (37.114)$$

Si nous prenons $Y = \mathbb{1}_{\{Z \neq Z'\}}$, nous avons

$$0 = \int_{\Omega} (Z - Z') \mathbb{1}_{Z \neq Z'} = \int_{Z \neq Z'} Z - Z', \quad (37.115)$$

d'où le fait que $P(Z \neq Z') = 0$. □

Si X est une variable aléatoire dont la tribu engendrée est indépendante de la tribu \mathcal{F} , nous voudrions que la connaissance de \mathcal{F} n'influence pas la connaissance de X , c'est-à-dire que

$$E(X|\mathcal{F}) = E(X). \quad (37.116)$$

Ce que nous avons est même mieux. Nous avons le lemme suivant.

Lemme 37.38 ([236]).

Les tribus \mathcal{F}_1 et \mathcal{F}_2 sont indépendantes si et seulement si

$$E(U|\mathcal{F}_1) = E(U) \quad (37.117)$$

pour toute variable aléatoire U étant \mathcal{F}_1 -mesurable.

Ici, par $E(U)$ nous entendons la variable aléatoire constante prenant la valeur numérique $E(U)$ en tout point de Ω .

Démonstration. Si \mathcal{F}_1 et \mathcal{F}_2 sont indépendantes, alors pour tout $B \in \mathcal{F}_2$ nous avons

$$\int_B U dP = E(U \mathbb{1}_B) \quad (37.118a)$$

$$= E(U)E(\mathbb{1}_B) \quad (37.118b)$$

$$= E(U) \int_{\Omega} \mathbb{1}_B dP \quad (37.118c)$$

$$= \int_B E(U) dP. \quad (37.118d)$$

Justifications.

- L'intégrale $\int_B UdP$ a un sens même si $B \in \mathcal{F}_2$ alors que U est \mathcal{F}_1 -mesurable. Le supremum (15.424) définissant l'intégrale est tout de même bien défini, en particulier, l'ensemble sur lequel on prend le supremum est non vide.
- Pour (37.118b), la variable aléatoire U est \mathcal{F}_1 -mesurable (donc la tribu engendrée par U est dans \mathcal{F}_1) alors que $\mathbb{1}_B$ est \mathcal{F}_2 -mesurable. Les tribus engendrées étant indépendantes, les variables aléatoires le sont et nous pouvons décomposer l'espérance.

Ce que montre le calcul (37.118) est que $E(U)$ est une variable aléatoire \mathcal{F}_2 -mesurable (parce que constante) dont l'intégrale sur chaque élément de \mathcal{F}_2 vaut l'intégrale de U . Par la partie unicité du théorème 37.31, nous déduisons que $E(U) = E(U|\mathcal{F}_2)$. \square

Corollaire 37.39.

Si X est une variable aléatoire et si \mathcal{F} est une tribu, alors

$$E(E(X|\mathcal{F})) = E(X). \quad (37.119)$$

Démonstration. Il suffit d'appliquer la définition (37.83) à $B = \Omega$:

$$E(E(X|\mathcal{F})) = \int_{\Omega} E(X|\mathcal{F})(\omega)dP(\omega) = \int_{\Omega} X(\omega)dP(\omega) = E(X). \quad (37.120)$$

\square

Exemple 37.40

Soient X_1, X_2 deux variables aléatoires à valeurs dans $\{0, 1\}$ avec probabilité 1/2 et indépendantes. Nous considérons $S = X_1 + X_2$. La situation est modélisée par l'espace

$$\Omega = \{(0, 0), (0, 1), (1, 0), (1, 1)\} \quad (37.121)$$

et les variables aléatoires

$$X_i(\omega_1, \omega_2) = \omega_i \quad (37.122a)$$

$$S(\omega_1, \omega_2) = \omega_1 + \omega_2. \quad (37.122b)$$

Pour vérifier que de cette manière nous avons bien que X_1 est indépendante de X_2 , nous commençons par voir les tribus associées. Un ouvert de \mathbb{R} soit contient 0 et 1, soit contient un seul des deux soit n'en contient aucun des deux. En appliquant X_1^{-1} à chacune de ces quatre situations nous voyons que la tribu $\sigma(X_1)$ est

$$\mathcal{F}_1 = \sigma(X_1) = \{\{(0, 0), (0, 1)\}, \{(1, 0), (1, 1)\}, \Omega, \emptyset\}. \quad (37.123)$$

De la même façon nous avons

$$\mathcal{F}_2 = \sigma(X_2) = \{\{(0, 0), (1, 0)\}, \{(0, 1), (1, 1)\}, \Omega, \emptyset\}. \quad (37.124)$$

Nous posons

$$A_0 = \{(0, 0), (0, 1)\} \quad (37.125a)$$

$$A_1 = \{(1, 0), (1, 1)\} \quad (37.125b)$$

$$B_0 = \{(0, 0), (1, 0)\} \quad (37.125c)$$

$$B_1 = \{(0, 1), (1, 1)\}. \quad (37.125d)$$

Étant donné que $A_i \cap B_j = (i, j)$, nous avons toujours que $P(A_i \cap B_j) = \frac{1}{4} = P(A_i)P(B_j)$. L'indépendance est donc assurée.

Calculons l'espérance conditionnelle $E(S|\mathcal{F}_1)$. Une fonction \mathcal{F}_1 -mesurable doit être constante sur A_0 et A_1 , donc l'espérance conditionnelle est une fonction constante sur A_0 et A_1 dont l'intégrale sur ces ensembles est égale à l'intégrale de S . Nous avons en particulier

$$\int_{A_0} E(S|\mathcal{F}_1) = \int_{A_0} S, \quad (37.126)$$

c'est-à-dire

$$E(S|\mathcal{F}_1)(0,0) + E(S|\mathcal{F}_1)(0,1) = S(0,0) + S(0,1) = 1. \quad (37.127)$$

Nous en concluons que $E(S|\mathcal{F}_1)(0,0) = E(S|\mathcal{F}_1)(0,1) = \frac{1}{2}$. Cela correspond à l'intuition que si on est au point $(0,1)$ ou au point $(0,0)$ en ne sachant que X_1 , nous ne savons que le premier zéro, et donc l'espérance de la somme est $\frac{1}{2}$.

Un calcul très similaire montre que

$$E(S|\mathcal{F}_1)(1,0) = E(S|\mathcal{F}_1)(1,1) = \frac{3}{2}. \quad (37.128)$$

Cela correspond au fait qu'en ces points, nous ne savons que le fait que le premier tirage a donné 1, et donc que l'espérance est $\frac{3}{2}$.

Complétons ce tour d'horizon en mentionnant que la tribu engendrée par X_1 et X_2 est la tribu des parties de Ω , de telle façon que l'espérance conditionnelle de S sachant X_1 et X_2 est égale à S . \triangle

Proposition 37.41 ([236]).

Soit (Ω, \mathcal{A}, P) un espace probabilisé et X, Y deux variables aléatoires sur Ω réelles. Soit \mathcal{B} une sous-tribu de \mathcal{A} . Nous supposons que $X \in L^1(\Omega, \mathcal{A}, P)$, que $Z \in L^\infty(\Omega, \mathcal{B}, P)$ et que $XZ \in L^1(\Omega, P)$. Alors

$$E(ZX|\mathcal{B}) = ZE(X|\mathcal{B}) \quad (37.129)$$

presque sûrement.

Démonstration. Nous commençons par prouver que

$$\int_{\Omega} ZE(X|\mathcal{B}) = \int_{\Omega} ZX. \quad (37.130)$$

Si $Z = \mathbb{1}_B$ pour un ensemble $B \in \mathcal{B}$, alors cette égalité est vraie par définition de l'espérance conditionnelle³. Donc cette égalité est correcte tant que Z est une variable aléatoire \mathcal{B} -mesurable et étagée. Nous considérons alors, grâce au lemme 15.103, une suite Z_n de variables aléatoires étagées et \mathcal{B} -mesurables avec $|Z_n| < Z$. Pour chaque n nous avons donc

$$\int_{\Omega} Z_n X = \int_{\Omega} ZE(X|\mathcal{B}). \quad (37.131)$$

Notre idée est de passer à la limite. Vu que Z et Z_n sont bornées (et donc intégrables sur Ω), pour chaque n nous avons $|Z_n X| \leq M|X|$ où M majore Z et donc tous les Z_n de façon uniforme vis-à-vis de n . Tout cela pour dire que le théorème de la convergence dominée fonctionne et que

$$\lim_{n \rightarrow \infty} \int_{\Omega} Z_n X = \int_{\Omega} ZX. \quad (37.132)$$

D'autre part vu que $X \in L^1$ et que $\Omega \in \mathcal{B}$ nous avons l'égalité $\int_{\Omega} E(X|\mathcal{B}) = \int_{\Omega} X$, ce qui prouve que $|E(X|\mathcal{B})|$ est intégrable. Cela nous permet d'utiliser encore la convergence dominée avec l'inégalité $|Z_n E(X|\mathcal{B})| \leq |E(X|\mathcal{B})|$ pour écrire

$$\lim_{n \rightarrow \infty} \int_{\Omega} Z_n E(X|\mathcal{B}) = \int_{\Omega} ZE(X|\mathcal{B}). \quad (37.133)$$

En passant à la limite des deux côtés de (37.131) nous avons donc

$$\int_{\Omega} ZE(X|\mathcal{B}) = \int_{\Omega} ZX. \quad (37.134)$$

L'égalité (37.130) est prouvée.

3. Théorème 37.31.

Nous passons maintenant à la preuve de l'égalité demandée : $E(EX|\mathcal{B}) = ZE(X|\mathcal{B})$. Pour cela il faut montrer que pour tout $B \in \mathcal{B}$ nous avons

$$\int_B ZE(X|\mathcal{B}) = \int_B ZX. \quad (37.135)$$

Cela n'est rien d'autre que l'égalité (37.130) que nous venons de prouver avec $Z\mathbb{1}_B$ au lieu de Z . \square

Proposition 37.42.

Soit une variable aléatoire réelle $X \in L^1(\Omega, \mathcal{A}, P)$. Pour toute variable aléatoire $Y: \Omega \rightarrow \mathbb{R}^d$, il existe une fonction borélienne \mathcal{A}_Y -mesurable $h: \mathbb{R}^d \rightarrow \mathbb{R}$ telle que

$$E(X|Y) = h \circ Y. \quad (37.136)$$

Démonstration. Nous utilisons le résultat de Doob (théorème 37.14). Par définition $E(X|Y)$ est une variable aléatoire réelle \mathcal{A}_Y -mesurable, et il existe une fonction borélienne $h: \mathbb{R}^d \rightarrow \mathbb{R}$ telle que $E(X|Y) = h \circ Y$. \square

Cette fonction $h: \mathbb{R}^d \rightarrow \mathbb{R}$ nous permet de définir

$$E(X|Z = z) = h(z). \quad (37.137)$$

Cela est l'espérance conditionnelle d'une variable aléatoire par rapport à une valeur donnée d'une autre variable aléatoire.

37.2.9 Probabilité conditionnelle : tribu

Soit un espace probabilisé (Ω, \mathcal{A}, P) .

Lemme 37.43.

Soit $(B_i)_{i \in \mathbb{N}}$ une partition de Ω en éléments de \mathcal{A} deux à deux disjoints tels que $P(B_i) \neq 0$. Soit \mathcal{F} la tribu engendrée par les B_i . Une variable aléatoire réelle est \mathcal{F} -mesurable si et seulement si elle est constante sur chaque B_i .

Proposition 37.44.

Soit $(B_i)_{i \in \mathbb{N}}$ une partition de Ω en éléments de \mathcal{A} deux à deux disjoints tels que $P(B_i) \neq 0$. Soit \mathcal{F} la tribu engendrée par les B_i . Soit une variable aléatoire X . Alors nous avons :

$$E(X|\mathcal{F}) = \sum_{i \in \mathbb{N}} \left(\frac{1}{P(B_i)} \int_{B_i} X dP \right) \mathbb{1}_{B_i}. \quad (37.138)$$

Démonstration. Si X est une variable aléatoire, alors la variable aléatoire $E(X|\mathcal{F})$ définie en 37.31 est une variable aléatoire \mathcal{F} -mesurable et elle est donc constante sur les ensembles B_i par le lemme 37.43 :

$$E(X|\mathcal{F}) = \sum_{i \in \mathbb{N}} a_i \mathbb{1}_{B_i}. \quad (37.139)$$

Étant donné que, par construction, B_i est \mathcal{F} -mesurable, nous avons

$$\int_{B_i} X dP = \int_{B_i} E(X|\mathcal{F}) = \sum_j a_j \int_{B_i} \mathbb{1}_{B_j} = \sum_j a_j \delta_{ij} P(B_j) = a_i P(B_i). \quad (37.140)$$

Par conséquent

$$a_i = \frac{1}{P(B_i)} \int_{B_i} X dP \quad (37.141)$$

et

$$E(X|\mathcal{F}) = \sum_{i \in \mathbb{N}} \left(\frac{1}{P(B_i)} \int_{B_i} X dP \right) \mathbb{1}_{B_i}, \quad (37.142)$$

ce qu'il fallait. \square

Notons que si $B \in \mathcal{A}$ alors la tribu engendrée par B est aussi celle engendrée par la partition $\{B, \complement B\}$ de Ω . Cette circonstance nous permet d'aller plus loin.

Proposition 37.45.

Soit un espace probabilisé (Ω, \mathcal{A}, P) et un événement $B \in \mathcal{A}$ avec sa tribu engendrée $\mathcal{F} = \sigma(B)$. Alors

$$E(\mathbb{1}_A | \mathcal{F}) = P(A|B)\mathbb{1}_B + P(A|\complement B)\mathbb{1}_{\complement B}. \quad (37.143)$$

Démonstration. Nous allons particulariser la formule (37.138). Si $B \in \mathcal{A}$ nous considérons la partition $\{B, \complement B\}$ de Ω et la tribu engendrée

$$\mathcal{F} = \{\emptyset, B, \complement B, \Omega\}. \quad (37.144)$$

La formule (37.138) devient

$$E(X | \mathcal{F}) = \left(\frac{1}{P(B)} \int_B X dP \right) \mathbb{1}_B + \left(\frac{1}{P(\complement B)} \int_{\complement B} X dP \right) \mathbb{1}_{\complement B}. \quad (37.145)$$

Si nous considérons $A \in \mathcal{A}$, nous écrivons cette égalité avec $X = \mathbb{1}_A$ pour obtenir

$$E(\mathbb{1}_A | \mathcal{F}) = \frac{P(A \cap B)}{P(B)} \mathbb{1}_B + \frac{P(A \cap \complement B)}{P(\complement B)} \mathbb{1}_{\complement B} = P(A|B)\mathbb{1}_B + P(A|\complement B)\mathbb{1}_{\complement B} \quad (37.146)$$

parce que nous avons reconnu la probabilité conditionnelle $P(A|B)$ de la définition 37.29. \square

Remarque 37.46.

Les nombres $P(A|\sigma(B)) = P(\mathbb{1}_A | \sigma(B))$ n'est pas la probabilité conditionnelle de A sachant B .

Il nous reste à définir la probabilité conditionnelle d'un événement relativement à une variable aléatoire.

Définition 37.47.

Si la variable aléatoire X est à valeurs discrètes, nous disons que $P(A|X)$ est la variable aléatoire de valeur

$$P(A|X)(\omega) = P(A|X = X(\omega)). \quad (37.147)$$

Dans le cas d'une variable aléatoire à valeurs continues, cette définition ne fonctionne pas parce que la condition $X = X(\omega)$ est souvent de probabilité nulle, tandis que c'est toujours une mauvaise idée de conditionner par rapport à un événement de probabilité nulle. C'est la base du **paradoxe de Borel**. La bonne définition du conditionnement de l'événement A par rapport à la variable aléatoire X est

Définition 37.48.

Si A est un événement et X une variable aléatoire à valeurs continues dans \mathbb{R} , nous définissons

$$P(A|X) = P(A|\sigma(X)) = E(\mathbb{1}_A | \sigma(X)). \quad (37.148)$$

La première égalité est une notation. La seconde est la définition.

Cette définition s'appuie également sur la définition 37.31.

Proposition 37.49.

Si X est une variable aléatoire et si A est un événement, alors

$$E(P(A|X)) = P(A). \quad (37.149)$$

Démonstration. Nous commençons par le cas discret, c'est-à-dire $X: \Omega \rightarrow \mathbb{N}$. Nous notons $p_k = P(X = k)$. En décomposant l'intégrale sur Ω par rapport à l'union disjointe

$$\Omega = \bigcup_{k \in \mathbb{N}} A_k = \bigcup_{k \in \mathbb{N}} \{\omega \in \Omega \text{ tel que } X(\omega) = k\}, \quad (37.150)$$

nous obtenons

$$E(P(A|X)) = \int_{\Omega} P(A|X)(\omega) dP(\omega) \quad (37.151a)$$

$$= \sum_{k=0}^{\infty} \int_{A_k} P(A|X = X(\omega)) dP(\omega) \quad (37.151b)$$

$$= \sum_k \int_{A_k} \frac{P(A \cap X = k)}{P(X = k)} dP(\omega) \quad \text{dans } A_k, X(\omega) = k \quad (37.151c)$$

$$= \sum_k \frac{1}{p_k} P(A \cap X = k) \underbrace{\int_{A_k} 1 dP(\omega)}_{=P(A_k)=p_k} \quad (37.151d)$$

$$= \sum_k P(A \cap X = k) \quad (37.151e)$$

$$= P(A). \quad (37.151f)$$

Nous devons maintenant prouver la propriété dans le cas où X prend des valeurs continues. Pour cela il suffit d'appliquer le corollaire 37.39 :

$$E(E(\mathbb{1}_A | \sigma(A))) = E(\mathbb{1}_A) = P(A). \quad (37.152)$$

□

37.2.10 Variables de Rademacher indépendantes

Une variable aléatoire de Rademacher est une variable aléatoire qui prend les valeurs 1 et -1 avec probabilité $\frac{1}{2}$. Nous pouvons en décrire une explicitement de la façon suivante. L'espace probabilité est à deux éléments : $\Omega = \{a, b\}$ avec la mesure $P(\{a\}) = P(\{b\}) = \frac{1}{2}$. La variable aléatoire est alors l'application $X: \Omega \rightarrow \mathbb{R}$ donnée par $X(a) = 1$ et $X(b) = -1$.

Soient X et Y deux variables aléatoires de Rademacher indépendantes. Cela donne $\Omega = \{a, b\}^2$ et

$$\begin{aligned} X(a, a) &= 1 & X(a, b) &= 1 & X(b, a) &= -1 & X(b, b) &= -1 \\ Y(a, a) &= 1 & Y(a, b) &= -1 & Y(b, a) &= 1 & Y(b, b) &= -1 \end{aligned} \quad (37.153)$$

Remarque 37.50.

Si une variable aléatoire d'un certain type est donnée par une application $X: \Omega \rightarrow \mathbb{R}$, pour construire des variables aléatoires indépendantes identiquement distribuées, il faut considérer les variables aléatoires sur (au moins) le produit $\Omega \times \Omega$ munie de la mesure produit.

Tribu du produit XY Quelle est la tribu de la variable aléatoire produit XY ? Le produit XY peut prendre les valeurs 1 et -1 . Nous avons

$$(XY)^{-1}(1) = \{(a, a), (b, b)\} \quad (XY)^{-1}(-1) = \{(a, b), (b, a)\} \quad (37.154)$$

La tribu est donc

$$\sigma(XY) = \{\Omega, \emptyset, A, B\} \quad (37.155)$$

avec $A = \{(a, a), (b, b)\}$ et $B = \{(a, b), (b, a)\}$.

Calcul de $E(X|XY)$ La définition de l'espérance à calculer est le théorème 37.31. Pour chaque élément B de $\sigma(XY)$ nous avons besoin de $\int_B X = \int_B E(X|XY)$. Nous notons $V = E(X|XY)$ pour alléger la notation. Nous avons

$$4 \int_A V = V(a, a) + V(b, b) \quad (37.156)$$

et

$$4 \int_A X = X(a, a) + X(b, b) = 0. \quad (37.157)$$

Pourquoi le facteur 4? Parce que sur Ω nous avons la mesure produit de celle que dont nous avons parlé sur $\{a, b\}$. C'est la mesure d'équiprobabilité et donc chaque singleton a mesure $1/4$. Pour plus de détails, il y a le théorème 15.207.

Nous en déduisons $V(a, a) + V(b, b) = 0$. Mais pour tout $t \in \mathbb{R}$ nous avons $V^{-1}(t) \in \sigma(XY)$ parce que la contrainte est que V soit XY -mesurable. En particulier

$$V^{-1}(V(a, a)) \quad (37.158)$$

est un mesurable qui contient (a, a) . C'est donc soit Ω , soit $\{(a, a), (b, b)\}$. Dans les deux cas nous avons $V(a, a) = V(b, b)$ et nous en déduisons $V(a, a) = V(b, b) = 0$.

En faisant de même avec $\int_B V = V(a, b) + V(b, a)$ nous déduisons $V(a, b) = V(b, a) = 0$ et au final nous avons

$$E(X|XY) = 0. \quad (37.159)$$

Cette égalité signifie $E(X|XY)(\omega) = 0$ pour tout $\omega \in \Omega$.

Calcul de $E(X|X+Y)$ Il ne faudrait pas croire que, seulement parce que X a une espérance nulle, nous trouverons une espérance nulle quel que soit le conditionnement. Juste pour le plaisir, nous calculons $E(X|X+Y)$.

La variable aléatoire $X+Y$ peut prendre trois valeurs : $-2, 0$ et 2 . La tribu engendrée par $X+Y$ doit en particulier contenir $A = \{(a, a)\}$, $B = \{(b, b)\}$ et $C = \{(a, b), (b, a)\}$.

Nous notons $V = E(X|\sigma(X+Y))$. Vu que

$$\int_A V = \int_A V, \quad (37.160)$$

nous avons $V(a, a) = X(a, a) = 1$. Même chose pour B qui donne $V(b, b) = X(b, b) = -1$. En ce qui concerne l'intégrale sur C nous avons

$$V(a, b) + V(b, a) = X(a, b) + X(b, a) = 0. \quad (37.161)$$

Par ailleurs l'ensemble $V^{-1}(V(a, b))$ est un ensemble mesurable qui doit au moins contenir (a, b) . Vu la tribu que nous avons, cela doit également contenir (b, a) , de telle sorte que $V(a, b) = V(b, a)$. La relation (37.161) nous permet alors de conclure que $V(a, b) = V(b, a) = 0$.

Quoi qu'il en soit, l'espérance conditionnelle $E(XY|X+Y)$ n'est pas nulle.

Calcul de $E(XY|\sigma(XY))$. Celle-là, elle est facile par 37.32 : c'est XY .

Nous aurions pu croire que si X et Y sont indépendantes, alors

$$E(XY|\mathcal{A}) = E(X|\mathcal{A})E(Y|\mathcal{A}). \quad (37.162)$$

L'exemple que nous venons de faire montre qu'il n'en est rien.

Exemple 37.51 ([484])

Un autre exemple, peut-être plus simple, pour contredire l'équation (37.162). Soient X et Y des expériences indépendantes de pile ou face non truquées. Les résultats sont représentés par 0 et 1. Nous notons \mathcal{A} la tribu engendrée par l'événement « les résultats des deux lancers sont différents » ; c'est-à-dire la tribu engendrée par l'événement $A = (1, 0), (0, 1)$. La variable aléatoire X et la tribu \mathcal{A} sont indépendants (définition 37.6), donc, donc $E(X|\mathcal{A}) = E(X) = 1/2$. Pareil pour Y . En revanche, le produit XY est nul sur A donc $E(XY|\mathcal{A})$ aussi. Ça ne peut donc être égal à la constante $1/4 = (1/2)^2$. \triangle

37.2.11 Un petit paradoxe

Attention : ce qui est écrit ici est ma réflexion personnelle sur le sujet. Merci de me dire si je me trompe.

Soit une famille dont vous savez seulement qu'il y a exactement deux enfants. Trois situations :

- (1) Vous frappez, une fille ouvre la porte et dit « Bonjour, je suis l'aînée ». Quelle est la probabilité que l'autre enfant soit une fille ?
- (2) Vous frappez, une fille ouvre la porte et dit « Bonjour ». Quelle est la probabilité que l'autre enfant soit une fille ?
- (3) Vous demandez aux parents s'il y a au moins une fille, ils répondent « oui ». Quelle est la probabilité que les deux enfants soient des filles ?

Dans les trois cas l'intuition dit que la probabilité est $1/2$. Il semble que de plus la (2) et la (3) soient les mêmes parce que l'on sait qu'il y a une fille et on se demande quelle est la probabilité qu'il y ait deux filles.

Nous allons voir ça de plus près.

37.2.11.1 « Bonjour, je suis l'aînée »

Résolution Si nous notons X_0 et X_1 les variables aléatoires donnant le sexe des deux enfants, ce sont des variables aléatoires indépendantes et identiquement distribuées, avec $P(X_i = f) = \frac{1}{2}$. La formule (37.73) de la probabilité conditionnelle ainsi que l'indépendance donnent :

$$P(X_1 = f | X_2 = f) = \frac{P(X_1 = f, X_2 = f)}{P(X_2 = f)}. \quad (37.163)$$

Le numérateur vaut $\frac{1}{4}$ et le dénominateur vaut $\frac{1}{2}$; le résultat vaut $\frac{1}{2}$. Fin de l'histoire.

Simulation Voici un petit programme qui simule la situation. Il retourne clairement $1/2$.

```

1  #! /usr/bin/python3
2  # -*- coding: utf8 -*-
3
4  """
5  Vous frappez à la porte d'une famille qui a deux enfants. Une ←
6  fille ouvre la porte et vous dit "Je suis l'aînée".
7  Quelle est la probabilité que l'autre soit une fille ?
8  """
9
10 import random
11
12 def famille():
13     """
14     return a pair of 'f' and 'g'.
15     """
16     a=[random.choice( ['g', 'f'] )]
17     a.append(random.choice( ['g', 'f'] ))
18     return a
19
20 def toctoc():
21     """
22     - Create a family with two children.
23     - Pick the second one, the elder.
```

```

24     - if it is a 'g', return None.
25     - if it is a 'f', return the other one.
26     """
27     F=famille()
28     if F[1] != 'f':
29         return None
30     else :
31         if F[0]=='f':
32             return 1
33         else :
34             return 0
35
36     N_girl_opens=0
37     N_girl_other=0
38     for k in range(1,10000):
39         res=toctoc()
40         if res is not None:
41             N_girl_opens = N_girl_opens+1
42             N_girl_other = N_girl_other + res
43
44     proba=N_girl_other/N_girl_opens
45     print(proba)           # ~0.5, intuitively correct.

```

tex/sage/simul_famille_aine.py

37.2.11.2 « Bonjour »

Nous frappons à la porte, une fille ouvre en disant « bonjour », sans préciser si elle est la première ou la seconde. Quelle est la probabilité que l'autre soit une fille? Naïvement on croirait que la probabilité est également $\frac{1}{2}$.

Un raisonnement moins naïf montre le contraire.

Et nous allons voir qu'un raisonnement encore moins naïf montre que la probabilité est bien $\frac{1}{2}$.

Premier raisonnement (incorrect) Voici le raisonnement qui est, à mon avis, faux. Vu que l'enfant qui ouvre la porte est une fille, la famille a une des compositions suivantes : fg , ff ou gf . Le cas où une fille ouvre la porte *et* que l'autre est également une fille est seulement le cas ff dont la probabilité est $\frac{1}{3}$.

Pour justifier cela nous considérons le couple de variables aléatoires (X_1, X_2) et le conditionnement $A = \{X_1 = f\} \cup \{X_2 = f\}$: évidemment $P(A) = \frac{3}{4}$. Nous calculons facilement la loi du couple (X_1, X_2) conditionné à A :

$$P(X_1 = f, X_2 = f|A) = \frac{P(\{X_1 = f, X_2 = f\} \cap A)}{P(A)} = \frac{1/4}{3/4} = \frac{1}{3}. \quad (37.164)$$

Donc sachant A , la probabilité que la famille soit constituée de deux filles est $\frac{1}{3}$.

Comment faire mieux? Ce calcul semble être correct, mais il ne l'est pas. Ce raisonnement fait l'hypothèse implicite que l'espace probabilisé décrivant la situation contient deux variables aléatoires X_1 et X_2 représentant les deux enfants. Or nous avons bien trois événements aléatoires dans l'histoire : le sexe des deux enfants et le *choix* de l'enfant qui ouvre la porte.

Certes, nous pouvons penser que cette troisième variable aléatoire ne change rien. Oui oui, on peut le penser. Mais ici, on ne doit pas penser, on doit *démontrer*.

Nous allons donc rédiger un calcul complet, en introduisant toutes les variables aléatoires, et en décrivant correctement l'espace probabilisé Ω et la mesure de probabilité P .

Peut-être que ça ne changera rien. Ou peut-être pas. Mais au moins nous serons sûrs d'avoir résolu le problème correctement.

La vraie réponse Nous considérons les variables aléatoires $X_0, X_1: \Omega_E \rightarrow \{f, g\}$ avec probabilité $\frac{1}{2}$. De plus nous considérons une nouvelle variable aléatoire qui donne le numéro de l'enfant qui ouvre la porte :

$$\sigma: \Omega_C \rightarrow \{1, 2\}. \quad (37.165)$$

Notre espace de probabilité est donc l'ensemble $\Omega = \{f, g\} \times \{f, g\} \times 0, 1$ sur lequel nous considérons la mesure d'équiprobabilité⁴.

Nous introduisons les variables aléatoires⁵

$$\begin{aligned} X_1: \Omega &\rightarrow \{f, g\} \\ (s_1, s_2, n) &\mapsto s_1 \end{aligned} \quad (37.166)$$

et

$$\begin{aligned} X_2: \Omega &\rightarrow \{f, g\} \\ (s_1, s_2, n) &\mapsto s_2 \end{aligned} \quad (37.167)$$

et

$$\begin{aligned} \sigma: \Omega &\rightarrow \{1, 2\} \\ (s_1, s_2, n) &\mapsto n \end{aligned} \quad (37.168)$$

Nous devons calculer

$$P(X_{1-\sigma} = f | X_\sigma = f) = \frac{P(X_{1-\sigma} = f, X_\sigma = f)}{P(X_\sigma = f)}. \quad (37.169)$$

Pour être explicite jusqu'au bout, nous énumérons tous les éléments de Ω :

- | | | | |
|---------------|---------------|---------------|-----------------|
| (1) $g, g, 0$ | (3) $g, f, 0$ | (5) $f, g, 0$ | (7) $f, f, 0$ |
| (2) $g, g, 1$ | (4) $g, f, 1$ | (6) $f, g, 1$ | (8) $f, f, 1$. |

Et tant qu'à être explicite, l'événement vulgairement noté $\{X_\sigma = f\}$ est la partie

$$\{X_\sigma = f\} = \{\omega \in \Omega \text{ tel que } X_{\sigma(\omega)}(\omega) = f\} \quad (37.170a)$$

$$= \{(s_1, s_2, n) \text{ tel que } X_n(s_1, s_2, n) = f\} \quad (37.170b)$$

$$= \{(s_1, s_2, n) \text{ tel que } s_n = f\}. \quad (37.170c)$$

Méditez la dernière égalité; elle n'est pas totalement indispensable au raisonnement, mais elle est cool.

Nous avons

$$\{X_\sigma = f\} = \{(g, f, 1), (f, g, 0), (f, f, 0), (f, f, 1)\}. \quad (37.171)$$

et

$$\{X_{1-\sigma} = f\} \cap \{X_\sigma = f\} = \{(f, f, 0), (f, f, 1)\}. \quad (37.172)$$

Donc

$$P(X_{1-\sigma} = f, X_\sigma = f) = \frac{2}{8} = \frac{1}{4} \quad (37.173)$$

et

$$P(X_\sigma = f) = \frac{4}{8} = \frac{1}{2}. \quad (37.174)$$

Au final,

$$P(X_{1-\sigma} = f | X_\sigma = f) = \frac{1/4}{1/2} = \frac{2}{4} = \frac{1}{2}. \quad (37.175)$$

4. C'est une hypothèse forte faisant appel d'un part ce que l'on sait de la reproduction humaine, et d'autre part ce que l'on sait de la sociologie de deux enfants qui entendent une sonnette.

5. Sur Ω , sur $\{f, g\}$ et sur $\{0, 1\}$ nous mettons la tribu des parties. Vérifiez que X_1, X_2 et σ sont mesurables.

Simulation Vous avez encore un doute ? Faites tourner la simulation suivante :

```

1  #! /usr/bin/python3
2  # -*- coding: utf8 -*-
3
4  """
5  Vous frappez à la porte d'une famille qui a deux enfants. Une ←
6  fille ouvre la porte.
7  Quelle est la probabilité que l'autre soit une fille ?
8  """
9
10 import random
11
12 def famille():
13     """
14     return a pair of 'f' and 'g'
15     """
16     a=[random.choice( ['g','f'] )]
17     a.append(random.choice( ['g','f'] ))
18     return a
19
20 def toctoc():
21     """
22     - Create a family with two children.
23     - Choose one (the one who opens the door)
24     - if it is a 'g', return None.
25     - if it is a 'f', return the other one.
26     """
27     F=famille()
28     s=random.choice([0,1])
29     if F[s] != 'f':
30         return None
31     else :
32         t=(s+1)%2
33         if F[t]=='f':
34             return 1
35         else :
36             return 0
37
38 N_girl_opens=0
39 N_girl_other=0
40 for k in range(1,10000):
41     res=toctoc()
42     if res is not None:
43         N_girl_opens = N_girl_opens+1
44         N_girl_other = N_girl_other + res
45
46 proba=N_girl_other/N_girl_opens
47 print(proba)          # ~0.5, intuitively correct.

```

tex/sage/simul_famille_simple.py

Le faisant tourner, la réponse est sans appel : la fréquence observée est beaucoup plus proche

de 0.5 que de 0.33 ou 0.66.

37.2.11.3 Le parent qui répond aux questions

Nous avons une famille de deux enfants dont nous savons qu'au moins un des deux est une fille. Quelle est la probabilité que la famille contienne deux filles ? Cela est a priori la même question que celle où une fille ouvre la porte sans dire si elle est l'aînée ou non.

Simulation Commençons par la simulation :

```
1  #!/usr/bin/python3
2  # -*- coding: utf8 -*-
3
4  """
5  Vous savez qu'une famille a deux enfants.
6  Vous demandez à un parent si il y a une fille.
7  Réponse : Oui.
8  Question : quelle est la probabilité que ce soient deux filles ?
9  """
10
11
12 import random
13
14 def famille():
15     """
16     return a pair of 'f' and 'g'.
17     """
18     a=[random.choice( ['g','f'] )]
19     a.append(random.choice( ['g','f'] ))
20     return a
21
22 def toctoc():
23     """
24     - Create a family with two children.
25     - If 'gg', there are no girls -> return None
26     - If 'gf', there is a girl but the other is a boy -> 0
27     - If 'fg', there is a girl but the other is a boy -> 0
28     - If 'ff', there is a girl and the other is a girl -> 1
29     """
30     F=famille()
31     if F==['g','g'] :
32         return None
33     if F==['g','f']:
34         return 0
35     if F==['f','g']:
36         return 0
37     if F==['f','f']:
38         return 1
39
40 N_at_least_one_girl=0
41 N_two_girls=0
42 for k in range(1,10000):
43     res=toctoc()
```

```

44     if res is not None:
45         N_at_least_one_girl=N_at_least_one_girl+1
46         N_two_girls=N_two_girls+res
47
48 proba=N_two_girls/N_at_least_one_girl
49 print(proba)           # ~0.333. Beware !

```

tex/sage/simul_famille_une_fille.py

Et là, bing, la réponse est clairement plutôt 0.33 que 0.5.

Résolution Nous avons les variables aléatoires X_1 et X_2 qui valent 0 ou 1 suivant que l'enfant soit une fille ou un garçon ; ce sont des variables aléatoires indépendantes et identiquement distribuées. Nous définissons la variable aléatoire somme

$$S = X_1 + X_2 \quad (37.176)$$

qui compte le nombre de filles. La question est de calculer

$$P(S = 2 | S \geq 1) = \frac{P(S = 2 \cap S \geq 1)}{P(S \geq 1)} = \frac{P(S = 2)}{P(S \geq 1)}. \quad (37.177)$$

L'événement $S = 2$ est réduit au singleton $\{ff\}$ et sa probabilité est $\frac{1}{4}$. Au contraire l'événement $S \geq 1$ est l'ensemble $\{fg, gf, ff\}$ et sa probabilité est $\frac{3}{4}$. Nous avons donc

$$P(S = 2 | S \geq 1) = \frac{1/4}{3/4} = \frac{1}{3}. \quad (37.178)$$

Et là, la réponse est $1/3$ et non $1/2$ comme d'aucuns auraient pu le croire.

Précision Notons que l'événement $S \geq 1$ n'est pas le même que l'événement $X_\sigma = f$. En effet

$$S \geq 1 = \{(ff, 1), (ff, 2), (fg, 1), (fg, 2), (gf, 1), (gf, 2)\} \quad (37.179)$$

tandis que

$$\{X_\sigma = f\} = \{(ff, 1), (ff, 2), (fg, 1), (gf, 2)\}. \quad (37.180)$$

37.2.11.4 Conclusion

L'internet regorge de sites discutant du paradoxe des deux enfants⁶.

Beaucoup insistent sur le fait que non seulement certaines informations apparemment anodines sont importantes, mais en plus *la façon* dont on obtient l'information est importante. Dans la situation « une fille ouvre », nous obtenons l'information « il y a au moins une fille » en voyant une ; dans la situation « la parent dit qu'il y a au moins une fille », nous obtenons l'information « il y a au moins une fille » de façon plus « pure ».

Personnellement je ne souscris pas vraiment à cette façon de penser. Le fait est que la formule

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad (37.181)$$

n'est pas seulement une formule dans laquelle il faut remplacer A par « la question » et B par « ce qu'on sait ». Il faut également remplacer P par « la bonne » mesure de probabilité.

Il est important de construire le bon espace de probabilité, avec la bonne mesure. Et pour cela, il faut bien s'assurer d'introduire une variable aléatoire pour chaque événement aléatoire se produisant dans l'histoire.

6. Par exemple [485].

37.2.11.5 À propos des simulations

Si vous lisez ces lignes avec l'intention de passer l'agrégation en utilisant Sage à l'épreuve de modélisation, vous devez être capable de refaire les trois simulations. Les bouts de code donnés ici sont écrits pour python3 alors que Sage utilise Python2. Je ne vous dit pas si ça change quelque chose.

Allez oui, je vous dit. Si vous changez dans `simul_famille_une_fille.py` la première ligne pour utiliser python2 au lieu de python3, le résultat affiché sera 0 et non 0.333. La raison est que dans Python2, l'opérateur / entre deux entiers est une **division entière**. Autrement dit : le résultat 0.33 est arrondi à zéro.

Solution : forcer python à interpréter le / comme une vraie division. Pour Sage, ça donne ceci comme début de programme :

```

1 #! /usr/bin/sage
2 # -*- coding: utf8 -*-
3
4 from __future__ import division

```

tex/frido/codeSnip_3.py

Importez toujours `division` de `__future__` .

Ah oui, et dernière remarque : pour autant que je le sache, le jour de l'oral, vous n'aurez que Sage en mode notebook. Je ne sais pas si l'import fonctionne aussi bien.

Sinon vous pouvez forcer la division dans les `float` de la façon suivante : `a/float(b)`.

37.2.12 Inégalité de Jensen

Proposition 37.52 (Inégalité de Jensen).

Soit g une fonction convexe⁷ sur \mathbb{R} et une variable aléatoire $Y \in L^1(\Omega, \mathcal{A}, P)$ telle que $g \circ Y$ soit également L^1 . Alors

$$g(E(Y|\mathcal{F})) \leq E((g \circ Y)|\mathcal{F}). \quad (37.182)$$

Démonstration. La convexité de g et la proposition 18.89 nous donnent deux suites (a_n) et (b_n) dans \mathbb{R} telles que pour tout $x \in \mathbb{R}$,

$$g(x) = \sup_{n \in \mathbb{N}} (a_n x + b_n). \quad (37.183)$$

Nous avons alors

$$a_n E(Y|\mathcal{F}) \stackrel{p.s.}{=} E(a_n Y + b_n|\mathcal{F}) \leq E(g \circ Y|\mathcal{F}). \quad (37.184)$$

L'inégalité est due au fait que $g \circ Y$ est le supremum sur les n de $a_n Y + b_n$. Pour chaque n , l'inégalité (37.184) est fautive sur un ensemble de mesure nulle $R_n \subset \Omega$. L'union

$$R = \bigcup_{n \in \mathbb{N}} R_n \quad (37.185)$$

est encore de mesure nulle. Sur $\Omega \setminus R$, nous avons

$$a_n E(Y|\mathcal{F}) + b_n \leq E(g \circ Y|\mathcal{F}). \quad (37.186)$$

Vu que cela est vrai presque partout et pour tout n nous passons à supremum et nous avons encore presque partout l'inégalité

$$\sup_{n \in \mathbb{N}} (a_n E(Y|\mathcal{F}) + b_n) \leq E(g \circ Y|\mathcal{F}). \quad (37.187)$$

□

Si nous ne nous intéressons pas à $E(Y|\mathcal{F})$ mais seulement à $E(Y)$, alors une démonstration plus simple est donnée sur Wikipédia[486].

7. Définition 18.76.

37.2.13 Fonction de répartition

Définition 37.53.

Si X est une variable aléatoire réelle, nous définissons sa **fonction de répartition** par

$$\begin{aligned} F_X: \mathbb{R} &\rightarrow [0, 1] \\ F_X(x) &= P(X \leq x). \end{aligned} \quad (37.188)$$

Remarque 37.54.

La fonction de répartition est discontinue en a si $P(X = a) > 0$. En particulier nous ne pouvons pas dire

$$P(X \geq a) = 1 - F_X(a). \quad (37.189)$$

37.2.14 Fonction caractéristique

Définition 37.55.

La **fonction caractéristique** de la variable aléatoire $X: \Omega \rightarrow \mathbb{R}$ est la fonction réelle définie par

$$\Phi_X(t) = E(e^{itX}). \quad (37.190)$$

Une autre façon d'écrire la définition est

$$\Phi_X(t) = \int_{\mathbb{R}} e^{itx} dP_X(x), \quad (37.191)$$

ou encore, si X a une densité f_X ,

$$\Phi_X(t) = \int_{\mathbb{R}} e^{itx} f_X(x) dx \quad (37.192)$$

Nous reconnaissons la transformée de Fourier :

$$\Phi_X(t) = \hat{f}_X(-t/2\pi). \quad (37.193)$$

La proposition suivante se déduit en utilisant le théorème de dérivation sous l'intégrale [18.21](#).

Proposition 37.56.

Soit X une variable aléatoire qui accepte un moment d'ordre $r \geq 1$. Alors la fonction caractéristique Φ_X est r fois continument dérivable et

$$\Phi_X^{(r)}(t) = E((iX)^r e^{itX}). \quad (37.194)$$

Démonstration. Nous étudions la fonction

$$\Phi(t) = \int_{\Omega} e^{itX(\omega)} dP(\omega). \quad (37.195)$$

Nous considérons la fonction

$$\begin{aligned} f: \mathbb{R} \times \Omega &\rightarrow \mathbb{R} \\ (t, \omega) &\mapsto e^{itX(\omega)}. \end{aligned} \quad (37.196)$$

et nous regardons si ce contexte vérifie les hypothèses du théorème [18.21](#).

- (1) Étant donné que X est mesurable, f sera mesurable.
- (2) La fonction $t \mapsto e^{itX(\omega)}$ est absolument continue pour chaque ω .

(3) Note : par rapport aux notations du théorème 18.21, nous avons ici $A = \mathbb{R}$. Prenons donc un intervalle (compact) $[a, b] \subset \mathbb{R}$ et calculons

$$\frac{\partial f}{\partial t}(t, \omega) = iX(\omega)e^{itX}, \quad (37.197)$$

et

$$\int_a^b \int_{\Omega} |iX(\omega)e^{itX(\omega)}| d\omega dt = \int_a^b \int_{\Omega} |X(\omega)| d\omega dt. \quad (37.198)$$

Par hypothèse X accepte un moment d'ordre 1, de sorte que l'intégrale par rapport à ω converge vers un nombre qui ne dépend pas de t . L'intégrale sur t ne pose alors aucun problèmes.

Par conséquent nous pouvons effectuer la première dérivation :

$$\Phi'(t) = \frac{d\Phi}{dt}(t) = \int_{\Omega} iX(\omega)e^{itX(\omega)} d\omega = E(iXe^{itX}) \quad (37.199)$$

et la fonction Φ' est absolument continue. Ce dernier point est important parce que c'est lui qui permet de faire la récurrence et passer à l'ordre deux.

Le résultat ressort alors en dérivant successivement l'expression (37.197). \square

Exemple 37.57

Sachant la fonction caractéristique de X , nous pouvons calculer les moments. Par exemple

$$E(X^2) = \Phi_X''(0). \quad (37.200)$$

\triangle

Théorème 37.58.

Si $\Phi_X = \Phi_Y$, alors $P_X = P_Y$.

Notons que cela n'implique pas que $X = Y$. En effet X et Y peuvent même être définis sur des espaces probabilisés différents.

Dans le cas d'une variable aléatoire vectorielle, nous définissons $\Phi_X: \mathbb{R}^d \rightarrow \mathbb{R}$ par

$$\Phi_X(v) = E(e^{i\langle v, X \rangle}) \quad (37.201)$$

37.2.15 Fonction génératrice des moments, transformée de Laplace

Soit X une variable aléatoire. Sa **transformée de Laplace** ou **fonction génératrice des moments** est la fonction

$$M_X(t) = E(e^{tX}) \quad (37.202)$$

pour chaque t tel que cette espérance existe.

Théorème 37.59 ([487]).

Soit X une variable aléatoire réelle et

$$I_X = \{t \in \mathbb{R} \text{ tel que } E(e^{tX}) \text{ existe}\}. \quad (37.203)$$

La fonction

$$\begin{aligned} M_X: I_X &\rightarrow \mathbb{R} \\ t &\mapsto E(e^{tX}) \end{aligned} \quad (37.204)$$

est la **transformée de Laplace** de X .

(1) I_X est un intervalle contenant 0.

(2) Si I_X n'est pas réduit à $\{0\}$ alors M_X se développe en série entière

$$M_X(t) = \sum_{n=0}^{\infty} \frac{E(X^n)}{n!} t^n. \quad (37.205)$$

(3) Si X et Y sont des variables aléatoires indépendantes, alors $I_{X+Y} = I_X \cap I_Y$ et

$$M_{X+Y} = M_X M_Y \quad (37.206)$$

sur I_{X+Y} .

Démonstration. Le fait que 0 soit dans I_X est évident : $E(1) = 1$. Pour montrer que I_X est un intervalle nous prenons $z \in I_X$ et $0 < s < z$ ou $z < s < 0$, puis nous montrons que $s \in I_X$. Il faut remarquer que dans tous les cas,

$$e^{sX} \leq 1 + e^{zX}. \quad (37.207)$$

En effet soit sX et zX sont tous deux à gauche de zéro et alors ils sont tous deux plus petit que 1 ; soit ils sont tous deux à droite de 0 et alors $e^{zX} > e^{sX}$ par croissance de l'exponentielle. Nous avons donc dans tous les cas que

$$E(e^{sX}) = \int_{\mathbb{R}} f_X(x) e^{sX} dx \leq \int_{\mathbb{R}} f_X(x) (1 + e^{zx}) = 1 + E(e^{zX}). \quad (37.208)$$

Soit maintenant $a > 0$ tel que $[-a, a] \in I_X$. Étant donné que $e^{a|X|} < e^{aX} e^{-aX}$, l'espérance $E(e^{a|X|})$ existe toujours pour $|t|$. Nous avons

$$\left| M_X(t) - \sum_{n=0}^N \frac{E(X^n)}{n!} t^n \right| = \left| E \left(e^{tX} - \sum_{n=0}^N \frac{X^n}{n!} t^n \right) \right| \quad (37.209a)$$

$$= \left| E \left(\sum_{n=N+1}^{\infty} \frac{X^n}{n!} t^n \right) \right| \quad (37.209b)$$

$$\leq E \left(\sum_{n=N+1}^{\infty} \frac{|tX|^n}{n!} \right). \quad (37.209c)$$

Maintenant le but est de prendre la limite $N \rightarrow \infty$ en inversant la limite et l'espérance par le théorème de la convergence dominée (15.184). L'intégrale à traiter est

$$\lim_{N \rightarrow \infty} \int_{\Omega} \sum_{n=N+1}^{\infty} \frac{|tX(\omega)|^n}{n!} dP(\omega). \quad (37.210)$$

L'intégrande est uniformément borné (en N) par $e^{tX(\omega)}$, qui est intégrable par hypothèse (choix de t). Du coup

$$\lim_{N \rightarrow \infty} \int_{\Omega} \sum_{n=N+1}^{\infty} \frac{|tX(\omega)|^n}{n!} dP(\omega) = E \left(\lim_{N \rightarrow \infty} \sum_{n=N+1}^{\infty} \frac{|tX|^n}{n!} \right) = 0. \quad (37.211)$$

□

37.2.16 Loi d'une variable aléatoire

La loi de la variable aléatoire X , notée P_X est la mesure image de P par X , c'est-à-dire

$$P_X(B) = P(X \in B) \quad (37.212)$$

pour tout borélien $B \subset \mathbb{R}^d$. Note :

$$P(X \in B) = P(\{\omega \in \Omega \text{ tel que } X(\omega) \in B\}) = P(X^{-1}(B)). \quad (37.213)$$

En particulier P_X est une mesure de probabilité sur \mathbb{R}^d parce que

$$P_X(\mathbb{R}^d) = P(\Omega) = 1. \quad (37.214)$$

Si Q est une mesure de probabilité sur \mathbb{R}^d , nous notons $X \sim Q$ si $P_X = Q$. Nous disons alors que « X suit la loi Q ».

La proposition suivante permet de calculer en pratique les intégrales qui définissent par exemple l'espérance mathématique d'une variable aléatoire.

Proposition 37.60 (Théorème de transfert[488]).

Si X est une variable aléatoire, alors

$$E(f \circ X) = \int_{\Omega} f(X(\omega))dP(\omega) = \int_{\mathbb{R}^d} f(x)dP_X(x) \quad (37.215)$$

dès que $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ est telle qu'une des deux intégrales existe. En particulier, ça marche si f est borélienne.

En utilisant cette proposition nous trouvons une formule pratique pour l'espérance d'une variable aléatoire réelle :

$$E(X) = \int_{\Omega} X(\omega)dP(\omega) = \int_{\mathbb{R}} xdP_X(x), \quad (37.216)$$

en vertu de la proposition 37.60 appliquée à la fonction $f(x) = x$.

Proposition 37.61.

Une variable aléatoire réelle X est intégrable si et seulement si $P(x = \pm\infty) = 0$ et

$$\int_{\mathbb{R}} |x|dP_X(x) < \infty. \quad (37.217)$$

Le lien entre la densité f_X de la variable aléatoire X et sa loi est

$$P_X(A) = \int_A f_X(x)dx \quad (37.218)$$

pour tout ensemble mesurable $A \subset \mathbb{R}$. Le lien entre la mesure de Lebesgue et celle de la loi de X est alors donné par

$$dP_X(x) = f_X(x)dx. \quad (37.219)$$

En particulier l'espérance de X peut être calculée à partir de sa densité via la formule

$$E(X) = \int_{\mathbb{R}} xdP_X(x) = \int_{\mathbb{R}} xf_X(x)dx. \quad (37.220)$$

37.2.17 Changement de variables

Théorème 37.62.

Soit \mathcal{O} , un ouvert de \mathbb{R}^n et \mathcal{O}' un ouvert de \mathbb{R}^m ainsi qu'un difféomorphisme C^1 $\varphi: \mathcal{O} \rightarrow \mathcal{O}'$. Soit $X: \Omega \rightarrow \mathbb{R}^n$ une variable aléatoire prenant presque sûrement ses valeurs dans \mathcal{O} . Si nous supposons que X a la densité f_X , alors la variable aléatoire $Y = \varphi(X)$ accepte la densité $f_Y: \mathcal{O}' \rightarrow \mathbb{R}$ donnée par

$$f_Y(v) = f_X(\varphi^{-1}(v))|J_{\varphi^{-1}}(v)|. \quad (37.221)$$

Démonstration. Nous devons vérifier la relation

$$P(Y \in B) = \int_B f_Y(v)dv \quad (37.222)$$

pour tout borélien $B \subset \mathcal{O}'$. Nous avons

$$P(Y \in B) = \int_{\mathbb{R}^m} \mathbb{1}_B(v) dP_Y(v) \quad (37.223a)$$

$$= E(\mathbb{1}_B \circ Y) \quad (37.223b)$$

$$= E((\mathbb{1}_B \circ \varphi) \circ X) \quad (37.223c)$$

$$= \int_{\mathbb{R}^n} (\mathbb{1}_B \circ \varphi)(u) dP_X(u) \quad (37.223d)$$

$$= \int_{\mathbb{R}^n} (\mathbb{1}_B \circ \varphi)(u) f_X(u) du. \quad (37.223e)$$

À ce niveau, nous utilisons la formule de changement de variables du théorème 15.252. Nous trouvons alors

$$P(Y \in B) = \int_{\mathbb{R}^m} \mathbb{1}_B(\varphi^{-1}(v)) f_X(\varphi^{-1}(v)) |J_{\varphi^{-1}}(v)| dv. \quad (37.224)$$

□

37.3 Convergence

Soient X_i des variables aléatoires réelles définies sur le même espace de probabilité (Ω, \mathcal{A}, P) . Nous disons que X_i converge **presque sûrement** vers la variable aléatoire X et nous notons

$$X_n \xrightarrow{p.s.} X \quad (37.225)$$

si

$$P(\{\omega \in \Omega \text{ tel que } X_n(\omega) \rightarrow X(\omega)\}) = 1 \quad (37.226)$$

où la convergence $X_n(\omega) \rightarrow X(\omega)$ est la convergence usuelle dans \mathbb{R} .

Lemme 37.63.

Nous avons $X_n \xrightarrow{p.s.} X$ si et seulement s'il existe un événement $A \in \mathcal{A}$ tel que $P(A) = 1$ et tel que $X_n(\omega) \rightarrow X(\omega)$ pour tout $\omega \in A$.

Lemme 37.64.

Si X est une variable aléatoire à valeurs dans $\mathbb{R} \cup \{\infty\}$, alors

$$X \wedge n \xrightarrow{p.s.} X \quad (37.227)$$

Démonstration. Si $\omega \in \{X = \infty\}$ alors $(X \wedge n)(\omega) = n$ et d'accord. Si par contre $\omega \in \{X < \infty\}$ alors il existe N tel que si $n \geq N$ alors $n \geq T(\omega)$ et pour ces grandes valeurs de n nous avons $(X \wedge n)(\omega) = T(\omega)$. □

Définition 37.65 ([489]).

Nous disons que les variables aléatoires réelles X_n convergent **en probabilité** vers la variable aléatoire X si pour tout $\eta > 0$, on a

$$P(|X_n - X| \geq \eta) \rightarrow 0, \quad (37.228)$$

et on note

$$X_n \xrightarrow{P} X. \quad (37.229)$$

Définition 37.66 (Convergence en loi).

Nous disons que X_n converge vers X **en loi** vers la variable aléatoire X et nous notons

$$X_n \xrightarrow{\mathcal{L}} X \quad (37.230)$$

si pour toute fonction continue et bornée g nous avons

$$E(g(X_n)) \rightarrow E(g(X)) = \int g dP_X. \quad (37.231)$$

Proposition 37.67.

Deux autres caractérisations de la convergence en loi.

(1) Nous avons $X_n \xrightarrow{\mathcal{L}} X$ si et seulement si

$$\Phi_{X_n}(v) \rightarrow \Phi_X(v) \quad (37.232)$$

pour tout $v \in \mathbb{R}^d$. Ici Φ_X est la fonction caractéristique de X .

(2) Dans la définition de la convergence en loi nous pouvons indifféremment utiliser les fonctions continues et bornées, les fonctions continues à support compact ou les fonctions bornées uniformément continues.

Proposition 37.68.

Les types de convergence sont reliées par les implications suivantes :

$$\text{presque sure} \Rightarrow \text{en probabilité} \Rightarrow \text{en loi.} \quad (37.233)$$

La convergence en loi n'implique pas la convergence en probabilité, et par conséquent pas non plus la convergence presque certaine.

Dans le cas particulier $d = 1$ nous avons quelques critères supplémentaires.

Proposition 37.69.

Supposons que les variables aléatoires X_n soient réelles, et notons F_n la fonction de répartition de X_n . Si $F_n(x) \rightarrow F(x)$ pour tout x dans l'ensemble des points de continuité de F , alors $X_n \xrightarrow{\mathcal{L}} X$.

Proposition 37.70.

Si les X_n sont des variables aléatoires réelles positives, et si X est une variable aléatoire positive, alors $X_n \xrightarrow{\mathcal{L}} X$ si les transformées de Laplace des fonctions de répartition convergent ponctuellement, c'est-à-dire si

$$E(e^{-\alpha X_n}) \rightarrow E(e^{-\alpha X}) \quad (37.234)$$

pour tout $\alpha \geq 0$.

Proposition 37.71.

Si les X_n et X sont des variables aléatoires réelles discrètes à valeurs dans $\{x_0, x_1, \dots\}$ alors $X_n \xrightarrow{\mathcal{L}} X$ si et seulement si

$$P(X_n = x_k) \rightarrow P(X = x_k) \quad (37.235)$$

pour tout $k \in \mathbb{N}$.

Proposition 37.72 ([481]).

Soient X_n et X des variables aléatoires réelles. Nous avons

$$X_n \xrightarrow{\mathcal{L}} X \quad (37.236)$$

si et seulement si pour tout t où F_X est continue,

$$\lim_{n \rightarrow \infty} F_{X_n}(t) = F_X(t). \quad (37.237)$$

Proposition 37.73 ([481]).

Soit X_n une suite de variables aléatoires $\Omega \rightarrow \mathbb{R}^d$ et $a \in \mathbb{R}^d$. Si $X_n \xrightarrow{\mathcal{L}} a$, alors

$$X_n \xrightarrow{P} a. \quad (37.238)$$

Démonstration. Quitte à passer aux composantes, nous pouvons supposer que $d = 1$. Soit $\eta > 0$; nous savons que l'inégalité $|x| > a$ a pour solution $x > a$ ou $x < -a$. Dans notre cas,

$$P(|X_n - a| > \eta) = P(X_n - a > \eta) + P(X_n - a < -\eta) \quad (37.239a)$$

$$= P(X_n > \eta + a) + P(X_n < a - \eta) \quad (37.239b)$$

$$= 1 - P(X_n \leq \eta + a) + P(X_n \leq a - \eta) - P(X_n = a - \eta) \quad (37.239c)$$

$$\leq 1 - F_{X_n}(\eta + a) + F_{X_n}(a - \eta) \quad (37.239d)$$

où la majoration est l'oubli du terme $P(X_n = a - \eta)$, lequel est positif ou nul et F_{X_n} est la fonction de répartition de X_n , définition 37.53. Nous allons utiliser la proposition 37.72. La fonction de répartition de la variable aléatoire constante $X = a$ est donnée par

$$F_a(t) = P(a \leq t) = \mathbb{1}_{[0, \infty[}(t - a). \quad (37.240)$$

Par conséquent, la convergence en loi $X_n \xrightarrow{\mathcal{L}} a$ nous montre que

$$F_{X_n}(t) \rightarrow \mathbb{1}_{[0, \infty[}(t - a) \quad (37.241)$$

pour tout $t \neq a$ parce que $t = 0$ est un point de discontinuité de $\mathbb{1}_{[0, \infty[}$. Nous avons par conséquent

$$P(|X_n - a| > \eta) = 1 - \mathbb{1}_{[0, \infty[}(\eta) + \mathbb{1}_{[0, \infty[}(-\eta) = 1 - 1 + 0 = 0 \quad (37.242)$$

parce que $\eta > 0$. □

Le lemme de Slutsky sera utilisé en combinaison avec la proposition 37.75 pour calculer des intervalles de confiance, voir par exemple ce qui se passe autour de l'équation (38.126).

Lemme 37.74 (Slutsky[490]).

Soient X_n et Y_n des suites de variables aléatoires réelles telles que

$$\begin{aligned} X_n &\xrightarrow{\mathcal{L}} X \\ Y_n &\xrightarrow{P} a \in \mathbb{R}. \end{aligned} \quad (37.243)$$

Alors $(X_n, Y_n) \xrightarrow{\mathcal{L}} (X, a)$.

Démonstration. Étant donné que $Y_n \xrightarrow{\mathcal{L}} a$, nous avons $Y_n \xrightarrow{P} a$ par la proposition 37.73. Soit une fonction $f: \mathbb{R}^2 \rightarrow \mathbb{R}^2$; nous devons prouver que

$$E(f(X_n, Y_n)) \rightarrow E(f(X, a)). \quad (37.244)$$

Soit $\epsilon > 0$. Nous avons

$$E(\|f(X_n, Y_n) - f(X, a)\|) \leq E(\|f(X_n, Y_n) - f(X_n, a)\|) + E(\|f(X_n, a) - f(X, a)\|). \quad (37.245)$$

La fonction $g(t) = f(t, a)$ étant continue et bornée, la convergence en loi $X_n \xrightarrow{\mathcal{L}} X$ donne

$$E(\|f(X_n, a) - f(X, a)\|) \rightarrow 0. \quad (37.246)$$

Étudions à présent le premier terme du membre de droite de (37.245). Pour tout $\eta > 0$ et toute variables aléatoires Z et Z' nous avons

$$E(Z) = E(Z \mathbb{1}_{|Z'| < \eta}) + E(Z \mathbb{1}_{|Z'| \geq \eta}). \quad (37.247)$$

Nous décomposons donc le premier terme de (37.245) en

$$\begin{aligned} E(\|f(X_n, Y_n) - f(X_n, a)\|) &= E(\|f(X_n, Y_n) - f(X_n, a)\| \mathbb{1}_{|Y_n - a| < \eta}) \\ &\quad + E(\|f(X_n, Y_n) - f(X_n, a)\| \mathbb{1}_{|Y_n - a| \geq \eta}). \end{aligned} \quad (37.248)$$

Choisissons maintenant une valeur de η telle que

$$|(x, y) - (x', y')| < \eta \Rightarrow |f(x, y) - f(x', y')| \leq \epsilon. \quad (37.249)$$

Un tel η existe par l'uniforme continuité de f . Dans le premier terme, $|Y_n - a| < \eta$, par conséquent

$$\|(X_n, Y_n) - (X_n, a)\| = |Y_n - a| < \eta \quad (37.250)$$

et donc

$$\|f(X_n, Y_n) - f(X_n, a)\| \leq \epsilon. \quad (37.251)$$

Le premier terme devient donc

$$E(\|f(X_n, Y_n) - f(X_n, a)\| \mathbb{1}_{|Y_n - a| < \eta}) \leq \epsilon E(\mathbb{1}_{|Y_n - a| < \eta}) \leq \epsilon \quad (37.252)$$

parce que $E(\mathbb{1}_A) = P(A) \leq 1$. Pour le second terme de (37.248) nous effectuons la majoration

$$\|f(X_n, Y_n) - f(X_n, a)\| \leq 2\|f\|_\infty \quad (37.253)$$

tandis que la convergence $Y_n \xrightarrow{P} a$ entraîne

$$P(|Y_n - a| \geq \eta). \quad (37.254)$$

□

Proposition 37.75 ([491]).

Soient X_i des variables aléatoires telles que

$$X_i \xrightarrow{\mathcal{L}} X \quad (37.255)$$

et h , une fonction mesurable sur l'espace d'arrivée de X_i . Soit C l'ensemble des points de continuité de h au sens

$$C = \{\omega \in \Omega \text{ tel que } h \text{ est continue en } X_i(\omega)\}. \quad (37.256)$$

Alors si $P(X \in C) = 1$, nous avons

$$h(X_i) \xrightarrow{\mathcal{L}} h(X). \quad (37.257)$$

Une conséquence de cette proposition couplée au lemme de Slutsky est le résultats suivant, qui est donné sous le nom de **théorème de Slutsky** sur wikipédia.

Corollaire 37.76.

En reprenant les notations du lemme de Slutsky, si

$$X_n \xrightarrow{\mathcal{L}} X \quad (37.258a)$$

$$Y_n \xrightarrow{P} a, \quad (37.258b)$$

alors

$$X_n + Y_n \xrightarrow{\mathcal{L}} X + a \quad (37.259a)$$

$$X_n Y_n \xrightarrow{\mathcal{L}} aX \quad (37.259b)$$

$$Y_n^{-1} X_n \xrightarrow{\mathcal{L}} a^{-1} X \quad (37.259c)$$

$$(37.259d)$$

pourvu que a soit inversible.

Lemme 37.77 (Borel-Cantelli).

Soit (A_n) une suite d'événements (avec $A_n \in \mathcal{A}$ pour tout n).

(1) Si $\sum_{n=0}^{\infty} P(A_n)$ converge, alors

$$P(A_n \text{ i.s.}) = 0. \quad (37.260)$$

(2) Si la somme $\sum_n P(A_n)$ diverge, et si de plus les A_i sont indépendants, alors

$$P(A_n \text{ i.s.}) = 1. \quad (37.261)$$

La notation $P(A_n \text{ i.s.})$ signifie « infiniment souvent », c'est-à-dire

$$P(A_n \text{ i.s.}) = P\left(\bigcap_{N \in \mathbb{N}} \bigcup_{k \geq N} A_k\right) = P(\limsup A_n) \quad (37.262)$$

Une façon de paraphraser le lemme de Borel-Cantelli est que nous avons l'alternative

$$P(\limsup A_n) = \begin{cases} 0 & \text{si } \sum_{n \geq 0} P(A_n) < \infty \\ 1 & \text{sinon.} \end{cases} \quad (37.263)$$

Proposition 37.78.

Soit X_n , une suite de variables aléatoires et X une variable aléatoires. Si

$$\sum_n P(\|X_n - X\| > \eta) < \infty \quad (37.264)$$

pour tout ϵ , alors $X_n \xrightarrow{p.s.} X$.

Démonstration. Fixons ϵ et considérons les événements $A_n = \|X_n - X\| > \epsilon$. L'hypothèse dit que

$$\sum_n P(A_{n,\epsilon}) < \infty \quad (37.265)$$

et le lemme de Borel-Cantelli implique que

$$P(\limsup \|X_n - X\| > \epsilon) = 0. \quad (37.266)$$

Un élément ω est dans $\limsup A_n$ s'il est contenu dans tous les A_n , par conséquent, pour chaque ϵ nous avons l'inclusion

$$\{\omega \in \Omega \text{ tel que } X_n(\omega) \rightarrow X(\omega)\} \subset \complement \limsup A_n. \quad (37.267)$$

Nous pouvons aller plus loin et écrire

$$\{\omega \in \Omega \text{ tel que } X_n(\omega) \rightarrow X(\omega)\} = \complement \{\omega \in \Omega \text{ tel que } \|X_n - X\| > \epsilon, \forall \epsilon > 0\}. \quad (37.268)$$

Or la probabilité de l'ensemble

$$\{\omega \in \Omega \text{ tel que } \|X_n - X\| > \epsilon\} \quad (37.269)$$

est 0 pour chaque ϵ , et par conséquent la probabilité du membre de droite de (37.268) est 1. \square

Exemple 37.79

Considérons une suite de 0 et de 1 dans laquelle le 1 arrive avec une probabilité p et le 0 avec une probabilité $1 - p$. Une telle suite est modélisée par une suite de variables aléatoires de Bernoulli $(X_n)_{n \in \mathbb{N}}$ indépendantes de paramètre p .

Question : une telle suite contient elle une infinité de 1 ? Considérons les événements indépendants $A_n = \{X_n = 1\}$. Nous avons

$$\sum_n P(A_n) = \sum_n P(X_n = 1) = \sum_n p = \infty. \quad (37.270)$$

Par Borel-Cantelli et son expression (37.263), nous avons alors

$$P(\limsup A_n) = 1. \quad (37.271)$$

Donc une infinité d'événements A_n se produisent, et nous avons bien une infinité de 1 dans la suite.

Remarque : dans ce raisonnement nous pouvons considérer une probabilité non constante p_n tant que la série $\sum_n p_n$ diverge. \triangle

Exemple 37.80

À propos de maximum. La fonction $h: \mathbb{R}^d \rightarrow \mathbb{R}$ donnée par $h(x_1, \dots, x_d) = \max_i \{x_i\}$ est une fonction continue. Nous voudrions prouver que si on a une famille (finie en $i = 1, \dots, l$) de suites (en n) variables aléatoires $X_n^{(i)} \xrightarrow{p.s.} a$ convergeant toutes vers la même limite a , alors

$$M_n = \max_i \{X_n^{(i)}\} \xrightarrow{p.s.} a. \quad (37.272)$$

D'abord si nous avons l suite numériques $(x_n^{(i)})$, alors la suite

$$M_n = \max_i x_n^{(i)} \quad (37.273)$$

converge vers la même limite. En effet si $\epsilon > 0$ est donné, il suffit de prendre N_i l'entier tel que $|x_n^{(i)} - a| \leq \epsilon$ pour tout $n > N_i$. Et ensuite on prend $N > \max\{N_i\}$.

Si maintenant au lieu de suites numériques, nous avons des variables aléatoires, le résultat reste valable. Nous cherchons à prouver que

$$P\left(\{\omega \in \Omega \text{ tel que } \max\{X_n^{(i)}(\omega)\} \rightarrow a\}\right) = 1. \quad (37.274)$$

Par ce que nous venons de dire sur les suites numériques, un élément ω n'est pas dans cet ensemble seulement s'il y a un ∞ pour lequel $X_n^{(i)}$ ne converge pas vers a . Or cela, pour chaque i est un événement de probabilité zéro.

Les ω qui ne fonctionneront pas dans l'équation (37.274) sont ceux de la réunion d'un ensemble fini d'ensembles de probabilité nulle. C'est donc de probabilité nulle. \triangle

37.4 Loi des grands nombres, théorème central limite

37.4.1 Loi des grands nombres

Lemme 37.81 (Inégalité de Markov[492]).

Soit une variable aléatoire $X \in L^p$ et $\epsilon > 0$. Nous avons

$$P(|X| \geq \epsilon) \leq \frac{1}{\epsilon^r} E(|X|^r). \quad (37.275)$$

Démonstration. Nous avons

$$E(|X|^r) \geq \int_{|X| \geq \epsilon} |X|^r dP(\omega) \geq \epsilon^r \int_{|X| \geq \epsilon} dP = \epsilon^r P(|X| \geq \epsilon). \quad (37.276)$$

□

Corollaire 37.82 ([492]).

Soit ϕ , une fonction croissante et positive ou nulle sur l'intervalle I . Soit aussi une variable aléatoire $Y: \Omega \rightarrow \mathbb{R}$ telle que $P(Y \in I) = 1$. Alors pour tout $b \in I$ tel que $\phi(b) > 0$ nous avons

$$P(Y \geq b) \leq \frac{E[\phi(Y)]}{\phi(b)}. \quad (37.277)$$

Théorème 37.83 (Loi forte des grands nombres).

Soit (X_n) une suite de variables aléatoires réelles

- (1) indépendantes et identiquement distribuées,
- (2) intégrables (c'est-à-dire dans L^1),

alors

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p.s.} E(X_1). \quad (37.278)$$

Note : étant donné que les variables aléatoires sont identiquement distribuées, nous avons évidemment $E(X_1) = E(X_2) = \dots$

Problèmes et choses à faire

Est-ce que les variables aléatoires doivent vraiment être réelles ?

Corollaire 37.84.

Si les variables aléatoires réelles X_n sont

- (1) indépendantes et identiquement distribuées,
- (2) dans L^2

alors

$$\bar{X}_n \xrightarrow{P} E(X_1). \quad (37.279)$$

Démonstration. Nous voulons prouver que pour tout $\eta > 0$,

$$P(|\bar{X}_n - E(X_1)| > \eta) \rightarrow 0. \quad (37.280)$$

Remarquons d'abord que les variables aléatoires X_n étant identiquement distribuées, $E(\bar{X}_n) = E(X_1)$ parce que $E(X_i) = E(X_1)$ pour tout i . L'inégalité de Markov avec $r = 2$ nous donne

$$P(|\bar{X}_n - E(\bar{X}_n)| > \eta) \leq \frac{1}{\eta^2} E(|\bar{X}_n - E(\bar{X}_n)|^2) \quad (37.281a)$$

$$= \frac{1}{\eta^2} \text{Var}(\bar{X}_n) \quad (37.281b)$$

$$= \frac{1}{n\eta^2} \text{Var}(X_1) \quad (37.281c)$$

où nous avons utilisé la proposition 37.27 : $\text{Var}(\bar{X}_n) = \text{Var}(X_1)/n$. Au final nous avons prouvé que

$$P(|\bar{X}_n - E(\bar{X}_n)| > \frac{1}{n\eta^2} \text{Var}(X_1), \quad (37.282)$$

qui tend vers zéro lorsque $n \rightarrow \infty$. □

Proposition 37.85.

Soient X_n des variables aléatoires indépendantes et identiquement distribuées avec $X_n \geq 0$. Nous acceptons $E(X_1) = \infty$, c'est-à-dire que nous relaxons la condition $X_n \in L^1$ par rapport à la loi des grands nombres.

Alors

$$\bar{X}_n \xrightarrow{p.s.} E(X_1) \in [0, \infty]. \quad (37.283)$$

Démonstration. Si $E(X_1) < \infty$, nous sommes dans le cas de la loi des grands nombres. Pour chaque $N \in \mathbb{N}$ nous considérons la suite de variables aléatoires

$$X_n^{(N)} = \min(X_n, N). \quad (37.284)$$

Nous avons évidemment $\bar{X}_n^{(N)} \leq \bar{X}_n$. Les variables aléatoires $X_n^{(N)}$ étant bornées par N , elles vérifient la loi des grands nombres pour chaque N séparément. Par conséquent nous avons pour chaque N la limite

$$\bar{X}_n^{(N)} \rightarrow E(X_1^{(N)}) \quad (37.285)$$

Nous supposons que $E(X_1) = \infty$, par conséquent $\lim_{N \rightarrow \infty} E(X_1^{(N)}) = \infty$. Soit $\eta > 0$ et choisissons N de telle manière à avoir

$$E(X_1^{(N)}) > \eta + 1. \quad (37.286)$$

La limite (37.285) nous permet de trouver n_0 tel que pour tout $n > n_0$ nous ayons $\bar{X}_n^{(N)} > \eta$. Au final,

$$\eta < \bar{X}_n^{(N)} \leq \bar{X}_n, \quad (37.287)$$

ce qui montre que $\bar{X}_n \rightarrow \infty$. □

Exemple 37.86

La loi des grands nombres justifie la pratique courante d'approximer une grandeur physique par la moyenne empirique d'un grand nombre de mesures. \triangle

Exemple 37.87

Citons ici le dernier paragraphe de *Le mystère de Marie Roget* par Edgar Allan Poe, traduit par Charles Baudelaire⁸.

Rien, par exemple, n'est plus difficile que de convaincre le lecteur non spécialiste que, si un joueur de dés a amené les six deux fois coup sur coup, ce fait est une raison suffisante de parier gros que le troisième coup ne ramènera pas les six. Une opinion de ce genre est généralement rejetée tout d'abord par l'intelligence. On ne comprend pas comment les deux coups déjà joués, et qui sont maintenant complètement enfouis dans le Passé, peuvent avoir de l'influence sur le coup qui n'existe que dans le Futur. La chance pour amener les six semble être précisément ce qu'elle était à n'importe quel moment, c'est-à-dire soumise seulement à l'influence de tous les coups divers que peuvent amener les dés. Et c'est là une réflexion qui semble si parfaitement évidente, que tout effort pour la controverser est plus souvent accueilli par un sourire moqueur que par une condescendance attentive.

Dans le cours de la nouvelle, Edgar Poe cite et utilise la théorie des probabilités avec une justesse inaccoutumée dans la littérature. Mais dans ce paragraphe final, Poe montre de façon la plus formelle qu'il n'a *rien* compris à la loi des grands nombres. \triangle

37.4.2 Théorème central limite**Lemme 37.88.**

Soit $z_n \rightarrow z$ une suite convergente dans \mathbb{C} . Alors

$$\left(1 + \frac{z_n}{n}\right)^n \rightarrow e^z. \quad (37.288)$$

Théorème 37.89.

Si les variables aléatoires X_n sont

- (1) indépendantes et identiquement distribuées de loi parente X ,
- (2) $X_1 \in L^2(\Omega, \mathcal{A})$,

alors nous notons $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$, $m = E(X_1)$ et $\sigma^2 = \text{Var}(X_1)$ et nous avons

$$\frac{\bar{X}_n - E(X)}{\sqrt{\text{Var}(\bar{X}_n)}} = \frac{\bar{X}_n - m}{\sigma/\sqrt{n}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1). \quad (37.289)$$

Démonstration. Nous allons écrire la démonstration dans le cas de variables aléatoires réelles. La proposition 37.67 dit que la suite X_n converge en loi vers X si et seulement si les fonctions caractéristiques convergent ponctuellement. Nous devons donc prouver, pour chaque⁹ $t \in \mathbb{R}$, que

$$\Phi_{\frac{S_n - nm}{\sigma\sqrt{n}}}(t) \rightarrow \Phi_{\mathcal{N}(0,1)}(t). \quad (37.290)$$

8. Disponible sur https://fr.wikisource.org/wiki/Le_Mystère_de_Marie_Roget

9. Chuck Norris peut *vraiment* le faire pour *chaque* $t \in \mathbb{R}$

Supposons dans un premier temps que $E(X_i) = 0$ et $\sigma(X_i) = 1$. Dans ce cas nous considérons la fonction

$$\Phi_{\frac{S_n}{\sqrt{n}}} = E\left(e^{i\frac{t}{\sqrt{n}}\sum_{k=1}^n X_k}\right) \quad (37.291a)$$

$$= \prod_{k=1}^n E\left(e^{i\frac{t}{\sqrt{n}}X_k}\right) \quad (37.291b)$$

$$= \prod_{k=1}^n \Phi_{X_1}\left(\frac{t}{\sqrt{n}}\right) \quad (37.291c)$$

$$= \Phi_{X_1}\left(\frac{t}{\sqrt{n}}\right)^n. \quad (37.291d)$$

Cette quantité est a priori complexe; nous ne pouvons donc pas immédiatement passer au logarithme. Nous pouvons par contre utiliser un développement en puissances de t en nous servant de la proposition 37.56 et de l'hypothèse comme quoi $X_1 \in L^2$:

$$\Phi_{X_1}(t) = \Phi_{X_1}(0) + \Phi'_{X_1}(0)t + \Phi''_{X_1}(0)\frac{t^2}{2} + \alpha(t)t^2 \quad (37.292)$$

où α est une fonction qui a la propriété $\lim_{x \rightarrow 0} \alpha(x) = 0$.

En utilisant les hypothèses et la formule de dérivation de la fonction caractéristique,

$$\Phi_{X_1}(0) = 1 \quad (37.293a)$$

$$\Phi'_{X_1}(0) = E(iX) = 0 \quad (37.293b)$$

$$\Phi''_{X_1}(0) = E(-X^2) = -\text{Var}(X_1) = -1. \quad (37.293c)$$

Nous avons donc

$$\Phi_{X_1}\left(\frac{t}{\sqrt{n}}\right) = \underbrace{1 - \frac{1}{2}\frac{t^2}{n}}_{\in \mathbb{R}} + \underbrace{\frac{t^2}{n}\alpha\left(\frac{t}{\sqrt{n}}\right)}_{\in \mathbb{C}}, \quad (37.294)$$

de telle sorte que, en considérant une valeur fixée de t ,

$$\Phi_{\frac{S_n}{\sqrt{n}}}(t) = \left(1 - \frac{\frac{t^2}{2} + \beta_n}{n}\right)^n \quad (37.295)$$

où $\beta_n = t^2\alpha(t/\sqrt{n})$. Nous avons bien entendu $\lim_{n \rightarrow \infty} \beta_n = 0$.

Nous pouvons appliquer le lemme 37.88 pour obtenir la limite

$$\lim_{n \rightarrow \infty} \Phi_{\frac{S_n}{\sqrt{n}}}(t) = e^{-t^2/2}. \quad (37.296)$$

La convergence (37.290) est par conséquent prouvée dans le cas où $E(X_i) = 0$ et $\text{Var}(X_i) = 1$.

Considérons maintenant des variables aléatoires avec $E(X_i) = m$ et $\text{Var}(X_i) = \sigma^2$. Elles peuvent être écrites sous la forme

$$X_i = \sigma X'_i + m \quad (37.297)$$

où X'_i est d'espérance nulle et de variance un. Nous avons alors

$$S_n = \sigma \sum_{i=1}^n X'_i + nm, \quad (37.298)$$

et

$$\frac{S_n - nm}{\sigma\sqrt{n}} = \frac{S'_n}{\sqrt{n}} \quad (37.299)$$

où $S'_n = \sum_i X'_i$. L'étude de la variable aléatoire

$$\frac{S_n - nm}{\sigma\sqrt{n}} \quad (37.300)$$

revient donc à celle de S'_n/\sqrt{n} qui vient d'être effectuée. \square

37.90.

À propos du théorème central limite [37.89](#). Si pour une certaine variable aléatoire X on a $E(X) = m$, alors nous n'avons pas forcément $P(X = m + a) = P(X = m - a)$. Est-ce que le théorème central limite permet cependant d'affirmer que dans un certaine mesure nous avons

$$P(\bar{X}_n = m + a) = P(\bar{X}_n = m - a) \quad (37.301)$$

lorsque n est grand ?

Tel quelle, l'équation [\(37.301\)](#) est en général fautive pour chaque n parce qu'il existe des distributions non symétriques. Mais bien entendu les deux membres tendent vers zéro pour $n \rightarrow \infty$. Mais cela n'est pas lié à la symétrie de la distribution gaussienne. C'est seulement le fait que la gaussienne n'a pas de masses ponctuelles.

Par contre, il y a effectivement une assurance de symétrie pour \bar{X}_n lorsque $n \rightarrow \infty$. Le fait est que $X_n \xrightarrow{\mathcal{L}} X$ où X est une gaussienne. La fonction de répartition de X est continue partout et la proposition [37.72](#) nous dit que pour tout $x \in \mathbb{R}$,

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x). \quad (37.302)$$

Vu que le nombre $P(\bar{X}_n \in B(m + a, \delta))$ peut être exprimé avec des sommes et différences F_{X_n} , nous avons

$$\lim_{n \rightarrow \infty} P(\bar{X}_n \in B(m + a, \delta)) = P(X \in B(m + a, \delta)). \quad (37.303)$$

Par symétrie de la gaussienne le membre de droite est égal à $P(X \in B(m - a, \delta))$ et nous avons bien

$$\lim_{n \rightarrow \infty} P(\bar{X}_n \in B(m + a, \delta)) = \lim_{n \rightarrow \infty} P(\bar{X}_n \in B(m - a, \delta)). \quad (37.304)$$

Remarque 37.91.

Le théorème central limite s'applique quelle que soit la distribution des variables aléatoires X_i (dans les limites de hypothèses); en particulier il ne dit rien sur la moyenne des X_i . Il dit seulement que l'écart de la moyenne « mesurée » à la moyenne « théorique » est une variable aléatoire gaussienne si on a mesuré assez de fois.

Autrement dit, si la durée d'attente à la poste est de 5 minutes, et si j'y vais 2000 fois, alors la probabilité que ma moyenne d'attente soit de 4 minutes est la même que la probabilité qu'elle soit de 6 minutes¹⁰.

Problèmes et choses à faire

La remarque [37.91](#) est une interprétation personnelle. J'aimerais avoir l'avis de quelqu'un de plus compétent.

Remarque 37.92.

Nous pouvons obtenir la limite [\(37.296\)](#) d'une façon alternative. Nous considérons la détermination du logarithme complexe sur $\mathbb{C} \setminus \mathbb{R}^-$; cela est une fonction analytique (théorème [27.59](#)) vérifiant l'équation

$$e^{\ln(z)} = z \quad (37.305)$$

pour tout $z \in \mathbb{C} \setminus \mathbb{R}^-$ et le développement

$$\ln(1 + z) = \sum_{n=1}^{\infty} (-1)^{n+1} \frac{z^n}{n}. \quad (37.306)$$

En particulier, $\ln(1 + z) = z + z\alpha(z)$ où $\lim_{z \rightarrow 0} \alpha(z) = 0$. Nous reprenons à l'équation [\(37.295\)](#) en

¹⁰. Et en l'occurrence, cette probabilité est nulle parce qu'on est en train de parler de variable aléatoire continue, mais vous voyez l'idée.

fixant t . Nous avons

$$\Phi_{\frac{S_n}{\sqrt{n}}}(t) = \exp \left[\ln \left(\Phi_{\frac{S_n}{\sqrt{n}}}(t) \right) \right] \quad (37.307a)$$

$$= \exp \left[n \ln \left(1 + \frac{-\frac{t^2}{2} - \beta_n}{n} \right) \right] \quad (37.307b)$$

$$= \exp \left[-\frac{t^2}{2} - \beta_n + \left(-\frac{t^2}{2} - \beta_n \right) \alpha \left(\frac{-\frac{t^2}{2} - \beta_n}{n} \right) \right]. \quad (37.307c)$$

À la limite $n \rightarrow \infty$ nous tombons sur $e^{-t^2/2}$.

Remarque 37.93.

Étant donné que la variable aléatoire

$$\frac{S_n - nm}{\sigma\sqrt{n}} \quad (37.308)$$

converge en loi vers $\mathcal{N}(0, 1)$, nous avons la convergence des fonctions de répartition partout où la fonction de répartition de la normale est continue¹¹ (donc sur tout \mathbb{R}). En particulier,

$$\left| P \left(\frac{S_n - nm}{\sigma\sqrt{n}} \leq x \right) - \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy \right| \rightarrow 0. \quad (37.309)$$

Nous avons la borne de **Berry-Esséen** qui donne une estimation de la vitesse de convergence : si $X \in L^3$, alors il existe une constante C , indépendante de x , des X_i et de n telle que

$$\left| P \left(\frac{S_n - nm}{\sigma\sqrt{n}} \leq x \right) - \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy \right| \leq \frac{X\mu_3}{\sigma^3\sqrt{n}} \quad (37.310)$$

où $\mu_3 = E(|X_1 - m|^3)$ est le moment d'ordre 3 de X . La chose à retenir est que la convergence est à la vitesse de $1/\sqrt{n}$.

En dimension $d > 1$, nous avons encore un théorème central limite.

Théorème 37.94.

Si $d > 1$, et si nous avons des variables aléatoires X_n à valeurs dans \mathbb{R}^d avec

- (1) les X_n sont indépendantes et identiquement distribuées
- (2) les X_n sont dans L^2 .

Alors nous notons $X_1 = (X_1^{(1)}, \dots, X_1^{(d)})$. Nous avons

$$\frac{S_n - nm}{\sqrt{n}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, \Sigma) \quad (37.311)$$

où Σ est ma matrice de covariance du vecteur aléatoire X_1 :

$$\Sigma = \left(\text{Cov}(X_1^{(1)}, \dots, X_1^{(d)}) \right)_{i,j=1, \dots, d}. \quad (37.312)$$

37.4.3 Marche aléatoire

Nous considérons un mobile qui se déplace sur l'axe \mathbb{Z} . À chaque pas de temps, nous supposons qu'il va faire un pas à gauche avec une probabilité p et un pas à droite avec une probabilité $(1 - p)$. Nous nous demandons quel est le mouvement du mobile sur le long terme.

La position S_n du mobile à l'instant n est donnée par

$$S_n = \sum_{i=1}^n X_i \quad (37.313)$$

11. Proposition 37.72.

où X_i est le pas effectué à l'instant i . Ce sont des variables de Bernoulli indépendantes avec

$$X_i \stackrel{\mathcal{L}}{=} p\delta_{-1} + (1-p)\delta_1 \quad (37.314)$$

c'est-à-dire

$$P(X_i = -1) = p \quad (37.315a)$$

$$P(X_i = 1) = 1 - p. \quad (37.315b)$$

Ces variables vérifient les hypothèses de la loi des grands nombres :

- (1) elles sont indépendantes et identiquement distribuées,
- (2) elles sont intégrables.

Pour le second point, le calcul est

$$\int_{\Omega} |X_i| dP = \int_{\mathbb{R}} |x| dP_X = \int_{\mathbb{R}} |x|(p\delta_{-1} + (1-p)\delta_1) = |1-p| + |p| = 1. \quad (37.316)$$

Nous avons par conséquent

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p.s.} E(X_1) = (1-2p) \quad (37.317)$$

et

$$\frac{S_n}{n} \rightarrow (1-2p). \quad (37.318)$$

Si $p \neq 1/2$ nous pouvons conclure que

- (1) si $p > 1/2$, alors $S_n \xrightarrow{p.s.} -\infty$
- (2) si $p < 1/2$, alors $S_n \xrightarrow{p.s.} \infty$.

De plus nous connaissons la vitesse de divergence : elle est linéaire. Le mobile suit essentiellement l'équation $p(n) = (1-2p)n$.

Remarque 37.95.

Cela ne traite pas le cas $p = 1/2$. Dans ce cas, nous pouvons simplement dire que $S_n = o(n)$.

37.5 Les lois usuelles

37.5.1 Loi de Bernoulli

Une expérience de Bernoulli consiste à tirer au hasard un 0 ou un 1 avec une probabilité p de tomber sur 1 et $1-p$ de tomber sur zéro. Il s'agit donc d'une expérience qui réussit ou qui rate.

Le cas typique est une urne avec des boules indiscernables blanches ou noires. La probabilité p est la proportion de blanches dans l'urne (avec remise entre les tirages). Dans ce cas, nous avons l'espace de probabilité (Ω, \mathcal{A}, P) où Ω représente l'ensemble des boules, \mathcal{A} est l'ensemble des parties de Ω et P est l'équiprobabilité sur Ω . Une variable aléatoire est une application

$$X: \Omega \rightarrow \{0, 1\} \quad (37.319)$$

$\omega \mapsto$ couleur de la boule ω .

Nous notons $\mathcal{B}(1, p)$ la loi de Bernoulli. Elle a une expression très simple :

$$\mathcal{B}(0, 1)(\{1\}) = p \quad (37.320a)$$

$$\mathcal{B}(0, 1)(\{0\}) = 1 - p \quad (37.320b)$$

Une variable aléatoire réelle est de **Bernoulli** de paramètre p ($0 < p < 1$) si

$$X: \Omega \rightarrow \mathbb{R} \quad (37.321)$$

avec $P(x = 1) = p$ et $P(X = 0) = 1 - p$. En tant que mesure sur \mathbb{R} , nous avons

$$P_X = p\delta_1 + (1 - p)\delta_0. \quad (37.322)$$

Une fonction h qui réalise le supremum de la formule (15.424) est par exemple une fonction en escalier qui vaut en x le plus petit entier plus grand ou égal à x . L'espérance d'une loi de Bernoulli est alors

$$E(x) = p. \quad (37.323)$$

Étant donné que la variable aléatoire X prend seulement les valeurs 0 et 1, nous avons pour tout ensemble mesurable B

$$P_{X^2}(B) = P(X^2 \in B) = P(X \in B), \quad (37.324)$$

et par conséquent $P_{X^2} = P_X$ et $E(X^2) = E(X)$. Nous trouvons donc la variance

$$\text{Var}(X) = E(X^2) - E(X)^2 = p - p^2 = p(1 - p). \quad (37.325)$$

37.5.2 Loi binomiale

Une expérience binomiale consiste à répéter n expériences de Bernoulli de paramètre p et de compter le nombre de réussites. Une telle expérience peut être réalisée selon la procédure suivante.

Soit une urne contenant N boules dont une proportion p de 1 et $1 - p$ de 0. Une expérience binomiale de paramètres n et p consistera à prendre n boules *avec remise* et à compter le nombre de 1 obtenus.

En termes d'espaces probabilisés, nous avons Ω qui est l'ensemble des tuples de taille n à valeurs dans $\{0, 1\}$, la tribu \mathcal{A} est l'ensemble des parties de Ω , et la probabilité P est l'équiprobabilité :

$$P(\omega) = \frac{1}{N^n} \quad (37.326)$$

si il y a N boules dans l'urne. Nous construisons alors la variable aléatoire

$$X(\omega) = \sum_{i=1}^n \omega_i \quad (37.327)$$

où ω est une suite de taille n de 0 et de 1.

Calculons $P(X = k)$. Il s'agit de considérer tous les sous-ensembles de taille n de Ω contenant exactement k fois 1. Il y a $\binom{n}{k}$ manière de décider lesquelles des n boules seront blanches. Ensuite, chaque boule blanche peut être choisie parmi les m boules disponibles, et chaque boule noire peut être choisie parmi les $(N - m)$ disponibles. Nous avons donc

$$P(X = k) = \binom{n}{k} \frac{m^k (N - m)^{n-k}}{N^n}. \quad (37.328)$$

En effet la mesure de probabilité sur Ω est la mesure de comptage renormalisée par le cardinal de Ω qui vaut N^n . Étant donné que $p = m/N$, nous transformons facilement (37.328) en

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}. \quad (37.329)$$

Une variable aléatoire de loi binomiale étant une somme de variables aléatoires de Bernoulli indépendantes, l'espérance¹² et la variance¹³ s'obtiennent en sommant les espérances et variances termes à terme :

$$E(X) = np \quad (37.330)$$

et

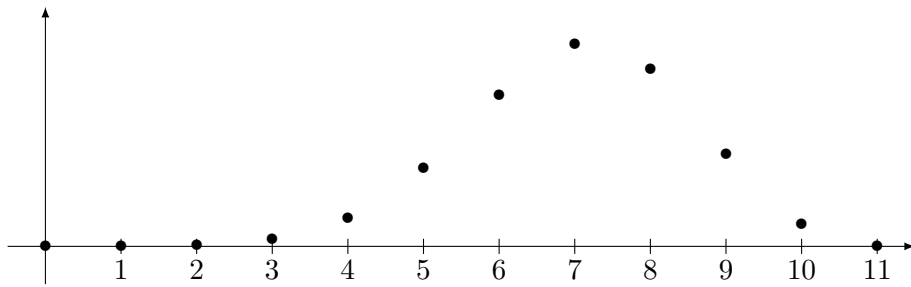
$$\text{Var}(X) = \sum_{i=1}^n \text{Var}(X_i) = np(1 - p) \quad (37.331)$$

en vertu de (37.325).

La loi binomiale lorsque $p = 0.7$ et $n = 10$.

12. Valable même sans indépendance, proposition 37.25

13. Lemme 37.23.



37.5.3 Loi multinomiale

La loi multinomiale $\mathcal{M}(n; k; p_1, \dots, p_k)$ consiste à effectuer n épreuves d'une démarche aléatoire qui peut avoir k issues différentes avec probabilités p_1, \dots, p_k . Les variables aléatoires multinomiales sont N_i avec les contraintes

$$\sum_{i=1}^k N_i = n \quad (37.332a)$$

$$\sum_{i=1}^k p_i = 1. \quad (37.332b)$$

La fonction de probabilité multinomiale est

$$P(N_1 = n_1, \dots, N_k = n_k) = \frac{n!}{n_1! \dots n_k!} p_1^{n_1} \dots p_k^{n_k}. \quad (37.333)$$

Chacune des N_i est une binomiale de probabilité p_i .

37.5.4 Loi géométrique

Soit (X_n) une suite indépendante et identiquement distribuée de lois de Bernoulli de paramètre p . Alors la variable aléatoire

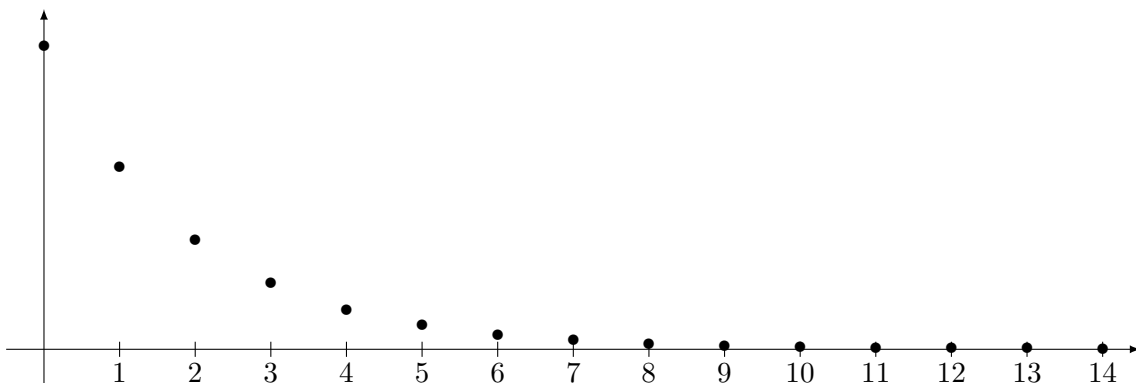
$$Z = \inf\{n \geq 1 \text{ tel que } X_n = 1\} \quad (37.334)$$

est une loi géométrique de paramètre p .

La loi géométrique compte donc le nombre d'expériences de Bernoulli à effectuer avant que le premier succès soit au rendez-vous. Nous avons

$$P(Z = k) = P(X_k = 1)P(X_1, \dots, X_{k-1} = 0) = p(1-p)^{k-1} \quad (37.335)$$

La loi géométrique de paramètre $p = 0.2$.



Note : si p est trop grand, on ne voit vite plus rien parce que la probabilité d'attendre longtemps est vite très faible.

37.5.5 Loi de Poisson

Une variable aléatoire Z suit une **loi de Poisson** de paramètre λ , notée $\mathcal{P}(\lambda)$ si

$$P(Z = k) = e^{-\lambda} \frac{\lambda^k}{k!} \quad (37.336)$$

pour tout $k \in \mathbb{N}$.

La **loi de Poisson** est une loi de probabilité discrète qui décrit le comportement du nombre d'évènements se produisant dans un laps de temps fixé, si ces évènements se produisent avec une fréquence moyenne connue et indépendamment du temps écoulé depuis l'évènement précédent.

Si un évènement se produit en moyenne p fois par seconde, la probabilité d'observer l'évènement k fois durant n secondes est donnée par $P(Z = k)$ où Z est une loi de Poisson de paramètre $\lambda = pn$.

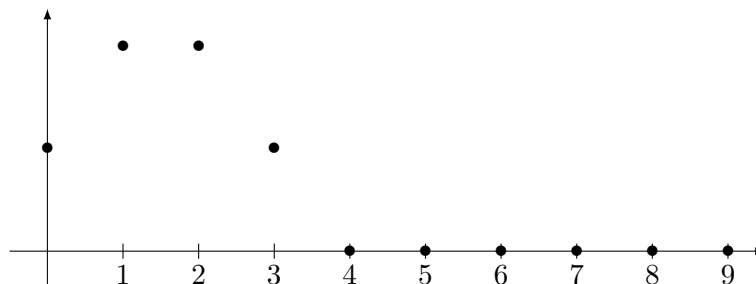
Théorème 37.96 ([wikipedia](#)).

L'espérance et la variance d'une variable aléatoire de Poisson sont λ :

$$E(X) = \lambda \quad (37.337a)$$

$$\text{Var}(X) = \lambda. \quad (37.337b)$$

La loi de Poisson de paramètre $\lambda = 2$.



37.5.6 Loi exponentielle

La loi de exponentielle représente le temps qu'il faut attendre pour qu'une particule se désintègre si elle a en permanence une probabilité¹⁴ λdt de se désintégrer entre t et $t + dt$. L'espérance est donc $1/\lambda$.

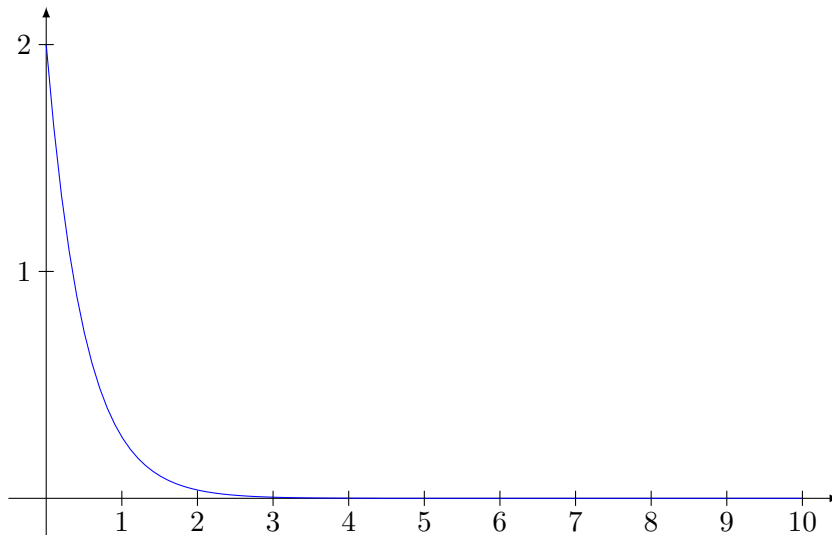
Il se passe donc en moyenne λ évènements par seconde. La proposition [37.103](#) nous montrera que le nombre d'évènements se produisant en une seconde suit une loi de Poisson de paramètre λ .

Plus formellement, la loi exponentielle de paramètre λ , notée $\mathcal{E}(\lambda)$ est la loi de densité

$$f_X : x \mapsto \begin{cases} \lambda e^{-\lambda x} & \text{si } x \geq 0 \\ 0 & \text{sinon.} \end{cases} \quad (37.338)$$

Densité de la loi exponentielle de paramètre $\lambda = 2$.

¹⁴. Étant donné que λ n'est pas limité à 1, en réalité ce n'est pas une probabilité. Je suis preneur d'une bonne interprétation physique de ce paramètre.

**Proposition 37.97.**

Si $X \sim \mathcal{E}(\lambda)$, alors la fonction de répartition de X est donnée par

$$F(x) = P(X \leq x) = \begin{cases} 1 - e^{-\lambda x} & \text{si } x \geq 0 \\ 0 & \text{sinon.} \end{cases} \quad (37.339)$$

La fonction caractéristique est donnée par

$$E(e^{itX}) = \frac{\lambda}{\lambda - it}. \quad (37.340)$$

L'espérance et la variance sont données par

$$\begin{aligned} E(X) &= \frac{1}{\lambda} \\ \text{Var}(X) &= \frac{1}{\lambda^2}. \end{aligned} \quad (37.341)$$

Démonstration. Pour la fonction caractéristique,

$$E(e^{itX}) = \int_{-\infty}^{\infty} e^{itx} \lambda e^{-\lambda x} \mathbb{1}_{[0, \infty[}(x) dx \quad (37.342a)$$

$$= \lambda \int_0^{\infty} e^{x(-\lambda + it)} dx \quad (37.342b)$$

$$= \lim_{A \rightarrow \infty} \left[\frac{e^{x(-\lambda + it)}}{-\lambda + it} \right]_{x=0}^{x=A} \quad (37.342c)$$

$$= \lim_{A \rightarrow \infty} \frac{e^{A(it - \lambda)}}{it - \lambda} - \frac{1}{it - \lambda}. \quad (37.342d)$$

Le premier terme est nul parce que si on prend la norme,

$$\left| \frac{e^{A(-\lambda + it)}}{-\lambda + it} \right| = \frac{e^{-\lambda A}}{|it - \lambda|} \rightarrow 0. \quad (37.343)$$

En ce qui concerne l'espérance nous faisons le calcul suivant :

$$E(X) = \int_{\mathbb{R}} x f_X(x) dx = \lambda \int_{\mathbb{R}^+} x e^{-\lambda x} dx = \frac{1}{\lambda}. \quad (37.344)$$

Pour la variance, nous utilisons la formule (37.58). Nous avons

$$E(X^2) = \int_{\mathbb{R}} x^2 f_X(x) dx = \int_0^{\infty} x^2 \lambda e^{-\lambda x} dx = \frac{2}{\lambda^2}. \quad (37.345)$$

Donc $\text{Var}(X) = E(X^2) - E(X)^2 = \frac{1}{\lambda^2}$. □

La loi exponentielle est une loi **sans mémoire** en ce sens que

$$P(X > x + y | X > y) = P(X > x). \quad (37.346)$$

En effet nous utilisons la règle de la probabilité conditionnelle

$$P(A|B) = \frac{P(A \cap B)}{P(B)}. \quad (37.347)$$

Ici,

$$P(X > x + y | X > y) = \frac{P(X > x + y)}{P(X > y)} = e^{-\lambda x}. \quad (37.348)$$

La proposition suivante montre que la loi exponentielle est à peu près la seule à être sans mémoire. D'où son importance dans l'étude des machines dont les pièces ne subissent pas d'usure.

Proposition 37.98.

Soit X , une variable aléatoire admettant une densité continue f par rapport à la mesure de Lebesgue. Si elle est sans mémoire, alors elle est exponentielle.

Démonstration. Nous posons $\varphi(x) = P(X \geq x)$. Cela est la fonction de répartition de X (à part que cette dernière est $1 - \varphi(x)$ mais c'est pas grave), et est donnée par

$$\varphi(x) = 1 - \int_0^x f(t) dt. \quad (37.349)$$

Cette dernière intégrale vérifie les hypothèses du théorème 37.98, de telle sorte que φ soit une fonction dérivable et $\varphi'(x) = f(x)$.

D'autre part en utilisant la définition de la probabilité conditionnelle la propriété de ne pas avoir de mémoire donne

$$\varphi(x + y) = \varphi(x)\varphi(y) \quad (37.350)$$

et de plus $\varphi(0) = 1$. Calculons la dérivée de φ :

$$\varphi'(x) = \lim_{\epsilon \rightarrow 0} \frac{\varphi(x + \epsilon) - \varphi(x)}{\epsilon} = \varphi(x) \lim_{\epsilon \rightarrow 0} \frac{\varphi(\epsilon) - 1}{\epsilon} = \varphi(x)\varphi'(0). \quad (37.351)$$

Donc φ vérifie l'équation différentielle de l'exponentielle. □

Exemple 37.99

Une machine a une durée de vie représentée par une variable aléatoire suivant une loi de Poisson de paramètre λ . Soit T_y la variable aléatoire qui représente le temps de vie restant sachant que la machine a déjà vécu un temps y . Nous voulons trouver la fonction de répartition de T_y . Nous avons

$$P(T_y > x) = P(X > x + y | X > y) = P(X > x) = e^{-\lambda x}. \quad (37.352)$$

Dans ce cas, la loi de T_y ne dépend pas de y . Cela signifie que la machine ne vieillit pas et surtout que le modèle n'est pas réaliste. △

Proposition 37.100.

Si $X \sim \mathcal{E}(\lambda)$ et $Y \sim \mathcal{E}(\mu)$ sont indépendantes, alors

(1) $P(X < Y) = \frac{\lambda}{\lambda + \mu}$

(2) $P(X > Y) = \frac{\mu}{\lambda + \mu}$

(3) $P(X = Y) = 0$

(4) $\min(X, Y) \sim \mathcal{E}(\lambda + \mu)$.

De plus les variables aléatoires exponentielles ont une propriété d'absence de mémoire :

$$P(X > t + s | X > s) = P(X > t) = e^{-\lambda t}. \quad (37.353)$$

Démonstration. Étant donné que X et Y sont indépendantes, la densité conjointe est le produit des densités (37.20). Nous avons donc

$$P(X > Y) = \int_D \lambda e^{-\lambda x} \mu e^{-\mu y} dx dy \quad (37.354)$$

où D est le domaine $D = \{(x, y) \in \mathbb{R}^2 \text{ tel que } x, y > 0, x > y\}$. Nous avons donc

$$P(X > Y) = \lambda \mu \int_0^\infty dx \int_0^x dy e^{-\lambda x} e^{-\mu y} = \frac{\mu}{\lambda + \mu}. \quad (37.355)$$

```
sage: var('a,b')
(a, b)
sage: f(x,y)=exp(-a*x)*exp(-b*y)
sage: assume(a>0)
sage: assume(b>0)
sage: a*b*f.integrate(y,0,x).integrate(x,0,oo)
(x, y) |--> a*b/(a^2 + a*b)
```

Pour trouver la loi de $\min(X, Y)$, nous écrivons

$$P(\min(X, Y) > t) = P(X > t, Y > t) \quad (37.356a)$$

$$= P(X > t)P(Y > t) \quad \text{par indépendance} \quad (37.356b)$$

$$= (1 - F_X(t))(1 - F_Y(t)) \quad (37.356c)$$

$$= e^{-(\lambda+\mu)t} \mathbb{1}_{[0, \infty[}(t) \quad (37.356d)$$

$$= 1 - F_Z(t) \quad (37.356e)$$

où $Z \sim \mathcal{E}(\lambda + \mu)$. □

37.5.7 Approximation de la binomiale par une Poisson

Proposition 37.101.

Soit (X_n) une suite de variables aléatoires avec $X_n \sim \mathcal{B}(n, p_n)$ telle que np_n converge vers une constante $\lambda > 0$. Alors $X_n \xrightarrow{\mathcal{L}} \mathcal{P}(\lambda)$.

Démonstration. Commençons par écrire la loi binomiale sous une forme plus adaptée au passage à la limite :

$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k} = \frac{n(n-1)\dots(n-k+1)}{k!} p^k (1-p)^{n-k}. \quad (37.357)$$

Le produit au numérateur contient k termes dans lesquels nous mettons n en évidence. Nous trouvons

$$P(X = k) = \frac{(np)^k \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \dots \left(1 - \frac{k-1}{n}\right)}{k!} p^k (n-p)^{n-k}. \quad (37.358)$$

Lorsque nous passons à la limite, tous les facteurs du type $1 - l/n$ tendent vers 1 ainsi que $(1-p_n)^{-k}$. Les facteurs dont la limite n'est pas 1 sont donc

$$P(X_n = k) \simeq \frac{(np_n)^k}{k!} (1-p_n)^k. \quad (37.359)$$

Nous avons

$$\lim_{n \rightarrow \infty} (1-p_n)^n = \lim_{n \rightarrow \infty} \left(1 - \frac{np_n}{n}\right)^n = e^{-\lambda}. \quad (37.360)$$

La thèse est alors obtenue en remettant les morceaux ensemble. □

Exemple 37.102

Considérons un serveur informatique qui reçoit des requêtes. Toutes les 10^{-3} s il reçoit une requête avec une probabilité $p = 0.05$. La variable aléatoire qui consiste à donner le nombre de requêtes effectivement effectuées en une seconde suit une loi binomiale $\mathcal{P}(1000, p)$.

Déterminons la probabilité que le serveur reçoive 20 requêtes en une seconde. Nous approximations $\mathcal{B}(1000, 0.05)$ par $\mathcal{P}(50)$, et la réponse est

$$e^{-50} \frac{50^{20}}{20!} \simeq 7 \cdot 10^{-7}. \quad (37.361)$$

△

37.5.8 Loi de Poisson et loi exponentielle

Soient X_1, \dots, X_n des variables aléatoires réelles indépendantes de loi exponentielle de paramètre λ . En utilisant le produit de convolution, nous pouvons trouver la fonction de densité de la somme (voir point 37.2.3). Commençons avec deux variables aléatoires X et Y . Les densités sont

$$f_X(x) = \mathbb{1}_{[x \geq 0]} \lambda e^{-\lambda x} \quad (37.362a)$$

$$f_Y(y) = \mathbb{1}_{[y \geq 0]} \lambda e^{-\lambda y}, \quad (37.362b)$$

et la densité conjointe est alors

$$f_{X+Y}(x) = \int_{\mathbb{R}} \mathbb{1}_{[x-t \geq 0]} \lambda e^{-\lambda(x-t)} \mathbb{1}_{[t \geq 0]} \lambda e^{-\lambda t} dt \quad (37.363a)$$

$$= \lambda^2 e^{-\lambda x} \int_0^x 1 dt \quad (37.363b)$$

$$= x \lambda^2 e^{-\lambda x}. \quad (37.363c)$$

Par récurrence si $S = X_1 + \dots + X_n$ nous trouvons

$$f_S(x) = x^{n-1} \lambda^n e^{-\lambda x}. \quad (37.364)$$

Proposition 37.103 ([493]).

Soit $(T_k)_{k \in \mathbb{N}}$ une suite de variables aléatoires indépendantes de loi $\mathcal{E}(\lambda)$. Nous considérons la variable aléatoire $S_n = \sum_{i=1}^n X_i$ et pour chaque $t \in \mathbb{R}^+$ nous considérons

$$N_t = \max\{n \geq 1 \text{ tel que } \sum_{i=1}^n T_i \leq t\}. \quad (37.365)$$

Alors $N_t \sim \mathcal{P}(\lambda t)$.

Démonstration. Ce que nous devons calculer est

$$P(N_t = k) = P(S_n \leq t \leq S_{n+1}). \quad (37.366)$$

Nous introduisons la variable aléatoire $V_{n+1} = (X_1, \dots, X_{n+1})$ ainsi que l'ensemble

$$A_{n+1} = \{x \in \mathbb{R}^{n+1} \text{ tel que } x_1 + \dots + x_n \leq t \leq x_1 + \dots + x_n + x_{n+1}\}. \quad (37.367)$$

Le problème est donc de calculer

$$P(S_n \leq t \leq S_{n+1}) = P(V_{n+1} \in A_{n+1}^+) = \int_{A_{n+1}^+} f_{n+1}(x) dx \quad (37.368)$$

où A_{n+1}^+ est la partie de A_{n+1} dans laquelle $x_i \geq 0$ pour tout i et f_{n+1} est la fonction de densité conjointe des variables aléatoires X_i . Nous effectuons le changement de variables

$$s_k = \sum_{i \leq k} x_i \quad (37.369a)$$

$$x_k = s_k - s_{k-1} \quad (37.369b)$$

dont le déterminant vaut 1. D'autre part par indépendance des variables aléatoires X_i , la fonction de partition jointe f_{n+1} s'exprime sous la forme

$$f_{n+1}(x_1, \dots, x_{n+1}) = f_{X_1}(x_1) \dots f_{X_{n+1}}(x_{n+1}) \quad (37.370a)$$

$$= \lambda^{n+1} e^{-\lambda(x_1 + \dots + x_{n+1})} \quad (37.370b)$$

$$= \lambda^{n+1} e^{-\lambda s_{n+1}}. \quad (37.370c)$$

En ce qui concerne les bornes de l'intégrale dans les variables s_i , nous voulons que tous les x_i soient positifs, par conséquent $s_1 \geq 0$ et ensuite l'équation $x_k = s_k - s_{k-1}$ demande $s_k \geq s_{k-1}$. Les bornes sont donc données par l'ensemble

$$0 \leq s_1 \leq s_2 \leq \dots \leq s_n \leq t \leq s_{n+1}, \quad (37.371)$$

c'est-à-dire $B_n \times]t, \infty[$ où

$$B_n = \{(s_1, \dots, s_n) \text{ tel que } 0 \leq s_1 \leq s_2 \leq \dots \leq s_n \leq t\}. \quad (37.372)$$

Le théorème de Fubini nous permet de décomposer l'intégrale :

$$P(S_n \leq t \leq S_{n+1}) = \int_{B_n \times]t, \infty[} \lambda^{n+1} e^{-\lambda s_{n+1}} ds_1 \dots ds_{n+1} \quad (37.373a)$$

$$= \lambda^{n+1} \left(\int_{B_n} ds_1 \dots ds_n \right) \underbrace{\left(\int_t^\infty e^{-\lambda s_{n+1}} ds_{n+1} \right)}_{= \lambda^{-1} e^{-\lambda t}} \quad (37.373b)$$

$$= \lambda^n e^{-\lambda t} \text{Vol}(B_n) \quad (37.373c)$$

où $\text{Vol}(B_n)$ est le volume de B_n qui reste à calculer. L'ensemble $C^n = [0, t]^n$ se décompose en cellules disjointes (à ensemble de mesure nulle près) de la forme

$$C_\sigma = \{0 \leq s_{\sigma(1)} \leq s_{\sigma(2)} \leq \dots \leq s_{\sigma(n)} \leq t\} \quad (37.374)$$

pour chaque permutation $\sigma \in S_n$. Il y a exactement $n!$ telles cellules dans C^n . Par conséquent

$$t^n = \text{Vol}(C^n) = n! \text{Vol}(C_\sigma) = n! \text{Vol}(B_n) \quad (37.375)$$

et $\text{Vol}(B_n) = \frac{t^n}{n!}$. Finalement nous avons

$$P(n_t = n) = P(S_n \leq t \leq S_{n+1}) = \frac{(\lambda t)^n}{n!} e^{-(\lambda t)}. \quad (37.376)$$

□

37.5.9 Loi normale

La loi normale de paramètres m et $\sigma > 0$, notée $\mathcal{N}(m, \sigma^2)$ est la loi donnée par la densité

$$\gamma_{m, \sigma^2}(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left[-\frac{1}{2} \left(\frac{x - m}{\sigma} \right)^2 \right]. \quad (37.377)$$

Proposition 37.104.

Si la variable aléatoire réelle X suit une loi normale $\mathcal{N}(m, \sigma^2)$, alors nous avons $E(X) = m$ et $\text{Var}(X) = \sigma^2$.

Démonstration. L'espérance d'une variable aléatoire se calcule à partir de la formule (37.220) :

$$E(X) = \frac{1}{\sigma\sqrt{2\pi}} \int_{\mathbb{R}} x \exp \left[-\frac{1}{2} \left(\frac{x-m}{\sigma} \right)^2 \right] dx \quad (37.378a)$$

$$= \frac{1}{\sigma\sqrt{2\pi}} \sigma \int_{\mathbb{R}} (\sigma u + m) e^{-u^2/2} du \quad (37.378b)$$

où nous avons effectué le changement de variable $u = (x-m)/\sigma$. Nous utilisons ensuite l'intégrale remarquable

$$\int_{\mathbb{R}} e^{-u^2/2} du = \sqrt{2\pi}. \quad (37.379)$$

En ce qui concerne la variance, nous avons le même genre de calculs. \square

La **loi normale réduite** est la densité

$$\gamma(x) = \gamma_{0,1}(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}. \quad (37.380)$$

La variable aléatoire X suit la loi $\mathcal{N}(m, \sigma^2)$ si et seulement si la variable aléatoire $Z = \frac{X-m}{\sigma}$ suit la loi normale réduite $\mathcal{N}(0, 1)$.

Proposition 37.105.

La fonction caractéristique de la distribution normale $\mathcal{N}(m, \sigma^2)$ est

$$\Phi_{\mathcal{N}(m, \sigma^2)}(t) = \exp \left(itm - \frac{\sigma^2 t^2}{2} \right). \quad (37.381)$$

Démonstration. En suivant la formule (37.192), l'intégrale à calculer est

$$\Phi_{\mathcal{N}(m, \sigma^2)}(t) = \frac{1}{\sigma\sqrt{2\pi}} \int_{\mathbb{R}} e^{itx} e^{-\frac{1}{2} \left(\frac{x-m}{\sigma} \right)^2} dx. \quad (37.382)$$

Nous reconnaissons une transformée de Fourier. Afin de la calculer sans encombres, nous passons par les fonctions intermédiaires suivantes :

$$\begin{aligned} g(x) &= e^{-\frac{1}{2}x^2} \\ h(x) &= g\left(\frac{x}{\sigma}\right) \\ k(x) &= h(x-m). \end{aligned} \quad (37.383)$$

La fonction caractéristique que nous cherchons est $\frac{1}{\sigma\sqrt{2\pi}} \hat{k}(t)$. Les formules liées à la transformée de Fourier nous donnent

$$\hat{k}(t) = \hat{h}(t) e^{itm} \quad (37.384a)$$

$$\hat{h}(t) = \sigma \hat{g}(\sigma t) \quad (37.384b)$$

$$\hat{g}(t) = \int_{\mathbb{R}} e^{-itx} e^{-\frac{1}{2}x^2} dx = \sqrt{2\pi} e^{-t^2/2}. \quad (37.384c)$$

Attention : l'intégrale à calculer est une transformée de Fourier *inverse*, d'où la formule (37.384a) qui a un signe de différence avec la formule usuelle. En recombinaison toutes ces expressions nous trouvons

$$\Phi_{\mathcal{N}(m, \sigma^2)} = e^{-\sigma^2 t^2/2} e^{itm}, \quad (37.385)$$

ce qu'il nous fallait. \square

Exemple 37.106

Une espérance qui sert de temps en temps est celle de $X = e^{\beta Z}$ lorsque $Z \sim \mathcal{N}(0, 1)$. Elle se calcule en remarquant que $x^2 - 2\beta x = (x - \beta)^2 - \beta^2$, donc

$$E(e^{\beta Z}) = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{\beta x} e^{-x^2/2} dx \quad (37.386a)$$

$$= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{1}{2}(x-\beta)^2} e^{\beta^2/2} dx \quad (37.386b)$$

$$= \frac{e^{\beta^2/2}}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-y^2/2} dy \quad (37.386c)$$

$$= e^{\beta^2/2}. \quad (37.386d)$$

△

37.5.10 Vecteurs gaussiens

Source : [236, 494].

Définition 37.107.

Un vecteur aléatoire $X: \Omega \rightarrow \mathbb{R}^d$ est un **vecteur gaussien** si toutes les combinaisons linéaires de ses composantes sont des variables aléatoires normales. En d'autres termes, X est un vecteur gaussien si pour tout vecteur u , la variable aléatoire $u \cdot X$ est gaussienne.

Le vecteur moyenne d'un vecteur gaussien est $E(X) = (E(X_1), \dots, E(X_n))$ et sa matrice de **variance-covariance** est la matrice

$$K_X = \text{Var}(X) = E\left[(X - E(X)) \otimes (X - E(X))\right] \quad (37.387)$$

où l'opération \otimes est celle introduite autour de l'équation (12.295). Cela n'est rien d'autre que la matrice de covariance de la variable aléatoire $X: \Omega \rightarrow \mathbb{R}^d$.

Lemme 37.108.

Si les variables aléatoires réelles X_1, \dots, X_n sont des variables aléatoires réelles gaussiennes indépendantes, alors le vecteur $X = (X_1, \dots, X_n)$ est un vecteur gaussien.

Démonstration. Nous devons montrer que si X et Y sont des variables aléatoires gaussiennes indépendantes, alors $X+Y$ est encore gaussienne. L'indépendance nous assure les égalités suivantes pour la fonction caractéristique :

$$\Phi_{X+Y}(t) = E(e^{it(X+Y)}) = E(e^{itX})E(e^{itY}). \quad (37.388)$$

Dans le cas où X et Y sont gaussiens nous trouvons

$$\Phi_{X+Y}(t) = \exp\left(im_1 t - \frac{\sigma_1^2 t^2}{2}\right) \exp\left(im_2 t - \frac{\sigma_2^2 t^2}{2}\right) = \exp\left(i(m_1 + m_2)t - \frac{(\sigma_1^2 + \sigma_2^2)t^2}{2}\right). \quad (37.389)$$

Étant donné que la loi d'une variable aléatoire est entièrement déterminée par sa fonction caractéristique (théorème 37.58), nous déduisons que $X+Y$ est une normale de moyenne $m_1 + m_2$ et de variance $\sigma = \sigma_1^2 + \sigma_2^2$. □

Proposition 37.109.

La fonction caractéristique d'un vecteur gaussien est donnée par

$$\Phi_X(u) = \exp\left(iu \cdot E(X) - \frac{1}{2}u \cdot K_X u\right) \quad (37.390)$$

où K_X est la matrice de covariance de X .

Démonstration. Nous considérons la variable aléatoire réelle gaussienne $u \cdot X$. Son espérance $m = E(u \cdot X) = u \cdot E(X)$. Nous commençons par établir la formule suivante :

$$u^t K_X u = u \cdot K_X u = E\left([X - E(X)] \cdot u\right)^2. \quad (37.391)$$

Utilisant la linéarité de l'espérance,

$$\sum_{kl} E(A_{kl}) u_k u_l = \sum_{kl} E(A_{kl} u_k u_l), \quad (37.392)$$

nous trouvons

$$K_X(u, u) = E\left([X - E(X)] \otimes [X - E(X)]\right)(u, u) \quad (37.393a)$$

$$= E\left([X - E(X)] \otimes [X - E(X)]\right)(u, u) \quad (37.393b)$$

$$= E\left([X - E(X)] \cdot u\right)\left([X - E(X)] \cdot u\right) \quad (37.393c)$$

$$= E\left([X - E(X)] \cdot u\right)^2. \quad (37.393d)$$

Par la linéarité du produit scalaire et de l'espérance,

$$[X - E(X)] \cdot u = u \cdot X - E(u \cdot X), \quad (37.394)$$

ce qui nous ramène à la variable aléatoire $u \cdot X$. Nous avons alors

$$K_X(u, u) = E\left([X - E(X)] \cdot u\right)^2 = \text{Var}(u \cdot X). \quad (37.395a)$$

Nous avons donc obtenu une forme pour la variance de la variable aléatoire $u \cdot X$. Étant donné que $u \cdot X$ est gaussienne de moyenne $m = u \cdot E(X)$ et de variance $\sigma^2 = K_X(u, u)$, nous avons

$$\Phi_{u \cdot X}(t) = \exp\left(itm - \frac{1}{2}\sigma^2 t^2\right). \quad (37.396)$$

Par ailleurs nous avons $\Phi_X(u) = \Phi_{u \cdot X}(1)$ parce que

$$\Phi_X(u) = E(e^{iu \cdot X}) = \Phi_{u \cdot X}(1). \quad (37.397)$$

En utilisant la forme (37.396) pour $\Phi_{u \cdot X}$ nous trouvons

$$\Phi_X(u) = \Phi_{u \cdot X}(1) = \exp\left(im - \frac{1}{2}\sigma^2\right) = \exp\left(iE(u \cdot X) - \frac{1}{2}u \cdot K_X u\right). \quad (37.398)$$

□

Théorème 37.110.

Soit $X = (X_1, \dots, X_d)$, un vecteur gaussien. Les composantes sont indépendantes si et seulement si elles sont non corrélées.

Démonstration. Nous savons que les variables aléatoires indépendantes sont non corrélées. Nous devons donc surtout prouver le contraire. Le fait que les variables aléatoires X_i soient non corrélées signifie que la matrice de covariance est

$$K_X = \begin{pmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_d^2 \end{pmatrix} \quad (37.399)$$

où $\sigma_k^2 = \text{Var}(X_k)$. Notons $m_k = E(X_k)$. Si $u \in \mathbb{R}^d$, nous avons en vertu de la proposition 37.109 que

$$\Phi_X(u) = \exp\left(i(u_1 m_1 + \dots + u_d m_d) - \frac{1}{2}(\sigma_1^2 u_1^2 + \dots + \sigma_d^2 u_d^2)\right) \quad (37.400a)$$

$$= \exp\left(iu_1 m_1 - \frac{1}{2}\sigma_1^2 u_1^2\right) \dots \exp\left(iu_d m_d - \frac{1}{2}\sigma_d^2 u_d^2\right) \quad (37.400b)$$

$$= \Phi_{X_1}(u_1) \dots \Phi_{X_d}(u_d). \quad (37.400c)$$

Les variables aléatoires X_i sont donc indépendantes parce que la fonction caractéristique se factorise. \square

Exemple 37.111

Nous donnons à présent un exemple de deux variables aléatoires gaussiennes qui ne forment pas un vecteur gaussien. Pour ce faire nous devons chercher des variables aléatoires non indépendantes. Soit $Y \sim \mathcal{N}(0, 1)$ et ϵ une variable aléatoire (indépendante de Y) donnée par $P(\epsilon = -1) = \frac{1}{2}$, $P(\epsilon = 1) = \frac{1}{2}$. Nous considérons le vecteur $(Y, \epsilon Y)$.

D'abord montrons que ϵY est une variable aléatoire gaussienne. Soit A un borélien de \mathbb{R} . Nous avons

$$P(\epsilon Y \in A) = P(\epsilon = -1, Y \in A) + P(\epsilon = -1, -Y \in A). \quad (37.401)$$

Par indépendance et par symétrie de Y ¹⁵ nous trouvons

$$P(\epsilon Y \in A) = P(\epsilon = 1)P(Y \in A) + P(\epsilon = -1)P(-Y \in A) \quad (37.402a)$$

$$= \frac{1}{2}P(Y \in A) + \frac{1}{2}P(-Y \in A) \quad (37.402b)$$

$$= P(Y \in A). \quad (37.402c)$$

Nous avons donc $Y \sim \epsilon Y$, et donc ϵY est gaussienne.

En ce qui concerne la covariance, nous savons que $E(Y) = E(\epsilon) = 0$, donc

$$\text{Cov}(Y, \epsilon Y) = E(Y \cdot \epsilon Y) = E(\epsilon Y^2) = E(\epsilon)E(Y^2) = 0. \quad (37.403)$$

Note : $E(Y^2) = 1$.

Les variables aléatoires Y et ϵY ne sont pas indépendantes. En effet si elles l'étaient, Y serait aussi indépendante de $(\epsilon Y)^2 = Y^2$, alors que Y et Y^2 ne sont pas indépendantes. Donc $X = (Y, \epsilon Y)$ n'est pas un vecteur gaussien. \triangle

Théorème 37.112 ([236]).

Soit $m \in \mathbb{R}^d$ et $K \in \mathbb{M}(d, \mathbb{R})$ une matrice symétrique et positive. Alors il existe un vecteur gaussien de moyenne m et de matrice de covariance K .

Démonstration. Nous effectuons la preuve avec $m = 0$. Nous choisissons l'espace probablisé $(\Omega, \mathcal{A}) = (\mathbb{R}^r, \mathcal{B}(\mathbb{R}^r))$ où r est le rang de K muni de la probabilité de densité

$$\gamma(u) = \frac{1}{(2\pi)^{r/2}} \exp\left(-\frac{1}{2}\|u\|^2\right). \quad (37.404)$$

Nous considérons la variable aléatoire

$$Y_0 : \Omega \rightarrow \mathbb{R}^r \quad (37.405)$$

$$u \mapsto u.$$

C'est une variable aléatoire gaussienne de loi $P_{Y_0} = \gamma \lambda_d$ où λ_d est la mesure de Lebesgue sur \mathbb{R}^d . Sa densité (37.404) s'écrit comme le produit de r gaussiennes indépendantes ; sa matrice de covariance est donc $\mathbb{1}_{r \times r}$.

15. C'est-à-dire que $P(Y \in A) = P(Y \in -A) = P(-Y \in A)$.

Étant donné que K est symétrique et positive, il existe une matrice $d \times r$ telle que $K = AA^t$. Pour voir cela, remarquons qu'il existe une matrice $d \times d$ qui fait le travail. En effet K se diagonalise par une orthogonale (théorème 11.189) :

$$K = ADA^t = A\sqrt{D}\sqrt{D}A^t \tag{37.406}$$

où D est une matrice diagonale contenant $d - r$ zéros et \sqrt{D} est la matrice que l'on imagine. Donc la matrice $L = A\sqrt{D}$ est une matrice telle que $LL^t = K$. Maintenant, étant donné que les $d - r$ dernières lignes de D sont vides, les $d - r$ dernières lignes de L n'ont pas d'importance et peuvent être choisies nulles, voire même ne pas exister. La matrice $A \in \mathbb{M}_{d \times r, \mathbb{R}}$ qui réalise $AA^t = K$ est la « troncature » de L à ses r premières lignes.

Nous considérons la variable aléatoire $Y = AY_0: \Omega \rightarrow \mathbb{R}^d$. Étant donné que AY_0 est une transformation linéaire d'un vecteur gaussien, c'est un vecteur gaussien. Nous avons encore $m(Y) = 0$ et

$$\text{Var}(Y) = E(Y \otimes Y) = E((AY_0) \otimes (AY_0)) = AE(Y_0 \otimes Y_0)A^t = AA^t = K \tag{37.407}$$

parce que $E(Y_0 \otimes Y_0) = \mathbb{1}$. Nous avons utilisé les formules du produit tensoriel introduit en (12.295), et en particulier la formule (12.298). □

Lorsqu'une matrice symétrique et positive K est donné, nous avons créé un vecteur gaussien de covariance K en créant $X = AY$ où Y est le vecteur gaussien « le plus simple » et où A est donné par $AA^t = K$. La proposition suivante montre l'inverse : un vecteur gaussien X peut se réduire au vecteur gaussien « le plus simple » en utilisant la transformation $Y = A^{-1}X$ où A est encore donnée par $AA^t = K$. Ce résultat nous permettra de voir les vecteurs gaussiens généraux comme des « changements de coordonnées » par rapport au vecteur gaussien simple $Y = (Y_1, \dots, Y_d)$ avec $Y_i \sim \mathcal{N}(0, 1)$.

Proposition 37.113.

Soit $Y = (Y_1, \dots, Y_n)$ le vecteur gaussien formé des variables aléatoires indépendantes $Y_i \sim \mathcal{N}(0, 1)$. Soit K une matrice symétrique et positive et la matrice A telle que $AA^t = K$. Alors le vecteur $X = AY$ est gaussien de covariance K .

Démonstration. Nous montrons que la covariance de $X = AY$ est donnée par K . Nous avons

$$K_X = E([X - E(X)] \otimes [X - E(X)]) \tag{37.408a}$$

$$= E(A(Y - E(Y)) \otimes A(Y - E(Y))) \tag{37.408b}$$

$$= E\left(A([Y - E(Y)] \otimes [Y - E(Y)])A^t\right) \tag{37.408c}$$

$$= E(AA^t) \tag{37.408d}$$

$$= K_X. \tag{37.408e}$$

Nous avons utilisé le lemme 12.132 ainsi que le fait que

$$[Y - E(Y)] \otimes [Y - E(Y)] = K_Y = \mathbb{1}. \tag{37.409}$$

Notons que $E(Y) = 0$, mais cela ne joue pas ici. □

Théorème 37.114.

Un vecteur gaussien $X: \Omega \rightarrow \mathbb{R}^d$ possède une densité par rapport à la mesure de Lebesgue si et seulement si sa matrice de covariance est inversible. Dans ce cas nous avons la densité

$$f_X(x) = \frac{1}{(2\pi)^{d/2}} \frac{1}{\sqrt{|\det(K_X)|}} \exp\left(-\frac{1}{2}[X - E(X)] \cdot K_X^{-1}[X - E(X)]\right). \tag{37.410}$$

Démonstration. Nous supposons que $E(X) = 0$. En utilisant la proposition 37.113, nous posons $X = AY$ où Y est un vecteur gaussien $Y \sim \mathcal{N}(0, \mathbb{1})$ et $AA^t = K_X$. Si K n'est pas inversible, alors

A n'est pas inversible non plus. Notons $r < d$ le rang de A . Étant donné que $X = AY$, la variable aléatoire X prend presque sûrement ses valeurs dans l'image de A , c'est-à-dire dans un sous-espace de dimension $r < d$ de \mathbb{R}^d . Ce sous-espace est de mesure nulle pour la mesure de Lebesgue, mais de mesure 1 pour la mesure P_X . La mesure P_X ne peut donc pas avoir de densité par rapport à celle de Lebesgue.

Supposons maintenant que K_X soit inversible. La matrice A l'est aussi. Nous anticipons l'utilisation du théorème de changement de variable 15.252. Ici le changement de variable sera la transformation linéaire A dont le jacobien vaut

$$|\det(A^{-1})| = \frac{1}{\det A} = \frac{1}{\sqrt{|\det K_X|}}. \quad (37.411)$$

Soit γ la densité de Y . Nous posons

$$f_X(x) = \frac{1}{\sqrt{|\det K_X|}} \gamma(A^{-1}x) \quad (37.412)$$

où γ est le produit des densités de d gaussiennes usuelles $\mathcal{N}(0, 1)$. Nous allons d'abord montrer que cette formule est bien la fonction (37.410) et ensuite que X admet f_X comme densité. Nous avons

$$\gamma(A^{-1}x) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2}\|A^{-1}x\|^2\right), \quad (37.413)$$

et

$$\|A^{-1}x\|^2 = \langle A^{-1}x, A^{-1}x \rangle \quad (37.414a)$$

$$= \langle (A^{-1})^t A^{-1}x, x \rangle \quad (37.414b)$$

$$= \langle (AA^t)^{-1}x, x \rangle \quad (37.414c)$$

$$= x \cdot K_X^{-1}x, \quad (37.414d)$$

ce qui nous donne bien la formule (37.410). Nous vérifions maintenant que f_X est bien une densité pour X . Soit B un borélien de \mathbb{R}^d . Nous avons d'abord

$$P(X \in B) = P(AY \in B) = P(Y \in A^{-1}B). \quad (37.415)$$

Ici nous avons utilisé le fait que A était bijectif. Nous avons ensuite

$$P(X \in B) = \int_{A^{-1}B} \gamma(t) dy = \int_B \gamma(A^{-1}x) |J_{A^{-1}}(x)| dx = \int_B f_X(x) dx. \quad (37.416)$$

C'est ici que nous avons utilisé le théorème de changement de variable 15.252. \square

37.5.11 Variable aléatoire de Rademacher

Une variable aléatoire **de Rademacher** est une variable aléatoire de loi

$$\epsilon \sim \delta_0 - \delta_1, \quad (37.417)$$

sur l'ensemble $\Omega = \{0, 1\}$. C'est la variable aléatoire qui prend valeur 1 ou -1 avec probabilité $\frac{1}{2}$.

Parmi les propriétés évidentes de cette variable aléatoire nous avons $E(\epsilon) = 0$ et $E(\epsilon^2) = 1$.

Proposition 37.115 (Inégalité de Khintchine[43]).

Soient r_1, \dots, r_n des variables aléatoires indépendantes et identiquement distribuées de Rademacher et une combinaison linéaire

$$X = \sum_{i=1}^n a_i r_i \quad (37.418)$$

avec $a_i \in \mathbb{R}$. Alors

$$\|X\|_2 \leq \sqrt{e} E(|X|). \quad (37.419)$$

Démonstration. Pour rappel la définition est que

$$\|X\|_2 = \sqrt{E(|X|^2)}. \quad (37.420)$$

Vu que c'est de la quantité $E(|X|)$ que nous voulons parler, nous notons l'inégalité

$$E(|X|) = \int_{\Omega} |X(\omega)| dP(\omega) \geq \left| \int_{\Omega} X(\omega) dP(\omega) \right| = |E(X)|. \quad (37.421)$$

Nous supposons que $\sum_{j=1}^n a_j^2 = 1$; sinon nous multiplions les a_j parce qu'il faut pour l'avoir et ce facteur sortira des deux côtés de (37.419).

Nous allons passer par la variable aléatoire intermédiaire

$$Y = \prod_{j=1}^n (1 + ia_j r_j) \quad (37.422)$$

et pour presque tout $\omega \in \Omega$ nous avons

$$|Y(\omega)| = \prod_{j=1}^n |1 + ia_j r_j(\omega)| \quad (37.423a)$$

$$= \prod_{j=1}^n \sqrt{1 + a_j^2 \underbrace{r_j(\omega)^2}_{=1}} \quad (37.423b)$$

$$= \prod_{j=1}^n \sqrt{1 + a_j^2} \quad (37.423c)$$

$$\leq \prod_{j=1}^n \sqrt{e^{a_j^2}} \quad (37.423d)$$

$$= \sqrt{\prod_{j=1}^n e^{a_j^2}} \quad (37.423e)$$

$$= \sqrt{e^{\sum a_j^2}} \quad (37.423f)$$

$$= \sqrt{e}. \quad (37.423g)$$

Donc $\|Y\|_{\infty} \leq \sqrt{e}$ où $\|\cdot\|_{\infty}$ est la norme supremum sur Ω . Cela nous permet de donner une première inégalité à propos de $E(|X|)$. D'abord

$$|E(XY)| = \left| \int_{\Omega} X(\omega)Y(\omega) dP(\omega) \right| \leq \int_{\Omega} |X(\omega)| \|Y\|_{\infty} dP(\omega) = \|Y\|_{\infty} E(|X|), \quad (37.424)$$

ensuite en remplaçant $\|Y\|_{\infty}$ par la majoration que nous venons de donner de $|Y(\omega)|$,

$$\sqrt{e} E(|X|) \geq |E(XY)|. \quad (37.425)$$

Il nous reste à prouver que $|E(XY)| \geq \|X\|_2$.

Pour ce faire nous commençons par noter que pour chaque j nous avons $E(r_j) = 0$ et $E(ia_j r_j^2) = ia_j$; en utilisant l'indépendance des r_j et le lemme 37.22 nous avons alors

$$E(r_j Y) = E\left(r_j (1 + ia_j r_j) \prod_{k \neq j} (1 + ia_k r_k)\right) \quad (37.426a)$$

$$= E(r_j (1 + ia_j r_j)) \prod_{k \neq j} E(1 + ia_k r_k) \quad (37.426b)$$

$$= ia_j \prod_{k \neq j} (1 + ia_k E(r_k)) \quad (37.426c)$$

$$= ia_j. \quad (37.426d)$$

Par conséquent, en utilisant la proposition 37.25 dans le cas non indépendant,

$$E(XY) = \sum_{j=1}^n a_j E(r_j Y) = \sum_j a_j i a_j = i \sum_j a_j^2 = i. \quad (37.427)$$

Nous pouvons compléter l'équation (37.425) en

$$\sqrt{e} E(|X|) \geq |E(XY)| = 1, \quad (37.428)$$

et nous nous empressons de montrer que $\|X\|_2 = 1$. En effet $\|X\|_2 = \sqrt{E(|X|^2)}$ alors que

$$\|X\|_2^2 = E\left(\sum_i a_i r_i\right)^2 \quad (37.429a)$$

$$= E\left(\sum_k a_k^2 r_k^2 + 2 \sum_{i \neq j} a_i a_j r_i r_j\right) \quad (37.429b)$$

$$= \sum_k a_k^2 + 2 \sum_{i \neq j} a_i a_j E(r_i r_j) \quad (37.429c)$$

$$= 1 \quad (37.429d)$$

□

37.5.12 Loi de Student

Définition 37.116.

La loi χ^2 à d degrés de liberté est la loi de la variable aléatoire $Y_1^2 + \dots + Y_n^2$ si les (Y_i) sont des variables aléatoires normales indépendantes centrées et réduites.

La loi de **Student** à d degrés de liberté est la loi de la variable aléatoire

$$\frac{X}{\sqrt{K/d}} \quad (37.430)$$

où $X \sim \mathcal{N}(0, 1)$ et $K \sim \chi^2(d)$ sont des variables aléatoires indépendantes. Cette loi est notée $\mathcal{T}(d)$

Nous avons une illustration de la densité de la loi $\chi^2(10)$ à la figure 37.1.

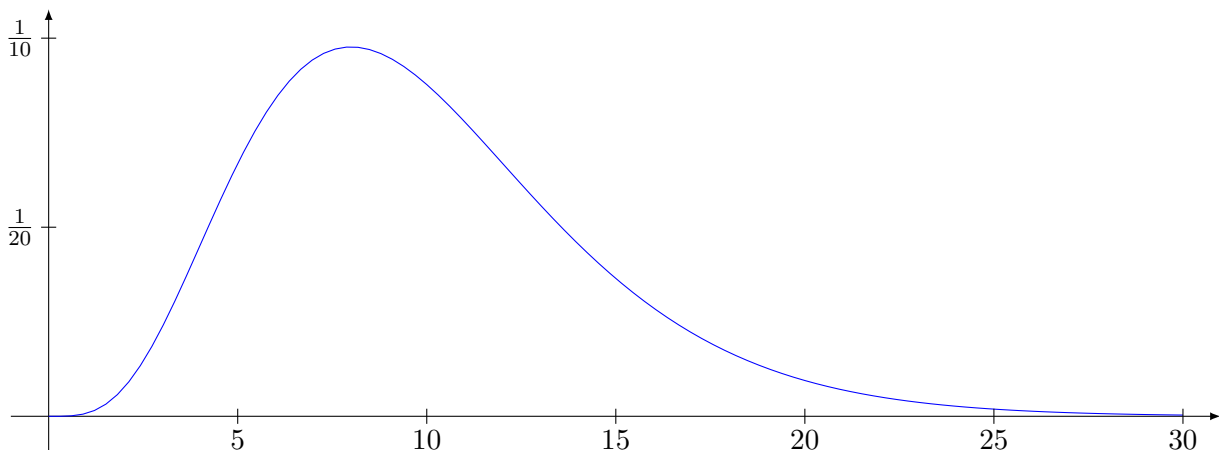


FIGURE 37.1 – La densité de $\chi^2(10)$.

L'importance de cette loi sera dans le théorème de Cochran 38.15.

37.5.13 Indépendance, covariance et variance de somme

Si X et Y sont des variables aléatoires réelles, nous avons défini la covariance par (37.67) :

$$\text{Cov}(X, Y) = E\left[(X - E(X))(Y - E(Y))\right] \quad (37.431)$$

Une clef est le lemme 37.22 qui dit que lorsque X et Y sont indépendantes, alors $E(XY) = E(X)E(Y)$. Nous avons les liens suivants.

- (1) Si X et Y sont indépendantes, alors $\text{Cov}(X, Y) = 0$.
- (2) La réciproque n'est pas vraie par l'exemple 37.24.
- (3) Si Z est un vecteur gaussien, les composantes Z_i sont indépendantes si et seulement si la matrice de covariance est diagonale, c'est-à-dire que les Z_i sont deux à deux non corrélées; c'est le théorème 37.110.
- (4) En ce qui concerne la variance d'une somme,

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2 \text{Cov}(X, Y). \quad (37.432)$$

Donc lorsque X et Y ne sont pas corrélées, la variance est sympa avec la somme. En particulier lorsqu'elles sont indépendantes, mais pas seulement.

37.6 Estimation des grands écarts

Si S_n est une somme de variables aléatoires de Bernoulli X_i indépendantes de probabilité p , l'espérance de S_n/n est p , et nous voudrions savoir quelle est la probabilité d'avoir un taux de succès un peu plus important :

$$P\left(\frac{S_n}{n} \geq p + \epsilon\right). \quad (37.433)$$

La loi des grands nombres nous permet de dire que ça ne va pas être très grand. En effet si nous posons

$$Y_i = X_i - p - \epsilon \quad (37.434)$$

et

$$Z_n = \frac{1}{n} \sum_{i=1}^n Y_i = \frac{S_n}{n} - p - \epsilon, \quad (37.435)$$

la loi des grands nombres (théorème 37.83) nous indique que

$$Z_n \xrightarrow{p.s.} E(Y_1) = -\epsilon. \quad (37.436)$$

La proposition 37.68 sur le lien entre les types de convergence nous donne immédiatement la convergence en loi. C'est-à-dire que pour tout $\eta > 0$,

$$P\left(|Z_n + \epsilon| \geq \eta\right) \rightarrow 0. \quad (37.437)$$

En prenant $\eta = \frac{\epsilon}{2}$,

$$P(Z_n = 0) \leq P(|Z_n + \epsilon| \geq \frac{\epsilon}{2}) \rightarrow 0. \quad (37.438)$$

Tout cela montre que

$$\lim_{n \rightarrow \infty} P\left(\frac{S_n}{n} \geq p + \epsilon\right) = 0. \quad (37.439)$$

Autrement dit, la probabilité de tomber à une distance fixée de la moyenne tend vers zéro lorsque le nombre d'essais augmente. Rien d'étonnant.

Le théorème suivant nous indique la vitesse de convergence. Elle est exponentielle et le coefficient est donné en fonction de p et de ϵ .

Théorème 37.117 ([43]).

Soient des variables aléatoires $X_i \sim \mathcal{B}(p)$ et

$$S_n = \sum_{i=1}^n X_i \sim \mathcal{B}(n, p). \quad (37.440)$$

Pour $\epsilon \in]0, 1 - p[$, nous définissons

$$h_+(\epsilon) = (p + \epsilon) \ln \left(\frac{p + \epsilon}{p} \right) + (1 - p - \epsilon) \ln \left(\frac{1 - p - \epsilon}{1 - p} \right). \quad (37.441)$$

Alors

(1) $h_+(\epsilon) > 0$.

(2) Pour tout $n \geq 1$, nous avons

$$P \left(\frac{S_n}{n} \geq p + \epsilon \right) \leq e^{-nh_+(\epsilon)}. \quad (37.442)$$

(3) L'estimation (37.442) est optimale au sens que

$$\lim_{n \rightarrow \infty} \frac{1}{n} \ln \left(P \left(\frac{S_n}{n} \geq p + \epsilon \right) \right) = -h_+(\epsilon). \quad (37.443)$$

Démonstration. Pour tout $t \geq 0$ nous avons :

$$P \left(\frac{S_n}{n} \geq p + \epsilon \right) = P(S_n \geq np + n\epsilon) \quad (37.444a)$$

$$\leq E \left(e^{t(S_n - np - n\epsilon)} \right) \quad (37.444b)$$

$$= e^{-nt(p+\epsilon)} E \left(e^{tS_n} \right) \quad (37.444c)$$

$$= e^{-nt(p+\epsilon)} \sum_{k=0}^{\infty} e^{tk} P(S_n = k) \quad (37.444d)$$

$$= e^{-nt(p+\epsilon)} \sum_{k=0}^n e^{tk} \binom{n}{k} p^k (1-p)^{n-k} \quad (37.444e)$$

$$= e^{-nt(p+\epsilon)} \left((1-p) + pe^t \right)^n \quad (37.444f)$$

$$= \exp \left(-n(t(p+\epsilon) - \ln(1-p+pe^t)) \right). \quad (37.444g)$$

Justifications :

— L'inégalité (37.444b) est l'inégalité de Markov (corollaire 37.82) avec $\phi(x) = e^{tx}$.

— La ligne (37.444d) est une utilisation du théorème de transfert 37.60.

Nous posons maintenant

$$h = \sup_{t>0} (t(p+\epsilon) - \ln(1-p+pe^t)). \quad (37.445)$$

Pour cette valeur de h nous avons

$$P \left(\frac{S_n}{n} \geq p + \epsilon \right) \leq e^{-nh}. \quad (37.446)$$

Nous considérons la fonction

$$\begin{aligned} g: \mathbb{R}^+ &\rightarrow \mathbb{R} \\ t &\mapsto t(p+\epsilon) - \ln(1-p+pe^t). \end{aligned} \quad (37.447)$$

Elle vérifie $g(0) = 0$ et

$$g'(t) = (p + \epsilon) - \frac{pe^t}{1 - p + pe^t} \tag{37.448a}$$

$$g'(0) = (p + \epsilon) - \frac{p}{1 - p + p} = \epsilon > 0. \tag{37.448b}$$

Par conséquent sur un voisinage de $t = 0$ la fonction g est strictement croissante et nous concluons que g prend (au moins) quelques valeurs strictement positives. Du coup nous avons

$$h = \|g\|_\infty > 0. \tag{37.449}$$

Nous cherchons maintenant pour quelle valeur de t est réalisé le maximum de g . D'abord résoudre $g'(t) = 0$ donne

$$t_0 = \ln \left(\frac{(p + \epsilon)(1 - p)}{p(1 - p - \epsilon)} \right) \tag{37.450}$$

Vu que $g'(t) \rightarrow p + \epsilon - 1 < 0$, nous avons $\lim_{t \rightarrow \infty} g(t) = -\infty$ et donc le t_0 trouvé est bien un maximum et non un minimum.

Il est maintenant loisible de calculer une valeur pour h : il suffit de calculer $g(t_0)$. La calcul n'est pas très compliqué et donne

$$h = g(t_0) = (p + \epsilon) \ln \left(\frac{p + \epsilon}{p} \right) + (p + \epsilon - 1) \ln \left(\frac{1 - p}{1 - p - \epsilon} \right), \tag{37.451}$$

ce qui est bien $h = h_+(\epsilon)$. Cela démontre les points (1) et (2).

Nous montrons à présent l'aspect optimal de l'estimation. Nous savons déjà que

$$\frac{1}{n} \ln \left(P \left(\frac{S_n}{n} \geq p + \epsilon \right) \right) \leq h_+(\epsilon). \tag{37.452}$$

Nous posons $k_n = \lceil n(p + \epsilon) \rceil$. Vu que $p + \epsilon < 1$ et que $S_n \leq n$, nous avons

$$P \left(\frac{S_n}{n} \geq p + \epsilon \right) = P(S_n \geq n(p + \epsilon)) \geq P(S_n = k_n) = \binom{n}{k_n} p^{k_n} (1 - p)^{n - k_n}. \tag{37.453}$$

C'est maintenant que nous utilisons la formule de Stirling (lemme 21.204) pour chacune des factorielles intervenant dans le coefficient binomial. Nous trouvons :

$$P \left(\frac{S_n}{n} \geq p + \epsilon \right) \geq P(S_n = k_n) = \clubsuit = \frac{\left(\frac{n}{e}\right)^n \sqrt{2\pi n} \alpha(n) p^{k_n} (1 - p)^{n - k_n}}{\left(\frac{k_n}{e}\right)^{k_n} \sqrt{2\pi k_n} \alpha(k_n) \left(\frac{n - k_n}{e}\right)^{n - k_n} \sqrt{2\pi(n - k_n)}}. \tag{37.454}$$

Nous savons que $k_n = n(p + \epsilon) + \sigma(n)$ avec σ borné par 1. Par conséquent $n - k_n \rightarrow \infty$ et nous pouvons regrouper les coefficients en α en

$$\beta(n) = \frac{\alpha(n)}{\alpha(k_n)\alpha(n - k_n)} \rightarrow 1. \tag{37.455}$$

Nous remarquons aussi que les e se simplifient. Nous récrivons \clubsuit sous la forme

$$\clubsuit = \frac{1}{\sqrt{2\pi}} \underbrace{\sqrt{\frac{n}{k_n(n - k_n)}}}_{=A(n)} \beta(n) \frac{n^n p^{k_n} (1 - p)^{n - k_n}}{k_n^{k_n} (n - k_n)^{n - k_n}} \tag{37.456a}$$

$$= A(n) n^n \left(\frac{p}{k_n}\right)^{k_n} \left(\frac{1 - p}{n - k_n}\right)^{n - k_n} \tag{37.456b}$$

$$= A(n) \left(\frac{np}{k_n}\right)^{k_n} \left(\frac{n(1 - p)}{n - k_n}\right)^{n - k_n}. \tag{37.456c}$$

Nous passons au logarithme et nous étudions $\frac{1}{n} \ln (P(S_n = k_n))$. Nous avons les termes suivants à étudier :

$$\begin{aligned} \frac{1}{n} \ln (P(S_n = k_n)) &= -\frac{1}{2n} \ln(2\pi) + \frac{1}{2n} \ln \left(\frac{n}{k_n(n - k_n)} \right) \\ &+ k_n \ln \left(\frac{np}{k_n} \right) + (n - k_n) \ln \left(\frac{n(1 - p)}{n - k_n} \right) + \frac{1}{n} \ln (\alpha(n)). \end{aligned} \quad (37.457)$$

Nous étudions terme à terme la limite de cela lorsque $n \rightarrow \infty$.

- (1) Le terme $\frac{1}{2n} \ln(2\pi)$ ne pose pas de problèmes. Il tend vers zéro.
- (2) Si nous remplaçons k_n par $n(p + \epsilon) + \sigma(n)$ nous voyons que ce qui est dans le logarithme est majoré par $\frac{1}{P(n)}$ pour un certain polynôme P . Ce terme est dans le cas $\frac{\ln(P(n))}{n}$ qui tend vers zéro lorsque $n \rightarrow \infty$.
- (3) Pour ce terme nous remplaçons k_n par $n(p + \epsilon) + k_n - n(p + \epsilon)$. Nous devons alors étudier la limite de

$$(p + \epsilon) \ln \left(\frac{np}{k_n} \right) + \frac{k_n - n(p + \epsilon)}{n} \ln \left(\frac{np}{k_n} \right). \quad (37.458)$$

Ce qui est dans les logarithmes est encadré de la façon suivante :

$$\frac{n(p + \epsilon)}{np} \leq \frac{k_n}{np} \leq \frac{n(p + \epsilon) + 1}{np}. \quad (37.459)$$

Donc la limite de k_n/np est $(p + \epsilon)/p$. Les logarithmes restent bornés. Pour le second terme de (37.458), le numérateur du coefficient est borné par 1. Donc le second terme tend vers zéro et le tout tend vers

$$(p + \epsilon) \ln \left(\frac{p}{p + \epsilon} \right). \quad (37.460)$$

- (4) Nous devons enfin étudier le dernier terme. La combinaison $\frac{n - k_n}{n}$ s'étudie de la façon suivante :

$$\frac{n - k_n}{n} = \frac{n - n(p + \epsilon) + n(p + \epsilon) - k_n}{n} = \frac{n(1 - p - \epsilon) + n(p + \epsilon) - k_n}{n} \rightarrow 1 - p - \epsilon \quad (37.461)$$

parce que $n(p + \epsilon) - k_n$ est borné par 1. Sachant cela, notre terme a pour limite

$$\frac{n - k_n}{n} \ln \left(\frac{n(1 - p)}{n - k_n} \right) \rightarrow (1 - p - \epsilon) \ln \left(\frac{1 - p}{1 - p - \epsilon} \right). \quad (37.462)$$

En remettant tous les morceaux bouts à bout,

$$\frac{1}{n} P(S_n = k_n) \rightarrow (1 + \epsilon) \ln \left(\frac{p}{p + \epsilon} \right) + (1 - p - \epsilon) \ln \left(\frac{1 - p}{1 - p - \epsilon} \right) = -h_+(\epsilon). \quad (37.463)$$

Étant donné que nous avons déjà prouvé que $P\left(\frac{S_n}{n} \geq p + \epsilon\right) \geq P(S_n = k_n)$, nous avons

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \ln \left(P\left(\frac{S_n}{n} \geq p + \epsilon\right) \right) \geq \lim_{n \rightarrow \infty} \frac{1}{n} P(S_n = k_n) = -h_+(\epsilon). \quad (37.464)$$

En combinant avec (37.452) nous trouvons que

$$\lim_{n \rightarrow \infty} \frac{1}{n} \ln \left(P\left(\frac{S_n}{n} \geq p + \epsilon\right) \right) = h_+(\epsilon). \quad (37.465)$$

Pour cette dernière déduction nous utilisons le fait que si (a_n) est une suite telle que $a_n \leq l$ et $\liminf_{n \rightarrow \infty} a_n = l$, alors (a_n) admet une limite qui vaut l . \square

37.7 Simulations de réalisations de variables aléatoires

Le générateur de base que possède un système informatique est un générateur de nombres pseudo-aléatoires de nombres entiers entre 0 et $m - 1$ généré par une suite du type

$$x_{n+1} = (ax_n + b) \pmod{m}. \quad (37.466)$$

37.7.1 Générateur uniforme

37.7.1.1 Première méthode

Une première façon de générer une variable aléatoire de loi uniforme sur $]0, 1[$ est de diviser par $m - 1$ la suite de x_n . En effet nous avons la proposition suivante.

Proposition 37.118.

Si (Y_n) est une suite de variables aléatoires indépendantes et identiquement distribuées de loi uniforme sur $\{0, \dots, m - 1\}$. Alors

$$\frac{Y_n}{m} \xrightarrow{\mathcal{L}} \mathcal{U}[0, 1]. \quad (37.467)$$

Démonstration. Nous prouvons la convergence en loi en passant par la fonction de répartition et la proposition 37.69. La fonction de répartition de la densité $\mathcal{U}[0, 1]$ est

$$F_X(x) = x \mathbb{1}_{[0,1]}. \quad (37.468)$$

La fonction de répartition de la variable aléatoire (discrète) $X_n = \frac{Y_n}{m}$ est

$$P\left(\frac{Y_n}{m} \leq x\right) = P(Y_n \leq mx) = \frac{\lfloor mx \rfloor}{m} \quad (37.469)$$

où $\lfloor a \rfloor$ est le plus grand entier inférieur à a . Nous avons évidemment

$$\lim_{n \rightarrow \infty} \frac{\lfloor mx \rfloor}{m} = x, \quad (37.470)$$

ce qui montre la convergence des fonctions de répartitions et donc la convergence en loi qui nous intéresse. \square

37.7.1.2 Seconde méthode

Soit (Y_n) une suite de variables aléatoires indépendantes et identiquement distribuées selon la loi $Y_n \sim \mathcal{U}\{0, \dots, m - 1\}$. Alors la série de variables aléatoires

$$Z = \sum_{k=0}^{\infty} \frac{Y_k}{m^{k+1}} \quad (37.471)$$

est une série qui converge presque sûrement parce que Y_k est borné par m . Avec probabilité zéro nous avons $Z = \sum_k 1/m^k$ qui converge. Nous avons

$$Z \sim \mathcal{U}[0, 1]. \quad (37.472)$$

L'argument pour montrer cette loi est qu'en base m , la variable aléatoire Z a un développement décimal $Z = 0.Y_1Y_2Y_3Y_4 \dots$.

37.7.2 Simulation par inversion

Nous cherchons maintenant à simuler une loi X de fonction de répartition F .

Définition 37.119.

Soit $f: \mathbb{R} \rightarrow [0, 1]$ une fonction croissante, continue à droite et telle que

$$\lim_{x \rightarrow -\infty} f(x) = 0 \quad (37.473a)$$

$$\lim_{x \rightarrow \infty} f(x) = 1. \quad (37.473b)$$

L'inverse généralisé de f , notée f^{-1} est la fonction définie par

$$f^{-1}(t) = \inf\{x \text{ tel que } f(x) \geq t\}. \quad (37.474)$$

Remarque 37.120.

L'inverse généralisé d'une fonction bijective est la vraie fonction réciproque usuelle.

Proposition 37.121.

Soit f une fonction admettant un inverse généralisé f^{-1} . Alors nous avons $f^{-1}(t) \leq a$ si et seulement si $t \leq f(a)$.

La continuité à droite joue pour démontrer cette proposition.

Proposition 37.122.

Si F est la fonction de répartition de la variable aléatoire X et si V est une variable aléatoire de loi uniforme $\mathcal{U}[0, 1]$, alors $F^{-1}(U)$ a la même loi que X .

Démonstration. Nous montrons que les fonctions de répartition de X et de $F^{-1}(U)$ sont identiques. En utilisant la proposition 37.121, nous avons

$$P(F^{-1}(U) \leq y) = P(U \leq F(y)) \quad (37.475a)$$

$$= F(y) \quad (37.475b)$$

$$= P(X \leq y). \quad (37.475c)$$

Donc $F^{-1}(U)$ est la fonction de répartition de X . □

La difficulté de la méthode par inversion est qu'il faut être capable de calculer l'inverse de la fonction de répartition de la loi à simuler.

37.7.2.1 Loi exponentielle

La loi exponentielle est une loi qui peut être simulée par inversion. La fonction de répartition vaut

$$F(x) = 1 - e^{-\lambda x}, \quad (37.476)$$

et l'inverse vaut

$$F^{-1}(x) = -\frac{1}{\lambda} \ln(1 - y). \quad (37.477)$$

Par conséquent, une bonne formule pour simuler une loi exponentielle est

$$-\frac{1}{\lambda} \ln(1 - U). \quad (37.478)$$

Notez que U étant uniforme, nous pouvons tout autant prendre $-\ln(U)/\lambda$.

37.7.3 Algorithme de Box-Muller

Il s'agit de simuler une loi gaussienne. La proposition est la suivante.

Proposition 37.123.

Si U et V sont des variables aléatoires indépendantes de même loi uniforme sur $[0, 1]$, alors le couple

$$(X, Y) = (\sqrt{-2 \ln(U)} \cos(2\pi V), \sqrt{-2 \ln(U)} \sin(2\pi V)) \quad (37.479)$$

vérifie

- (1) X est indépendante de Y
 (2) X et Y sont de loi $\mathcal{N}(0, 1)$.

Démonstration. Nous allons montrer la proposition en utilisant les fonctions tests. Soit donc $\varphi: \mathbb{R}^2 \rightarrow \mathbb{R}$ une fonction bornée et mesurable. Soient Z et W , deux variables aléatoires indépendantes de loi $\mathcal{N}(0, 1)$. Nous allons montrer que

$$F(\varphi(X, Y)) = E\left[\varphi(\sqrt{-2\ln(U)} \cos(2\pi V), \sqrt{-2\ln(U)} \sin(2\pi V))\right]. \quad (37.480)$$

Par indépendance de U et V , la densité du couple est le produit des densités, donc en passant aux coordonnées polaires,

$$\diamond = E[\varphi(Z, W)] \quad (37.481a)$$

$$= \frac{1}{2\pi} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \varphi(x, y) e^{-x^2/2} e^{-y^2/2} dx dy \quad (37.481b)$$

$$= \frac{1}{2\pi} \int_0^{2\pi} d\theta \int_0^{\infty} \varphi(r \cos \theta, r \sin \theta) e^{-r^2/2} r dr. \quad (37.481c)$$

Nous posons $u = e^{-r^2/2}$ et $v = \frac{\theta}{2\pi}$. En particulier $r = \sqrt{-2\ln(u)}$ et

$$\diamond = \int_0^1 \int_0^1 \varphi(\sqrt{-2\ln(u)} \cos(2\pi v), \sqrt{-2\ln(u)} \sin(2\pi v)) dudv \quad (37.482a)$$

$$= E\left(\varphi(\sqrt{-2\ln(U)} \cos(2\pi V), \sqrt{-2\ln(U)} \sin(2\pi V))\right) \quad (37.482b)$$

parce que mesure $dudv$ est la densité de la loi uniforme. \square

En pratique, la formule

$$(x, y) \mapsto (\sqrt{-2\ln x} \cos(2\pi y), \sqrt{-2\ln x} \sin(2\pi y)) \quad (37.483)$$

est une façon d'obtenir deux gaussiennes à partir de deux variables uniformes.

37.7.4 Méthode du rejet

La méthode du rejet permet de simuler des lois à densité. Soit f la densité de la loi à simuler. Nous faisons les hypothèses suivantes.

- (1) Il existe une densité g d'une variable aléatoire facile à simuler.
 (2) Il existe un $k \geq 0$ tel que $f(x) \leq kg(x)$.

Remarque 37.124.

Le k de la seconde hypothèse est nécessairement plus grand que 1. En effet,

$$1 = \int f \leq k \int g = k \quad (37.484)$$

parce que f et g sont des densités et ont donc une intégrale égale à 1.

Proposition 37.125.

Soient (X_n) et (U_n) des suites de variables aléatoires indépendantes au sens où non seulement les X_i et U_k sont indépendants entre eux, mais de plus X_i est indépendant de U_j pour tout i et j . Nous supposons que les X_i sont indépendantes et identiquement distribuées, de densité g et que les U_i sont indépendantes et identiquement distribuées de loi uniforme.

Nous introduisons la variable aléatoire à valeurs dans \mathbb{N}

$$p(\omega) = \inf\{n \geq 0 \text{ tel que } \alpha(X_n(\omega)) \geq U_n(\omega)\} \quad (37.485)$$

où α est la fonction définie par

$$\alpha(x) = \begin{cases} \frac{f(x)}{kg(x)} & \text{si } g(x) \neq 0 \\ 0 & \text{si } g(x) = 0. \end{cases} \quad (37.486)$$

Alors la variable aléatoire Y définie par

$$Y(\omega) = X_{p(\omega)}(\omega) \quad (37.487)$$

admet f pour densité.

Démonstration. D'abord étant donné que $f(x) \leq kg(x)$ nous avons $\alpha(x) \in [0, 1]$. Nous pouvons a priori avoir $p(\omega) = \infty$, ce qui rendrait caduque la définition de $Y(\omega)$. Montrons donc pour commencer que $P(p = \infty) = 0$. En utilisant l'indépendance nous avons

$$P(\alpha(X_n) < U_n, \forall n) = \lim_{N \rightarrow \infty} \prod_{i=1}^N P(\alpha(X_i) < U_i) \quad (37.488a)$$

$$= \lim_{N \rightarrow \infty} P(\alpha(X_1) < U_1)^N. \quad (37.488b)$$

Pour conclure nous devons prouver que $P(\alpha(X_1) < U_1) < 1$. Pour cela nous calculons

$$P(\alpha(X_1) < U_1) = \int_{\mathbb{R}} dx \int_0^1 du \mathbb{1}_{\alpha(x) < u} g(x) \quad (37.489a)$$

$$= \int_{\mathbb{R}} g(x)(1 - \alpha(x)) dx \quad (37.489b)$$

$$= \int_{\mathbb{R}} \left(g(x) - \frac{f(x)}{k} \right) \quad (37.489c)$$

$$= 1 - \frac{1}{k} \quad (37.489d)$$

$$< 1. \quad (37.489e)$$

L'équation (37.488) nous permet donc de conclure que $P(\alpha(X_n) < U_n, \forall n) = 0$. Par conséquent la variable aléatoire $Y(\omega) = X_{p(\omega)}(\omega)$ a un sens.

Nous devons maintenant prouver que Y a bien f pour densité. Pour cela nous considérons un ensemble mesurable $A \in \mathcal{B}(\mathbb{R})$ et nous montrons que $P(Y \in A) = \int_A f(x) dx$. Nous avons

$$P(Y \in A) = P(X_p \in A) = \sum_{j=1}^{\infty} P(X_j \in A, p = j). \quad (37.490)$$

Par ailleurs nous avons

$$\begin{aligned} P(X_j \in A, p = j) &= P(X_j \in A, \alpha(X_j) \geq U_j, \alpha(X_m) < U_m \forall m \leq j-1) \\ &= P(X_j \in A, \alpha(X_j) \geq U_j) P(\alpha(X_1) < U_1)^{j-1} \\ &= P(X_j \in A, \alpha(X_j) \geq U_j) \left(1 - \frac{1}{k}\right)^{j-1}. \end{aligned} \quad (37.491)$$

Étant donné que $g(x) dx$ est la densité de X_j et que du est la densité de U , nous avons

$$P(X_j \in A, \alpha(X_j) \geq U_j) = \int_{\mathbb{R}} \int_0^1 \mathbb{1}_{x \in A} \mathbb{1}_{\alpha(x) \geq u} g(x) du dx \quad (37.492a)$$

$$= \int_{\mathbb{R}} g(x) \mathbb{1}_{x \in A} \underbrace{\int_0^1 \mathbb{1}_{\alpha(x) \geq u} du}_{\alpha(x)} dx. \quad (37.492b)$$

$$= \int_{\mathbb{R}} \mathbb{1}_{x \in A} \frac{f(x)}{k} dx \quad (37.492c)$$

$$= \frac{1}{k} P(X \in A). \quad (37.492d)$$

En remplaçant dans l'équation (37.491) nous trouvons

$$P(X_j \in A, p = j) = \frac{1}{k} P(X \in A) \left(1 - \frac{1}{k}\right)^{j-1}. \quad (37.493)$$

Et enfin, l'équation (37.490) donne

$$P(Y \in A) = \frac{1}{k} P(X \in A) \sum_{i=1}^{\infty} \left(1 - \frac{1}{k}\right)^{i-1} = P(X \in A). \quad (37.494)$$

□

37.7.5 Simuler une loi géométrique à l'ordinateur

Si (X_n) est une suite de variables aléatoires indépendantes et identiquement distribuées avec $X_n \sim \mathcal{B}(p)$, alors

$$Z = \min\{k \geq 1 \text{ tel que } X_k = 1\} \sim \mathcal{G}(p). \quad (37.495)$$

Nous avons alors $P(Z = k) = (1 - p)^k p$.

Si nous avons un générateur de lois de Bernoulli de paramètre p , alors nous on simulons jusqu'à obtenir 1 et nous comptons combien de simulations ont été nécessaires.

37.7.6 Simuler une loi exponentielle à l'ordinateur

Nous pouvons utiliser la méthode de l'inversion. Étant donné que la fonction de répartition de la loi exponentielle est $F(x) = 1 - e^{-\lambda x}$, nous avons $F^{-1}(y) = \frac{1}{\lambda} \ln(1 - y)$. Par conséquent à partir d'un générateur uniforme U , nous pouvons calculer

$$F^{-1}(U) = \frac{1}{\lambda} \ln(U) \quad (37.496)$$

qui suivra une loi exponentielle d'espérance $1/\lambda$.

37.7.7 Simuler une loi de Poisson à l'ordinateur

Nous savons du point 37.5.8 que si les T_i sont des variables aléatoires indépendantes et identiquement distribuées de loi $\mathcal{E}(\lambda)$, alors nous avons

$$\max\{n \geq 1 \text{ tel que } \sum_i T_i \leq 1\} \sim \mathcal{P}(\lambda). \quad (37.497)$$

La façon usuelle pour créer une loi exponentielle est d'avoir un générateur de loi uniforme U_i et d'écrire que

$$-\frac{1}{\lambda} \ln(U_i) \sim \mathcal{E}(\lambda). \quad (37.498)$$

Nous devons donc faire la somme de telles variables aléatoires et voir à partir de quel moment la somme dépasse 1. Le calcul est le suivant :

$$-\sum_{i=1}^n \frac{1}{\lambda} \ln(U_i) \leq 1 \quad (37.499)$$

implique

$$\prod_{i=1}^n U_i \leq e^{-\lambda}. \quad (37.500)$$

En pratique, la variable aléatoire qui se comporte comme une loi de Poisson de paramètre λ est

$$N = \max\{n \geq 1 \text{ tel que } \prod_{i=1}^n U_i \geq e^{-\lambda}\}. \quad (37.501)$$

Nous générons donc des nombres aléatoires entre 1 et 1 et nous effectuons le produit jusqu'à ce qu'il passe en dessous de $e^{-\lambda}$. À ce moment, nous retournons le nombre de nombres qu'il a fallu générer.

37.8 Sage

Nous allons montrer maintenant quelques trucs importants dans l'utilisation de Sage pour réaliser des petits graphiques.

Remarque 37.126.

Dans ce qui suit, nous allons parler de « Sage », mais en réalité nous allons surtout parler du module `scipy` qui fait partie des modules hyper-usuels de Python. Les remerciements vont donc au moins autant du côté de l'équipe de `scipy` que vers celle de Sage.

37.8.1 Loi exponentielle

Il faut savoir que la définition d'une loi continue retourne automatiquement la loi centrée réduite. Pour avoir une loi exponentielle de moyenne donnée, il faut donc préciser de façon plus maligne que ce que l'on croit.

```

1 from scipy import stats
2
3 X=stats.expon(scale=5)
4 print(X.mean())      # retourne 5
5
6 P=plot( X.pdf ,x,0,10 )
7 show(P)              # Affiche le graphique

```

tex/frido/code_sage1.py

37.8.2 Inverser des lois

Pour trouver des intervalles de confiance, il faut souvent calculer des inverses de loi. Bien entendu Sage le fait. Ce que sage connaît, c'est l'inverse de la fonction de survie. Autrement dit si X est une variable aléatoire, $X.sf$ est la fonction $x \mapsto 1 - P(X < x)$ et $X.isf$ en est l'inverse. Pour résoudre $P(X < \xi) = \alpha$, il faut résoudre $F(\xi) = \alpha$, c'est-à-dire

$$1 - F(\xi) = 1 - \alpha, \quad (37.502)$$

ce qui se fait de la façon suivante : le programme suivant donne pour une loi normale centrée réduite la valeur de ξ pour laquelle $P(N < \xi) = 0.05$:

```

1 from scipy import stats
2
3 N=stats.norm
4 print N.mean()      # 0
5 print N.var()       # 1
6
7 xi = N.isf(0.95)
8 print xi            # -1.64485
9
10 N.cdf(xi)          # Vérification : 0.05
11
12 # Graphiques de la fonction de densité et la cumulative.
13 P=plot(N.cdf ,x,-10,10)
14 Q=plot(N.pdf ,x,-10,10,color="red")
15 show(P+Q)

```

tex/frido/code_sage2.py

37.9 Monte-Carlo

Nous voudrions calculer une valeur approchée de l'intégrale

$$I = \int_a^b f(x)dx. \quad (37.503)$$

Les méthodes classiques consistent à discrétiser l'intervalle $[a, b]$ et en calculant une somme de la forme $\sum_i w_i f(x_i)$.

L'idée de Monte Carlo est de remplacer le découpage déterministe x_i par des variables aléatoires X_i en trois étapes.

- (1) Pour cela nous commençons par écrire l'intégrale comme une espérance : $I = E(X)$ où X est une variable aléatoire définie sur un espace probabilisé (Ω, \mathcal{F}, P) à déterminer. Une contrainte est évidemment d'avoir $X \in L^1(\Omega, \mathcal{F}, P)$.
- (2) Nous générons une suite indépendantes et identiquement distribuée de variables aléatoires (X_n) de même loi que X et la loi (forte) des grands nombres implique que

$$\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k \xrightarrow{p.s.} E(X) = I. \quad (37.504)$$

- (3) Le dernier point sera de donner un intervalle de confiance.

Exemple 37.127

Nous voudrions déterminer de façon approchée l'intégrale $I = \int_0^1 f(x)dx$. Si $U \sim \mathcal{U}[0, 1]$, alors

$$I = E(f(U)) \quad (37.505)$$

et il suffit de faire

$$I \simeq \frac{1}{n} \sum_{k=1}^n f(U_k). \quad (37.506)$$

où mes U_i sont indépendantes et identiquement distribuées de loi $\mathcal{U}[0, 1]$. △

Exemple 37.128

Supposons que la fonction à intégrer se présente sous la forme $f(x) = h(x)g(x)$ avec $g \geq 0$ et telle que l'intégrale $\int_{\mathbb{R}} g$ existe. Notons

$$c = \int_{\mathbb{R}} g \quad (37.507)$$

et

$$I = \int_{\mathbb{R}} h(x)c \frac{g(x)}{c} dx. \quad (37.508)$$

Nous avons alors $I = E(ch(Y))$ où Y admet la densité $g(x)c$. △

Passons au cas de plusieurs variables et considérons l'intégrale

$$I = \int_{[0,1]^d} f(x_1, \dots, x_d) dx_1 \dots dx_d. \quad (37.509)$$

Nous écrivons

$$I = E(f(U_1, \dots, U_d)) \quad (37.510)$$

où les U_i sont de loi uniformes sur $[0, 1]$. En pratique, nous générons une suite de variables aléatoires de (Z_k) de lois uniformes que nous regroupons par paquets :

$$V_k = (Z_{dk}, Z_{dk+1}, \dots, Z_{d(k+1)-1}). \quad (37.511)$$

Ces variables aléatoires V_k sont indépendantes et identiquement distribuées de loi $\mathcal{U}[0, 1]^d$. Ensuite la loi des grands nombres nous indique que

$$I \sim \frac{1}{n} \sum_{k=1}^n f(V_k). \quad (37.512)$$

37.9.1 Intervalle de confiance

37.9.1.1 Principe

Nous supposons que nous travaillons sur une approximation de Monte-Carlo telle que la variable aléatoire choisie soit dans L^2 . La loi des grands nombres nous dit que $\bar{X}_n \sim I$ tandis que le théorème central limite nous enseigne que

$$\frac{\bar{X}_n - E(X)}{\sigma/\sqrt{n}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1). \quad (37.513)$$

Par conséquent

$$P\left(\frac{\bar{X}_n - E(X)}{\sigma/\sqrt{n}} \in [-u, u]\right) \simeq P(-u \leq Z \leq u) \quad (37.514)$$

où $Z \sim \mathcal{N}(0, 1)$. En remplaçant $E(X)$ par I et en effectuant les manipulations usuelles, nous trouvons que $P(I \in J_\alpha) = 1 - \alpha$ si

$$J_\alpha = \left[\bar{X}_n - u_\alpha \frac{\sigma}{\sqrt{n}}, \bar{X}_n + u_\alpha \frac{\sigma}{\sqrt{n}}\right] \quad (37.515)$$

où σ^2 est la variance de X . Si σ n'est pas connue, alors nous le remplaçons par un estimateur

$$S'_n = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2 \quad (37.516)$$

et nous considérons l'intervalle

$$J'_\alpha = \left[\bar{X}_n - u_\alpha \frac{S'_n}{\sqrt{n}}, \bar{X}_n + u_\alpha \frac{S'_n}{\sqrt{n}}\right]. \quad (37.517)$$

Il y a deux façons de faire diminuer la longueur de l'intervalle de confiance : augmenter n ou diminuer σ . Pour le second point, le choix de X dans $I = E(X)$ est essentiel.

Exemple 37.129

Soit à calculer

$$I = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{\beta z} e^{-z^2/2} dz \quad (37.518)$$

avec $\beta > 0$. Nous introduisons la variable aléatoire $X = e^{\beta Z}$ avec $Z \sim \mathcal{N}(0, 1)$. Nous avons alors

$$I = E(X). \quad (37.519)$$

Par ailleurs l'intégrale demandée vaut $e^{\beta^2/2}$. En appliquant les formules vues plus haut nous trouvons

$$J_\alpha = \left[\bar{X}_n - u_\alpha \frac{\sigma}{\sqrt{n}}, \bar{X}_n + u_\alpha \frac{\sigma}{\sqrt{n}}\right] \quad (37.520)$$

où

$$\sigma^2 = \text{Var}(X) = E(e^{2\beta Z}) - E(e^{\beta Z})^2 = e^{2\beta^2} - e^{\beta^2}. \quad (37.521)$$

Nous avons utilisé la formule (37.386). Si nous choisissons $\beta = 2$, nous trouvons $\sigma^2 \simeq 2926$. Donc si nous voulons une longueur de J_α plus petite que 10^{-2} tout en demandant $\alpha = 0.05$ (ce qui implique $u_\alpha = 1.96$), nous devons avoir

$$1.96 \frac{2973}{\sqrt{n}} < 10^{-2}, \quad (37.522)$$

c'est-à-dire environ $n = 10^{11}$, ce qui soit dit en passant est très largement au delà des capacités de la commande `rand` de scilab. \triangle

Nous allons maintenant voir quelques méthodes pour réduire la variance.

37.9.1.2 Échantillonnage préférentiel

Nous devons calculer $I = \int_{\mathbb{R}} f(x)dx$. Pour cela nous introduisons artificiellement une densité g et nous écrivons

$$I = \int_{\mathbb{R}} \frac{f(x)}{g(x)} g(x) dx = E \left(\frac{f(Y)}{g(Y)} \right) \quad (37.523)$$

où Y est de densité g . Il faut essayer de trouver g de telle sorte à ce que

$$\text{Var} \left(\frac{f(Y)}{g(Y)} \right) \quad (37.524)$$

soit la plus petite possible.

37.9.1.3 Méthode de la variable de contrôle

Soit $I = E(X)$. Nous introduisons une variable aléatoire Z et nous écrivons

$$I = E(X - Z) + E(Z). \quad (37.525)$$

Il faut alors choisir Z de telle sorte que $E(Z)$ soit calculable et que $X - Z$ ait une variance plus faible. En particulier, Z ne peut pas être indépendante de X .

37.9.1.4 Variables antithétiques

Soit $I = \int_0^1 f(x)dx$. La première idée (exemple 37.127) est d'écrire

$$I = E(f(U)) \quad (37.526)$$

où $U \sim \mathcal{U}[0, 1]$, mais nous n'avons pas de garanties sur la variance de $f(U)$. Nous pouvons écrire

$$I = E(f(U)) = E \left[\frac{1}{2} (f(U) - f(1 - U)) \right]. \quad (37.527)$$

Ici $1 - U$ est encore une variable aléatoire uniforme sur $[0, 1]$, mais il se fait que la variable aléatoire

$$Z = \frac{1}{2} (f(U) + f(1 - U)) \quad (37.528)$$

a une variance inférieure à $\text{Var}(f(U))$. En effet, $f(U)$ et $f(1 - U)$ ne sont pas indépendantes, par conséquent le résultat du lemme 37.23 n'est pas valide, par contre la proposition 37.28 reste vraie et nous avons

$$\text{Var}(Z) = \frac{1}{4} \text{Var}(f(U)) + \frac{1}{4} \text{Var}(f(1 - U)) + \frac{1}{2} \text{Cov}(f(U), f(1 - U)). \quad (37.529)$$

Nous avons $\text{Var}(f(1 - U)) = \text{Var}(f(U))$. En ce qui concerne le terme avec la covariance, nous lui appliquons l'équation (37.56) :

$$\text{Cov}(f(U), f(1 - U)) = E \left((f(U) - I)(f(1 - U) - I) \right) \quad (37.530a)$$

$$\leq E((f(U) - I)^2)^{1/2} E((f(1 - U) - I)^2)^{1/2} \quad (37.530b)$$

$$= \text{Var}(f(U)) \quad (37.530c)$$

où nous avons utilisé le fait que $E(f(U)) = E(f(1 - U)) = I$. Au final nous avons bien obtenu

$$\text{Var}(Z) \leq \text{Var}(f(U)). \quad (37.531)$$

37.10 Résultats qui se démontrent avec des variables aléatoires

37.10.1 Nombres normaux

Tout nombre $x \in [0, 1[$ admet un unique¹⁶ développement en base $b \geq 2$:

$$x = \sum_{n=1}^{\infty} \frac{\epsilon_n(x)}{b^n} \quad (37.532)$$

avec $\epsilon_n(x) \in \mathcal{A} = \{0, \dots, b-1\}$.

Soit $k \geq 1$ et $r \in \mathcal{A}^k$; nous posons

$$N_x(r, n) = \text{Card} \{i \in \{1, \dots, n-k+1\} \text{ tel que } \epsilon_1(x) = r_1, \dots, \epsilon_{i+k-1}(x) = r_k\}. \quad (37.533)$$

C'est le nombre d'occurrences du motif r (de longueur k) dans les n premières décimales de x .

Définition 37.130.

Un nombre $x \in [0, 1[$ est **normal** en base b si pour tout $r \in \{0, \dots, b-1\}^k$ nous avons

$$\frac{N_x(b, n)}{n} \rightarrow \frac{1}{b^k}. \quad (37.534)$$

Un nombre est normal s'il est normal en toute base.

Proposition 37.131 ([495, 43]).

Au sens de la mesure de Lebesgue, presque tous les nombres de $[0, 1[$ sont normaux.

Démonstration. Pour $x \in [0, 1[$, nous notons $\epsilon_n(x)$ son développement en base b . Cela nous donne des variables aléatoires $\epsilon_i: [0, 1[\rightarrow \mathcal{A}$ dont la loi de probabilité est donnée par

$$P(\epsilon_1 = d) = P\left(\left[\frac{d}{b}, \frac{d+1}{b}\right[\right) = \frac{1}{b} \quad (37.535)$$

parce que l'intervalle $\left[\frac{d}{b}, \frac{d+1}{b}\right[$ est l'ensemble des nombres de $[0, 1[$ dont la première décimale est d . Pour la loi des ϵ_i , il faut un peu plus découper, mais ça donne le même résultat : $P(\epsilon_i = d) = 1/b$. Ces variables aléatoires sont indépendantes et identiquement distribuées. Nous considérons aussi la variable aléatoire

$$\begin{aligned} N(r, n): [0, 1[&\rightarrow \mathbb{N} \\ x &\mapsto N_x(r, n) \end{aligned} \quad (37.536)$$

Pour un $r \in \mathcal{A}$ fixé, nous définissons encore la variable aléatoire

$$\begin{aligned} X_j: \mathcal{A} &\rightarrow \{0, 1\} \\ x &\mapsto \begin{cases} 1 & \text{si } \epsilon_j(x) = b \\ 0 & \text{sinon.} \end{cases} \end{aligned} \quad (37.537)$$

Les variables aléatoires X_j sont des variables aléatoires de Bernoulli indépendantes et identiquement distribuées de paramètre $E(X_j) = P(X_j = 1) = P(\epsilon_1 = b) = \frac{1}{b}$. Nous pouvons utiliser dessus la loi forte des grands nombres (théorème 37.83). Pour dire que

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p.s.} E(X_1) = \frac{1}{b}. \quad (37.538)$$

Mais en réalité nous avons aussi $\sum_{j=1}^n X_j = N(r, n)$ parce que en appliquant à $x \in [0, 1[$:

$$\sum_{j=1}^n \begin{cases} 1 & \text{si } \epsilon_j(x) = r \\ 0 & \text{sinon} \end{cases} = \text{Card} \{i \in \{1, \dots, n\} \text{ tel que } \epsilon_i(x) = r\} = N_x(r, n), \quad (37.539)$$

16. Nous excluons 1 parce que son développement en puissances négatives de b est zéro.

de sorte que l'équation (37.538) nous dit exactement que pour tout $r \in \mathcal{A}$,

$$\lim_{n \rightarrow \infty} \frac{N_x(r, n)}{n} = \frac{1}{b} \quad (37.540)$$

pour presque tout $x \in [0, 1[$.

Il reste à prouver la même chose pour tout $r \in \mathcal{A}^k$. Voyons avec $k = 2$ et $r = (u, v) \in \mathcal{A}^2$. Nous posons

$$Y_j = \mathbb{1}_{\{\epsilon_j = u, \epsilon_{j+1} = v\}}, \quad (37.541)$$

et $N(r, n) = \sum_{j=1}^{n-1} Y_j$. Les Y_i sont encore des binomiales de paramètre $\frac{1}{b^2}$, mais elles ne sont pas indépendantes. En effet pour avoir $Y_1(x) = Y_2(x) = 1$, il faut que les trois premières décimales de x soit en même temps de la forme uv . et $.uv$, donc

$$P(Y_1, Y_2 = 1) = b^3 \delta_{u,v} \quad (37.542)$$

alors que $P(Y_1 = 1)P(Y_2 = 2) = b^4$. Nous pouvons contourner ce problème en remarquant que les ϵ_i , eux, sont indépendants. Donc le lemme de regroupement 37.16 nous dit que la famille $\{Y_{2n}\}$ est une famille de variables aléatoires indépendantes (et idem pour la famille Y_{2n-1}). En effet, les variables aléatoires Y_{2n} correspondent à la partition 23, 45, 67, etc.

Nous appliquons la loi des grands nombres sur les deux familles indépendamment :

$$\frac{1}{n} \sum_{j=1}^n Y_{2j-1} \xrightarrow{p.s.} \frac{1}{b^2} \quad (37.543)$$

et

$$\frac{1}{n} \sum_{j=1}^n Y_{2j} \xrightarrow{p.s.} \frac{1}{b^2} \quad (37.544)$$

Pour rappel, le but pour l'instant est d'établir la limite $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n Y_j = \frac{1}{b^2}$. Nous allons l'établir séparément pour les termes pairs et impairs de la suite. Pour les pairs :

$$\frac{1}{2n} \sum_{j=1}^{2n} Y_j = \frac{1}{2} \left(\frac{1}{n} \sum_{j=1}^n Y_{2j-1} \right) + \frac{1}{2} \left(\frac{1}{n} \sum_{j=1}^n Y_{2j} \right) \xrightarrow{p.s.} \frac{1}{b^2} \quad (37.545)$$

Pour les impairs¹⁷,

$$\frac{1}{2n-1} \sum_{j=1}^{2n-1} Y_j = \frac{n}{2n-1} \left(\frac{1}{n} \sum_{j=1}^n Y_{2j-1} \right) + \frac{n}{2n-1} \left(\frac{1}{n} \sum_{j=1}^n Y_{2j} \right) \rightarrow \frac{1}{b^2} \quad (37.546)$$

parce que les deux parenthèses convergent vers $\frac{1}{b^2}$ alors que les coefficients devant convergent vers $\frac{1}{2}$.

Au final nous avons bien

$$\frac{N(r, n)}{n} = \frac{1}{n} \sum_{j=1}^{n-1} Y_j = \frac{n-1}{n} \left(\frac{1}{n-1} \sum_{j=1}^{n-1} Y_j \right) \rightarrow \frac{1}{b^2} \quad (37.547)$$

tant que $r \in \mathcal{A}^2$.

Pour prouver la même chose avec $r \in \mathcal{A}^k$, il suffit de faire le même raisonnement en divisant en plus de paquets : $\{Y_{kj+m}\}_{m=1, \dots, k-1}$ sont indépendants et nous utilisons k fois la loi des grands nombres.

Donc pour toute base b nous savons que les nombres normaux en base b forment un ensemble de mesure nulle dans $[0, 1[$. Il reste à voir que leur union reste de mesure nulle. Cela est vrai parce que nous avons une union dénombrable et qu'une union dénombrable d'ensembles de mesure nulle est de mesure nulle par le lemme 15.27. \square

17. Ici dans [43], la seconde somme va jusqu'à $n-1$ et je ne comprends pas pourquoi.

Remarque 37.132.

Un nombre x est normal en base b si et seulement si la suite $u_k = xb^k$ est équirépartie modulo 1 sur $[0, 1]$ (c'est-à-dire que la suite des parties fractionnelles des u_k est équirépartie). Pour le nombre 0.2357873..., nous parlons de la suite 0.2357873...; 0.357873...; 0.57897... etc. C'est la suite des queues de suites de la suite de ses décimales¹⁸.

37.10.2 Théorème de Bernstein

Théorème 37.133 (Théorème de Bernstein[43]).

Soit $f \in C^0([0, 1], \mathbb{C})$ et son module de continuité

$$\begin{aligned} \omega: [0, 1] &\rightarrow \mathbb{R} \\ h &\mapsto \sup\{|f(u) - f(v)| \text{ tel que } |u - v| < h\}. \end{aligned} \quad (37.548)$$

Pour $n \geq 0$ nous définissons le n^e **polynôme de Bernstein** de f par

$$B_n(f)(x) = \sum_{k=0}^n \binom{n}{k} x^k (1-x)^{n-k} f\left(\frac{k}{n}\right). \quad (37.549)$$

Alors il existe C tel que pour tout $n \geq 1$:

$$(1) \quad \|f - B_n(f)\|_\infty \leq C\omega\left(\frac{1}{\sqrt{n}}\right). \quad (37.550)$$

$$(2) \quad B_n(f) \xrightarrow{\text{unif}} f \quad (37.551)$$

sur $[0, 1]$.

(3) L'inégalité (37.550) est optimale : il existe une fonction $g \in C^0([0, 1], \mathbb{C})$ et $\delta > 0$ tels que pour tout $N \geq 1$, $\|g - B_N(g)\|_\infty \geq \frac{\delta}{\sqrt{N}}$. Cette fonction peut être choisie Lipschitzienne. Une telle fonction est donnée par exemple par $g(x) = |x - \frac{1}{2}|$.

(4) Les polynômes forment une partie dense dans $(C^0([0, 1]), \|\cdot\|_\infty)$.

Démonstration. Soit $x \in [0, 1]$ et une suite de variables aléatoires de Bernoulli indépendantes¹⁹ et identiquement distribuées $(X_i)_{i \geq 1}$ de paramètre x . Nous notons $S_n = \sum_{k=1}^n X_k$.

(1) Pour cette histoire de convergence, il faut majorer la quantité $|f(x) - B_n(f)(x)|$. Pour cela il y a trois astuces. La première est de se souvenir que $E(f(x)) = f(x)$, et la seconde est que le théorème de transfert 37.60 appliqué à $x \mapsto f(x/n)$ donne²⁰

$$E\left(f\left(\frac{S_n}{n}\right)\right) = \sum_{k=0}^n f\left(\frac{k}{n}\right) P(S_n = k) = \sum_{k=0}^n f\left(\frac{k}{n}\right) \binom{n}{k} x^k (1-x)^{n-k}, \quad (37.552)$$

c'est-à-dire que

$$B_n(f)(x) = E\left(f\left(\frac{S_n}{n}\right)\right). \quad (37.553)$$

Et enfin la troisième astuce est d'utiliser le lemme 12.171 pour avoir

$$\omega\left(|x - \frac{S_n}{n}|\right) = \omega\left(\frac{1}{\sqrt{n}}\left|\sqrt{n} - \frac{S_n}{\sqrt{n}}\right|\right) \leq \left(\sqrt{n}\left|x - \frac{S_n}{n}\right| + 1\right) \omega\left(\frac{1}{\sqrt{n}}\right). \quad (37.554)$$

18. C'est pas trop bien dit, mais on se comprend, non ?

19. Définition 37.9.

20. Nous avons aussi utilisé la formule de l'espérance pour les variables aléatoires discrètes.

À partir de là nous pouvons un peu calculer :

$$|f(x) - B_n(f)(x)| = \left| E \left(f(x) - f\left(\frac{S_n}{n}\right) \right) \right| \quad (37.555a)$$

$$\leq E \left(\left| f(x) - f\left(\frac{S_n}{n}\right) \right| \right) \quad (37.555b)$$

$$\leq E \left(\omega \left(\left| x - \frac{S_n}{n} \right| \right) \right) \quad (37.555c)$$

$$\leq \omega \left(\frac{1}{\sqrt{n}} \right) E \left(\left| \sqrt{n}x - \frac{S_n}{n} \right| + 1 \right). \quad (37.555d)$$

Le dernier facteur peut être réécrit sous la forme

$$E \left(\left| \sqrt{n}x - \frac{S_n}{n} \right| + 1 \right) = \sqrt{n} E \left(\left| x - \frac{S_n}{n} \right| \right) + 1, \quad (37.556)$$

et c'est là que nous pouvons utiliser l'inégalité de Hölder **28.33** :

$$E(|X|) = \|X\|_1 \leq \|X\|_2 \quad (37.557)$$

où $\|X\|_2$ désigne

$$\|X\|_2 = \sqrt{E(|X|^2)}. \quad (37.558)$$

Nous pouvons donc écrire

$$|f(x) - B_n(f)(x)| \leq \omega \left(\frac{1}{\sqrt{n}} \right) \left(\sqrt{n} \|x - \frac{S_n}{n}\|_2 + 1 \right). \quad (37.559)$$

Nous étudions maintenant de plus près la quantité $\|x - \frac{S_n}{n}\|_2$. D'abord

$$E \left(\left| x - \frac{S_n}{n} \right|^2 \right) = x^2 - 2 \frac{x}{n} E(S_n) + \frac{1}{n^2} E(S_n^2). \quad (37.560)$$

Ensuite nous savons l'espérance de S_n (qui vaut $E(S_n) = nx$) par **(37.330)** et le lemme **37.22** nous permet de calculer $E(S_n^2)$ par indépendance des X_i qui composent S_n . Nous avons alors

$$E \left(\left| x - \frac{S_n}{n} \right|^2 \right) = x^2 - 2x^2 + \frac{1}{n^2} \sum_{1 \leq i \neq j \leq n} E(X_i)E(X_j) + \frac{1}{n^2} \sum_{i=1}^n E(X_i^2) \quad (37.561a)$$

$$= -x^2 + \frac{n^2 - n}{n^2} x^2 + \frac{nx}{n^2} \quad (37.561b)$$

$$= \frac{x(1-x)}{n}. \quad (37.561c)$$

Quelques justifications :

- $E(X_i) = E(X_i^2) = x$ parce que X_i est une variable aléatoire de Bernoulli de paramètre x .
- La première somme contient tous les couples (i, j) sauf les diagonaux ; il y en a donc $n^2 - n$.

En recombinaison le tout,

$$|f(x) - B_n(f)(x)| \leq \omega \left(\frac{1}{\sqrt{n}} \right) \left(\sqrt{n} \sqrt{\frac{x(1-x)}{n}} + 1 \right) \quad (37.562a)$$

$$= \omega \left(\frac{1}{\sqrt{n}} \right) (\sqrt{x(1-x)} + 1) \quad (37.562b)$$

$$\leq \frac{3}{2} \omega \left(\frac{1}{\sqrt{n}} \right). \quad (37.562c)$$

La dernière majoration est une rapide étude de la fonction $x(1-x)$.

Étant donné que les majorations (37.562) sont valables pour tout x , en passant au supremum nous avons

$$\|f - B_n(f)\|_\infty \leq \frac{3}{2}\omega\left(\frac{1}{\sqrt{n}}\right) \rightarrow 0. \quad (37.563)$$

Ceci prouve les deux premiers points du théorème.

(2) Fait.

(3) Nous considérons la fonction

$$g(x) = \|x - \frac{1}{2}j\| \quad (37.564)$$

et nous vérifions qu'elle vérifie toutes les conditions. D'abord si $u, v \in [0, 1]$ alors

$$|g(u) - g(v)| \leq |u - v| \quad (37.565)$$

et donc $\omega(h) \leq h$, ce qui signifie que g est 1-Lipschitz. Le principe de cette partie est de montrer que $\|g - B_n(g)\|_\infty$ est plus grand que d'autres trucs (et non plus petit que d'autres trucs comme d'habitude). Nous commençons par

$$\|g - B_n(g)\|_\infty \geq g\left(\frac{1}{2}\right) - B_n(g)\left(\frac{1}{2}\right). \quad (37.566)$$

Très vite nous nous rendons compte que $g(1/2) = 0$. Ensuite nous nous souvenons que

$$B_n(g)\left(\frac{1}{2}\right) = E\left(g\left(\frac{S_n}{n}\right)\right) = E\left(\left|\frac{S_n}{n} - \frac{1}{2}\right|\right) = \frac{1}{2n}E(|2S_n - n|). \quad (37.567)$$

si nous posons $\epsilon_i = 2X_i - 1$, alors les ϵ_i sont des variables aléatoires de Rademacher indépendantes et identiquement distribuées qui satisfont à $2S_n - n = \sum_{i=1}^n \epsilon_i$. Nous utilisons la proposition 37.115 :

$$\|f - B_n(f)\|_\infty \geq \frac{1}{2n}E\left(|\sum_i \epsilon_i|\right) \geq \frac{1}{2n\sqrt{e}}\left\|\sum_{i=1}^n \epsilon_j\right\|_2. \quad (37.568)$$

Calculons ce qui est dans la norme :

$$\left\|\sum_{j=1}^n \epsilon_j\right\|_2^2 = E\left(\left(\sum_{j=1}^n \epsilon_j\right)^2\right) = \sum_{1 \leq i \neq j \leq n} E(\epsilon_i)E(\epsilon_j) + \sum_{i=1}^n E(\epsilon_i^2) = 0 + n = n. \quad (37.569)$$

Nous finissons alors notre travail de majoration :

$$\|f - B_n(f)\|_\infty \geq \frac{1}{2n\sqrt{e}}\left\|\sum_{i=1}^n \epsilon_j\right\|_2 \geq \frac{1}{2\sqrt{n}\sqrt{e}} \geq \frac{1}{2\sqrt{e}}\omega\left(\frac{1}{\sqrt{n}}\right). \quad (37.570)$$

(4) Nous avons trouvé une suite de polynômes qui converge uniformément vers un élément arbitraire de $L^2([0, 1])$. Cela prouve la densité. □

Corollaire 37.134.

Dans \mathbb{R} , si $I = [a, b]$ alors les polynômes forment une partie dense dans $(C^0(I), \|\cdot\|_\infty)$.

Démonstration. Nous supposons que $b > a$. Le cas $a = b$ est assez facile parce que l'espace des fonctions sur $\{a\}$ est de dimension 1.

Nous considérons une bijection affine $\varphi: [0, 1] \rightarrow [a, b]$ telle que $\varphi(0) = a$ et $\varphi(1) = b$. Soit $f \in C^0(I)$.

Si $g = f \circ \varphi$, alors le théorème de Bernstein 37.133 nous donne une suite de polynômes g_k sur $[0, 1]$ tels que

$$g_k \xrightarrow{\text{unif}} g. \quad (37.571)$$

Nous considérons $f_k = g_k \circ \varphi^{-1}$ qui est encore un polynôme parce que φ^{-1} est affine. Étant donné que φ^{-1} est une bijection, si h est une fonction sur $[0, 1]$, nous avons

$$\sup_{x \in [a, b]} \|(h \circ \varphi^{-1})(x)\| = \sup_{y \in [0, 1]} \|h(y)\|. \quad (37.572)$$

Cela nous permet le calcul suivant :

$$\|f_k - f\|_\infty = \|g_k \circ \varphi^{-1} - g \circ \varphi^{-1}\| \quad (37.573a)$$

$$= \|(g_k - g) \circ \varphi^{-1}\| \quad (37.573b)$$

$$= \sup_{x \in [a, b]} \|(g_k - g)(\varphi^{-1}(x))\| \quad (37.573c)$$

$$= \sup_{y \in [0, 1]} \|(g_k - g)(y)\| \quad (37.573d)$$

$$= \|g_k - g\|_\infty. \quad (37.573e)$$

Nous avons donc

$$\lim_{k \rightarrow \infty} \|f_k - f\|_\infty = 0, \quad (37.574)$$

ce qui prouve la densité. □

Chapitre 38

Statistiques

38.1 Notations et hypothèses

Nous notons X le caractère à étudier, et Ω l'ensemble des individus. Le caractère à étudier est vu comme une fonction sur Ω :

$$X: \Omega \rightarrow \mathbb{R}, \mathbb{N}, \{0, 1\}, \dots \quad (38.1)$$

Les **statistiques descriptives** sont les techniques pour présenter et résumer les données : diagrammes, graphiques, indicateurs numériques : moyenne, écart-type, médiane, ...

Nous faisons les hypothèses suivantes :

- (1) Chaque observation x_i est la réalisation de la variable aléatoire X qui sera de loi inconnue μ .
- (2) Le n -uple (x_1, \dots, x_n) est la réalisation de (X_1, \dots, X_n) qui est l'échantillon de taille n .
- (3) Les variables aléatoires X_i sont indépendantes et identiquement distribuées, de loi commune μ . La loi μ est la **loi parente** de l'échantillon.

Exemple 38.1

Un échantillon de taille 1 consisterait à tirer au sort une personne dans une population et mesurer sa taille. △

Exemple 38.2

Une échantillon de taille n consisterait à tirer au sort n personnes dans une population et de mesurer leurs tailles. △

L'**inférence statistique** est l'art de dégager des informations sur la population à partir d'informations partielles : intervalles de confiance, estimateurs, test d'hypothèses, ...

En théorie des probabilités, nous connaissons la loi de la variable aléatoire X et nous en déduisons des informations sur les réalisations de X : valeur la plus probable, moyenne, intervalle dans lequel $X(\omega)$ a le plus de chance d'appartenir. En statistique, au contraire, la loi est inconnue et nous cherchons des informations sur la loi à partir d'un échantillon de données numériques observées.

38.2 Modèle statistique

Un **modèle statistique** est un triplet

$$\mathcal{S} = \left[(\Omega, \mathcal{F}, P), (X_\theta)_{\theta \in \Theta}, (\mu_\theta)_{\theta \in \Theta} \right] \quad (38.2)$$

où (Ω, \mathcal{F}, P) est un espace probabilisé, (X_θ) est une famille de variables aléatoires définies sur Ω et telles que pour tout $\theta \in \Theta$, la variable aléatoire X_θ suit la loi μ_θ . Les μ_θ sont des mesures sur les boréliens de \mathbb{R} et pour tout $B \in \mathcal{B}(\mathbb{R})$ nous avons

$$P(X_\theta \in B) = \mu_\theta(B). \quad (38.3)$$

Remarque 38.3.

D'une certaine manière, l'introduction de μ_θ dans la définition est redondante parce que ces mesures sont déjà contenues dans la données des variables aléatoires X_θ .

Exemple 38.4(Modèle statistique gaussien)

Si nous savons que les variables aléatoires X_i suivent une loi gaussienne, alors nous considérons $\Theta = \mathbb{R} \times \mathbb{R}^+$ et $\theta = (m, \sigma^2 j)$. Dans ce cas, $\mu_\theta = \mathcal{N}(m, \sigma^2)$ et le but de la statistique est de déterminer la valeur de θ qui correspond à une population en partant de l'observation d'un échantillon. \triangle

Définition 38.5.

Si $\Theta \subset \mathbb{R}^k$, nous disons que le modèle statistique est un modèle **paramétrique**.

Le modèle gaussien est un modèle paramétrique : dès que m et σ^2 sont déterminés, la loi du phénomène X est connue.

Définition 38.6.

Pour chaque $\theta \in \Theta$, nous disons qu'un **échantillon** de taille n associé à un modèle statistique $[(\Omega, \mathcal{F}, P), (X_\theta), (\mu_\theta)]$ est un vecteur $(X_{\theta,1}, \dots, X_{\theta,n})$ de taille n de variables aléatoires indépendantes et identiquement distribuées de la même loi que la variable aléatoire X_θ . La loi μ_θ est la **loi parente** de l'échantillon.

Définition 38.7.

Un **modèle d'échantillonnage** sur le modèle statistique \mathcal{S} est une famille $(X_{\theta,1}, \dots, X_{\theta,n})_{\theta \in \Theta}$ d'échantillons de taille $n \geq 1$.

Nous noterons souvent (X_1, \dots, X_n) à la place de $(X_{\theta,1}, \dots, X_{\theta,n})$ un échantillon, mais il faut se souvenir que les X_i suivent toujours la même loi donnée par θ . La loi du vecteur (X_1, \dots, X_n) est $\mu_\theta \otimes \dots \otimes \mu_\theta$ et est définie sur l'espace $(\Omega^n, \mathcal{F} \otimes \dots \otimes \mathcal{F}, P^{\otimes n})$.

Remarque 38.8.

Le travail du statisticien est de proposer un modèle statistique \mathcal{S} a priori. Si nous étudions la taille d'une population, nous allons choisir un modèle gaussien. Plus le modèle est précis, plus l'espace Θ est petit mais plus il y a de risques que le vérité soit hors de l'ensemble considéré.

Exemple 38.9

Soit X une variable aléatoire de carré intégrable que l'on sait simuler. Afin d'évaluer la moyenne μ de X , nous pouvons considérer la moyenne empirique des simulations : $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ où les variables aléatoires X_i sont indépendantes, identiquement distribuées et de même loi que X .

La loi des grands nombres nous enseigne que $\bar{X}_n \rightarrow \mu$. De plus,

$$\lim_{n \rightarrow \infty} P \left(\mu \in \left[\bar{X}_n - \frac{a\sigma}{\sqrt{n}}, \bar{X}_n + \frac{a\sigma}{\sqrt{n}} \right] \right) = \int_{-a}^a e^{-x^2/2} \frac{dx}{\sqrt{2\pi}}. \quad (38.4)$$

En effet, la condition sur μ est équivalente à

$$-a \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq a, \quad (38.5)$$

tandis que le théorème central limite nous enseigne que la variable aléatoire $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$ se comporte comme $\mathcal{N}(0, 1)$ lorsque n est grand. Dans ce cas, nous avons que

$$P \left(\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \in [-a, a] \right) = \frac{1}{\sqrt{2\pi}} \int_{-a}^a e^{-x^2/2} dx. \quad (38.6)$$

Notons que dans ce calcul nous avons utilisé le fait que $\mu = E(X_1)$.

Montrons que la suite

$$\sigma_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2 \quad (38.7)$$

converge presque sûrement vers σ^2 . Le théorème central limite implique que

$$\frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow{p.s.} E(X_1^2) \quad (38.8)$$

et que

$$\left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2 \xrightarrow{p.s.} E(X_1)^2. \quad (38.9)$$

La différence converge donc presque sûrement vers σ^2 en vertu de la proposition 37.26.

Nous avons également $E(\sigma_n^2) = \sigma^2$. En effet, sachant que $E(X_i) = E(X_1) = \mu$ et que $E(X_i^2) = E(X_1^2) = \mu^2 + \sigma^2$,

$$E(\sigma_n^2) = \frac{1}{n} \sum_{i=1}^n E(X_i^2) - \frac{1}{n^2} \left(\sum_{i=1}^n E(X_i^2) \right) + \sum_{i \neq j} E(X_i X_j) \quad (38.10a)$$

$$= \sigma^2 + \mu^2 - \frac{1}{n^2} (n(\sigma^2 + \mu^2) + (n^2 - n)\mu^2) \quad (38.10b)$$

$$= \sigma^2 - \frac{1}{n} \sigma^2, \quad (38.10c)$$

dont la limite $n \rightarrow \infty$ donne bien σ^2 .

Nous voudrions à présent montrer que

$$\frac{\bar{X}_n - \mu}{\sigma_n/\sqrt{n}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1). \quad (38.11)$$

Vu que le théorème central limite donne une convergence en loi, nous pouvons utiliser le lemme de Slutsky pour montrer que

$$\left(\frac{\bar{X}_n - \mu}{1/\sqrt{n}}, \sigma_n^2 \right) \xrightarrow{\mathcal{L}} (\sigma Z, \sigma^2) \quad (38.12)$$

où $Z \sim \mathcal{N}(0, 1)$. En vertu de la proposition 37.75 appliqué à la fonction $f: \mathbb{R}^2 \rightarrow \mathbb{R}$,

$$f(x, y) = \begin{cases} \frac{x}{\sqrt{y}} & \text{si } y \neq 0 \\ 0 & \text{sinon} \end{cases} \quad (38.13)$$

nous avons la convergence en loi

$$f \left(\frac{\bar{X}_n - \mu}{1/\sqrt{n}}, \sigma_n^2 \right) \xrightarrow{\mathcal{L}} f(\sigma Z, \sigma^2), \quad (38.14)$$

c'est-à-dire

$$\frac{\bar{X}_n - \mu}{\sigma_n/\sqrt{n}} \xrightarrow{\mathcal{L}} Z. \quad (38.15)$$

Afin d'être complet, précisons que

$$P((\sigma Z, \sigma) \in \mathbb{R} \times \{0\}) = 0. \quad (38.16)$$

△

Proposition 38.10.

Soit X_i des variables aléatoires indépendantes et identiquement distribuées de loi parente $\mathcal{B}(n, p)$. Alors si \bar{X}_n désigne la moyenne empirique,

$$\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{\bar{X}_n(1 - \bar{X}_n)}} \xrightarrow{\mathcal{L}} Z \sim \mathcal{N}(0, 1). \quad (38.17)$$

Démonstration. Cela est une application de la loi des grands nombres, du théorème central limite, du lemme de Slutsky et de la proposition 37.75.

D'abord, la loi des grands nombres nous indique que $\bar{X}_n \rightarrow p$ parce que p est l'espérance de Bernoulli. Ensuite nous avons

$$\frac{\bar{X}_n - E(X)}{\sqrt{\text{Var}(X)}/\sqrt{n}} = \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \quad (38.18)$$

parce que la variance d'une loi de Bernoulli est $p(1-p)$. Le théorème central limite nous indique par conséquent que

$$\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \xrightarrow{\mathcal{L}} Z \sim \mathcal{N}(0, 1). \quad (38.19)$$

Le lemme de Slutsky implique alors la convergence du couple :

$$\left(\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}}, \bar{X}_n \right) \xrightarrow{\mathcal{L}} (Z, p). \quad (38.20)$$

Nous appliquons maintenant la proposition 37.75 avec la fonction

$$f(x, y) = \frac{\sqrt{p(1-p)}x}{\sqrt{y(1-y)}} \quad (38.21)$$

qui est une fonction dont l'ensemble des points de discontinuité est $C = \{0\}$. Étant donné que $P(\bar{X}_n = 0) = 0$, la proposition s'applique et nous avons

$$f \left(\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}}, \bar{X}_n \right) \xrightarrow{\mathcal{L}} f(Z, p), \quad (38.22)$$

c'est-à-dire

$$\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{\bar{X}_n(1 - \bar{X}_n)}} \xrightarrow{\mathcal{L}} Z \sim \mathcal{N}(0, 1). \quad (38.23)$$

□

38.3 Modèles d'échantillonnages

Soit X , une variable aléatoire sur (Ω, \mathcal{F}, P) . Un **échantillon** de taille n pour X est une suite de n variables aléatoires (X_1, \dots, X_n) définies sur (Ω, \mathcal{F}, P) indépendantes et de même loi que X . Nous disons que la loi de X est la **loi parente** de la suite X_i .

Définition 38.11.

Soit

$$\mathcal{S} = \left[(\Omega, \mathcal{F}, P), (X_\theta), (\mu_\theta) \right]_{\theta \in \Theta}, \quad (38.24)$$

un modèle statistique. Un **modèle d'échantillonnage** de taille n associée au modèle statistique \mathcal{S} est la donnée d'une famille de n -échantillons $(X_{\theta,1}), \dots, X_{\theta,n}$ telle que pour tout $\theta \in \Theta$, l'échantillon $(X_{\theta,i})$ soit de variable parente X_θ .

La **moyenne empirique** du n -échantillon (X_i) est la variable aléatoire

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i. \quad (38.25)$$

La proposition suivante signifie que la moyenne empirique est une « bonne » façon d'approcher la variable aléatoire.

Proposition 38.12.

Soit X , une variable aléatoire dans $L^2(\Omega)$ (c'est-à-dire $E(X^2) < \infty$) d'espérance m et de variance σ^2 . Alors

(1) $E(\bar{X}_n) = m$ et $\text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$.

(2) Nous avons les convergences

$$\bar{X}_n \xrightarrow{p.s.} m \quad (38.26a)$$

$$\frac{\bar{X}_n - m}{\sigma/\sqrt{n}} \xrightarrow{\mathcal{L}} Z \sim \mathcal{N}(0, 1). \quad (38.26b)$$

(3) Si X est de loi $\mathcal{N}(m, \sigma^2)$, alors $\bar{X}_n \sim \mathcal{N}(m, \frac{\sigma^2}{n})$, c'est-à-dire

$$\frac{\bar{X}_n - m}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1). \quad (38.27)$$

Remarque 38.13.

L'intérêt de cette proposition en statistique descriptive expérimentale est le suivant. La taille moyenne des français est un nombre m qui existe, mais qui est largement hors de portée de l'expérience (mesurer $65 \cdot 10^6$ personnes risque de prendre un sacré temps). Si on mesure seulement n personnes dont les tailles sont $(x_i)_{i=1, \dots, n}$ (ici x_i est un *nombre expérimental*, pas une variable aléatoire), alors on peut calculer la moyenne \bar{x}_n de ces n personnes-là. La proposition indique que si n est assez grand, alors \bar{x}_n donne une bonne idée de m .

Ne pas confondre X_n qui est une variable aléatoire, c'est-à-dire une application mesurable, qui nous sert à démontrer des théorèmes en mathématique, avec x_n qui est un nombre mesuré sur le terrain, qui a une existence *physique* bien définie, mais aucun status mathématique.

Si on croit que toute cette histoire de variables aléatoires, de tribu et de mesures décrit effectivement la réalité, alors on peut croire que le comportement de la suite \bar{X}_n décrit bien le comportement de la suite \bar{x}_n (cette dernière n'étant même pas une suite parce qu'on n'a jamais qu'un nombre fini de mesures expérimentales).

La **variance empirique** d'un échantillon est la variable aléatoire

$$V_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2. \quad (38.28)$$

La **variance empirique corrigée** est la variable aléatoire

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2. \quad (38.29)$$

Lemme 38.14.

La variance corrigée et la variance empirique ont comme espérances :

$$E(V_n^2) = \frac{n-1}{n} \sigma^2 \quad (38.30a)$$

$$E(S_n^2) = \sigma^2. \quad (38.30b)$$

Démonstration. Nous commençons par calculer l'espérance de la variance non corrigée. La première étape est de la réécrire sous la forme

$$\begin{aligned}
 V_n^2 &= \frac{1}{n} \sum_k (X_k^2 - 2X_k\bar{X}_n + \bar{X}_n^2) \\
 &= \frac{1}{n} \sum_k X_k^2 - \frac{2}{n} \bar{X}_n \underbrace{\sum_k X_k}_{=n\bar{X}_n} + \frac{1}{n} \underbrace{\sum_k \bar{X}_n^2}_{=n\bar{X}_n^2} \\
 &= \frac{1}{n} \sum_k X_k^2 - 2\bar{X}_n^2 + \bar{X}_n^2 \\
 &= \frac{1}{n} \sum_k X_k^2 - \bar{X}_n^2.
 \end{aligned} \tag{38.31}$$

Nous calculons séparément l'espérance de ces deux termes. Si X est la loi parente des X_i , en utilisant l'indépendance des X_i nous trouvons

$$\begin{aligned}
 E\left(\frac{1}{n} \sum_k X_k^2\right) &= \frac{1}{n} \sum_{k=1}^n E(X_k^2) \\
 &= E(X^2) \\
 &= E(X)^2 + \text{Var}(X).
 \end{aligned} \tag{38.32}$$

Nous devons à présent calculer l'espérance de \bar{X}_n^2 :

$$E(\bar{X}_n^2) = E(\bar{X}_n)^2 + \text{Var}(\bar{X}_n). \tag{38.33}$$

En utilisant le lemme 37.23,

$$\text{Var}(\bar{X}_n) = \text{Var}\left(\frac{1}{n} \sum_k X_k\right) \tag{38.34a}$$

$$= \frac{1}{n^2} \sum_k \text{Var}(X_k) \tag{38.34b}$$

$$= \frac{1}{n} \text{Var}(X). \tag{38.34c}$$

Par conséquent

$$E(V_n^2) = \text{Var}(X) \left(1 - \frac{1}{n}\right) = \frac{n-1}{n} \text{Var}(X). \tag{38.35}$$

En ce qui concerne la variance corrigée,

$$S_n^2 = \frac{1}{n-1} \sum_k (X_k - \bar{X}_n)^2 = \frac{n}{n-1} V_n^2, \tag{38.36}$$

par conséquent $E(S_n^2) = \frac{n}{n-1} E(V_n^2) = \text{Var}(X)$. □

Théorème 38.15 (Théorème de Cochran[236]).

Soient (X_i) des variables aléatoires gaussiennes indépendantes de loi $\mathcal{N}(m, \sigma^2)$ avec $\sigma > 0$. Alors

- (1) $\bar{X}_n \sim \mathcal{N}(m, \frac{\sigma^2}{n})$,
- (2) $(\frac{n-1}{\sigma^2}) S_n^2 = (\frac{n}{\sigma^2}) \bar{V}_n^2 \sim \chi^2(n-1)$,
- (3) les variables aléatoires \bar{X}_n et \bar{V}_n sont indépendantes et

$$\frac{\bar{X}_n - m}{S_n/\sqrt{n}} = \frac{\bar{X}_n - m}{\sqrt{\bar{V}_n/(n-1)}} \sim \mathcal{T}(n-1). \tag{38.37}$$

Nous pouvons aussi écrire le dernier résultat en termes de la variance corrigée S_n , l'estimateur sans biais de la variance parce que

$$\sqrt{\frac{\bar{V}_n}{n-1}} = \frac{1}{\sqrt{n}} S_n \quad (38.38)$$

en vertu de la définition (38.29).

Proposition 38.16.

Soit X une variable aléatoire de variance $\text{Var}(X) = \sigma$. Si $E(X^4) < \infty$, alors

(1) $S_n^2 \xrightarrow{p.s.} \sigma^2$.

(2)

$$\frac{S_n^2 - \sigma^2}{\sqrt{\frac{\mu^4 - \sigma^4}{n}}} \xrightarrow{\mathcal{L}} Z \sim \mathcal{N}(0, 1) \quad (38.39)$$

où $\mu^4 = E(X^4)$ est le moment d'ordre 4 de X .

Théorème 38.17.

Si (X_1, \dots, X_n) est un n -échantillon de loi parente $\mathcal{N}(m, \sigma^2)$, alors

(1) Les variables aléatoires

$$\frac{\bar{X}_n - m}{\sigma/\sqrt{n}} \quad \text{et} \quad (n-1) \frac{S_n^2}{\sigma^2} \quad (38.40)$$

sont indépendantes.

(2) La loi de $(n-1) \frac{S_n^2}{\sigma^2}$ est $\chi^2(n-1)$.

Si un échantillon vérifie ces deux propriétés, alors les X_i sont de loi $\mathcal{N}(m, \sigma^2)$.

L'inégalité de Markov donne une borne supérieure à la probabilité qu'une variable aléatoire positive soit plus grande ou égale à une constante.

Théorème 38.18 (Inégalité de Markov).

Soit X une variable aléatoire à valeurs dans \mathbb{R}^d et $\varphi: \mathbb{R}^d \rightarrow [0, \infty[$. Alors

$$P(\varphi(X) \geq a) \leq \frac{E(\varphi(X))}{a} \quad (38.41)$$

pour tout $a > 0$.

Démonstration. Calculons le second membre :

$$\begin{aligned} \frac{E(\varphi(X))}{a} &= \int_{\Omega} \frac{\varphi(X)}{a} dP \\ &= \int_{\varphi(X) \geq a} \underbrace{\frac{\varphi(X)}{a}}_{\leq 1} dP + \int_{\varphi(X) < a} \frac{\varphi(X)}{a} dP \\ &\geq \int_{\varphi(X) \leq a} dP \\ &= P(\varphi(X) \leq a). \end{aligned} \quad (38.42)$$

D'où l'inégalité voulue. □

38.4 Estimation ponctuelle

Nous considérons un modèle statistique

$$\mathcal{S} = [(\Omega, \mathcal{F}, P), (X_\theta), (\mu_\theta)]_{\theta \in \Theta} \quad (38.43)$$

et pour tout θ nous notons $X = (X_{\theta,1}, \dots, X_{\theta,n})$ un échantillon de loi parente μ_θ . Tant que nous travaillerons avec un θ fixé, nous écrirons $X = (X_1, \dots, X_n)$ sans expliciter la paramètre θ . Nous noterons

$$E_\theta(\varphi(X_1, \dots, X_n)) = \int_{\mathbb{R}^n} \varphi(x_1, \dots, x_n) d\mu_\theta^{\otimes n}(x_1, \dots, x_n). \quad (38.44)$$

Dans cette notation nous plaçons le θ sur l'espérance, tandis qu'en réalité le θ devrait être sur chaque X_1 . Tant qu'aucune confusion n'est possible nous ferons toujours cet abus d'écriture.

Le but de la théorie de l'estimation est de déduire la valeur de θ (et donc la loi μ_θ) à partir d'un échantillon de loi parente θ .

Nous posons les hypothèses suivantes.

- (1) Le modèle statistique \mathcal{S} est paramétré, c'est-à-dire que $\Theta \subset \mathbb{R}^d$ avec le plus souvent $d = 1, 2$. Typiquement les paramètres seront la moyenne et la variance.
- (2) Le modèle statistique est **identifiable**, c'est-à-dire que pour tout couple $(\theta_1, \theta_2) \in \Theta^2$, si $\theta_1 \neq \theta_2$, alors $\mu_{\theta_1} \neq \mu_{\theta_2}$.
- (3) Le modèle \mathcal{S} est **dominé** par la mesure de Lebesgue si les lois μ_θ sont continues et par la mesure de comptage si les lois μ_θ sont discrètes.

Exemple 38.19

Quelques familles identifiables :

- La famille des lois exponentielles $(\mathcal{E}(\lambda))_{\lambda > 0}$ est identifiable.
- Les lois gaussiennes $(\mathcal{N}(m, \sigma^2))_{m \in \mathbb{R}, \sigma^2 > 0}$ sont également identifiables.

En réalité il est assez compliqué de trouver un exemple de modèle non identifiable à moins de la faire exprès. Par exemple en paramétrant les lois exponentielles de la façon suivante : $(\mathcal{E}(\sin(\lambda)))_{\lambda \in \mathbb{R}}$. Cette famille n'est pas identifiable. \triangle

Le corollaire 15.197 ainsi que l'hypothèse de modèle dominé implique que les lois ont des densités. Si la loi μ_θ est discrète, nous notons

$$p(x, \theta) = \mu_\theta(\{x\}) \quad (38.45)$$

la densité de μ_θ par rapport à la mesure de comptage. Si la loi μ_θ est continue, nous notons

$$p(x, \theta) = f_\theta(x) \quad (38.46)$$

la densité par rapport à la mesure de Lebesgue.

Si μ_θ est une loi discrète et si (X_1, \dots, X_n) est un échantillon de taille n , alors pour tout $(x_1, \dots, x_n) \in \mathbb{R}^n$ nous avons

$$p_n(x_1, \dots, x_n; \theta) = \mu_\theta^{\otimes n}(\{x_1, \dots, x_n\}) = \mu_\theta(\{x_1\}) \dots \mu_\theta(\{x_n\}) = p(x_1, \theta) \dots p(x_n, \theta). \quad (38.47)$$

La première et la dernière égalité sont des notations ; la seconde est une conséquence de l'indépendance des X_i contenues dans l'échantillon. Pour une loi continue, nous adoptons la même notation. Le vecteur aléatoire (X_1, \dots, X_n) admet la densité

$$(x_1, \dots, x_n) \mapsto p_n(x_1, \dots, x_n; \theta) = f_\theta(x_1) \dots f_\theta(x_n) = p(x_1, \theta) \dots p(x_n, \theta). \quad (38.48)$$

Exemple 38.20

Soit (X_1, \dots, X_n) un n -échantillon de loi $\mathcal{B}(1, \theta)$ avec $\theta \in]0, 1[$. C'est une loi discrète portée par l'ensemble $\{0, 1\}$. Nous avons

$$p(x, \theta) = \begin{cases} 0 & \text{si } 0 \neq x \neq 1 \\ 1 - \theta & \text{si } x = 0 \\ \theta & \text{si } x = 1. \end{cases} \quad (38.49)$$

De façon plus condensée nous pouvons écrire

$$p(x, \theta) = \theta^x (1 - \theta)^{1-x} \mathbb{1}_{\{0,1\}}(x). \quad (38.50)$$

Pour tout $(x_1, \dots, x_n) \in \mathbb{R}^p$, la densité du n -échantillon est donnée par

$$p_n(x_1, \dots, x_n; \theta) = \theta^{x_1 + \dots + x_n} (1 - \theta)^{1 - \sum_i (1 - x_i)} \mathbb{1}_{\{0,1\}^n}(x_1, \dots, x_n). \quad (38.51)$$

△

38.5 Statistiques et estimateurs

Définition 38.21.

Une **statistique** sur un modèle d'échantillonnage est une variable aléatoire fonction de l'échantillon (X_1, \dots, X_n) ne dépendant pas de θ ¹. C'est-à-dire une application borélienne $T: \mathbb{R}^n \rightarrow \mathbb{R}$ ne dépendant pas de θ . La statistique associée à cette application est $S = T(X_1, \dots, X_n)$.

Les fonctions $T(X_1, \dots, X_n)$ données par $\sum_i X_i$, $e^{\sum_i X_i}$ sont des statistiques. La constante $\frac{1}{2}$ est également une statistique (mais elle est moins intéressante).

Un **estimateur** est une statistique qui prend ses valeurs dans Θ . Nous la noterons

$$\hat{\theta}_n = \theta(X_1, \dots, X_n). \quad (38.52)$$

La fonction $\hat{\theta}_n$ est borélienne à valeurs dans Θ .

Exemple 38.22

Soit un n -échantillon de loi $\mathcal{B}(1, \theta)_{\theta \in [0,1]}$. Les fonctions $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$ et $\hat{\varphi}_n = \frac{1}{2}$ sont des estimateurs. Cependant nous devinons que la première va être plus intéressante que la seconde.

△

Pour la suite, nous travaillerons avec des estimateurs de carré intégrable, c'est-à-dire que

$$E_\theta(|\hat{\theta}_n(X_1, \dots, X_n)|^2) < \infty \quad (38.53)$$

pour tout $\theta \in \Theta$.

38.5.1 Qualité des estimateurs

Définition 38.23.

Une estimateur est **convergent** ou **consistant** si pour tout $\theta \in \Theta$, la suite de variables aléatoires $\hat{\theta}_n(X_1, \dots, X_n)$ converge en probabilité vers θ .

En d'autres termes, l'estimateur $\hat{\theta}_n$ est convergent si pour tout $\theta \in \Theta$ et pour tout $\eta > 0$,

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n(X_1, \dots, X_n) - \theta| > \eta) = 0. \quad (38.54)$$

La probabilité dans le membre de gauche est donnée par

$$\mu_\theta^{\otimes n} \left(\{ (x_1, \dots, x_n) \in \mathbb{R}^n \text{ tel que } |\hat{\theta}_n(x_1, \dots, x_n) - \theta| > \eta \} \right). \quad (38.55)$$

Soit $\hat{\theta}_n$ un estimateur pour θ . Nous cherchons à minimiser l'erreur commise en remplaçant θ par $\hat{\theta}_n$. Nous introduisons donc le **risque quadratique** de l'estimateur $\hat{\theta}_n$ par

$$R(\hat{\theta}_n, \theta) = E_\theta((\hat{\theta}_n - \theta)^2). \quad (38.56)$$

Nous disons qu'un estimateur $\hat{\theta}_{n,1}$ est préférable à $\hat{\theta}_{n,2}$ si pour tout $\theta \in \Theta$ nous avons

$$R(\hat{\theta}_{n,1}, \theta) < R(\hat{\theta}_{n,2}, \theta). \quad (38.57)$$

1. Parce que d'habitude c'est ce qu'on cherche à estimer.

Lemme 38.24.

Une formule alternative pour le risque quadratique :

$$R(\hat{\theta}_n, \theta) = \text{Var}(\hat{\theta}_n) + (E_\theta(\hat{\theta}_n) - \theta)^2 \quad (38.58)$$

Démonstration. Nous avons

$$E_\theta \left((\hat{\theta}_n - \theta)^2 \right) = \text{Var}(\hat{\theta}_n - \theta) + E_\theta(\hat{\theta}_n - \theta)^2. \quad (38.59)$$

D'une part $\text{Var}(\hat{\theta}_n - \theta) = \text{Var}(\hat{\theta}_n)$ et d'autre part $E_\theta(\hat{\theta}_n - \theta)^2 = [E(\hat{\theta}_n) - \theta]^2$. Par conséquent

$$R(\hat{\theta}_n, \theta) = \text{Var}(\hat{\theta}_n) + (E(\hat{\theta}_n) - \theta)^2. \quad (38.60)$$

□

Le **biais** de l'estimateur $\hat{\theta}_n$ est la quantité

$$E_\theta(\hat{\theta}_n) - \theta. \quad (38.61)$$

À ce niveau, nous rappelons que nous écrivons E_θ l'espérance calculée en supposant la valeur θ pour le paramètre des différentes variables aléatoires entrant dans le calcul. Voir la discussion autour de la définition (38.44).

Exemple 38.25

Dans le cadre de la proposition 37.103, nous voulons savoir si $\frac{N_t}{t}$ est un estimateur sans biais de λ . Pour ce faire nous calculons

$$E_\lambda \left(\frac{N_t}{t} \right) = \frac{1}{t} E_\lambda(N_t) = \lambda \quad (38.62)$$

parce que $E(N_t) = \lambda t$. Ici nous avons calculé $E(N_t)$ en prenant λ pour valeur du paramètre du processus de Poisson, alors que en principe c'est justement le paramètre que nous voulons estimer.

△

Exemple 38.26

La moyenne empirique \bar{X}_n est un estimateur sans biais de la moyenne. L'estimateur

$$\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \quad (38.63)$$

est un estimateur sans biais de la variance.

△

Un estimateur sans biais n'est pas toujours de meilleure qualité qu'un estimateur avec biais. En effet ce que nous voulons est de se donner un (petit) intervalle I autour de la bonne valeur de θ et de maximiser $P(\hat{\theta}_n \in I)$. Sur la figure 38.1, c'est l'estimateur biaisé rouge tombe plus souvent sur le bon intervalle que l'estimateur non biaisé bleu.

Nous allons maintenant étudier quelques manières de construire des estimateurs convergents². Il vont évidemment s'appuyer sur la loi des grands nombres.

38.5.2 Méthode des moments

Sans surprises, un bon estimateur pour la moyenne est

$$\hat{\theta}_n(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n X_i. \quad (38.64)$$

2. Là je ne suis vraiment pas sûr de l'orthographe. Comme partout où j'utilise le verbe « converger » d'ailleurs.

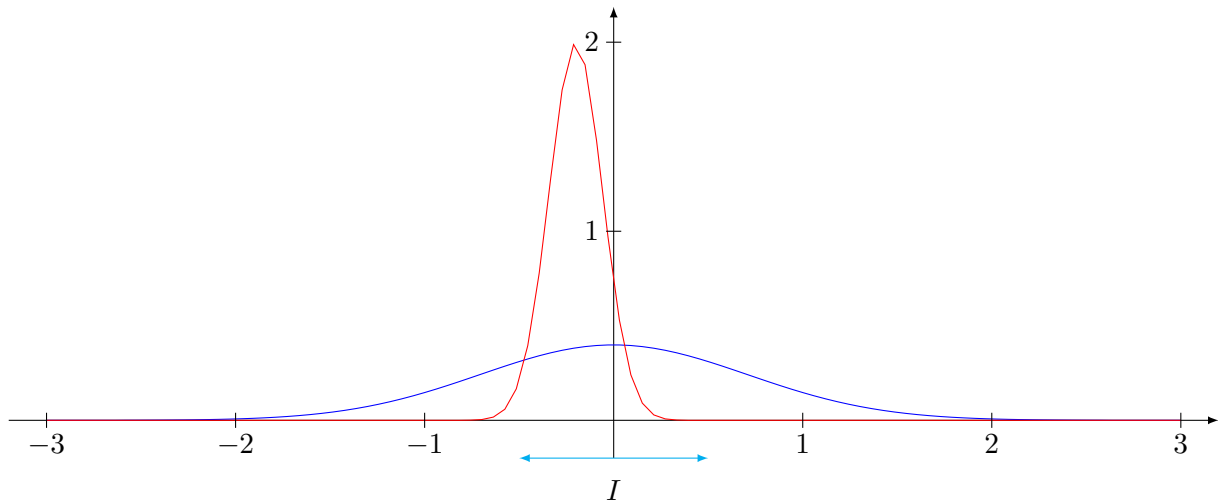


FIGURE 38.1 – Un estimateur sans biais et un avec biais.

Plus généralement, nous supposons qu'il existe une fonction borélienne³ $M: \mathbb{R} \rightarrow \mathbb{R}$ telle que

$$E_\theta(|M(X)|) < \infty \tag{38.65}$$

où X est la loi parente de l'échantillon. Supposons également que la fonction

$$h(\theta) = E_\theta(M(X)) \tag{38.66}$$

soit inversible et continue sur Θ . Dans ce cas, pour estimer le paramètre θ , nous considérons l'estimateur

$$\hat{\theta}_n = h^{-1}\left(\frac{1}{n} \sum_{i=1}^n M(X_i)\right). \tag{38.67}$$

Cela est un estimateur convergent. En effet, la loi des grands nombres dit que

$$\frac{1}{n} \sum_i M(X_i) \xrightarrow{p.s.} E_\theta(M(X)). \tag{38.68}$$

En composant avec la fonction h , nous avons

$$\hat{\theta}_n \xrightarrow{p.s.} h^{-1}\left(E_\theta(M(X))\right) = \theta. \tag{38.69}$$

Dans cette construction, $M(X)$ est le moment de X que l'on souhaite déterminer.

Exemple 38.27

Soit (X_1, \dots, X_n) un échantillon de loi $\mathcal{E}(\lambda)$. Construire $\hat{\lambda}_n$. Pour une loi exponentielle,

$$E(X) = \frac{1}{\lambda}. \tag{38.70}$$

Nous devons donc déterminer le moment d'ordre 1 de X (c'est-à-dire sa moyenne). Nous considérons donc la fonction $M(x) = x$; par conséquent

$$h(\lambda) = E(X) = \frac{1}{\lambda} \tag{38.71}$$

et $h^{-1}(\theta) = 1/\theta$. L'estimateur que nous considérons pour λ est finalement

$$\hat{\theta}_n = \frac{1}{\frac{1}{n} \sum_{i=1}^n X_i}. \tag{38.72}$$

△

3. Définition 15.42.

38.5.3 Méthode de substitution

Supposons que nous connaissions un estimateur convergeant $\hat{\theta}_n \rightarrow \theta$. Si $g: \mathbb{R}^d \rightarrow \mathbb{R}$ est une fonction continue, alors

$$g(\hat{\theta}_n) \rightarrow g(\theta). \quad (38.73)$$

38.5.4 Méthode du maximum de vraisemblance

Exemple 38.28

Nous désirons contrôler la qualité d'une chaîne de production ; pour cela nous prélevons un échantillon de 10 pièces, et nous en trouvons 3 défectueuses. Que dire de la proportion de pièces défectueuses ?

Évidemment, le plus probable est que la proportion de pièces défectueuses soit de $1/3$. Analysons en détail comment nous arrivons à ce résultat. Nous considérons que le fait de tirer 10 pièces revient à une expérience binomiale de paramètres 10 et de probabilité p inconnue. Dans ce cas, la probabilité d'observer exactement 3 pièces défectueuses est de

$$L(p) = P(X = 3) = \binom{10}{3} p^3 (1-p)^7. \quad (38.74)$$

Le maximum de $L(p)$ est $p = 3/10$.

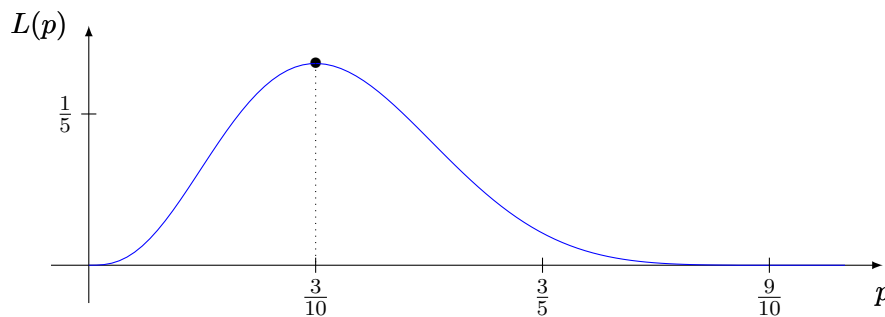


FIGURE 38.2 – La fonction de vraisemblance de l'exemple 38.28.

△

Soit (x_1, \dots, x_n) , une réalisation de l'échantillon (X_1, \dots, X_n) . L'application

$$\theta \mapsto p_n(x_1, \dots, x_n; \theta) = \prod_{i=1}^n p(x_i, \theta) \quad (38.75)$$

est la **vraisemblance** de l'échantillon. Nous définissons $\hat{\theta}_n$ par

$$p_n(x_1, \dots, x_n; \hat{\theta}_n) = \sup_{\theta \in \Theta} p_n(x_1, \dots, x_n; \theta). \quad (38.76)$$

Remarque 38.29.

Nous passons sous le silence le fait que la fonction sup soit une fonction mesurable, et que par conséquent $\hat{\theta}_n$ soit bien une variable aléatoire.

La variable aléatoire $\hat{\theta}_n = \hat{\theta}_n(X_1, \dots, X_n)$ est l'**estimateur de maximum de vraisemblance** de θ .

38.5.5 Exemples sous forme d'exercices

Exemple 38.30

Soit (X_1, \dots, X_n) un échantillon de loi $\mathcal{B}(1, \theta)$ avec $\theta \in [0, 1]$. Trouver l'estimateur de maximum de vraisemblance de θ .

En utilisant la densité de la loi multinomiale,

$$p_n(x_1, \dots, x_n; \theta) = \prod_i p(x_i, \theta) = \prod_i \theta^{x_i} (1 - \theta)^{1-x_i} \mathbb{1}_{\{0,1\}^n}(x_1, \dots, x_n). \quad (38.77)$$

En passant au logarithme, si $x_i \in \{0, 1\}$,

$$L_n(\theta) = \sum_i x_i \ln(\theta) + (1 - x_i) \ln(1 - \theta) \quad (38.78)$$

En passant à la dérivée, nous trouvons que l'extremum est donné par

$$\theta = \frac{1}{n} \sum_{i=1}^n x_i. \quad (38.79)$$

△

Exemple 38.31

Soit (X_1, \dots, X_n) un échantillon de loi $\mathcal{E}(\lambda)$ avec $\theta > 0$. Trouver l'estimateur de maximum de vraisemblance de θ .

Nous avons

$$p_n(x_1, \dots, x_n) = \prod_{i=1}^n \theta e^{-\theta x_i}. \quad (38.80)$$

En passant au logarithme la fonction à minimiser est

$$L(\theta) = n \ln(\theta) - \theta \sum_i x_i. \quad (38.81)$$

La minimisation donne

$$\theta = \frac{n}{\sum_i x_i}, \quad (38.82)$$

c'est-à-dire

$$\hat{\theta}_n(X_1, \dots, X_n) = \frac{n}{\sum_i X_i} = \frac{1}{\bar{X}_n}. \quad (38.83)$$

Ce résultat n'est pas étonnant vu que le paramètre λ de la loi exponentielle est l'inverse de la moyenne.

△

Exemple 38.32

Soit (X_1, \dots, X_n) , un échantillon de loi parente uniforme sur $[0, \theta]$ avec $\theta > 0$.

(1) Montrer que la fonction de vraisemblance est donnée par

$$L(x_1, \dots, x_n; \theta) = \frac{1}{\theta^n} \mathbb{1}_{[0, \infty[}(\min(x_1, \dots, x_n)) \mathbb{1}_{[\max(x_1, \dots, x_n), \infty[}(\theta). \quad (38.84)$$

(2) Déterminer l'estimateur du maximum de vraisemblance de θ .

Lorsque nous parlons de paramètres qui peuvent prendre un spectre continu de valeurs, il est inutile de calculer la *probabilité* parce qu'elle est nulle. Le système de maximum de vraisemblance fonctionne avec les densités. Dans notre cas, la fonction de vraisemblance est le produit des densités :

$$L(x_1, \dots, x_n; \theta) = \prod_{i=1}^n p(x_i; \theta) \quad (38.85a)$$

$$= \prod_i \frac{1}{\theta} \mathbb{1}_{[0, \theta]}(x_i) \quad (38.85b)$$

$$= \frac{1}{\theta^n} \mathbb{1}_{[0, \theta]^n}(x_1, \dots, x_n). \quad (38.85c)$$

De cette expression nous voyons que $\min\{x_1, \dots, x_n\}$ doit être positif en même temps que le maximum doit être plus petit que θ . Cette seconde condition peut s'écrire $\mathbb{1}_{[\max\{x_1, \dots, x_n\}, \infty]}(\theta)$. Au final nous avons

$$L(x_1, \dots, x_n; \theta) = \frac{1}{\theta^n} \mathbb{1}_{[0, \infty]}(\min(x_1, \dots, x_n)) \mathbb{1}_{[\max\{x_1, \dots, x_n\}, \infty]}(\theta). \quad (38.86)$$

Il n'est évidemment pas possible de dériver explicitement cette expression. Par contre pour que cette fonction soit non nulle, il faut obligatoirement $\theta \geq \max\{x_1, \dots, x_n\}$. Par conséquent elle prend son maximum pour $\theta = \max\{x_1, \dots, x_n\}$.

La conclusion est que l'estimateur de maximum de vraisemblance de θ est $\hat{\theta}_n = \max_i \{X_i\}$. \triangle

Exemple 38.33

Déterminer l'estimateur de maximum de vraisemblance de la moyenne pour un modèle statistique dans laquelle la famille de probabilités est

$$(\mu_\theta) = \{\mathcal{N}(\theta, \sigma^2) \text{ tel que } \theta \in \mathbb{R}\}. \quad (38.87)$$

Nous supposons que σ est connu.

La densité de la variable aléatoire conjointe (X_1, \dots, X_n) au point (x_1, \dots, x_n) est le produit des densités, donc

$$\prod_{i=1}^n \gamma_{m, \sigma}(x_1, \dots, x_n) = \frac{1}{\sigma^n (2\pi)^{n/2}} \exp -\frac{1}{2} \left(\sum_i \left(\frac{x_i - m}{\sigma} \right)^2 \right). \quad (38.88)$$

Étant donné que le but est de minimiser, nous pouvons oublier le facteur et passer au logarithme :

$$L_0(\bar{x}) = -\frac{1}{2} \sum_{i=0}^n \left(\frac{x_i - m}{\sigma} \right)^2. \quad (38.89)$$

Nous pouvons également supprimer le $\frac{1}{2}$ et le $1/\sigma^2$. La fonction à minimiser devient

$$L(x_1, \dots, x_n) = -\sum_i (x_i - m)^2, \quad (38.90)$$

dont la dérivée vaut $2nm - 2\sum_i x_i$. Par conséquent nous avons un minimum avec

$$m = \frac{1}{n} \sum_{i=1}^n x_i. \quad (38.91)$$

\triangle

38.5.6 Estimation d'une fonction de répartition

Théorème 38.34 (Glivenko-Cantelli[496]).

Soient X_i des variables aléatoires indépendantes et identiquement distribuées suivant une loi dont la fonction de distribution est F (inconnue). Nous définissons l'estimateur

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i \leq x}. \quad (38.92)$$

Alors

$$P\left(\lim_{n \rightarrow \infty} \|F_n - F\|_{\infty} = 0\right) = 1. \quad (38.93)$$

Autrement dit, pour presque tout $\omega \in \Omega$,

$$\lim_{n \rightarrow \infty} \|F_n(\omega, \cdot) - F\|_{\infty} = 0. \quad (38.94)$$

C'est-à-dire qu'il y a presque certainement convergence en probabilité.

Notons que de façon générale lorsqu'on parle d'estimateurs, partout où il y a un « n » dans une variable aléatoire, il y a une dépendance sous-entendue en ω .

Proposition 38.35.

Pour presque tout x , l'estimateur $F_n(x)$ est sans biais par rapport à $F(x)$:

$$E(F_n(x)) = F(x). \quad (38.95)$$

Démonstration. C'est juste un calcul :

$$E(F_n(x)) = \frac{1}{n} \sum_{i=1}^n E(\mathbb{1}_{\{X_i \leq x\}}) = \frac{1}{n} \sum_{i=1}^n F(x) = F(x). \quad (38.96)$$

□

38.5.7 Exemples sous forme d'exercices

Exemple 38.36

Dans cet exercice nous construisons un estimateur biaisé qui présente un risque quadratique inférieur à un estimateur non biaisé.

Soit un modèle statistique dont la famille de lois est $\{\mathcal{N}(m, \sigma^2)\}_{\theta \in]0, \infty[}$ où m est un paramètre réel connu. En ce qui concerne la variance nous considérons

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - m)^2. \quad (38.97)$$

- (1) Montrer que $\hat{\sigma}_n^2$ est un estimateur sans biais de θ .
- (2) Montrer que le risque quadratique de $\hat{\sigma}_n^2$ est

$$R(\hat{\sigma}_n^2, \theta) = \frac{2\theta^2}{n}. \quad (38.98)$$

- (3) Considérer la famille d'estimateurs $c\hat{\sigma}_n^2$ avec $c > 0$. Déterminer la valeur de c qui minimise le risque quadratique de l'estimateur.

- (1) Nous calculons

$$\begin{aligned} E(\hat{\sigma}_n^2) &= \frac{1}{n} \sum_i E[(X_i - m)^2] \\ &= \frac{1}{n} \sum_i \text{Var}(X_i - m) - \frac{1}{n} \sum_i E(X_i - m)^2 \\ &= \frac{1}{n} \sum_i \text{Var}(X_i) \\ &= \sigma^2. \end{aligned} \quad (38.99)$$

- (2) Par la formule (38.60) nous savons que le risque d'un estimateur sans biais est donné par sa variance :

$$R(\hat{\theta}_n, \theta) = \text{Var}(|\hat{\sigma}_n^2 - \theta|) = \text{Var}(\hat{\sigma}_n^2). \quad (38.100)$$

Étant donné que les variables aléatoires X_i sont indépendantes et identiquement distribuées, le lemme 37.23 nous enseigne que la variance de la somme est la somme des variances. Nous avons donc à calculer

$$\begin{aligned} \text{Var}(|\hat{\sigma}_n^2 - \theta|) &= \text{Var}(\hat{\sigma}_n^2) \\ &= \frac{1}{n^2} \sum_i \text{Var}[(X_i - m)^2] \\ &= \frac{1}{n^2} \sum_i E[(X_i - m)^4] - E[(X_i - m)^2]^2. \end{aligned} \quad (38.101)$$

Ces espérances ne sont pas très compliquées à calculer en utilisant la fonction caractéristique donnée par la proposition 37.105 :

$$\Phi_{X-m}(t) = E(e^{it(X-m)}) = e^{-itm} E(e^{itX}) = e^{-itm} \Phi_X(t) = \exp\left(-\frac{\sigma^2 t^2}{2}\right). \quad (38.102)$$

Nous avons $\Phi^{(4)}(0) = 3\sigma^4$ et $\Phi''(0) = -\sigma^2$. Par conséquent

$$\text{Var}(\hat{\sigma}_n^2) = \frac{1}{n^2} \sum_i 2\theta^2 = \frac{2\theta^2}{n}. \quad (38.103)$$

Notez ici que $\theta = \sigma^2$.

- (3) En tenant compte du fait que $\text{Var}(\hat{\sigma}_n^2) = \frac{2\theta^2}{n}$ et $E(\hat{\sigma}_n^2) = \theta$, nous avons

$$E((c\hat{\sigma}_n^2 - \theta)^2) = \text{Var}(c\hat{\sigma}_n^2 - \theta) + E(c\hat{\sigma}_n^2 - \theta)^2 \quad (38.104a)$$

$$= 2\frac{c^2\theta^2}{n} + (c-1)^2\theta^2. \quad (38.104b)$$

La dérivée (par rapport à c) de cela s'annule pour $c_0 = \frac{n}{n+2}$. Notons que nous n'avons pas tout à fait démontré que cela est bien un minimum. Calculons cependant le risque quadratique de notre estimateur pour cette valeur de c . Pour cela nous reportons $c = c_0$ dans l'expression (38.104b) :

$$E\left(\frac{n}{n+2}\hat{\sigma}_n^2\right) = \frac{2\theta^2}{n+2}. \quad (38.105)$$

Cela est effectivement plus petit que $R(\hat{\sigma}_n^2, \theta)$.

Nous avons ainsi construit un estimateur biaisé qui a un risque quadratique plus petit que l'estimateur non biaisé. \triangle

Exemple 38.37

Nous considérons la famille de probabilités $\mu_\theta = \mathcal{N}(\theta)$ où $\theta = (m, \sigma^2) \in \mathbb{R} \times]0, \infty[$. Déterminer l'estimateur de maximum de vraisemblance du paramètre $\theta = (m, \sigma^2)$.

Pour chaque observation x_i nous avons une densité gaussienne. Le produit donne

$$p(x_1, \dots, x_n; (m, \sigma^2)) = \frac{1}{\sigma^2(2\pi)^{n/2}} \exp\left[-\frac{1}{2\sigma^2} \sum_i (x_i - m)^2\right]. \quad (38.106)$$

En passant au logarithme et en supprimant des facteurs inutiles à la minimisation,

$$L(m, \sigma) = -n \ln(\sigma) - \frac{1}{2\sigma^2} \sum_i (x_i - m)^2. \quad (38.107)$$

L'annulation de la dérivée par rapport à m donne immédiatement $m = \frac{1}{n} \sum_i x_i$. L'annulation de la dérivée par rapport à σ donne

$$-n\sigma^2 + \sum_i (x_i - m)^2 = 0 \quad (38.108)$$

et donc

$$\sigma^2 = \frac{1}{n} \sum_i (x_i - \bar{x}_n). \quad (38.109)$$

L'estimateur de maximum de vraisemblance du couple $\theta = (m, \sigma^2)$ est donc

$$\hat{\theta}_n = \left(\frac{1}{n} \sum_i X_i, \frac{1}{n} \sum_i (X_i - \bar{X}_n)^2 \right). \quad (38.110)$$

△

38.5.8 Espérance et variance d'un estimateur

Soit $T_n = T(X_1, \dots, X_n)$ un estimateur du paramètre θ dans un modèle d'échantillonnage. Les moyennes et variances de l'estimateur sont les variables aléatoires

$$m_{\theta,n} = E_{\theta}(T_n) = \int_{\mathbb{R}^n} T_n(x_1, \dots, x_n) d\mu_{\theta}^{\otimes n}(x_1, \dots, x_n), \quad (38.111a)$$

$$\text{Var}_{\theta}(T_n) = \int_{\mathbb{R}^n} [T_n(x_1, \dots, x_n) - m_{\theta,n}]^2 d\mu_{\theta}^{\otimes n}(x_1, \dots, x_n), \quad (38.111b)$$

Lemme 38.38.

Si l'estimateur T_n satisfait

$$\lim_{n \rightarrow \infty} E_{\theta}(T_n) = \theta \quad (38.112a)$$

$$\lim_{n \rightarrow \infty} \text{Var}_{\theta}(T_n) = 0, \quad (38.112b)$$

alors il est convergent.

Démonstration. Nous utilisons l'inégalité de Markov (théorème 38.18) et nous introduisons l'espérance de l'estimateur :

$$P(|T_n(X_1, \dots, X_n) - \theta| > \epsilon) \leq \frac{1}{\epsilon} E(T_n(X_1, \dots, X_n) - \theta) \quad (38.113a)$$

$$\leq \frac{1}{\epsilon} E(|T_n - m_{n,\theta}|) + \frac{1}{\epsilon} E(|m_{n,\theta} - \theta|) \quad (38.113b)$$

$$(38.113c)$$

Le second terme est l'espérance d'une constante. Nous majorons le premier terme en utilisant le fait que $\|\cdot\|_1 \leq \|\cdot\|_2$ (voir la remarque 28.34 après l'inégalité de Hölder) :

$$P(|T_n(X_1, \dots, X_n) - \theta| > \epsilon) \leq \frac{1}{\epsilon} E(|T_n - m_{n,\theta}|^2)^{1/2} + \frac{1}{\epsilon} |E_{\theta}(T_n) - \theta| \quad (38.114a)$$

$$= \frac{1}{\epsilon} \text{Var}(T_n)^{1/2} + \frac{1}{\epsilon} |E_{\theta}(T_n) - \theta|. \quad (38.114b)$$

Les deux termes tendent séparément vers zéro par hypothèse. Nous avons par conséquent la convergence en probabilité $T_n \rightarrow \theta$. □

38.6 Estimation par intervalle de confiance

Nous voudrions estimer la proportion d'individus dans une population ayant un certain caractère déterminé par une variable booléenne : chaque individu a ou non le caractère étudié. L'échantillon sera donc une suite de 0 et de 1.

Pour tout $i \in \{1, \dots, n\}$ nous notons

$$x_i = \begin{cases} 1 & \text{si le } i\text{ème individu a la caractère} \\ 0 & \text{sinon.} \end{cases} \quad (38.115)$$

et $\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$. Notre modèle statistique sera

$$\mathcal{S} = \left[(\Omega, \mathcal{F}, P), (X_\theta), B(1, \theta) \right] \quad (38.116)$$

où Ω est l'ensemble des individus étudiés, P est la manière de choisir les individus lors du sondage (essentiellement c'est une loi uniforme) et X_θ est la variable aléatoire

$$X(\omega) = \begin{cases} 1 & \text{si } \omega \text{ a le caractère} \\ 0 & \text{sinon} \end{cases} \quad (38.117)$$

Cela est une variable aléatoire de distribution $\mathcal{B}(1, p)$ où p est inconnu. Ici, $\Theta = [0, 1]$ est l'ensemble des p possibles.

Remarque 38.39.

Nous supposons que Ω est la population entière et que la variable aléatoire est l'opinion de la personne ω . En cela, nous considérons que le tirage de l'échantillon est sans remise. Le fait que nous modélisons par une variable aléatoire de Bernoulli signifie que nous considérons l'approximation dans laquelle la population globale est grande.

Nous supposons que nous ayons un échantillon (X_1, \dots, X_n) dont nous avons observé une réalisation (x_1, \dots, x_n) de fréquence \bar{x}_n . Nous voudrions déterminer un intervalle dans lequel \bar{X}_n a de fortes chances de se trouver. Plus précisément nous considérons un petit α et nous cherchons ϵ tel que

$$P(p \in [\bar{X}_n - \epsilon, \bar{X}_n + \epsilon]) = 1 - \alpha. \quad (38.118)$$

Typiquement, $\alpha = 5\%$. Le nombre α est le **niveau de confiance** que nous nous fixons a priori.

Si nous trouvons un intervalle I tel que $P(p \in I) = 1 - \alpha$, nous disons que l'intervalle est **exact**, si nous avons $P(p \in I) \geq 1 - \alpha$, nous disons que l'intervalle est **par excès**.

Il y a deux points de départ pour trouver l'intervalle. Le plus simple est d'utiliser le théorème central limite et considérer

$$\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1). \quad (38.119)$$

La seconde est d'utiliser la loi exacte : $n\bar{X}_n = \sum_i X_i \sim \mathcal{B}(n, p)$.

Bien entendu la seconde donne lieu à des calculs plus compliquée.

Remarque 38.40.

Dans certaines vraies vies (par exemple en médecine), la taille des échantillons est très réduite. Dans ce cas le théorème central limite n'a aucun sens et les calculs exact s'imposent.

De plus dans de nombreux cas de la vraie vie, nous avons un ordinateur à disposition pour calculer avec la loi exacte. L'utilisation du théorème central limite dans le but de produire un intervalle de confiance semble donc de plus en plus être une survivance du passé.

Dans la suite, nous allons supposer que n est suffisamment grand pour justifier l'approximation normale du théorème central limite 37.89. Si Z est une variable aléatoire normale centrée réduite,

notre premier essai est de faire

$$1 - \alpha = P(p \in [\bar{X}_n - \epsilon, \bar{X}_n + \epsilon]) \quad (38.120a)$$

$$= P\left(\frac{-\epsilon\sqrt{n}}{\sqrt{p(1-p)}} \leq \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \leq \frac{\epsilon\sqrt{n}}{\sqrt{p(1-p)}}\right) \quad (38.120b)$$

$$\simeq P\left(-\frac{\epsilon\sqrt{n}}{\sqrt{p(1-p)}} \leq Z \leq \frac{\epsilon\sqrt{n}}{\sqrt{p(1-p)}}\right) \quad (38.120c)$$

$$= 2P\left(Z \leq \frac{\epsilon\sqrt{n}}{\sqrt{p(1-p)}}\right) - 1. \quad (38.120d)$$

La dernière ligne utilise la symétrie de la distribution $\mathcal{N}(0, 1)$. Le nombre ϵ que nous cherchons vérifie donc

$$P\left(Z \leq \frac{\epsilon\sqrt{n}}{\sqrt{p(1-p)}}\right) = 1 - \frac{\alpha}{2}. \quad (38.121)$$

De nos jours, les ordinateurs donnent la loi de répartition inverse des normales. Cela nous fournit un nombre t_α tel que

$$\frac{\epsilon\sqrt{n}}{\sqrt{p(1-p)}} = t_\alpha \quad (38.122)$$

où t_α est le nombre tel que $P(Z \leq t_\alpha) = 1 - \alpha/2$.

Conclusions de ce premier essai :

- (1) Le problème est que nous ne pouvons pas déduire ϵ de l'équation (38.122) parce que p est inconnu.
- (2) Cela ruine notre premier essai et nous demande de trouver mieux.
- (3) L'astuce est évidemment de remplacer p par \bar{X}_n , mais il faut le justifier.

Méthode Slutsky Le point de départ du premier essai infructueux était la convergence

$$\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \xrightarrow{\mathcal{L}} N \sim \mathcal{N}(0, 1) \quad (38.123)$$

donnée par le théorème central limite 37.89. Ce que nous voudrions en réalité est la convergence

$$\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{\bar{X}_n(1 - \bar{X}_n)}} \xrightarrow{\mathcal{L}} N \quad (38.124)$$

La loi des grands nombres nous donne

$$\bar{X}_n(1 - \bar{X}_n) \xrightarrow{p.s.} p(1 - p). \quad (38.125)$$

Par conséquent le lemme de Slutsky implique la convergence en loi du couple :

$$\left(\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}}, \sqrt{\bar{X}_n(1 - \bar{X}_n)}\right) \xrightarrow{\mathcal{L}} (N, \sqrt{p(1-p)}). \quad (38.126)$$

À ce point des opérations nous pouvons utiliser la proposition 37.75 au couple avec la fonction

$$f(x, y) = \frac{x}{y} \sqrt{p(1-p)} \quad (38.127)$$

dont la probabilité d'être continue est 1 ($y = 0$ est de mesure nulle dans \mathbb{R}^2). La conclusion du théorème est que

$$\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{\bar{X}_n(1 - \bar{X}_n)}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1). \quad (38.128)$$

C'est à partir de là que nous pouvons construire notre intervalle de confiance :

$$1 - \alpha = P(\bar{X}_n - \epsilon \leq p \leq \bar{X}_n + \epsilon) \quad (38.129a)$$

$$= P\left(\frac{\sqrt{n}\epsilon}{\sqrt{\bar{X}_n(1-\bar{X}_n)}} \geq \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{\bar{X}_n(1-\bar{X}_n)}} \geq \frac{-\sqrt{n}\epsilon}{\sqrt{\bar{X}_n(1-\bar{X}_n)}}\right) \quad (38.129b)$$

$$\simeq P\left(\frac{\sqrt{n}\epsilon}{\sqrt{\bar{X}_n(1-\bar{X}_n)}} \geq N \geq \frac{-\sqrt{n}\epsilon}{\sqrt{\bar{X}_n(1-\bar{X}_n)}}\right). \quad (38.129c)$$

Nous cherchons maintenant dans les tables le ξ qui fait

$$P(-\xi \leq N \leq \xi) = 1 - \alpha \quad (38.130)$$

puis nous cherchons ϵ de telle sorte à avoir

$$\frac{\sqrt{n}\epsilon}{\sqrt{\bar{X}_n(1-\bar{X}_n)}} = \xi. \quad (38.131)$$

Dans cette équation tout est connu à part le ϵ qui se découvre.

Méthode piétonne Nous remarquons que l'événement

$$-\frac{\epsilon\sqrt{n}}{\sqrt{p(1-p)}} \leq \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \leq \frac{\epsilon\sqrt{n}}{\sqrt{p(1-p)}} \quad (38.132)$$

est le même que l'événement

$$\left| \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \right| \leq t_\alpha. \quad (38.133)$$

Vu que t_α est positif, cela est encore le même événement que

$$n \frac{(\bar{X}_n - p)^2}{p(1-p)} \leq t_\alpha^2 \quad (38.134)$$

ou encore

$$p^2(n + t_\alpha^2) - p(2n\bar{X}_n + t_\alpha^2) + n\bar{X}_n^2 \leq 0. \quad (38.135)$$

Les racines du polynôme du membre de gauche sont

$$p_\pm = \frac{2n\bar{X}_n + t_\alpha^2 \pm \sqrt{(2n\bar{X}_n + t_\alpha^2)^2 - 4(n + t_\alpha^2)n\bar{X}_n^2}}{2(n + t_\alpha^2)}. \quad (38.136)$$

Le but étant d'effectuer une limite $n \rightarrow \infty$, nous factorisons d'abord n . Après simplification

$$p_\pm = \frac{\bar{X}_n + \frac{t_\alpha^2}{2n} \pm t_\alpha \sqrt{\frac{t_\alpha^2}{4n^2} + \frac{\bar{X}_n(1-\bar{X}_n)}{n}}}{1 + \frac{t_\alpha^2}{n}}. \quad (38.137)$$

Étant donné que nous considérons que n est grand, nous allons négliger les termes en $\frac{1}{n}$ en faisant attention à ce que le terme en $\frac{1}{n}$ sous la racine est en réalité $1/\sqrt{n}$ et ne doit pas être négligé. Nous trouvons, à cette approximation, que

$$p \in \left[\bar{X}_n - t_\alpha \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}, \bar{X}_n + t_\alpha \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}} \right] \quad (38.138)$$

avec une probabilité $1 - \alpha$.

38.6.1 Région de confiance

Soit un n -échantillon (X_1, \dots, X_n) de loi parente μ_θ . Nous supposons $\Theta \subset \mathbb{R}$ avec Θ ouvert. Soit $\alpha \in [0, 1]$ un intervalle de confiance et une application mesurable

$$\begin{aligned} \Lambda: \mathbb{R}^n &\rightarrow \mathcal{B}(\Theta) \\ (x_1, \dots, x_n) &\mapsto \Lambda(x_1, \dots, x_n). \end{aligned} \quad (38.139)$$

On appelle **région de confiance exact** au niveau de confiance $1 - \alpha$ une région aléatoire $\Lambda(x_1, \dots, x_n)$ telle que

$$P(\theta \in \Lambda(x_1, \dots, x_n)) = 1 - \alpha. \quad (38.140)$$

Si $d = 1$, la région $\Lambda(x_1, \dots, x_n)$ est un intervalle.

38.6.2 Fonction pivotale

Soit $\hat{\theta}_n$, un estimateur de θ . Une fonction v sur $\Theta \times \Theta$ est **pivotale** pour θ si la loi de la variable aléatoire $v(\hat{\theta}_n, \theta)$ ne dépend pas de θ . Elle est **asymptotiquement pivotale** si

$$v(\hat{\theta}_n, \theta) \xrightarrow{\mathcal{L}} \xi \quad (38.141)$$

où ξ est une variable aléatoire indépendante de θ .

En pratique, nous essayons de faire apparaître une variable aléatoire de loi connue qui ne dépend pas du paramètre que l'on recherche. Si la variance est connue et si l'échantillon est grand, le théorème central limite nous permet d'introduire une loi normale centrée réduite.

Exemple 38.41

Soit X_1, \dots, X_n des variables aléatoires de loi parente $\mathcal{N}(m, \sigma^2)$. Une fonction asymptotiquement pivotale pour m est

$$v(z_1, z_2) = \frac{z_1 - z_2}{\sigma/\sqrt{n}} \quad (38.142)$$

parce que la variable aléatoire

$$v(\bar{X}_n, m) = \frac{\bar{X}_n - m}{\sigma/\sqrt{n}} \quad (38.143)$$

tend vers $\mathcal{N}(0, 1)$ qui ne dépend pas de m . △

Exemple 38.42

Si (X_n) est une suite de variables aléatoires indépendantes et identiquement distribuées de moyenne m et d'écart type σ que nous supposons inconnus. Le fonction suivante est asymptotiquement pivotale pour m :

$$v(\bar{X}_n, m) = \frac{\bar{X}_n - m}{\sigma/\sqrt{n}}. \quad (38.144)$$

△

Soit (X_1, \dots, X_n) un échantillon de loi $\mathcal{N}(m, \sigma_0^2)$ avec σ_0^2 connu. Nous cherchons un intervalle de confiance $1 - \alpha$ pour m . Pour cela nous allons utiliser une fonction asymptotiquement pivotale, à savoir

$$\frac{\bar{X}_n - m}{\sigma_0^2/\sqrt{n}} \sim \mathcal{N}(0, 1). \quad (38.145)$$

Nous devrions chercher des valeurs z_+ de z_- telles que

$$P\left(z_- \leq \frac{\bar{X}_n - m}{\sigma_0/\sqrt{n}} \leq z_+\right) = 1 - \alpha. \quad (38.146)$$

Pour des raisons de symétries (de la courbe gaussienne), nous allons chercher un intervalle symétrique : $z_- = -z_+$. Le nombre à chercher est donc le z_α tel que

$$P(|Z| \leq z_\alpha) = 1 - \alpha. \quad (38.147)$$

Si nous demandons $\alpha = 5\%$, la réponse est $z_\alpha = 1.96$, c'est-à-dire que

$$P\left(-1.96 \leq \frac{\bar{X}_n - m}{\sigma/\sqrt{n}} \leq 1.96\right) = 0.95. \quad (38.148)$$

Nous avons donc

$$P\left(m \in \left[\bar{X}_n - \frac{1.96\sigma}{\sqrt{n}}, \bar{X}_n + \frac{1.96\sigma}{\sqrt{n}}\right]\right) = 0.95. \quad (38.149)$$

Supposons maintenant que nous avons observé 100 valeurs numériques avec $\bar{x}_n = 12$ et $\sigma = 1$. La réalisation de l'intervalle de confiance pour m au niveau de confiance 0.95 est :

$$[12 - 0.196, 12 + 0.196]. \quad (38.150)$$

Cet intervalle est à interpréter de la façon suivante : si nous recommençons un grand nombre de fois le sondage, la moyenne tombera 95% des fois dans l'intervalle ainsi calculé. Mais il faut bien comprendre que la probabilité

$$P(m \in [12 - 0.196, 12 + 0.196]) \quad (38.151)$$

vaut zéro ou un.

Exemple 38.43

Soit un échantillon (X_1, \dots, X_n) de loi parente $\mathcal{N}(m, \sigma^2)$ où m et σ^2 sont inconnus. Déterminer un intervalle de confiance exact symétrique au niveau de confiance $1 - \alpha$.

Déterminer un intervalle de confiance pour σ^2 .

Nous savons que la moyenne empirique est un estimateur de la moyenne :

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i. \quad (38.152)$$

Nous cherchons un intervalle du type $I = [\bar{X}_n - \epsilon, \bar{X}_n + \epsilon]$ pour lequel $P(m \in I) = 1 - \alpha$. Nous savons que la variable aléatoire

$$\frac{\bar{X}_n - m}{\sigma/\sqrt{n}} \quad (38.153)$$

suit une loi $\mathcal{N}(0, 1)$, mais la variance est inconnue. La subtilité à savoir est que la variable aléatoire

$$Z = \frac{\bar{X}_n - m}{S_n/\sqrt{\sigma}} \quad (38.154)$$

où $S_n^2 = \sum_i (X_i - \bar{X}_n)^2 / (n - 1)$ suit une loi de Student à n degrés de liberté $\mathcal{T}(n - 1)$ en vertu du théorème de Cochran 38.15. Comme il est usuel de le faire, nous inversons l'intervalle :

$$1 - \alpha = P(-\epsilon \leq \bar{X}_n - m \leq \epsilon) \quad (38.155a)$$

$$= P\left(-\frac{\epsilon\sqrt{n}}{S_n} \leq Z \leq \frac{\epsilon\sqrt{n}}{S_n}\right). \quad (38.155b)$$

Les valeurs se trouvent dans des tables ; par exemple pour $n = 10$ et $\alpha = 5\%$ nous trouvons

$$\frac{\epsilon\sqrt{n}}{S_n} = 2.262. \quad (38.156)$$

Plus généralement nous notons $t_{n-1,1-\frac{\alpha}{2}}$ le quantile d'ordre $1 - \frac{\alpha}{2}$ de la loi $\mathcal{T}(n-1)$, c'est-à-dire le nombre tel que

$$P(Z \leq t_{n-1,1-\frac{\alpha}{2}}) = 1 - \frac{\alpha}{2} \tag{38.157}$$

si $Z \sim \mathcal{T}(n-1)$. L'intervalle de confiance est alors donné par

$$I = \left[\bar{X}_n - \frac{t_{n-1,1-\frac{\alpha}{2}} S_n}{\sqrt{n}}, \bar{X}_n + \frac{t_{n-1,1-\frac{\alpha}{2}} S_n}{\sqrt{n}} \right]. \tag{38.158}$$

Cela est un intervalle exact pour m au niveau de confiance $1 - \alpha$.

Nous pouvons aussi trouver un intervalle asymptotique en utilisant le théorème central limite :

$$\frac{\bar{X}_n - m}{S_n/\sqrt{n}} \xrightarrow{\mathcal{L}} T \tag{38.159}$$

avec $T \sim \mathcal{N}(0, 1)$.

Rappel : dire que I_n est un **intervalle de confiance asymptotique** signifie que

$$\lim_{n \rightarrow \infty} P(m \in I_n) = 1 - \alpha. \tag{38.160}$$

En ce qui concerne la variance σ^2 , l'intervalle de confiance se construit en utilisant la partie (2) du théorème de Cochran 38.15. Nous introduisons la variable aléatoire pivot

$$Z = (n-1) \frac{S_n^2}{\sigma^2} \tag{38.161}$$

qui suit une loi $\chi^2(n-1)$. Cette loi n'étant pas symétrique (voir figure 37.1), nous n'allons pas chercher un intervalle de confiance symétrique. Nous cherchons c_1 et c_2 tels que

$$\begin{cases} P(\sigma^2 \in [c_1, c_2]) = 1 - \alpha & (38.162a) \\ P(\sigma^2 \in [0, c_1]) = \frac{\alpha}{2} & (38.162b) \\ P(\sigma^2 \in [c_2, \infty]) = \frac{\alpha}{2} & (38.162c) \end{cases}$$

La situation est représentée à la figure 38.3.

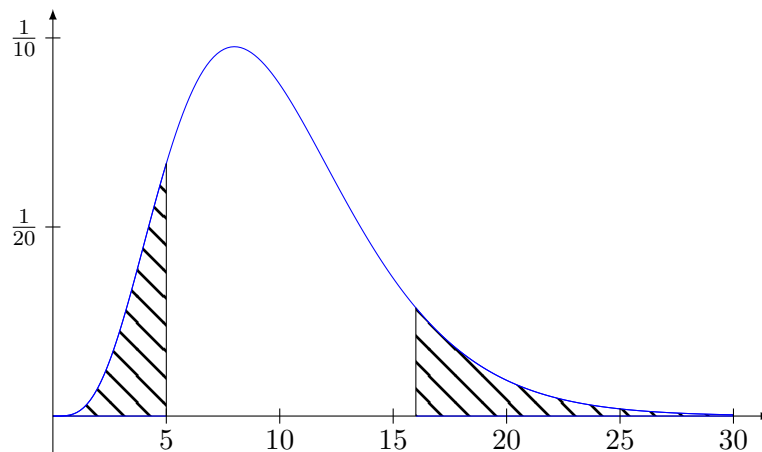


FIGURE 38.3 – L'intervalle de confiance pour une χ^2 .

Le construction des nombres c_1 et c_2 passe par la relation

$$P(c_1 \leq \sigma^2 \leq c_2) = P\left(\frac{(n-1)S_n^2}{c_2} \leq Z \leq \frac{(n-1)S_n^2}{c_1}\right). \tag{38.163}$$

Nous notons $t_{n-1, \frac{\alpha}{2}}$ et $t_{n-1, 1-\frac{\alpha}{2}}$ les quantiles donnés sur la figure 38.3, c'est-à-dire

$$P(Z \leq t_{n-1, \frac{\alpha}{2}}) = \frac{\alpha}{2} P(Z \geq t_{n-1, 1-\frac{\alpha}{2}}) = 1 - \frac{\alpha}{2}. \quad (38.164a)$$

Ce que nous obtenons est

$$\frac{(n-1)S_n^2}{c_2} = t_{n-1, \frac{\alpha}{2}} \quad (38.165a)$$

$$\frac{(n-1)S_n^2}{c_1} = t_{n-1, 1-\frac{\alpha}{2}}, \quad (38.165b)$$

et par conséquent l'intervalle de confiance pour σ^2 est

$$I = \left[\frac{(n-1)S_n^2}{t_{n-1, 1-\frac{\alpha}{2}}}, \frac{(n-1)S_n^2}{t_{n-1, \frac{\alpha}{2}}} \right] \quad (38.166)$$

avec $P(\sigma^2 \in I) = 95\%$.

Remarque 38.44.

Il n'est pas clair a priori que la longueur de l'intervalle I décroisse avec n parce qu'il y a n dans les t au numérateur.

△

38.6.3 Sondage de proportion

Une utilisation classique des statistiques est d'interpréter une proportion donnée par un sondage. Nous considérons une élection avec deux candidats A et B . Nous avons interrogés $n = 2500$ personnes et nous avons obtenus 51% pour le candidat A et 49% pour le candidat B . Que peut on dire ?

La modélisation de cette situation est que nous avons des variables aléatoires $X_i \sim \mathcal{B}(p_A)$ et que nous en avons observés n avec une moyenne

$$\bar{x}_n = 0.51. \quad (38.167)$$

La loi de \bar{X}_n est une binomiale. Sa densité n'est pas symétrique, mais si n est grand, elle le devient. Nous cherchons un intervalle

$$I = [\bar{X}_n - \epsilon, \bar{X}_n + \epsilon] \quad (38.168)$$

tel que $P(p_A \in I) = 1 - \alpha$. Pour cela nous considérons le fait que $n = 2500$ est grand et nous utilisons la limite de la proposition 38.10 :

$$Z_n = \sqrt{n} \frac{\bar{X}_n - p_A}{\sqrt{\bar{X}_n(1 - \bar{X}_n)}} \xrightarrow{\mathcal{L}} Z \sim \mathcal{N}(0, 1). \quad (38.169)$$

La variable aléatoire Z_n est asymptotiquement pivotale et normale centrée réduite. Nous cherchons donc un intervalle symétrique pour $\bar{X}_n - p_A$:

$$1 - \alpha = P(-\epsilon \leq \bar{X}_n - p_A \leq \epsilon), \quad (38.170)$$

c'est-à-dire, si n est grand,

$$1 - \alpha = P\left(-\epsilon \frac{\sqrt{n}}{\sqrt{\bar{X}_n(1 - \bar{X}_n)}} \leq Z_n \leq \epsilon \frac{\sqrt{n}}{\sqrt{\bar{X}_n(1 - \bar{X}_n)}}\right) \quad (38.171)$$

où Z_n est une normale centrée réduite. Nous trouvons ainsi, via les tables que

$$\frac{\epsilon\sqrt{n}}{\sqrt{\bar{X}_n(1-\bar{X}_n)}} = 1.96 \quad (38.172)$$

si nous voulons un intervalle à 5%. Par conséquent nous avons $\epsilon = 1.96\frac{\sqrt{\bar{X}_n(1-\bar{X}_n)}}{\sqrt{n}}$ et l'intervalle de confiance est

$$I_C = \left[\bar{X}_n - \frac{1.96\sqrt{\bar{X}_n(1-\bar{X}_n)}}{\sqrt{n}}, \bar{X}_n + \frac{2.96\sqrt{\bar{X}_n(1-\bar{X}_n)}}{\sqrt{n}} \right]. \quad (38.173)$$

La propriété de cet intervalle est que

$$\lim_{n \rightarrow \infty} P(p_A \in I_C) = 1 - \alpha. \quad (38.174)$$

Remarque 38.45.

À quel moment avons nous fait une hypothèse sur la taille de la population globale ? En modélisant les sondés par des variables de Bernoulli et leur somme par une binomiale, nous supposons que le sondage est *avec remise*, sinon, elles ne seraient pas indépendantes. En supposant les sondés indépendants, nous avons donc fait comme si la population totale était infinie.

38.7 Estimer une densité lorsqu'on ne sait rien

Nous supposons avoir une série d'observations issues d'un processus complexe dont nous n'avons aucune idée de la loi parente, et nous voudrions nous faire une idée de la densité de cette loi inconnue.

Nous observons une suite de réalisations que nous modélisons comme étant des variables aléatoires (X_1, \dots, X_n) de loi parente (inconnue) μ . Notre but est de trouver un estimateur $\hat{\mu}$ de μ . Par simplicité nous allons nous restreindre aux lois admettant une densité par rapport à la mesure de Lebesgue. C'est-à-dire que nous allons estimer μ par une suite de lois $\hat{\mu}_n$ qui sont toutes des lois acceptant une densité par rapport à la mesure de Lebesgue.

38.7.1 Distance entre des mesures

Si ν_1 et ν_2 sont deux mesures de densité sur \mathbb{R} , la **distance** entre ν_1 et ν_2 est définie par

$$d(\nu_1, \nu_2) = \sup_{A \in \mathcal{B}(\mathbb{R})} |\nu_1(A) - \nu_2(A)| \quad (38.175)$$

où $\mathcal{B}(\mathbb{R})$ désigne l'ensemble des boréliens de \mathbb{R} .

Théorème 38.46 (de Scheffé[497]).

Si f_1 et f_2 sont les densités de ν_1 et ν_2 par rapport à la mesure de Lebesgue, alors

$$d(\nu_1, \nu_2) = \int_{\mathbb{R}} (f_1 - f_2)_+ = \frac{1}{2} \int_{\mathbb{R}} |f_1 - f_2| = \frac{1}{2} \|f_1 - f_2\|_1 \quad (38.176)$$

où f_+ est la partie positive de f (pour la décomposition $f(x) = f_+(x) - f_-(x)$).

Démonstration. La dernière égalité est simplement une notation usuelle ; nous devons seulement prouver les deux premières. Pour la première nous commençons par prouver que le borélien réalisant le supremum est

$$B = \{x \in \mathbb{R} \text{ tel que } f_1(x) \geq f_2(x)\}. \quad (38.177)$$

En effet si A est un borélien nous avons

$$\nu_1(A) - \nu_2(A) = \int_A f_1 - f_2 \leq \int_{A \cap B} f_1 - f_2 \leq \int_B f_1 - f_2 = \int_B (f_2 - f_2)_+ = \int_{\mathbb{R}} (f_1 - f_2)_+ \quad (38.178)$$

Justifications :

- $f_1 - f_2$ négative sur $A \cap \mathbb{C}B$.
- Vu que $f_1 - f_2 \geq 0$ sur B , l'intégrale augmente si on augmente le domaine.
- Sur B nous avons $f_1 - f_2 = (f_1 - f_2)_+$.

Donc pour tout borélien A nous avons

$$d(\nu_1, \nu_2) \leq \int_{\mathbb{R}} (f_1 - f_2)_+ \quad (38.179)$$

Mais pour $A = B$ nous avons égalité :

$$\nu_1(B) - \nu_2(B) = \int_B f_1 - f_2 = \int_B (f_1 - f_2)_+ = \int_{\mathbb{R}} (f_1 - f_2)_+ \quad (38.180)$$

Pour la seconde égalité nous savons que f_1 et f_2 s'intègrent toutes deux à 1 (parce que ce sont des densités de probabilité), donc

$$\int_{\mathbb{R}} f_1 - f_2 = 0. \quad (38.181)$$

En particulier nous avons

$$\int_{\mathbb{R}} (f_1 - f_2)_+ = \int_{\mathbb{R}} (f_1 - f_2)_-, \quad (38.182)$$

ce qui donne bien

$$\int_{\mathbb{R}} (f_1 - f_2)_+ = \frac{1}{2} \int_{\mathbb{R}} |f_1 - f_2|. \quad (38.183)$$

□

38.7.2 Estimateur par fenêtres glissantes

Nous considérons les estimations suivantes de la fonction de répartition :

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \leq x\}}, \quad (38.184)$$

et un nombre h_n qui sera la taille de la fenêtre glissante. Nous avons en tête de faire $\lim_{n \rightarrow \infty} h_n = 0$. Nous considérons ceci comme estimateur de la densité inconnue f des variables aléatoires X_i :

$$\hat{f}_n(x) = \frac{F_n(x + h_n) - F_n(x - h_n)}{2h_n}. \quad (38.185)$$

L'idée sous-jacente est de prendre la dérivée de la fonction de répartition comme densité.

Lemme 38.47 ([497]).

Pour tout $h_n > 0$, l'estimateur \hat{f}_n est une densité de probabilité.

Démonstration. D'abord \hat{f}_n est bien à valeurs positives ou nulle. Ensuite devons parler de son intégrale. Pour le numérateur nous avons

$$F_n(x + h_n) - F_n(x - h_n) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i \in B(x, h_n)}. \quad (38.186)$$

En réalité cette égalité est valable seulement presque partout parce qu'elle n'est pas valable en $x = x \pm h_n$, mais cela ne va pas nous ennuyer dans la mesure où nous avons dans l'idée d'intégrer cela sur \mathbb{R} . Avant de nous lancer dans l'intégrale nous remarquons que $X_i \in B(x, h_n)$ est la même chose que $x \in B(X_i, h_n)$, c'est-à-dire que

$$\{X_i \in B(x, h_n)\} = \{x \in B(X_i, h_n)\}. \quad (38.187)$$

Donc

$$\int_{\mathbb{R}} \hat{f}_n = \frac{1}{2h_n} \frac{1}{n} \sum_{i=1}^n \underbrace{\int_{\mathbb{R}} \mathbb{1}_{B(X_i, h_n)}}_{2h_n} = \frac{1}{2nh_n} \sum_{i=1}^n 2h_n = 1. \quad (38.188)$$

□

Lemme 38.48 ([497]).

L'estimateur \hat{f}_n est déjà pas mal parce que

$$\lim_{h_n \rightarrow 0} E(\hat{f}_n(x)) \rightarrow f(x) \quad (38.189)$$

pour presque tout $x \in \mathbb{R}$.

Démonstration. Nous commençons par nous rappeler le fait que $F_n(x)$ est un estimateur sans biais de $F(x)$ (proposition 38.35). Donc

$$E(\hat{f}_n(x)) = \frac{F(x + h_n) - F(x - h_n)}{2h_n}. \quad (38.190)$$

Nous devons prendre la limite de cela lorsque $h_n \rightarrow 0$, c'est-à-dire considérer la dérivée de F . Attention : rien ne dit que F soit dérivable, si ce n'est la proposition 18.24 qui indique qu'elle est dérivable presque partout avec f comme dérivée.

La limite $h_n \rightarrow 0$ dans (38.190) nous donne donc bien presque partout

$$\lim_{h_n \rightarrow 0} E(\hat{f}_n(x)) = f(x). \quad (38.191)$$

□

Proposition 38.49 ([497]).

Si la suite (h_n) est telle que $h_n \rightarrow 0$ et $nh_n \rightarrow \infty$, alors pour presque tout $x \in \mathbb{R}$ nous avons les convergences

$$\hat{f}_n(x) \xrightarrow{L^2(P)} f(x) \quad (38.192)$$

et

$$\hat{f}_n(x) \xrightarrow{P} f(x). \quad (38.193)$$

Démonstration. Il faut d'abord comprendre ce que signifie la convergence $L^2(P)$ pour presque tout x . Pour cela il faut comprendre que $\hat{f}_n(x)$ est en soi une variable aléatoire et est en réalité une fonction $\omega \mapsto \hat{f}_n(x, \omega)$. Ce que nous allons montrer est que pour presque tout x (maintenant fixé), cette variable aléatoire converge vers une constante (par rapport à ω) et que cette constante est $f(x)$.

La convergence $X_n \xrightarrow{L^2(P)} X$ signifie $E(|X_n - X|^2) \rightarrow 0$, c'est-à-dire

$$\int_{\Omega} |X_n(\omega) - X(\omega)|^2 dP(\omega) \rightarrow 0. \quad (38.194)$$

En faisant une décomposition biais-variance nous devons donc étudier

$$E\left[(\hat{f}_n(x) - f(x))^2\right] = E[\hat{f}_n(x) - f(x)]^2 + \text{Var}(\hat{f}_n(x) - f(x)) \quad (38.195)$$

Ici $f(x)$ doit être vue comme la variable aléatoire constante sur Ω . Par le lemme 38.48 et la proposition 37.25 le terme de biais converge vers zéro lorsque $n \rightarrow \infty$.

Pour traiter le terme de biais, nous savons déjà que

$$2nh_n \hat{f}_n(x) = \sum_{i=1}^n \mathbb{1}_{\{X_i \in B(x, h_n)\}}, \quad (38.196)$$

où le membre de droite (et donc aussi celui de gauche) est une variable aléatoire binomiale comptant le nombre de succès de l'expérience $X_i \in B(x, h_n)$ en n essais. Nous notons $p_{x,n} = P(X_i \in B(x, h_n))$. Si μ est la loi parente des X_i , alors

$$p_{n,x} = P(X_i \in B(x, h_n)) = \mu(B(x, h_n)) = F(x + h_n) - F(x - h_n) \quad (38.197)$$

où F est la fonction de répartition (parente) des X_i .

Alors la variance de ladite binomiale est donnée par (37.331), c'est-à-dire $np_{x,n}(1-p_{x,n})$. Nous avons alors

$$\text{Var}(2nh_n\hat{f}_n(x)) = bp_{x,n}(1-p_{x,n}) \quad (38.198)$$

et

$$\text{Var}(\hat{f}_n(x)) = \frac{1}{4n^2h_n^2}np_{x,n}(1-p_{x,n}). \quad (38.199)$$

Nous pouvons faire la majoration $t(1-t) \leq t$ qui est valable pour tout t et écrire

$$\text{Var}(\hat{f}_n(x)) \leq \frac{1}{4nh_n} \frac{p_{x,n}}{h_n}. \quad (38.200)$$

Le premier facteur tend vers zéro parce que nous avons supposé que $nh_n \rightarrow \infty$. Pour le second facteur, il faut remarquer que l'expression (38.197) nous donne presque partout

$$\lim_{h_n \rightarrow 0} \frac{p_{n,x}}{h_n} = 2f(x), \quad (38.201)$$

qui est constant et certainement borné.

Nous avons maintenant prouvé que pour presque tout x nous avons $\hat{f}_n(x) \xrightarrow{L^2(P)} f(x)$. Montrons que cela implique la convergence en loi, c'est-à-dire que pour tout $\eta > 0$, nous avons la limite

$$P(|\hat{f}_n(x) - f(x)| > \eta) \rightarrow 0. \quad (38.202)$$

Si cela n'était pas vrai, nous aurions un nombre $\eta_0 > 0$ et $\epsilon > 0$ tel que pour tout n à partir d'une certaine taille,

$$P(|\hat{f}_n(x) - f(x)|^2 > \eta_0) > \epsilon, \quad (38.203)$$

et en particulier en notant A l'événement $|\hat{f}_n(x) - f(x)|^2 > \eta_0^2$, $P(A) > \epsilon$. Alors

$$\int_{\Omega} |\hat{f}_n(x, \omega) - f(x, \omega)|^2 dP(\omega) \geq \int_A |\hat{f}_n(x, \omega) - f(x, \omega)|^2 dP(\omega) \geq \int_A \eta_0^2 = \eta_0^2 P(A). \quad (38.204)$$

Cela signifie que

$$\|\hat{f}_n(x) - f(x)\|_{L^2(P)} \geq \eta_0 P(A), \quad (38.205)$$

ce qui contredit la première convergence démontrée. \square

Note : l'hypothèse $nh_n \rightarrow \infty$ revient à dire que nous voulons que chaque boîte contienne de plus en plus d'observations. Si nous avons $nh_n \rightarrow 0$, alors avec n qui augmente, la majorité des boîtes deviendraient vides, ce qui reviendrait à une perte d'information.

38.8 Test d'hypothèses, prise de décision

38.8.1 Exemple : qualité des pièces d'usine

Une usine fabrique des composantes électronique garantis un an. Le constructeur ne veut pas accepter que plus de 5% des pièces tombent en panne avant un an.

Nous supposons que la durée de vie T d'une pièce soit une variable aléatoire suivant une loi exponentielle de paramètre λ (qui est l'inverse de la moyenne : $\theta = 1/\lambda$). Le fabricant veut donc s'assurer que

$$0.95 \leq P(T \geq 1), \quad (38.206)$$

ou encore

$$P(T \geq 1) = \int_0^{\infty} \frac{1}{\theta} e^{-x/\theta} dx = e^{-1/\theta} \geq 0.95, \quad (38.207)$$

donc le fabricant doit s'assurer que

$$\theta \geq \frac{1}{\ln\left(\frac{1}{0.95}\right)}. \quad (38.208)$$

Nous posons donc $\theta_0 = 19.5$ et nous prenons comme modèle de décision que si $\theta < \theta_0$, alors la chaîne de production doit être revue, et si $\theta > \theta_0$, alors l'usine peut continuer son travail.

Ce dont nous disposons n'est pas de θ , mais d'une estimation de θ à partir d'un échantillon. Cela étant il faudra aussi pouvoir estimer la probabilité de faire un mauvais choix.

38.8.2 Exemple : la résistance d'un fil

Un artisan a besoin d'un fil qui a une résistance à une traction de 100 g en moyenne. Si la résistance est trop faible, le fil casse ; si elle est trop grande, c'est trop rigide et ça ne convient pas.

Remarque 38.50.

Dans l'exemple précédent, avoir $\theta > \theta_0$ ne dérange pas. Si la durée de vie moyenne est de 2 ans, le directeur de l'usine ne sera pas malheureux. Ici par contre l'artisan cherche une valeur précise et a donc une borne vers le haut et vers le bas.

L'artisan reçoit un lot de fils et souhaite savoir s'il est conforme. Pour cela, il prend 4 fils au hasard et mesure une moyenne de 112 g. Est-ce que cela est cohérent avec une moyenne de 100 g ?

Nous faisons l'hypothèse que la résistance des fils suit une loi normale $\mathcal{N}(m, \sigma^2)$ avec m inconnu. Pour la simplicité nous supposons que σ est connu et vaut 10.

Nous devons prendre une décision entre deux hypothèses. L'hypothèse H_0 sera de dire que le lot a une résistance de 100 g et l'hypothèse alternative sera que le lot a une résistance différente.

Les 4 observations sont quatre variables aléatoires $(X_i)_{i=1,2,3,4}$, et le nombre 112 est une réalisation de la variable aléatoire

$$\bar{X}_4 = \frac{1}{4}(X_1 + X_2 + X_3 + X_4). \quad (38.209)$$

Nous supposons que H_0 est vraie, et nous calculons quelle est l'intervalle autour de $m = 100$ qui a 95% de chances de contenir la moyenne observée. Si 112 est dedans, nous acceptons H_0 et si 112 est hors de cet intervalle, nous refusons H_0 .

Compte tenu de l'hypothèse H_0 , nous avons

$$\frac{\bar{X}_4 - 100}{\frac{10}{\sqrt{4}}} = \frac{\bar{X}_4 - 100}{5} \sim \mathcal{N}(0, 1). \quad (38.210)$$

Nous commençons à connaître par cœur l'intervalle de confiance à 95% d'une loi normale centrée réduite ; le quantile est à 1.96, c'est-à-dire

$$P\left(-1.96 \leq \frac{\bar{X}_4 - 100}{5} \leq 1.96\right) = 0.95, \quad (38.211)$$

ou encore

$$P(\bar{X}_4 \in [90.2, 109.8]). \quad (38.212)$$

Il y a donc moins de 5% de chances que la moyenne de ces quatre fils tombent en dehors de l'intervalle $[90.2, 109.8]$. L'artisan doit donc rejeter l'hypothèse et considérer que le lot est mauvais.

La région

$$]-\infty, 90.2] \cup [109.8, \infty[\quad (38.213)$$

est la **région de rejet**, ou **région critique**.

Ici, le nombre 5% représente le risque de refuser H_0 alors qu'elle était vraie. Notons que nous ne pouvons pas calculer le risque d'accepter H_0 alors qu'elle est fautive. En effet, si H_0 est fautive, nous ne savons pas quelles sont les valeurs de \bar{X}_4 acceptables parce qu'il y a une infinité de possibilités pour m qui soient alternatifs à $m = 100$.

Évidemment si la vraie moyenne est $100 + 10^{-7}$, l'hypothèse H_0 sera acceptée, mais nous n'avons aucun moyen de savoir si elle est vraie ou non.

38.8.3 Vocabulaire et théorie

Nous avons un modèle d'échantillonnage paramétrique $(X_{\theta,1}, \dots, X_{\theta,n})$ de taille n et de paramètre inconnu θ , de loi parente μ_θ appartenant à une famille paramétrique de lois $(\mu_\theta)_{\theta \in \Theta}$.

Soient H_0 et H_1 deux ensembles disjoints tels que $\Theta = H_0 \cup H_1$. L'ensemble H_0 sera nommé **hypothèse nulle** et l'ensemble H_1 sera l'**hypothèse alternative**.

Pour l'exemple des fils, nous avons $H_0 = \{100\}$ et $H_1 = \mathbb{R} \setminus \{100\}$. Si une hypothèse est réduite à un singleton, nous parlons d'hypothèse **simple** et sinon c'est une hypothèse **composite** ou **multiple**. Faire un tests consiste à déterminer une région critique.

Définition 38.51.

Un **test** est une application mesurable δ qui à $(x_1, \dots, x_n) \in \mathbb{R}^n$ associe

$$\delta(x_1, \dots, x_n) \in \{0, 1\}. \quad (38.214)$$

Si $\delta(x_1, \dots, x_n) = 0$ on accepte l'hypothèse H_0 pour l'échantillon x_1, \dots, x_n , et si $\delta(x_1, \dots, x_n) = 1$, alors on rejette H_0 et on choisit H_1 . L'ensemble

$$W = \{(x_1, \dots, x_n) \in \mathbb{R}^n \text{ tel que } \delta(x_1, \dots, x_n) = 1\} \quad (38.215)$$

est la **région de rejet** ou la **région critique**.

L'ensemble $W = \delta^{-1}(1)$ est un borélien de \mathbb{R}^n parce que δ est mesurable. L'événement auquel nous sommes intéressés est l'événement

$$R = \{(X_{\theta,1}, \dots, X_{\theta,n}) \in W\}. \quad (38.216)$$

Exemple 38.52

Pour l'exemple de 38.8.2 nous avons

$$\delta(x_1, \dots, x_4) = \mathbb{1}_{\mathbb{C}[90.2, 109.8]}(\bar{x}_4). \quad (38.217)$$

△

38.8.4 Risque de première et seconde espèce

Le modèle de décision que nous avons introduit comprend deux façons de se tromper. Soit nous rejetons H_0 alors qu'elle est vraie (c'est le **risque de première espèce**), soit nous acceptons H_0 alors qu'elle est fautive (risque de **seconde espèce**). Nous pouvons formaliser ces concepts de la façon suivante.

Nous considérons un test de région critique W . Le risque de première espèce, noté α est la fonction

$$\begin{aligned} \alpha: H_0 &\rightarrow [0, 1] \\ \theta &\mapsto P((X_{\theta,1}, \dots, X_{\theta,n}) \in W). \end{aligned} \quad (38.218)$$

Il s'agit de la probabilité de rejeter H_0 alors qu'elle es vraie. Le risque de seconde espèce est la fonction

$$\begin{aligned} \beta: H_1 &\rightarrow [0, 1] \\ \theta &\mapsto P((X_{\theta,1}, \dots, X_{\theta,n}) \notin W). \end{aligned} \quad (38.219)$$

C'est la probabilité d'accepter H_0 alors qu'elle est fautive.

Définition 38.53.

Soit δ un test de région critique W . La **puissance** du test est la fonction

$$\begin{aligned} \eta: H_1 &\rightarrow [0, 1] \\ \theta &\mapsto P((X_{\theta,1}, \dots, X_{\theta,n}) \in W). \end{aligned} \quad (38.220)$$

La courbe d'efficacité du test est la fonction

$$\begin{aligned} h: \Theta &\rightarrow [0, 1] \\ \theta &\mapsto P((X_{\theta,1}, \dots, X_{\theta,n}) \notin W). \end{aligned} \quad (38.221)$$

La puissance d'un test est la probabilité de rejeter H_0 lorsque H_1 est vraie. Plus la puissance est grande, mieux c'est. La courbe d'efficacité du test est la probabilité d'accepter H_0 pour une certaine valeur de θ .

Soit un test δ . Une statistique $T_\theta = T_n(X_{\theta,1}, \dots, X_{\theta,n})$ est une **variable de décision** pour δ si $\mathbb{C}W$ peut s'écrire d'une des façons suivantes

$$\mathbb{C}W = \begin{cases} \{(x_1, \dots, x_n) \in \mathbb{R}^n \text{ tel que } T_n(x_1, \dots, x_n) < c\} & \text{test unilatéral à droite} \\ \{(x_1, \dots, x_n) \in \mathbb{R}^n \text{ tel que } T_n(x_1, \dots, x_n) > c\} & \text{test unilatéral à gauche} \\ \{(x_1, \dots, x_n) \in \mathbb{R}^n \text{ tel que } c_1 \leq T_n(x_1, \dots, x_n) < c_2\} & \text{test bilatéral.} \end{cases} \quad (38.222)$$

Le plus souvent la variable de décision sera la moyenne : $T_n(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i$. Les valeurs c, c_1, c_2 sont des **valeurs critiques**.

En ce qui concerne les notations, ici T_n représente la valeur mesurée sur un n -échantillon (d'où l'indice n) alors que T_θ est la valeur *théorique* de T lorsque θ est la vraie valeur du paramètre qu'on veut estimer.

Pour un test unilatéral à gauche, nous fixons la valeur critique c de telle manière à avoir

$$P(T_\theta > c) \leq \alpha. \quad (38.223)$$

Pour un test unilatéral à gauche, nous fixons c de telle manière à avoir

$$P(T_\theta < c) \leq \alpha \quad (38.224)$$

et pour un test bilatéral nous fixons c_1 et c_2 de telle façon à avoir

$$P(T_\theta > c_2) = P(T_\theta < c_1) \leq \frac{\alpha}{2}. \quad (38.225)$$

38.8.5 Modèle paramétrique de loi gaussienne

Soit un modèle statistique paramétrique de lois gaussiennes $\mathcal{N}(m, 1)$ de moyenne m inconnue avec $m \in \mathbb{R}^+$. Nous avons $\Theta = [0, \infty[$.

Nous observons un échantillon de taille $n = 36$. Avec un risque de première espèce de 5% nous voulons estimer l'hypothèse $H_0 = \{0\}$ contre l'hypothèse $H_1 =]0, \infty[$. De notre échantillon nous construisons la variable aléatoire

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad (38.226)$$

dans laquelle les X_i sont les éléments de l'échantillon, elles sont indépendantes et identiquement distribuées de loi $\mathcal{N}(m, 1)$ avec m inconnu.

Si \bar{X}_n est proche de zéro nous acceptons H_0 , sinon nous la rejetons. La région de rejet s'écrit donc sous la forme

$$W = \{(x_1, \dots, x_n) \in \mathbb{R}^n \text{ tel que } \bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i > u\} \quad (38.227)$$

dans lequel il faut fixer le u pour satisfaire au risque de première espèce de 5%. La contrainte est d'avoir

$$P(\bar{X}_n > u) = \alpha \quad (38.228)$$

si H_0 est vérifiée. Cela revient à dire que dans 5% des cas où H_0 est correcte, nous la rejetterons. Si H_0 est vraie alors \bar{X}_n est une moyenne de gaussiennes de moyennes m et nous avons

$$\frac{\bar{X}_n - m}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1) \quad (38.229)$$

avec $m = 0$ et $\sigma = 1$. L'équation (38.228) devient donc

$$\alpha = P\left(\frac{\bar{X}_n}{1/\sqrt{n}} > \frac{u}{1/\sqrt{n}}\right) = P(T > \sqrt{nu}) \quad (38.230)$$

où $T \sim \mathcal{N}(0, 1)$. Avec $n = 36$ et $\alpha = 5\%$ nous trouvons

$$u = \frac{1.645}{6} \simeq 0.274 \quad (38.231)$$

La règle de décision est donc la suivante : si $\bar{x}_n > 0.274$ alors nous rejetons H_0 , et sinon nous l'acceptons.

Calculons la puissance de ce test (définition 38.53). C'est la fonction donnée par

$$\begin{aligned} \eta: H_1 &\rightarrow \mathbb{R} \\ m &\mapsto P((X_{1,m}, \dots, X_{n,m}) \in W) = P\left(\frac{1}{n} \sum X_i > u\right). \end{aligned} \quad (38.232)$$

Dans ce calcul, les X_i sont d'une loi normale $\mathcal{N}(m, 1)$, et non $\mathcal{N}(0, 1)$. En retranchant m et en divisant par $1/\sqrt{n}$ nous trouvons

$$\eta(m) = P\left(\frac{\frac{1}{n} \sum X_i - m}{1/\sqrt{n}} > \frac{u - m}{1/\sqrt{n}}\right) \quad (38.233a)$$

$$= P(T > \sqrt{n}(u - m)) \quad (38.233b)$$

$$= P(T > 16.45 - 6m) \quad (38.233c)$$

$$= 1 - \Phi(16.45 - 6m) \quad (38.233d)$$

où Φ est la fonction de répartition de $\mathcal{N}(0, 1)$. La fonction η a les propriétés suivantes :

$$\lim_{m \rightarrow -\infty} \eta(m) = 0 \quad (38.234a)$$

$$\lim_{m \rightarrow \infty} \eta(m) = 1 \quad (38.234b)$$

$$\eta(0) = \frac{5}{100}. \quad (38.234c)$$

Remarque 38.54.

Si nous regardons $m = 0.001$, le risque de seconde espèce est quasiment de 90%. En effet le risque de seconde espèce est d'accepter H_0 alors qu'il est faux. Lorsque $m = 0.001$, l'hypothèse H_0 est fautive, mais la probabilité qu'on l'accepte est grande. D'ailleurs les conséquences de l'accepter à tort ne sont peut-être pas si grandes que cela.

38.9 Tests paramétriques

La proposition suivante montre le lien entre région de confiance et les tests.

Proposition 38.55.

Soit $\Lambda(X_1, \dots, X_n)$ une région de confiance par excès de niveau de confiance $1 - \alpha$. Alors il existe un tests pur de niveau α pour tester $H_0 = \{\theta_0\}$ de région de rejet

$$W_n = \{x = (x_1, \dots, x_n) \in \mathbb{R}^n \text{ tel que } \theta_0 \notin \Lambda(x_1, \dots, x_n)\}. \quad (38.235)$$

Démonstration. L'hypothèse sur Λ signifie qu'avec les observations (X_1, \dots, X_n) , il y a une forte probabilité (plus grande que $1 - \alpha$) que θ soit dans $\Lambda(X_1, \dots, X_n)$. Avec ou sans H_0 nous avons donc

$$P(\theta \in \Lambda) \geq 1 - \alpha. \quad (38.236)$$

Supposons maintenant l'hypothèse H_0 , alors

$$P((X_1, \dots, X_n) \in W_n) = P(\theta_0 \notin \Lambda(X_1, \dots, X_n)) \leq \alpha. \quad (38.237)$$

□

Remarque 38.56.

Soit W_n la région de confiance d'un test de niveau α pour tester $H_0 = \{\theta_0\}$. Alors

$$\Lambda = \{x \in \mathbb{R}^n \text{ tel que } x \notin W_n\} \quad (38.238)$$

est une région de confiance $1 - \alpha$ pour θ .

Exemple 38.57

Soit (X_1, \dots, X_n) un échantillon de loi parente $\mathcal{N}(\theta, 1)$ avec $\theta \in \{\theta_0, \theta_1\}$. Nous supposons $\theta_0 < \theta_1$. Nous voulons tester $H_0 = \{\theta_0\}$ contre $H_1 = \{\theta_1\}$. Nous proposons le test suivant. La variable de décision sera \bar{X}_n et la région de rejet sera

$$W = \{(x_1, \dots, x_n) \in \mathbb{R}^n \text{ tel que } \frac{1}{n} \sum_i x_i > \frac{\theta_1 + \theta_0}{2}\}. \quad (38.239)$$

- (1) Donner le risque de première espèce de ce test.
- (2) Soit $0 < \alpha < 1$. Pour quelle valeur de n le tests a-t-il un risque de première espèce égal à α ?
- (3) Donner la puissance du test.

Les réponses peuvent être exprimées en termes de la fonction de répartition F de la loi normale centrée réduite.

Le risque de première espèce est donné par

$$\alpha = P\left(\frac{1}{n} \sum_i X_i > \frac{\theta_0 + \theta_1}{2}\right) \quad (38.240)$$

où les X_i sont des variables aléatoires indépendantes et identiquement distribuées de loi parente $\mathcal{N}(\theta_0, 1)$. Cela est la probabilité d'être dans la région de rejet alors que l'hypothèse H_0 est vraie. La formule (38.240) se transforme en

$$\alpha = P\left(\frac{\frac{1}{n} \sum X_i - \theta_0}{1/\sqrt{n}} > \frac{\frac{\theta_0 + \theta_1}{2} - \theta_0}{1/\sqrt{n}}\right) \quad (38.241a)$$

$$= P(T > \sqrt{n} \frac{\theta_1 - \theta_0}{2}). \quad (38.241b)$$

En termes de la fonction de répartition nous avons alors

$$\alpha = 1 - F\left(\sqrt{n} \frac{\theta_1 - \theta_0}{2}\right) \quad (38.242)$$

Il s'agit maintenant de trouver le nombre n qui réalise cette égalité. Pour cela nous utilisons l'inverse F^{-1} de la fonction de répartition de la normale :

$$n = \left(\frac{2}{\theta_1 - \theta_0} F^{-1}(1 - \alpha)\right)^2. \quad (38.243)$$

Le risque de seconde espèce est la possibilité d'accepter H_0 lorsque H_1 est vraie, c'est-à-dire

$$\beta = P\left(\frac{1}{n} \sum_i X_i < \frac{\theta_0 + \theta_1}{2}\right) \quad (38.244)$$

sous l'hypothèse H_1 . Dans le calcul de (38.244) nous prenons donc $X_i \sim \mathcal{N}(\theta_1, 1)$. Le résultat est que

$$\beta = F\left(\sqrt{n} \frac{\theta_0 - \theta_1}{2}\right). \quad (38.245)$$

Remarque 38.58.

L'expression (38.243) diminue lorsque θ_0 et θ_1 s'éloignent, ce qui est normal : plus les nombres à discerner sont éloignés, moins l'échantillon à prendre pour réaliser le travail doit être grand.

Notons aussi que $\theta_0 - \theta_1 < 0$, par conséquent augmenter n diminue la valeur de

$$\beta = F\left(\sqrt{n}\frac{\theta_0 - \theta_1}{2}\right) \quad (38.246)$$

$\Delta < ++ >$

38.10 Tests d'adéquation

Soit X_1, \dots, X_n un échantillon de loi parente X finie prenant ses valeurs dans $\{a_1, \dots, a_k\}$. La loi de X est donnée par les nombres

$$p_i = P(X = a_i) \quad (38.247)$$

pour $i = 1, \dots, k$. Nous introduisons l'**effectif empirique**, la variable aléatoire N_i qui compte le nombre de fois que a_i est observé dans l'échantillon. La **fréquence empirique** est la variable aléatoire

$$F_i = \frac{N_i}{n}. \quad (38.248)$$

Nous savons que la loi de N_i est $\mathcal{B}(n, p_i)$, et la loi des grands nombres dit que

$$F_i \xrightarrow{p.s.} p_i \quad (38.249)$$

pour chaque i . Le théorème central limite nous indique de plus que

$$\frac{N_i - np_i}{\sqrt{np_i(1 - p_i)}} \xrightarrow{\mathcal{L}} T \sim \mathcal{N}(0, 1). \quad (38.250)$$

Nous considérons un cas où les p_i sont inconnus. Ils peuvent être approchés par $N_i \simeq np_i$. Le théorème de Pearson nous indique comment.

Théorème 38.59 (Théorème de Pearson).

Nous avons

$$\sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i} \xrightarrow{\mathcal{L}} K \sim \chi^2(k - 1) \quad (38.251)$$

où la distribution $\chi^2(l)$ est la somme des carrés de l gaussiennes centrées réduites indépendantes.

Démonstration. Nous commençons par le cas $k = 2$. Dans ce cas nous avons $N_2 = n - N_1$ et $p_1 + p_2 = 1$. La sommes que nous regardons est

$$\frac{(N_1 - np_1)^2}{np_1} + \frac{(N_2 - np_2)^2}{np_2} = \frac{(N_1 - np_1)^2}{np_1} + \frac{(N_1 - np_1)^2}{n(1 - p_1)} \quad (38.252a)$$

$$= \frac{(N_1 - np_1)^2}{np_1(1 - p_1)}. \quad (38.252b)$$

Étant donné que N_1 est une variable aléatoire binomiale nous avons

$$\frac{N_1 - np_1}{\sqrt{np_1(1 - p_1)}} \xrightarrow{\mathcal{L}} T \sim \mathcal{N}(0, 1). \quad (38.253)$$

Par conséquent la limite de (38.252b) est

$$\left(\frac{N_1 - np_1}{\sqrt{np_1(1 - p_1)}}\right)^2 \xrightarrow{\mathcal{L}} T^2 \simeq \chi^2(1). \quad (38.254)$$

Cela conclut le cas $k = 2$.

Passons à présent au cas général. Le k -uplet (N_1, \dots, N_k) est une variable aléatoire multinomiale de loi

$$\mathcal{M}(n; k; p_1, \dots, p_k). \quad (38.255)$$

Nous introduisons les variables aléatoires U_i données par $U_i: \Omega \rightarrow \mathbb{R}^k$ avec

$$P(U_i = (0, \dots, 1, \dots, 0)) = p_i; \quad (38.256)$$

c'est le vecteur aléatoire qui prend ses valeurs dans les vecteurs de la base canonique de \mathbb{R}^k et qui prend la valeur e_i avec probabilité p_i . Par construction nous avons

$$(N_1, \dots, N_k) = \sum_{i=1}^n U_i. \quad (38.257)$$

Nous allons étudier la fonction caractéristique de (N_1, \dots, N_k) définie par l'équation (37.201) :

$$\begin{aligned} \Phi_{(N_1, \dots, N_k)}: \mathbb{R}^k &\rightarrow \mathbb{C} \\ e_j &\mapsto E(e^{ie_j \cdot N}) = E(e^{iN_j}). \end{aligned} \quad (38.258)$$

Plus généralement,

$$\Phi_{(N_1, \dots, N_k)}(t_1, \dots, t_k) = E(e^{i\langle t, N \rangle_{\mathbb{R}^k}}). \quad (38.259)$$

Nous avons

$$e^{i\langle t, N \rangle} = e^{i\sum_j \langle t, U_j \rangle} = \prod_j e^{i\langle t, U_j \rangle} \quad (38.260)$$

et vu que les U_i sont indépendantes et identiquement distribuées nous pouvons écrire U_1 à la place de U_j de façon à avoir

$$\Phi_{(N_1, \dots, N_k)}(t_1, \dots, t_k) = \prod_j E(e^{i\langle t, U_1 \rangle}) \quad (38.261a)$$

$$= \prod_j \sum_l p_l e^{i\langle t, e_l \rangle} \quad (38.261b)$$

$$= \prod_j \sum_l p_l e^{it_l} \quad (38.261c)$$

$$= \left(\sum_{l=1}^k p_l e^{it_l} \right)^n. \quad (38.261d)$$

Nous allons montrer que

$$\lim_{n \rightarrow \infty} \Phi_{(N_1, \dots, N_n)}(t_1, \dots, t_k) = e^{-\frac{1}{2} \sum_{j=1}^k t_j^2} - \left(\sum_{j=1}^k t_j \sqrt{p_j} \right)^2. \quad (38.262)$$

Pour ce faire, nous allons effectuer un développement limité. D'abord nous introduisons les variables aléatoires

$$\alpha_j = \frac{N_j - np_j}{\sqrt{np_j}} \quad (38.263)$$

et nous calculons

$$\Phi_{(\alpha_1, \dots, \alpha_k)}(t_1, \dots, t_k) = E \left[\exp \left(i \left\langle t, \left(\frac{N_1 - np_1}{\sqrt{np_1}}, \dots, \frac{N_k - np_k}{\sqrt{np_k}} \right) \right\rangle \right) \right] \quad (38.264)$$

Étant donné que n et p_j sont des variables déterministes, nous pouvons les sortir de l'espérance. Nous avons alors

$$\Phi_{(\alpha_1, \dots, \alpha_k)}(t_1, \dots, t_k) = \exp \left(-i \sum_{j=1}^k t_j \sqrt{np_j} \right) \Phi_{(N_1, \dots, N_k)} \left(\frac{t_1}{\sqrt{np_1}}, \dots, \frac{t_k}{\sqrt{np_k}} \right) \quad (38.265)$$

parce que

$$E \left(e^{t_j \frac{N_j - np_j}{\sqrt{np_j}}} \right) = e^{\frac{-t_j np_j}{\sqrt{np_j}}} E \left(e^{t_j N_j / \sqrt{np_j}} \right) \quad (38.266)$$

En remplaçant (38.261) dans (38.265) nous trouvons

$$\Phi_{(\alpha_1, \dots, \alpha_k)}(t_1, \dots, t_k) = \exp \left(-i \sum_{j=1}^k t_j \sqrt{np_j} \right) \underbrace{\left(\sum_{j=1}^k p_j e^{i \frac{t_j}{\sqrt{np_j}}} \right)^n}_A \quad (38.267)$$

Nous analysons maintenant le terme A . Nous écrivons l'égalité $A = A + 1 - 1$ en tenant compte de $\sum_j p_j = 1$ sous la forme

$$A = \left(1 + \sum_{j=1}^k p_j (\exp(it_j / \sqrt{np_j}) - 1) \right)^n, \quad (38.268)$$

Nous avons alors

$$\ln(A) = n \ln \left[1 + \sum_{j=1}^{\infty} p_j (e^{it_j / \sqrt{np_j}} - 1) \right] \quad (38.269)$$

Nous développons l'exponentielle en

$$e^{it_j / \sqrt{np_j}} - 1 = \frac{it_j}{\sqrt{np_j}} - \frac{t_j^2}{2np_j} + \frac{1}{n} \epsilon(1/n) \quad (38.270)$$

et ensuite le logarithme selon la formule

$$\ln(1+x) = x - \frac{x^2}{2} + x^2 \alpha(x^2). \quad (38.271)$$

Nous avons

$$\ln(A) = n \ln \left[1 + \sum_j p_j \left(\frac{it_j}{\sqrt{np_j}} - \frac{t_j^2}{2np_j} + \frac{1}{n} \epsilon\left(\frac{1}{n}\right) \right) \right] \quad (38.272a)$$

$$= n \ln \left[1 + \sum_j p_j \left(it_j \sqrt{\frac{p_j}{n}} - \frac{t_j^2}{2n} + \frac{1}{n} \epsilon\left(\frac{1}{n}\right) \right) \right] \quad (38.272b)$$

$$= n \underbrace{\sum_{j=1}^k \left(it_j \sqrt{\frac{p_j}{n}} - \frac{t_j^2}{2n} + \frac{1}{n} \epsilon(1/n) \right)}_K \quad (38.272c)$$

$$- n \frac{1}{2} \left[\sum_j \left(it_j \sqrt{\frac{p_j}{n}} - \frac{t_j^2}{2n} + \frac{1}{n} \epsilon(1/n) \right)^2 \right] \quad (38.272d)$$

$$+ nK^2 \alpha(K) \quad (38.272e)$$

Nous introduisons dans ϵ tous les termes en $1/n^2$ et nous trouvons

$$\ln(A) = \sum_j \left(it_j \sqrt{p_j n} - \frac{t_j^2}{2} \right) - \frac{1}{2} \left(\sum_j it_j \sqrt{p_j} \right)^2 + \epsilon(1/n) + K^2 \alpha(K). \quad (38.273)$$

En remplaçant dans (38.267) et en passant à la limite pour $n \rightarrow \infty$,

$$\Phi_{(\alpha_1, \dots, \alpha_k)}(t_1, \dots, t_k) = \exp \left(-\frac{1}{2} \sum_j t_j^2 + \frac{1}{2} \left(\sum_j t_j \sqrt{p_j} \right)^2 \right). \quad (38.274)$$

Nous reconnaissons des lois gaussiennes dans le premier terme de l'exponentielle. Nous allons maintenant nous atteler à identifier le second terme.

Soit C une matrice orthogonale dont la dernière ligne est $(\sqrt{p_1}, \dots, \sqrt{p_k})$. Nous considérons les vecteurs

$$\alpha = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_k \end{pmatrix}, \quad t = \begin{pmatrix} t_1 \\ \vdots \\ t_k \end{pmatrix}. \quad (38.275)$$

et ensuite nous notons

$$U = Ct = \begin{pmatrix} u_1 \\ \vdots \\ u_k \end{pmatrix}. \quad (38.276)$$

Étant donné que C est orthogonale, nous avons $\sum_{i=1}^k \alpha_i^2 = \sum_{i=1}^k \beta_i^2$ et

$$\Phi_{(\beta_1, \dots, \beta_k)} = E(e^{i\langle u, \beta \rangle}) = E(e^{i\langle t, \alpha \rangle}) = \Phi_{(\alpha_1, \dots, \alpha_k)}(t_1, \dots, t_k). \quad (38.277)$$

Nous pouvons récrire l'argument de l'exponentielle (38.274) de la façon suivante :

$$\sum_{i=1}^k t_j^2 = \sum_j u_j^2 \quad (38.278a)$$

$$\sum_{j=1}^k t_j \sqrt{p_j} = (Ct)_k, \quad (38.278b)$$

Nous avons alors

$$\lim_{n \rightarrow \infty} \Phi_{(\beta_1, \dots, \beta_k)}(t_1, \dots, t_k) = \lim_{n \rightarrow \infty} \Phi_{(\alpha_1, \dots, \alpha_k)}(u_1, \dots, u_k) \quad (38.279a)$$

$$= \exp\left(\frac{1}{2} \sum_{j=1}^{k-1} u_j^2\right) \quad (38.279b)$$

$$= \Phi_{(Z_1, \dots, Z_{k-1}, 0)}(u_1, \dots, u_k) \quad (38.279c)$$

où les Z_i sont des variables aléatoires indépendantes et identiquement distribuées de distribution normale centrée réduite. Nous avons donc montré que

$$(\beta_1, \dots, \beta_k) \xrightarrow{\mathcal{L}} (Z_1, \dots, Z_{k-1}, 0). \quad (38.280)$$

Étant donné que l'application $x \mapsto \|x\|^2$ est continue, nous avons aussi

$$\|(\beta_1, \dots, \beta_k)\|^2 \xrightarrow{\mathcal{L}} \|(Z_1, \dots, Z_{k-1}, 0)\|^2, \quad (38.281)$$

et par conséquent

$$\|(\alpha_1, \dots, \alpha_k)\|^2 \xrightarrow{\mathcal{L}} \sum_{j=1}^{k-1} Z_j^2 \sim \chi^2(k-1). \quad (38.282)$$

D'après la définition (38.263) nous avons

$$\|(\alpha_1, \dots, \alpha_k)\|^2 = \sum_{j=1}^k \frac{(N_j - p_j)^2}{np_j}. \quad (38.283)$$

□

Chapitre 39

Chaînes de Markov à temps discret

Mets tes deux pieds en canard, c'est la chaîne de Markov qui se prépare.

39.1 Généralités

Les chaînes de Markov interviennent pour la description des systèmes dont l'évolution future ne dépend que de l'état présent.

Définition 39.1.

Soit E un ensemble au plus dénombrable et (Ω, \mathcal{F}, P) un espace probabilisé. Une **chaîne de Markov** à valeurs dans E est une famille $(X_n)_{n \in \mathbb{N}}$ de variables aléatoires telles que pour tout $x_0, \dots, x_{n+1} \in E$,

$$P(X_{n+1} = x_{n+1} | X_n = x_n, \dots, X_0 = x_0) = P(X_{n+1} = x_{n+1} | X_n = x_n). \quad (39.1)$$

Pour une chaîne de Markov, il n'est pas important de savoir l'historique pour prédire la futur : X_{n+1} est seulement déterminé par X_n .

Remarque 39.2.

Il existe une théorie des chaînes de Markov à temps continu ou avec E non dénombrable, mais ce n'est pas au programme.

Nous notons

$$p_n(x, y) = P(X_{n+1} = y | X_n = x) \quad (39.2)$$

la **probabilité de transition** de la chaîne à l'instant n . Si cette probabilité ne dépend pas de n , nous disons que la chaîne de Markov est **homogène**, et nous notons $p(x, y)$ au lieu de $p_n(x, y)$. Nous notons Q la matrice (éventuellement infinie) de transition

$$Q_{xy} = p(x, y). \quad (39.3)$$

Nous avons

$$\sum_{y \in E} p(x, y) = 1 \quad (39.4)$$

parce que c'est la somme de toutes les transitions possibles en partant de x . Notons aussi que $p(x, y) \geq 0$.

Définition 39.3.

Une matrice dont tous les éléments sont positifs ou nuls et donc la somme de toutes les lignes sont 1 est une **matrice stochastique**.

Notons que l'ensemble des matrices stochastiques est un fermé dans l'ensemble des matrices.

Lemme 39.4.

Si U est une matrice stochastique¹, alors il existe une chaîne de Markov dont la matrice de transition est U .

Remarque 39.5.

La somme $\sum_{x \in E} p(x, y)$ ne vaut pas spécialement 1. Si les états x_1 et x_2 arrivent tous les deux en y de façon certaine, alors nous avons $\sum_x p(x, y) \geq 2$. Il n'y a donc pas de limites aux sommes des colonnes.

Exemple 39.6

Nous considérons une fourmi qui se déplace dans un appartement à trois pièces A, B, C . Supposons qu'à chaque minute, elle a une probabilité $1/3$ de rester dans la pièce et une probabilité $2/3$ de se déplacer. Le plan de l'appartement est

$$A \longrightarrow B \longrightarrow C \quad (39.5)$$

De la pièce A est donc uniquement possible d'aller vers la pièce B ; de la B il est possible d'aller en A et en C et de la C il est uniquement possible d'aller en B .

La matrice de transition de cette chaîne de Markov est

$$Q = \begin{pmatrix} 1/3 & 2/3 & 0 \\ 1/3 & 1/3 & 1/3 \\ 0 & 2/3 & 1/3 \end{pmatrix} \quad (39.6)$$

△

Exemple 39.7

Si N_t est un processus de Poisson, alors les variables aléatoires $X_n = N_n$ forment une chaîne de Markov. △

39.2 Chaînes de Markov sur un ensemble fini

Une chaîne de Markov est **finie** si l'ensemble E dans lequel elle prend ses valeurs est fini.

Proposition 39.8 ([498]).

Si (X_n) est une chaîne de Markov irréductible sur un ensemble fini, alors pour tout ensemble $A \subset E$ nous avons

$$P(\tau_A < \infty) = \lim_{n \rightarrow \infty} P(\tau_A \leq n) = 1 \quad (39.7)$$

où $\tau_A = \min\{k \text{ tel que } X_k \in A\}$.

Les proposition à venir vont montrer que

- (1) Toute matrice stochastique admet un état stationnaire, proposition 39.9.
- (2) Si la chaîne de Markov est irréductible, alors il y a unicité de l'état stationnaire, proposition 39.10. Mais attention : cela ne veut pas encore dire que la chaîne converge effectivement vers cet état.
- (3) Si la chaîne est irréductible et apériodique, alors il y a convergence en loi vers l'unique loi invariante, théorème 39.13.

Proposition 39.9 ([499]).

Toute matrice stochastique admet un état stationnaire.

1. Définition 39.3.

Proposition 39.10 ([499]).

Soit une chaîne de Markov irréductible finie. Alors il existe une unique loi stationnaire π et de plus nous avons $\pi_i > 0$ pour tout état i de E .

Définition 39.11.

Une chaîne de Markov finie est **régulière** s'il existe un $n \in \mathbb{N}$ tel que P^n a uniquement des éléments strictement positifs.

Théorème 39.12 ([498]).

Soit P la matrice de transition d'une chaîne de Markov régulière sur un ensemble E de cardinal N . Alors il existe des nombres π_1, \dots, π_N tels que

- (1) $\pi_i > 0$ pour tout $i = 1, \dots, N$
- (2) $\pi_1 + \dots + \pi_N = 1$
- (3)

$$\lim_{n \rightarrow \infty} P^n = \Pi = \begin{pmatrix} \pi_1 & \pi_2 & \dots & \pi_N \\ \vdots & \vdots & & \vdots \\ \pi_1 & \pi_2 & \dots & \pi_N \end{pmatrix} \quad (39.8)$$

De plus le vecteur $\pi = (\pi_1, \dots, \pi_N)$ est l'unique solution de

$$\pi P = \pi. \quad (39.9)$$

Démonstration. Si la chaîne n'a qu'un seul état c'est évident parce que la probabilité de transition est toujours 1 ; fin de l'histoire.

Hypothèse Sinon nous supposons que P n'a que des éléments positifs, quitte à considérer P^n au lieu de P . Nous notons d le plus petit élément de P ; il vérifie $d \leq \frac{1}{2}$ parce que la somme des éléments d'une ligne de la matrice P doit être égale à 1.

Les suites min et max Soit x un vecteur quelconque (de composantes positives). Nous notons $m_0 = \min\{x_i\}$ et $M_0 = \max\{x_i\}$. Étant donné que les éléments du vecteur Px sont des moyennes pondérées des éléments de x , si nous posons

$$m_k = \min\{(P^k x)_i\}_{i=1, \dots, N} \quad (39.10a)$$

$$M_k = \max\{(P^k x)_i\}_{i=1, \dots, N}, \quad (39.10b)$$

la suite (m_k) est croissante et la suite (M_k) est décroissante.

Stricte croissance et décroissance Si $M_{k+1} = M_k$, alors toutes les composantes de $P^k x$ sont égales à M_k et le théorème est prouvé. Cela est encore une propriété de la moyenne. Même remarque pour la suite (m_k) .

Nous pouvons donc supposer que la suite (m_k) est strictement croissante et que la suite (M_k) est strictement décroissante. Elles sont toutes les deux bornées dans $[m_0, M_0]$. Le lemme 8.18 nous donne la convergence.

Égalité des limites Vu que les éléments de $P^k x$ ne sont pas tous les mêmes et s'évaluent de m_k à M_k , pour majorer M_{k+1} nous mettons le plus petit coefficient possible (c'est-à-dire d) devant m_k et nous supposons que toutes les autres composantes sont M_k ; nous avons alors

$$M_{k+1} \leq dm_k + (1-d)M_k \quad (39.11)$$

parce que tous les autres coefficients de la ligne contenant le d (dans P^k) sont plus petits ou égaux à $1-d$. De la même façon nous avons la minoration

$$m_{k+1} \geq dM_k + (1-d)m_k. \quad (39.12)$$

En faisant la différence, et en nous souvenant que $0 < 1-2d < 1$,

$$M_{k+1} - m_k \leq (1-2d)(M_k - m_k), \quad (39.13)$$

ce qui signifie que

$$M_{k+1} - m_k \leq (1 - 2d)^k (M_0 - m_0), \quad (39.14)$$

et donc que les deux limites sont égales.

Conclusion pour la limite Pour tout vecteur x , la suite $P^k x$ tend vers un vecteur dont toutes les composantes sont égales. En particulier pour le vecteur e_i de la base canonique,

$$P^k e_i \rightarrow \begin{pmatrix} \pi_1 \\ \vdots \\ \pi_1 \end{pmatrix}. \quad (39.15)$$

Mais $P^k e_i$ est la i^e colonne de la matrice P^k . Cela prouve la convergence annoncée $P^k \rightarrow \Pi$.

Réglons rapidement le cas des deux autres allégations du théorème. D'abord les matrices P^k sont toutes des matrices stochastiques; et l'ensemble des matrices stochastiques est fermé, donc la convergence se fait à l'intérieur de l'ensemble des matrices stochastiques. Cela prouve que $\pi_1 + \dots + \pi_N = 1$.

Ensuite la suite (m_k) étant strictement croissante et m_0 étant égal à 0 dans le cas de e_i nous avons toujours $\pi_i > 0$ (strictement). \square

Théorème 39.13 ([499]).

Si (X_n) est une chaîne de Markov finie, irréductible et apériodique de loi stationnaire π , alors

(1) La suite de matrices stochastiques P^k converge vers la matrice

$$P^k \rightarrow \Pi = \begin{pmatrix} \pi \\ \vdots \\ \pi \end{pmatrix}. \quad (39.16)$$

(2) Nous avons convergence des variables aléatoires au sens où

$$P(X_k = \mu P^k) \rightarrow \pi. \quad (39.17)$$

39.3 Marche aléatoire sur \mathbb{Z}

Soit (Y_n) une suite de variables aléatoires indépendantes et identiquement distribuées valant -1 avec une probabilité p et 1 avec une probabilité $(1 - p)$. La loi est

$$Y_n \sim p\delta_{-1} + (1 - p)\delta_1. \quad (39.18)$$

Nous considérons la variable aléatoire

$$X_n = X_0 + \sum_{i=1}^n Y_i \quad (39.19)$$

où X_0 est une variable aléatoire indépendante des Y_i à valeurs dans \mathbb{Z} . Nous vérifions à présent que X_n est une chaîne de Markov avec comme espace d'états $E = \mathbb{Z}$. Nous devons montrer que

$$P(X_{n+1} = x_{n+1} | X_n = x_n, \dots, X_0 = x_0) = P(X_{n+1} = x_{n+1} | X_n = x_n). \quad (39.20)$$

Pour ce faire nous allons exprimer tout cela en termes des Y_i au lieu des X_i . D'abord étant donné que nous avons égalité des événements

$$\{X_{n+1} = x_{n+1}\} \cap \{X_n = x_n, \dots, X_0 = x_0\} = \{Y_{n+1} = x_{n+1} - x_n\} \cap \{X_n = x_n, \dots, X_0 = x_0\}, \quad (39.21)$$

nous pouvons, en vertu du principe (37.76), remplacer $X_{n+1} = x_{n+1}$ par $Y_{n+1} = x_{n+1} - x_n$ dans le membre de gauche de (39.20). Nous avons donc déjà

$$P(X_{n+1} = x_{n+1} | X_n = x_n, \dots, X_0 = x_0) = P(\underbrace{Y_{n+1} = x_{n+1} - x_n}_A | \underbrace{X_n = x_n, \dots, X_0 = x_0}_B). \quad (39.22)$$

L'événement B est égal à l'événement

$$\{X_0 = x_0, Y_1 = x_1 - x_0, Y_2 = x_2 - x_1, \dots, Y_n = x_n - x_{n-1}\}, \quad (39.23)$$

qui n'est autre que l'ensemble

$$X_0^{-1}(x_0) \cap Y_1^{-1}(x_1 - x_0) \cap \dots \cap Y_n^{-1}(x_n - x_{n-1}) \quad (39.24)$$

qui est dans la tribu engendrée par les variables aléatoires $X_0, (Y_i)_{i=1, \dots, n}$. Le point délicat du raisonnement est de montrer que les événements A et B donnés par

$$A = \{Y_{n+1} = x_{n+1} - x_n\} \quad (39.25a)$$

$$B = \{X_0 = x_0\} \cap \bigcap_{i=1}^n \{Y_i = x_i - x_{i-1}\} \quad (39.25b)$$

sont indépendants. Nous ne pouvons pas montrer directement que $P(A \cap B) = P(A)P(B)$ parce que cela est la formule que nous voulons utiliser pour montrer que la chaîne est de Markov. Nous passons donc par les tribus :

$$A \in \sigma(Y_{n+1}) \quad (39.26a)$$

$$B \in \sigma(X_0, Y_1, \dots, Y_n). \quad (39.26b)$$

Nous utilisons maintenant l'hypothèse d'indépendance des variables aléatoires X_0 et Y_i pour conclure que les deux tribus des équations (39.26) sont indépendantes. Les événements A et B sont par conséquent indépendants.

L'événement A est indépendant de l'événement $\{X_n = x_n\}$. Nous avons donc successivement

$$P(X_{n+1} = x_{n+1} | X_n = x_n, \dots, X_0 = x_0) = P(Y_{n+1} = x_{n+1} - x_n | X_n = x_n, \dots, X_0 = x_0) \quad (39.27a)$$

$$= P(Y_{n+1} = x_{n+1} - x_n | Y_i = x_i - x_{i-1}, X_0 = x_0) \quad (39.27b)$$

$$= P(Y_{n+1} = x_{n+1} - x_n) \quad (39.27c)$$

$$= P(Y_{n+1} = x_{n+1} - x_n | X_n = x_n) \quad (39.27d)$$

$$= P(Y_{n+1} = x_{n+1} - X_n | X_n = x_n) \quad (39.27e)$$

$$= P(X_{n+1} = x_{n+1} | X_n = x_n). \quad (39.27f)$$

Justifications :

— (39.27c) parce que les tribus $\sigma(Y_{n+1})$ et $\sigma(Y_i, X_0)$ sont indépendantes.

— (39.27d) Nous avons

$$\{X_n = x_n\} \in \sigma(X_0, Y_1, \dots, Y_n) \quad (39.28)$$

tandis que

$$\{Y_{n+1} = x_{n+1} - x_n\} \in \sigma(Y_{n+1}); \quad (39.29)$$

ce sont donc deux événements issus de tribus indépendantes. Donc conditionner ou non l'événement $Y_{n+1} = x_{n+1} - x_n$ à l'événement $X_n = x_n$ ne change rien.

— (39.27e) est encore l'utilisation du fait que $P(A|B) = P(K|B)$ dès que $A \cap B = K \cap B$.

La chaîne est par conséquent de Markov.

La matrice de transition de cette chaîne de Markov est une matrice infinie « dans tous les sens » :

$$p(x, y) = \begin{cases} p & \text{si } y = x - 1 \\ (1 - p) & \text{si } y = x + 1 \\ 0 & \text{sinon.} \end{cases} \quad (39.30)$$

Remarque 39.14.

La plupart du temps lorsqu'il faut démontrer qu'une chaîne est de Markov, il faut suivre la procédure que nous venons de suivre pour la marche aléatoire sur \mathbb{Z} .

- Écrire tout en fonction des incréments.
- Dire que les incréments conditionnés sont indépendants des incréments qui conditionnent (via les tribus engendrées).
- Écrire que la probabilité cherchée est égale à l'événement conditionné dans lequel on a juste remplacé l'incrément par sa valeur.
- Conditionner à nouveau par rapport au dernier incrément qui est indépendant.
- Changer la valeur du dernier incrément par la variable aléatoire.

Dans ce raisonnement nous utilisons deux fois le fait que $P(A|B) = P(K|B)$ si $A \cap B = K \cap B$.

39.3.1 Chaînes de Markov homogènes

Proposition 39.15.

Voici quelques propriétés des chaînes de Markov homogènes.

(1) *La probabilité d'une trajectoire donnée est*

$$P(X_n = x_n, X_{n-1} = x_{n-1}, \dots, X_0 = x_0) = p(x_{n-1}, x_n) \dots p(x_0, x_1) P(X_0 = x_0). \quad (39.31)$$

(2) *La probabilité de transition « en n coups » est donnée par la puissance n ème de la matrice de transition :*

$$P(X_n = x_n | X_0 = x_0) = Q_{x_0, x_n}^n. \quad (39.32)$$

(3) *Si l'espace des états E est fini, l'espérance d'une fonction bornée $f: E \rightarrow \mathbb{R}$ de l'état est donnée par*

$$E(f(X_{n+1}) | X_n = x_n, \dots, X_0 = x_0) = E(f(X_{n+1}) | X_n = x_n) \quad (39.33a)$$

$$= \sum_{y \in E} f(y) p(x_n, y). \quad (39.33b)$$

Démonstration. Pour (1), étant donné que $P(A \cap B) = P(A|B)P(B)$, nous avons

$$P(X_n = x_n, \dots, X_0 = x_0) = P(X_n = x_n | X_{n-1} = x_{n-1}, \dots, X_0 = x_0) P(X_{n-1} = x_{n-1}, \dots, X_0 = x_0). \quad (39.34)$$

Par la propriété de Markov, le premier facteur est

$$P(X_n = x_n | X_{n-1} = x_{n-1}) = p(x_{n-1}, x_n). \quad (39.35)$$

Le reste est une récurrence sur n .

Pour (2). Montrons avec $n = 2$. En utilisant les divers points du théorème 37.30, nous avons

$$P(X_2 = x_2 | X_0 = x_0) = \sum_{y \in E} P(X_2 = x_2, X_1 = y | X_0 = x_0) \quad (39.36a)$$

$$= \sum_{y \in E} P(X_2 = x_2 | X_1 = y, X_0 = x_0) P(X_1 = y | X_0 = x_0) \quad (39.36b)$$

$$= \sum_{y \in E} P(X_2 = x_2 | X_1 = y) P(X_1 = y | X_0 = x_0) \quad (39.36c)$$

$$= \sum_{y \in E} p(x_2, y) p(y, x_0) \quad (39.36d)$$

$$= Q_{x_2, x_0}^2. \quad (39.36e)$$

Bien entendu ici la notion de produit matriciel doit être comprise de façon formelle lorsque E est infini.

Remarque 39.16.

Nous avons utilisé l'homogénéité de la chaîne de Markov au moment d'écrire l'expression (39.36d). En principe nous aurions dû écrire $p_2(y, x_2)p_1(x_0, y)$.

Pour (3), ...

□

39.3.2 Graphe de transition

Le **graphe de transition** d'une chaîne de Markov est le graphe dont les sommets sont les éléments de l'espace des états de la chaîne et dont les sommets sont reliés par des arrêtes pondérées par la probabilité de transition correspondante.

Définition 39.17.

Une chaîne de Markov est **irréductible** si pour tout $x, y \in E$, il existe n tel que $p^n(x, y) > 0$ où

$$p^n(x, y) = P(X_n = y | X_0 = x). \quad (39.37)$$

Le nombre n peut dépendre de x et y .

Lemme 39.18.

Une chaîne de Markov homogène est irréductible si et seulement si son graphe de transition est connexe.

Démonstration. Pour chaque couple $(x, y) \in E^2$ nous avons

$$\begin{aligned} p^n(x, y) &= \sum_{z_i \in E} P(X_n = y, X_{n-1} = z_{n-1}, \dots, X_1 = z_1, X_0 = x) \\ &= \sum_{z_i} p(z_{n-1}, y) p(z_{n-2}, z_{n-1}) \dots p(z_1, z_2) p(x, z_1). \end{aligned} \quad (39.38)$$

La positivité d'un des termes de la somme signifie que le graphe est connexe tandis que la positivité de $p^n(x, y)$ signifie que la chaîne est irréductible. □

39.3.3 Chaîne de Markov définie par récurrence

39.3.3.1 Le cas général

Proposition 39.19.

Soit X_0 une variable aléatoire à valeurs dans E , un ensemble au plus dénombrable. Soit (Y_n) une suite de variables aléatoires réelles indépendantes et identiquement distribuées indépendantes de X_0 .

Soit (X_n) la suite de variables aléatoires à valeurs dans E définie par récurrence selon la formule

$$X_{n+1} = G(X_n, Y_{n+1}) \quad (39.39)$$

où $G: E \times \mathbb{R} \rightarrow E$ est une fonction mesurable. Alors (X_n) est une chaîne de Markov.

Démonstration. Soient x_0, \dots, x_{n+1} des éléments de E . Nous devons calculer la valeur de

$$P(X_{n+1} = x_{n+1} | X_n = x_n, \dots, X_0 = x_0). \quad (39.40)$$

Commençons par préciser les espaces sur lesquels nos variable aléatoires sont définies. Nous avons

$$X_0: \Omega_0 \rightarrow E \quad (39.41)$$

et

$$Y_i: \Omega \rightarrow \mathbb{R}. \quad (39.42)$$

La variable aléatoire X_1 est donnée par

$$\begin{aligned} X_1: \Omega_0 \times \Omega &\rightarrow E \\ (\omega_0, \omega_1) &\mapsto G(X_0(\omega_0), Y_1(\omega_1)). \end{aligned} \quad (39.43)$$

La variable aléatoire X_2 est

$$\begin{aligned} X_2: \Omega_0 \times \Omega^2 &\rightarrow E \\ (\omega_0, \omega_1, \omega_2) &\mapsto G(X_1(\omega_0, \omega_1), Y_2(\omega_2)) \\ &= G(G(X_0(\omega_0), \omega_1), Y_2(\omega_2)) \end{aligned} \quad (39.44)$$

et ainsi de suite.

Considérons maintenant l'événement

$$\{X_1 = x_1, X_0 = x_0\} \subset \Omega_0 \times \Omega. \quad (39.45)$$

Il est donné explicitement par

$$\{X_1 = x_1, X_0 = x_0\} = \{(\omega_0, \omega_1) \text{ tel que } G(X_0(\omega_0), Y_1(\omega_1)) = x_1, X_0(\omega_0) = x_0\} \quad (39.46a)$$

$$= \{(\omega_0, \omega_1) \text{ tel que } G(x_0, Y_1(\omega_1)) = x_1, X_0(\omega_0) = x_0\} \quad (39.46b)$$

$$= \{\omega_0 \in \Omega_0 \text{ tel que } X_0(\omega_0) = x_0\} \times \{\omega_1 \in \Omega \text{ tel que } G(x_0, Y_1(\omega_1)) = x_1\}. \quad (39.46c)$$

Le premier terme du produit cartésien est dans $\sigma(X_0)$, tandis que le second est dans $\sigma(Y_1)$. Étant donné la définition des tribus produit (définition 15.112) nous avons

$$\{X_1 = x_1, X_0 = x_0\} \in \sigma(X_0, Y_1). \quad (39.47)$$

Ce raisonnement se généralise immédiatement et nous trouvons que

$$\{X_n = x_n, \dots, X_0 = x_0\} \in \sigma(X_0, Y_1, \dots, Y_n). \quad (39.48)$$

Nous sommes donc à calculer

$$\diamond = P(X_{n+1} = x_{n+1} | X_n = x_n, \dots, X_0 = x_0) \quad (39.49a)$$

$$= P\left(\underbrace{G(x_n, Y_{n+1}) = x_{n+1}}_{\in \sigma(Y_{n+1})} \mid \underbrace{X_n = x_n, \dots, X_0 = x_0}_{\in \sigma(X_0, Y_1, \dots, Y_n)}\right). \quad (39.49b)$$

Les tribus $\sigma(Y_{n+1})$ et $\sigma(X_0, Y_1, \dots, Y_n)$ étant indépendantes nous avons

$$\diamond = P(G(x_n, Y_{n+1}) = x_{n+1}) \quad (39.50a)$$

$$= P(G(x_n, Y_{n+1}) = x_{n+1} | X_n = x_n) \quad (39.50b)$$

$$= P(G(X_n, Y_{n+1}) = x_{n+1} | X_n = x_n) \quad (39.50c)$$

$$= P(X_{n+1} = x_{n+1} | X_n = x_n). \quad (39.50d)$$

Pour (39.50b) nous avons utilisé le fait que $\sigma(Y_{n+1})$ est indépendante de $\sigma(X_n)$. Nous avons prouvé que la chaîne était de Markov. \square

Les probabilités de transition de la chaîne de Markov définie dans la proposition 39.19 sont

$$P(X_1 = y | X_0 = x) = P(G(X_0, Y_1) = y | X_0 = x_0) = P(G(x_0, Y_1) = y). \quad (39.51)$$

39.3.3.2 Exemple : la file de réparation de machines à laver

Nous considérons un magasin de réparation d'électroménager. Durant le jour n , un nombre aléatoire Z_n de machines en panne arrivent au magasin. Une machine est réparée chaque jour (aucune si le magasin est vide). Nous supposons que les Z_n soient indépendantes et identiquement distribuées, et nous posons X_n , le nombre de machines en magasin le jour n .

La loi d'avancement de X_n est

$$X_{n+1} = \begin{cases} X_n + Z_n - 1 & \text{si } X_n \neq 0 \\ Z_n & \text{si } X_n = 0. \end{cases} \quad (39.52)$$

Cela est une chaîne de Markov en vertu de la proposition 39.19. Ici la fonction est

$$G(x, y) = x + y - \mathbb{1}_{x \neq 0}. \quad (39.53)$$

Les probabilités de transitions sont

$$p(x, y) = \begin{cases} 0 & \text{si } x \leq y - 2 \\ P(Z = 0) & \text{si } x = y - 1 \\ P(Z = k) & \text{si } x = y + k - 1 \end{cases} \quad (39.54)$$

pour $x \neq 0$.

Exemple 39.20

Soit (X_n) une chaîne de Markov de matrice de transition

$$P = \begin{pmatrix} 0.2 & 0.5 & 0.3 \\ 0.1 & 0.1 & 0.8 \\ 0.5 & 0.2 & 0.3 \end{pmatrix}. \quad (39.55)$$

Calculer $P(X_3 = 1 | X_0 = 1)$ et $P(X_7 = 0 | X_4 = 0)$.

Déterminer, s'il en existe, une loi stationnaire vers laquelle converge la chaîne.

Nous avons

$$P^3 = \begin{pmatrix} 0.344 & 0.251 & 0.405 \\ 0.283 & 0.307 & 0.41 \\ 0.287 & 0.248 & 0.465 \end{pmatrix}. \quad (39.56)$$

La probabilité d'aller de l'état 1 à l'état 1 en trois étapes est donc 0.307. La chaîne étant de Markov, sans mémoire, les probabilités entre les temps 4 et 7 sont les mêmes qu'entre 0 et 3. Nous avons alors

$$P(X_7 = 0 | X_4 = 0) = 0.344. \quad (39.57)$$

La chaîne est irréductible et n'a pas d'états absorbants.

△

39.4 Classification des états

Sauf mention expresse du contraire, nous considérons toujours une chaîne de Markov homogène.

Définition 39.21.

Un état $x \in E$ est **absorbant** pour la chaîne (X_n) si $p(x, x) = 1$.

Il n'est pas spécialement impossible d'arriver sur un état absorbant, mais il est impossible d'en sortir.

Si $x \in E$, nous notons

$$T(x) = \inf\{k \geq 1 \text{ tel que } X_k = x\}, \quad (39.58)$$

le **premier temps d'atteinte** de l'état x . Si $X_0 = x$, alors $T(x)$ est le **temps de retour** en x . Si $p \in \mathbb{N}$ nous notons

$$T_p(x) = \inf\{k \geq 1 \text{ tel que } X_{k+p} = x\}. \quad (39.59)$$

C'est le temps mis pour atteindre x à partir de l'instant p .

Proposition 39.22.

La loi de la variable aléatoire $[T_p(x)|X_p = x]$ est la même que celle de la variable aléatoire $[T(x)|X_0 = x]$.

Démonstration. Nous devons montrer que

$$P(T_p(x) = k | X_p = x) = P(T(x) = k | X_0 = x). \quad (39.60)$$

Cela est intuitivement évident du fait qu'une chaîne de Markov soit un processus sans mémoire. Afin de prouver, nous allons sommer sur tous les états intermédiaires possibles :

$$P(T_p(x) = k | X_0 = x) = P(X_{p+k} = x, X_{p+k-1} \neq x, \dots, X_{p+1} \neq x | X_p = x) \quad (39.61a)$$

$$= \sum_{z_i \neq x} P(X_{p+k} = x, X_{p+k-1} = z_{k-1}, \dots, X_{p+1} = z_1 | X_p = x) \quad (39.61b)$$

$$= \sum_{z_i} P(X_{p+k} = x, X_{p+k-i} = z_i | X_{p+1} = z_1, X_p = x) \underbrace{P(X_{p+1} = z_1 | X_p = x)}_{=p(x, z_1)} \quad (39.61c)$$

$$= \sum_{z_i} P(X_{p+k} = x, X_{p+k-i} = z_i | X_{p+2} = z_2, X_{p+1} = z_1, X_p = x) \quad (39.61d)$$

$$\underbrace{P(X_{p+2} = z_2 | X_{p+1} = z_1, X_p = x)}_{P(X_{p+2}=z_2 | X_{p+1}=z_1)=p(z_1, z_2)} p(x, z_1) \quad (39.61e)$$

$$= \dots \quad (39.61f)$$

$$= \sum_{z_i} p(x, z_1) p(z_1, z_2) \dots p(z_{k-1}, z_{k-1}) p(z_{k-1}, x). \quad (39.61g)$$

À ce point ci, nous avons éliminé toute référence à p grâce à l'homogénéité de la chaîne. Nous pouvons refaire le calcul à l'envers pour reconstituer l'expression de départ sans le p :

$$\sum_{z_i} p(x, z_1) p(z_1, z_2) \dots p(z_{k-1}, z_{k-1}) p(z_{k-1}, x) \quad (39.62a)$$

$$= P(x_k = x, X_{k-1} \neq x, \dots, X_1 \neq x | X_0 = x) \quad (39.62b)$$

$$= P(T(x) = k), \quad (39.62c)$$

ce qu'il fallait obtenir. □

Définition 39.23.

Un état x est **récurrent** si $P(T(x) = \infty | X_0 = x) = 0$, c'est-à-dire si la probabilité de ne jamais retourner en x lorsqu'on y est passé est nulle. L'état x est **transient** ou **transitoire** dans le cas contraire.

Si x est un état récurrent, et si $E(T(x) | X_0 = x) < \infty$, nous disons que x est **récurrent positif**. Si $E(T(x) | X_0 = x) = \infty$ alors nous disons que est **récurrent nul**.

Nous introduisons une variable aléatoire qui compte le nombre de fois que la chaîne de Markov passe par l'état x :

$$N_x = \sum_{k=0}^{\infty} \mathbb{1}_{\{X_k = x\}}. \quad (39.63)$$

C'est une variable aléatoire à valeurs dans $\mathbb{N} \cup \{\infty\}$.

Proposition 39.24.

Les deux propriétés suivantes sont équivalentes à dire que x est récurrent :

- (1) $P(N_x < \infty | X_0 = x) = 0$
- (2) $E(N_x | X_0 = x) = \infty$.

Les deux propriétés suivantes sont équivalentes à dire que x est transient :

- (1) $P(N_x < \infty | X_0 = x) = 1$
- (2) $E(N_x | X_0 = x) < \infty$.

Démonstration. En tant que événements, nous avons l'égalité

$$N_x < \infty = \bigcup_{n \in \mathbb{N}} \underbrace{\{X_n = x, X_{n+k} \neq x \forall k \geq 1\}}_{F_n}. \quad (39.64)$$

Nous avons donc

$$P(N_x < \infty | X_0 = x) = \sum_{n=0}^{\infty} P(F_n | X_0 = x), \quad (39.65)$$

et

$$P(F_n | X_0 = x) = P(X_{n+k} \neq x, \forall k \geq 1, X_n = x | X_0 = x) \quad (39.66a)$$

$$= P(X_{n+k} \neq x, k \geq 1 | X_n = x, X_0 = x) P(X_n = x | X_0 = x) \quad (39.66b)$$

$$= P(X_{n+k} \neq x, k \geq 1 | X_n = x) P(X_n = x | X_0 = x) \quad (39.66c)$$

$$= P(X_k \neq x, k \geq 1 | X_0 = x) P(X_n = x | X_0 = x) \quad (39.66d)$$

$$= P(T(x) = \infty | X_0 = x) P(X_n = x | X_0 = x) \quad (39.66e)$$

Justifications :

- (1) Pour (39.66c), nous utilisons le fait que la chaîne soit « sans mémoire ».
- (2) Pour (39.66d), nous utilisons le fait que la chaîne soit homogène.
- (3) Pour (39.66e), l'événement $X_k \neq x$ pour tout $k \geq 1$ est exactement l'événement $T(x) = \infty$.

En nous servant de la proposition 15.263 (théorème de Fubini et mesure de comptage), nous permutons l'espérance et la somme dans l'expression

$$\sum_{n=0}^{\infty} P(X_n = x | X_0 = x) = \sum_{n=0}^{\infty} E(\mathbb{1}_{\{X_n=x\}} | X_0 = x) \quad (39.67a)$$

$$= E\left(\sum_{n=0}^{\infty} \mathbb{1}_{\{X_n=x\}} | X_0 = x\right) \quad (39.67b)$$

$$= E(N_x | X_0 = x). \quad (39.67c)$$

Voyons ce passage plus en détail. D'abord, en général nous avons

$$E(Y | X = x_0) = \int_{\{X=x_0\}} Y(\omega) dP(\omega) = \int_{\Omega} \mathbb{1}_{\{X=x_0\}}(\omega) Y(\omega) dP(\omega). \quad (39.68)$$

Dans notre cas,

$$E(\mathbb{1}_{\{X_n=x\}} | X_0 = x) = \int_{\Omega} \mathbb{1}_{X_0=x}(\omega) \mathbb{1}_{\{X_n=x\}}(\omega) dP(\omega). \quad (39.69)$$

La fonction qui correspond à la proposition 15.263 est

$$f(n, \omega) = f_n(\omega) = \delta_{X_0(\omega), x} \delta_{X_n(\omega), x}, \quad (39.70)$$

qui est bien une fonction positive et mesurable.

Nous reprenons à présent le calcul (39.65) en remplaçant les éléments par leurs valeurs que nous avons calculées :

$$P(N_x < \infty | X_0 = x) = P(T(x) = \infty | X_0 = x)E(N_x | X_0 = x). \quad (39.71)$$

Si x est récurrent, nous avons $P(T(x) = \infty | X_0 = x) = 0$, mais la relation (39.71) ne permet pas de conclure que le membre de gauche est nul parce qu'il reste la possibilité que $E(N_x | X_0 = x) = \infty$. Nous devons donc faire un pas en arrière et écrire cette espérance comme la limite des sommes partielles :

$$P(N_x < \infty | X_0 = x) = \lim_{N \rightarrow \infty} \sum_{n=0}^N P(T(x) = \infty | X_0 = x)P(X_n = x | X_0 = x) = 0 \quad (39.72)$$

parce que tous les termes de la suite des sommes partielles sont nuls. Nous avons donc bien que $P(N_x < \infty | X_0 = x) = 0$. Il s'ensuit immédiatement que $E(N_x | X_0 = x) = 1$.

Nous devons maintenant démontrer l'implication inverse. Supposons que $P(N_x < \infty | X_0 = x) = 0$. Dans ce cas nous avons immédiatement $P(N_x = \infty | X_0 = x) = 1$ et $E(N_x | X_0 = x) = \infty$. L'équation (39.71) nous indique alors que

$$P(T(x) = \infty | X_0 = x) = 0, \quad (39.73)$$

c'est-à-dire que x est récurrent. □

39.4.1 Chaînes irréductibles

Proposition 39.25.

Soit (X_n) une chaîne de Markov irréductible.

- (1) Un état x est récurrent si et seulement si tous les états sont récurrents.
- (2) Un état x est transient si et seulement si tous les états sont transients.

Démonstration. Soient x et y des états de la chaîne de Markov. Nous devons tester la valeur de $P(X_n = y | X_0 = y)$. Afin d'exploiter l'hypothèse d'irréductibilité, nous considérons $r, s \in \mathbb{N}$ tels que

$$p^r(x, y) > 0 \quad (39.74a)$$

$$p^s(y, x) > 0 \quad (39.74b)$$

et nous calculons majorons en passant par quelques intermédiaires :

$$P(X_{n+r+s} = y | X_0 = y) \geq P(X_{n+r+s} = y, X_{n+s} = x, X_s = x | X_0 = y) \quad (39.75a)$$

$$= P(X_{n+r+s} = y | X_{n+s} = x, X_s = x, X_0 = y) \quad (39.75b)$$

$$P(X_{n+s} = x | X_s = x, X_0 = y)P(X_s = x | X_0 = y).$$

Les deux premiers facteurs se calculent en utilisant la propriété de Markov et l'homogénéité de la chaîne. Pour le premier,

$$P(X_{n+s} = x | X_s = x, X_0 = y) = P(X_{n+s} = x | X_s = x) = P(X_n = x | X_0 = x). \quad (39.76)$$

Nous avons donc

$$\sum_{n \in \mathbb{N}} P(X_{n+r+s} = y | X_0 = y) \geq p^r(x, y)p^s(y, x) \sum_{n \in \mathbb{N}} P(X_n = x | X_0 = x). \quad (39.77)$$

En réutilisant Fubini comme dans l'équation (39.67), nous avons

$$\sum_{n \in \mathbb{N}} P(X_{n+r+s} = y | X_0 = y) \geq KE(N_x | X_0 = x) \quad (39.78)$$

où K est une constante strictement positive, par hypothèse d'irréductibilité de la chaîne de Markov.

Si x est un état récurrent, alors le membre de gauche est infini par la proposition (39.24) et donc

$$\sum_{n \in \mathbb{N}} P(X_{n+r+s} = y | X_0 = y) = \infty. \quad (39.79)$$

Aux $r + s$ premiers termes près (qui ne changent pas la somme), nous avons

$$\sum_{n \in \mathbb{N}} P(X_n = y | X_0 = y) = \infty, \quad (39.80)$$

ce qui signifie que y est récurrent. □

Nous rappelons que $T(x)$ est le temps que première atteinte de l'état x . Nous notons

$$\pi(x) = \frac{1}{E(T(x) | X_0 = x)}. \quad (39.81)$$

Étant donné que $T(x)$ est un entier positif ou nul nous avons $E(T(x) | X_0 = x) \in [1, \infty]$ et donc $\pi(x) \in [0, 1]$.

Si x est un état transiet, alors $T(x) = \infty$ lorsque $X_0 = x$ et donc $E(T(x) | X_0 = x) = 0$ et $\pi(x) = 0$. Si x est récurrent par contre, $P(T(x) < \infty | X_0 = x) = 1$ et il n'y a pas de garanties sur la valeur de $E(T(x) | X_0 = x)$.

Corollaire 39.26.

Un état récurrent est récurrent positif si et seulement si $\pi(x) > 0$. Un état récurrent est récurrent nul si et seulement si $\pi(x) = 0$.

Démonstration. C'est la formule (39.81). □

Proposition 39.27.

Soit (X_n) est une chaîne de Markov irréductible.

- (1) *Si x est un état récurrent, alors $T(x) < \infty$ presque surement.*
- (2) *Nous avons une égalité entre les lois*

$$\mathcal{L}(X_{k+T(x)} | T(x) < \infty) = \mathcal{L}(X_k | X_0 = x). \quad (39.82)$$

39.4.2 Nombre de visites

La fonction

$$\frac{1}{n} \sum_{k=1}^n \mathbb{1}_{\{X_k = x\}} \quad (39.83)$$

est la **fréquence empirique** de la chaîne de Markov.

Soit x un état récurrent, c'est-à-dire que $P(T(x) < \infty | X_0 = x) = 1$. Nous classons les visites de la façon suivante :

$$T_1(x) = T(x) = \inf\{k \geq 1 \text{ tel que } X_k = x\} \quad (39.84a)$$

$$T_2(x) = \inf\{k \geq 1 \text{ tel que } X_{T_1(x)+k} = x\} \quad (39.84b)$$

$$\vdots \quad (39.84c)$$

$$T_n(x) = \inf\{k \geq 1 \text{ tel que } X_{T_{n-1}(x)+k} = x\} \quad (39.84d)$$

La variable aléatoire T_i représente le temps entre la visite numéro $i - 1$ et la visite numéro i (si $X_0 \neq x$, sinon il faut décaler). Nous définissons l'instant la n visite numéro n :

$$S_n = \sum_{k=1}^n T_k(x). \quad (39.85)$$

Lemme 39.28.

Les variables aléatoires T_i sont indépendantes.

Démonstration. Nous choisissons n des T_i et nous calculons la probabilité

$$\spadesuit = P(T_{i_1} = k_1, T_{i_2} = k_2, \dots, T_{i_n} = k_n) \quad (39.86)$$

où nous supposons $i_1 > i_2 > \dots > i_n$. Nous décomposons cette probabilité en sommant sur toutes les histoires de la chaîne de Markov compatibles avec les nombres k_i donnés :

$$\spadesuit = \sum_{\substack{\{z_j\} \\ \text{compatibles}}} P(X_j = z_j, j = 1, \dots, N). \quad (39.87)$$

Notons qu'ici, le numéro du dernier terme de la somme n'est pas certain parce que tous les T_i ne sont pas fixés. Nous l'avons noté N , mais en réalité il est différent d'un terme à l'autre de la somme. Il est certain que $z_N = x$ et $z_{N-k_1} = x$ et si $N - k_1 < j < N$, alors $z_j \neq x$. Cela est simplement le fait que nous demandions aux z_i de respecter les conditions données par les k_i . Nous avons

$$\spadesuit = \sum_{\{z_j\}} P(X_N = x, X_j = z_j, N - k_1 < j < N | X_j = z_j, j \leq N - k_1) P(X_j = z_j, j < N - k_1) \quad (39.88a)$$

$$= \sum_{\{z_j\}} P(X_N = x, X_j = z_j, N - k_1 < j < N | X_{N-k_1} = x) P(X_j = z_j, j < N - k_1) \quad (39.88b)$$

$$(39.88c)$$

Le premier facteur est $P(T_{i_1} = k_1)$ tandis que le second facteur est précisément $P(T_j = k_j, j > 1)$. Nous avons donc montré que

$$P(T_{i_1} = k_1, T_{i_2} = k_2, \dots, T_{i_n} = k_n) = P(T_{i_1} = k_1) P(T_j = k_j, j > 1), \quad (39.89)$$

et donc les T_i sont indépendants. \square

Proposition 39.29.

Si (X_n) est une chaîne de Markov irréductible et si $x \in E$ alors

$$\pi(x) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{\{X_k = x\}} \quad (39.90)$$

presque sûrement.

Démonstration. Étant donné que la chaîne est irréductible, les états sont soit tous transients soit tous récurrents par la proposition 39.25. Nous commençons par considérer que x est transient.

En comparant la définition (39.63) de N_x et le membre de droite de (39.90), nous avons pour chaque n l'inégalité

$$\frac{1}{n} \sum_{k=1}^n \mathbb{1}_{\{X_k = x\}} \leq \frac{1}{n} E(N_x). \quad (39.91)$$

Dans le cas d'un élément transient, nous avons $\pi(x) = 0$, donc il serait bon de montrer que $E(N_x) < \infty$, de sorte que prendre la limite $n \rightarrow \infty$ dans (39.91) donne zéro.

Nous décomposons le calcul en deux morceaux :

$$E(N_x) = E(N_x | T(x) = \infty) P(T(x) = \infty) + E(N_x | T(x) < \infty) P(T(x) < \infty). \quad (39.92)$$

Le fait que le premier terme soit fini découle immédiatement du fait que $T(x) = \infty$ implique $X_k \neq x$ pour tout $k \geq 1$. Dans ce cas l'espérance de N_x est évidemment finie.

Pour le second terme nous avons

$$E(N_x | T(x) < \infty) = E\left(\sum_{k=0}^{\infty} \mathbb{1}_{\{X_k=x\}} | T(x) < \infty\right) \quad (39.93a)$$

$$= \sum_{k=1}^{\infty} E(\mathbb{1}_{\{X_k=x\}} | T(x) < \infty). \quad (39.93b)$$

Pour inverser la somme et l'espérance, nous avons utilisé le théorème de Fubini-Tonelli qui est encore valable pour des fonctions qui prennent la valeur ∞ . Le fait d'inverser ne signifie pas que ni la somme ni l'intégrale soit finie. D'ailleurs c'est exactement ce que nous sommes en train de déterminer.

Étant donné que nous voulons seulement savoir si cette somme est finie ou non, nous pouvons nous restreindre à la somme depuis $k = 1$ ou oublier le premier terme. D'autre par nous avons

$$\sum_{k=1}^{\infty} \mathbb{1}_{\{X_k=x\}} = \sum_{j=0}^{\infty} \mathbb{1}_{\{X_{j+T(x)}=x\}} \quad (39.94)$$

parce que les $T(x)$ premiers termes sont par définition nuls. Nous regardons donc

$$\sum_{j=0}^{\infty} E(\mathbb{1}_{X_{j+T(x)}=x} | T(x) < \infty) = \sum_j P(X_{j+T(x)} = x | T(x) < \infty) \quad (39.95a)$$

$$= \sum_j P(X_j = x | X_0 = x) \quad (39.95b)$$

$$= \sum_j E(\mathbb{1}_{\{X_j=x\}} | X_0 = x) \quad (39.95c)$$

$$= E\left(\sum_j \mathbb{1}_{X_j=x} | X_0 = x\right) \quad (39.95d)$$

$$= E(N_x | X_0 = x) \quad (39.95e)$$

$$< \infty \quad \text{parce que } x \text{ est transient.} \quad (39.95f)$$

L'équation (39.95b) provient de la proposition 39.27 et plus précisément de l'égalité entre les lois (39.82). Nous avons terminé la preuve dans le cas où x est transient.

Nous passons maintenant au cas où x est récurrent, c'est-à-dire $P(T(x) < \infty | X_0 = x) = 1$. Les variables aléatoires T_i définies en (39.84) pour $i \geq 2$ sont indépendantes et identiquement distribuées et

$$\mathcal{L}(T_k(x)) \sim \mathcal{L}(T(X) | X_0 = x). \quad (39.96)$$

La loi des grands nombres nous indique que

$$\frac{S_n}{n} = \frac{T_1(x)}{n} + \frac{1}{n} \sum_{k=2}^n T_k(x) \xrightarrow{p.s.} E(T_2(x)) \quad (39.97a)$$

$$= E(T(x) | X_0 = x). \quad (39.97b)$$

Remarque 39.30.

La loi des grands nombres est encore vraie sans l'hypothèse de variables aléatoires dans L^1 pourvu qu'elles soient positives. Alors dans la conclusion de la loi nous devons accepter la possibilité que l'espérance soit infinie.

Nous posons pour $m \in \mathbb{N}$

$$n(m) = \sum_{j=1}^m \mathbb{1}_{\{X_j=x\}} \quad (39.98)$$

qui est le nombre de visites de x avant l'instant m . Nous avons évidemment $n(m) \leq m$. Mais S_n est l'instant de la n ème visite, par conséquent $S_{n(m)}$ est l'instant de la dernière visite avant le moment m . Pour tout m nous avons les inégalités

$$S_{n(m)} \leq m < S_{n(m)+1}. \quad (39.99)$$

Nous divisons par $n(m)$ et nous effectuons la limite $m \rightarrow \infty$:

$$\frac{S_{n(m)}}{n(m)} \leq \frac{m}{n(m)} \leq \frac{S_{n(m)} + 1}{n(m)} \quad (39.100)$$

En ce qui concerne la limite de $n(m)$, nous utilisons la définition (39.98) :

$$n(m) \rightarrow \sum_{j=1}^{\infty} \mathbb{1}_{\{X_j=x\}} = \quad (39.101)$$

heur...

$$\lim_{m \rightarrow \infty} n(m) = \lim_{m \rightarrow \infty} \sum_{n=1}^m \mathbb{1}_{\{X_n=x\}} \xrightarrow{p.s.} \infty \quad (39.102)$$

par la proposition (39.24). Plus précisément, la limite vaut N_x qui vaut presque sûrement ∞ dans le cas où x est récurrent. Par ailleurs la loi des grands nombres (39.97) nous enseigne en particulier que

$$\frac{S_{n(m)}}{n(m)} \xrightarrow{p.s.} E(T(x)|X_0 = x). \quad (39.103)$$

Le terme de droite dans (39.100) se traite de façon usuelle :

$$\frac{S_{n(m)+1}}{n(m)} = \frac{S_{n(m)+1}}{n(m)+1} \frac{n(m)+1}{n(m)}. \quad (39.104)$$

Le dernier facteur tend vers 1 et le tout a pour limite $E(T(x)|X_0 = x)$. Par conséquent nous avons

$$\frac{m}{n(m)} \xrightarrow{p.s.} E(T(x)|X_0 = x) \quad (39.105)$$

et

$$\frac{n(m)}{m} = \frac{1}{m} \sum_{j=1}^m \mathbb{1}_{\{X_j=x\}} \rightarrow \frac{1}{E(T(x)|X_0 = x)} = \pi(x). \quad (39.106)$$

□

Lemme 39.31.

Soit (X_k) une chaîne de Markov dont l'espace des états est noté E . Pour chaque $x \in E$ nous notons

$$T(x) = \inf\{k \geq 1 \text{ tel que } X_k = x\} \quad (39.107)$$

et

$$T_p(x) = \inf\{k \geq 1 \text{ tel que } X_{k+p} = x\} \quad (39.108)$$

Alors nous avons

$$P(T_p(x) = k | X_p = y) = P(T(x) = k | X_0 = y). \quad (39.109)$$

La proposition suivante nous permet de parler de chaîne de Markov **récence positive**.

Proposition 39.32.

Soit (x_n) une chaîne de Markov irréductible.

- (1) Un état x est transient si et seulement si tous les états sont transients.
- (2) Un état est récurrent positif si et seulement si tous les états sont récurrents positifs.

Démonstration. Nous rappelons (proposition 39.29) que si la chaîne est irréductible

$$\pi(x) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{[X_k=x]} \quad (39.110)$$

Notons aussi que

$$\sum_{k=1}^N \mathbb{1}_{X_k=x} = \begin{cases} 0 & \text{si } N < T(x) \\ \sum_{k=0}^{N-T(x)} \mathbb{1}_{X_{k+T(x)}=x} & \text{si } N > T(x) \end{cases} \quad (39.111)$$

où dans la seconde ligne nous avons effectué le changement de variable de sommation $k' = k + T(x)$. Dans la limite (39.110) nous sommes toujours dans le cas où N est assez grand. Nous pouvons donc écrire

$$\pi(x) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-T(x)} \mathbb{1}_{X_{k+T(x)}=x}. \quad (39.112)$$

Nous pouvons aussi écrire

$$\frac{1}{N - T(x)} \sum_{k=0}^{N-T(x)} \mathbb{1}_{X_{k+T(x)}=x} = \frac{N}{N - T(x)} \frac{1}{N} \sum_{k=0}^{N-T(x)} \mathbb{1}_{X_{k+T(x)}=x}. \quad (39.113)$$

Dans cette dernière égalité le membre de droite tend vers $\pi(x)$ et nous avons

$$\lim_{N \rightarrow \infty} \frac{1}{N - T(x)} \sum_{k=0}^{N-T(x)} \mathbb{1}_{X_{k+T(x)}=x} = \pi(x) \quad (39.114)$$

ou encore

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^N \mathbb{1}_{X_{k+T(x)}=x} = \pi(x) \quad (39.115)$$

Étant donné que $\pi(x)$ est une constante nous avons évidemment $E(\pi(x)) = \pi(x)$. Nous pouvons cependant considérer les variables aléatoires

$$Z_n = \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{X_{k+T(x)}=x} \quad (39.116)$$

et remarquer que $Z_n \xrightarrow{p.s.} \pi(x)$ avec $0 \leq Z_n \leq 1$. Le théorème de la convergence dominée (15.184) nous permet d'inverser la limite et l'espérance et écrire

$$\pi(x) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n E(\mathbb{1}_{X_{k+T(x)}=x}) \quad (39.117a)$$

$$= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n P(X_{k+T(x)} = x). \quad (39.117b)$$

Par le lemme 39.31 nous avons

$$P(X_{k+T(x)} = x) = P(X_k = k | X_0 = x) \quad (39.118)$$

et $\pi(x)$ prend la forme

$$\pi(x) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n P(X_k = x | X_0 = x). \quad (39.119)$$

Soit maintenant un état x positif récurrent et y , un autre état. Par définition 39.23 et par corollaire 39.26 nous avons $\pi(x) > 0$. Nous devons prouver que $\pi(y) > 0$.

Étant donné que la chaîne est irréductible il existe r et s tels que

$$\begin{cases} p^r(x, y) = P(X_r = y | X_0 = x) > 0 & (39.120a) \\ p^s(x, y) = P(X_s = x | X_0 = y) > 0 & (39.120b) \end{cases}$$

Nous reprenons l'équation (39.77) multipliée par $1/N$:

$$\frac{1}{N} \sum_{n=1}^N P(X_{r+s+n=y} | X_0=y) \geq \underbrace{p^r(x,y)p^s(y,x)}_{\geq 0} \underbrace{\frac{1}{N} \sum_{n=1}^N P(X_n = x | X_0 = x)}_{\rightarrow \pi(x)} \quad (39.121)$$

et nous prenons la limite lorsque $N \rightarrow \infty$. À $r+s$ termes près, nous trouvons à gauche l'expression (39.119) de $\pi(y)$. Par conséquent

$$\pi(y) \geq \lim_{N \rightarrow \infty} \frac{1}{n} \sum_{n=1}^N P(X_{r+s+n} = y | X_0 = y) \geq \alpha \pi(x) \quad (39.122)$$

où α est une constante positive. Le nombre $\pi(x)$ étant strictement positif par hypothèse nous avons montré que $\pi(y) > 0$, c'est-à-dire que y est récurrent positif. \square

39.5 Mesure invariante

Définition 39.33.

Une mesure de probabilité μ sur l'espace des états E d'une chaîne de Markov est **invariante** si pour tout $x \in E$

$$\mu(x) = \sum_{y \in E} p(y, x) \mu(y). \quad (39.123)$$

Remarque 39.34.

Une mesure invariante est une mesure de probabilité et nous noterons par abus $\mu(x)$ pour $\mu(\{x\})$. Si $A \subset E$ nous avons

$$\mu(A) = \sum_{x \in A} \mu(x). \quad (39.124)$$

Remarque 39.35.

Une loi invariante associée à une chaîne de Markov est une loi associée à la matrice de transition de la chaîne, mais pas à la loi de X_0 . Par conséquent nous pouvons tester si μ est une mesure invariante pour une certaine chaîne de Markov (X_k) en considérant la chaîne (Y_k) avec $Y_k = X_k$ pour $k > 0$ et Y_0 arbitraire.

L'adjectif *invariant* provient du lemme suivant.

Lemme 39.36.

Soit (X_n) une chaîne de Markov telle que $X_0 \sim \mu$ où μ est une mesure invariante sur l'espace des états. Alors $X_k \sim \mu$ pour tout k .

Démonstration. Par hypothèse, $P(X_0 = x) = \mu(x)$. Ensuite nous avons

$$P(X_1 = y) = \sum_{x \in E} P(X_1 = y | X_0 = x) P(X_0 = x) \quad (39.125a)$$

$$= \sum_x p(x, y) \mu(x) \quad (39.125b)$$

$$= \mu(y). \quad (39.125c)$$

Par conséquent X_1 suit également la loi μ . Par récurrence tous les états suivent cette même loi. \square

Si les états d'une chaîne de Markov ont comme loi une mesure invariante, alors nous disons que la chaîne est **stationnaire**.

Remarque 39.37.

Pour une chaîne de Markov stationnaire de loi invariante μ nous avons

$$\mu(x) = \sum_y p(y, x) \mu(y) \quad (39.126)$$

et si l'ensemble E est fini cette équation signifie

$$\mu = Q\mu \quad (39.127)$$

où Q est la matrice de transition de la chaîne de Markov.

Théorème 39.38 (Théorème ergodique).

Une chaîne de Markov irréductible est positive récurrente si et seulement si elle accepte une mesure invariante. Cette mesure est invariante est alors unique et vérifie $\mu = Q\mu$ où Q est la matrice de transition.

Démonstration. Nous allons seulement prouver le théorème ergodique dans le cas où E est fini. Soit (X_n) une chaîne de Markov récurrente positive; nous avons $\pi(x) > 0$ pour tout $x \in E$. Nous allons montrer que π est une mesure invariante.

Nous commençons par montrer que

$$\sum_{x \in E} \pi(x) = 1. \quad (39.128)$$

Pour cela nous reprenons la propriété de chaîne irréductible pour écrire

$$\pi(x) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \mathbb{1}_{X_k=x} \quad (39.129)$$

Étant donné que E est fini nous pouvons sommer sur $x \in E$ et permuter la somme avec la limite :

$$\sum_{x \in E} \pi(x) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \underbrace{\sum_{x \in E} \mathbb{1}_{X_k=x}}_{=1}. \quad (39.130)$$

Nous nous retrouvons donc avec $\lim_{N \rightarrow \infty} \frac{1}{N} N = 1$. La fonction π définit donc bien une mesure de probabilité sur E .

Nous montrons à présent que cette mesure est invariante, c'est-à-dire que

$$\pi(x) = \sum_{y \in E} p(y, x) \pi(y). \quad (39.131)$$

Pour cela nous utilisons encore le théorème de la convergence dominée pour permuter la limite et l'intégrale dans

$$\pi(x) = E(\pi(x)) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \underbrace{E(\mathbb{1}_{X_k=x})}_{P(X_k=x)} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N P(X_{k+1} = x). \quad (39.132)$$

La dernière égalité découle du fait que en divisant par N et en faisant tendre N vers l'infini, le fait d'enlever un terme à la somme ne change pas la valeur de la limite. Nous pouvons substituer dans (39.132) la valeur

$$P(X_{k+1} = x) = \sum_{y \in E} p(y, x) P(X_k = y). \quad (39.133)$$

Nous avons alors

$$\pi(x) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \sum_{y \in E} p(y, x) P(X_k = y) \quad (39.134a)$$

$$= \sum_{y \in E} p(y, x) \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N P(X_k = y) \quad (39.134b)$$

$$= \sum_{y \in E} p(y, x) \pi(y), \quad (39.134c)$$

ce qui signifie que π est une mesure invariante. Notons que nous avons encore utilisé le fait que E soit fini pour permuter avec la limite.

Il nous reste à montrer l'unicité de la mesure invariante sur la chaîne de Markov. Soit μ une mesure invariante pour la chaîne de Markov (X_k) . Comme indiqué dans la remarque 39.35 nous pouvons supposer que X_0 suit la loi μ . Par le lemme 39.36 nous avons $P(X_k = x) = \mu(x)$ pour tout k . Par conséquent

$$\pi(x) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N P(X_k = x) = \mu(x). \tag{39.135}$$

□

Théorème 39.39 (loi des grands nombres pour les chaîne de Markov).

Soit (X_n) une chaîne de Markov irréductible acceptant une mesure invariante. Soit $f: E \rightarrow \mathbb{R}$ une fonction dans $L^1(E, \mu)$. Alors nous avons

$$\frac{1}{N} \sum_{k=1}^N f(X_k) \xrightarrow{p.s.} \sum_{x \in E} f(x) \mu(x). \tag{39.136}$$

En ce qui concerne les notations, l'hypothèse $f \in L^1(E, \mu)$ signifie

$$\sum_{x \in E} |f(x)| \mu(x) = \int_E |f(x)| d\mu(x) < \infty. \tag{39.137}$$

Démonstration. Nous prouvons le théorème dans le cas où E est fini. Si nous écrivons

$$f(X_k) = \sum_{y \in E} f(y) \mathbb{1}_{X_k=y}, \tag{39.138}$$

alors

$$\frac{1}{N} \sum_{k=1}^N f(X_k) = \sum_{y \in E} f(y) \frac{1}{N} \sum_{k=1}^N \mathbb{1}_{X_k=y}. \tag{39.139}$$

Étant donné que E est fini nous pouvons permuter les sommes et prendre la limite $N \rightarrow \infty$:

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_k f(X_k) = \sum_{y \in E} f(y) \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \mathbb{1}_{X_k=y} = \sum_{y \in E} f(y) \pi(y). \tag{39.140}$$

□

39.6 Convergence vers l'équilibre

Nous voudrions savoir sous quelles conditions la variable aléatoire X_n converge en loi vers quelque chose lorsque $n \rightarrow \infty$. Une telle loi limite doit dépendre de la loi initiale² comme le montre l'exemple de la chaîne de Markov

$$1 \circlearrowleft A \xleftarrow{1/2} C \xrightarrow{1/2} B \circlearrowright 1 \tag{39.141}$$

Si $X_0 = C$, alors la loi limite est

$$\frac{1}{2}(\delta_A + \delta_B). \tag{39.142}$$

Si par contre $X_0 = B$, la loi limite est δ_B . Notons que la chaîne de Markov proposée ici est irréductible.

2. Lorsque la loi limite ne dépend pas de la loi initiale, nous disons que la chaîne de Markov est ergodique, nous y reviendrons.

Notons qu'il n'y a pas toujours de lois limite comme le montre l'exemple

$$A \begin{array}{c} \xrightarrow{1} \\ \xleftarrow{1} \end{array} B \quad (39.143)$$

avec $X_0 = A$. La loi en est

$$X_k = \begin{cases} \delta_A & \text{si } k \text{ est pair} \\ \delta_B & \text{si } k \text{ est impair.} \end{cases} \quad (39.144)$$

Lemme 39.40.

Si nous avons une loi limite

$$P(X_n = x) \rightarrow l(x), \quad (39.145)$$

et que la chaîne est irréductible, alors nous avons $l = \pi$.

Démonstration. D'après la proposition 39.29 nous avons

$$\frac{1}{n} \sum_{k=1}^n P(X_k = x) \rightarrow \pi(x). \quad (39.146)$$

Par le lemme 12.82 sur la moyenne de Cesaro et l'hypothèse (39.145), nous avons aussi

$$\frac{1}{n} \sum_{k=1}^n P(X_k = x) \rightarrow l(x). \quad (39.147)$$

Du coup $\pi(x) = l(x)$. □

Lemme 39.41 ([499]).

Si π est une loi stationnaire et si x est un état transitoire, alors $\pi(x) = 0$.

Ce lemme (qui peut être prouvé rigoureusement) est principalement dû au fait que la chaîne de Markov ne visite un état transitoire qu'un nombre fini de fois par la proposition 39.24(1).

Définition 39.42.

Un état $x \in E$ est **apériodique** si

$$\text{pgcd}\{n \geq 1 \text{ tel que } p^n(x, x) > 0\} = 1. \quad (39.148)$$

Mettons que tous les n tels que $p^n(x, x) > 0$ ont 2 comme diviseur. L'état n'est alors pas apériodique, mais on voit que si $X_0 = x$, alors les états impairs ne peuvent pas être sur x . Cela est une forme de périodicité.

Si un état est apériodique, il existe p et q premiers entre eux tels que $p^p(x, x)$ et $p^q(x, x)$ sont non nuls. En particulier pour tout $n \in p\mathbb{N} + q\mathbb{N}$, $P(X_n = x) \neq 0$. Par conséquent la proposition 3.15 nous indique qu'à partir d'un certain moment tous les X_k pourraient être x .

L'état C de la chaîne de Markov suivante est apériodique :

$$\begin{array}{ccc} A & \begin{array}{c} \xrightarrow{1} \\ \xleftarrow{2/3} \end{array} & B \\ & \swarrow 1 \quad \searrow 1/3 & \\ & C & \end{array} \quad (39.149)$$

En effet $p^3(C, C) \neq 0$ par le chemin $C \rightarrow A \rightarrow B \rightarrow C$ tandis que $p^5(C, C) \neq 0$ également par le chemin $C \rightarrow A \rightarrow B \rightarrow A \rightarrow B \rightarrow C$. Or $\text{pgcd}\{3, 5\} = 1$.

Proposition 39.43 ([499]).

Soit (X_n) , une chaîne de Markov irréductible. Un état x est apériodique si et seulement s'il existe N tel que

$$p^k(x, x) = P(X_k = x | X_0 = x) > 0 \quad (39.150)$$

pour tout $k \geq N$.

La proposition suivante va nous permettre de parler de **chaîne apériodique** .

Proposition 39.44.

Si une chaîne de Markov est irréductible, alors un état est apériodique si et seulement si tous les états sont apériodiques.

Démonstration. Soit x un état apériodique de la chaîne de Markov $(X_n)_{n \in \mathbb{N}}$. En vertu de la proposition 39.43 il existe N_x tel que $p^k(x, x) \neq 0$ pour tout $k \geq N_x$. Soit $y \in E$. Étant donné que la chaîne est irréductible, il existe r et s tels que $p^r(x, y) > 0$ et $p^s(y, x) > 0$. Nous avons

$$p^{k+r+s}(y, y) = P(X_{k+r+s} = y | X_0 = y) \geq p^s(x, y)P(X_k = x | X_0 = x)p^r(y, x). \tag{39.151}$$

Si k est assez grand, cette quantité est strictement positive. Donc il suffit de prendre $N_y = N_x + r + s$ pour savoir que y est également apériodique. □

Exemple 39.45

Quelle est la différence entre une chaîne irréductible et une chaîne apériodique? Une chaîne est irréductible lorsque aucune sous-chaîne ne peut piéger le système. Pour toute paire d'états $x, y \in E$, il existe un n tel qu'il soit possible d'aller de x à y en n pas. Une chaîne est apériodique lorsqu'après un temps suffisamment long, tous les états soient possibles en même temps.

Un exemple de chaîne irréductible non apériodique :

$$A \begin{matrix} \xrightarrow{1} \\ \xleftarrow{1} \end{matrix} B \tag{39.152}$$

Cette chaîne est irréductible parce que le graphe est connexe, par contre il n'est pas apériodique parce que si $X_0 = A$ il n'est pas possible d'être dans l'état A après un nombre impair de pas.

Plus formellement, $p^n(A, A) = 1$ dès que n est pair ; le PGCD de la définition 39.42 n'est donc certainement pas 1. △

Si E est fini et si la chaîne de Markov est irréductible, alors en posant $N = \max_{x \in E} N(x)$, la matrice P^k a des éléments non nuls sur toute la diagonale pour tout $k > N$. Ces éléments diagonaux ne sont autre que les $p^k(x, x)$.

Théorème 39.46 (Convergence en loi des chaîne de Markov).

Si (X_n) est

- (1) irréductible,
- (2) récurrente positive,
- (3) apériodique,

alors X_n converge en loi vers l'unique probabilité invariante π vérifiant

$$\pi(x) = \sum_{u \in E} p(y, x)\pi(y) = \frac{1}{E(T(x) | X_0 = x)}. \tag{39.153}$$

Cette convergence est indépendante de la loi de X_0 et on a

$$P(X_n = x | X_0 = y) \rightarrow_{n \rightarrow \infty} \pi(x). \tag{39.154}$$

39.7 Processus de Galton-Watson

Nous considérons une maladie et notons Z_n le nombre de malades à l'instant n . Nous posons $Z_0 = 1$ et

$$Z_{n+1} = \begin{cases} 0 & \text{si } Z_n = 0 \\ \sum_{i=1}^{Z_n} \xi_i^{(n)} & \text{sinon} \end{cases} \tag{39.155}$$

où $\xi_i^{(n)}$ est le nombre de personnes contaminées par le malade i à l'instant n . Nous supposons que ces variables aléatoires sont indépendantes et identiquement distribuées et admettent un moment d'ordre 1.

L'équation de propagation 39.155 signifie que nous supposons qu'une personne malade à l'instant n n'est plus malade à l'instant $n + 1$. Par ailleurs les hypothèses d'indépendance signifient qu'à chaque instant, le nombre de personnes contaminées par le malade i est indépendant du nombre de personnes contaminées par le malade j . De plus la façon dont la contamination se passe à l'instant n est indépendant de la façon dont la contamination se passe à l'instant m . Ces hypothèses sont raisonnables tant que le nombre de personnes non contaminées est grand. À partir du moment où presque tout le monde est malade, l'approximation de Galton-Watson ne fonctionne plus.

Nous notons ξ la loi parente des $\xi_i^{(n)}$. Ensuite nous considérons

$$G(s) = E(s^\xi) \tag{39.156a}$$

$$m = E(\xi) \tag{39.156b}$$

$$G_n(s) = E(s^{Z_n}). \tag{39.156c}$$

Par le théorème de transfert (proposition 37.60) avec $f(t) = s^t$. Ce que nous avons est

$$G_n(s) = E(f(Z_n)) = \int_{\mathbb{R}} s^x dP_{Z_n}(x) = \sum_{k=0}^{\infty} s^k P(Z_n = k) \tag{39.157}$$

où l'intégrale s'est transformée en somme parce que la loi de Z_n est discrète : dP_{Z_n} est une somme de masses de Dirac. En particulier nous avons

$$G_n(s) = \sum_{k=0}^{\infty} s^k P(Z_n = k) \tag{39.158a}$$

$$G(0) = P(Z_n = 0) \tag{39.158b}$$

et

$$\eta = \lim_{n \rightarrow \infty} P(Z_n = 0) = \lim_{n \rightarrow \infty} G_n(0). \tag{39.159}$$

D'où l'intérêt d'étudier G_n .

Lemme 39.47.

Pour tout $n \in \mathbb{N}^*$ et pour tout $s \in [0, 1]$, nous avons

$$G_n(s) = \underbrace{G \circ G \circ \dots \circ G(s)}_{n \text{ fois}}. \tag{39.160}$$

Démonstration. Pour $n = 1$, nous avons $Z_1 = \xi_1^{(1)}$ et donc

$$G_1(s) = E(s^{\xi}) = G(s), \tag{39.161}$$

comme il se doit.

Si $n \neq 1$ nous écrivons

$$G_n(s) = E(s^{Z_n}) \tag{39.162a}$$

$$= E\left(s^{\sum_{i=1}^{Z_{n-1}} \xi_i^{(n-1)}}\right) \tag{39.162b}$$

$$= E\left(\sum_{k=0}^{\infty} \mathbb{1}_{\{Z_{n-1}=k\}} s^{\sum_{i=1}^k \xi_i^{(n-1)}}\right). \tag{39.162c}$$

À ce niveau, nous voulons permuter la somme et l'espérance. Étant donné que le lemme est facile à vérifier pour $s = 1$, nous supposons $s < 1$. Du coup

$$s^{\sum_{i=1}^k \xi_i^{(n-1)}} < 1 \tag{39.163}$$

et ce qui se trouve dans l'espérance est majoré par

$$\sum_{k=0}^{\infty} \mathbb{1}_{Z_{n-1}=k} = 1. \quad (39.164)$$

La fonction constante 1 est intégrable sur Ω (ici nous utilisons à fond le fait que l'espace Ω soit un espace de probabilité) et nous pouvons utiliser le théorème de convergence dominée de Lebesgue 16.2 pour permuter la somme et l'intégrale. Nous continuons donc le calcul (39.162) :

$$G_n(s) = \sum_{k=0}^{\infty} E \left(\mathbb{1}_{\{Z_{n-1}=k\}} s^{\sum_{i=1}^k \xi_i^{(n-1)}} \right). \quad (39.165)$$

La tribu engendrée par la variable aléatoire $\mathbb{1}_{\{Z_{n-1}=k\}}$ est une fonction des variables aléatoires $\xi_i^{(m)}$ avec $m \leq n-2$ tandis que la variable aléatoire $s^{\sum_{i=1}^k \xi_i^{(n-1)}}$ est une fonction des variables aléatoires $\xi_i^{(n-1)}$. Par conséquent le lemme de regroupement 37.16 nous dit que ces variables aléatoires sont indépendantes, donc

$$G_n(s) = \sum_{k=0}^{\infty} \underbrace{E(\mathbb{1}_{\{Z_{n-1}=k\}})}_{=P(Z_{n-1}=k)} E(s^{\sum_{i=1}^k \xi_i^{(n-1)}}). \quad (39.166)$$

Nous avons utilisé le fait que l'espérance d'une fonction indicatrice est la probabilité de l'événement.

En ce qui concerne la puissance de s , les événements ξ_i^{n-1} sont indépendants et suivent tous la même loi ξ , donc

$$s^{\sum_{i=1}^k \xi_i^{(n-1)}} = \prod_{i=1}^k s^{\xi_i^{(n-1)}} \quad (39.167)$$

et

$$E\left(\prod_{i=1}^k s^{\xi_i}\right) = E(s^{\xi})^k = G(s)^k. \quad (39.168)$$

En mettant tout bout à bout,

$$G_n(s) = \sum_{k=1}^{\infty} P(Z_{n-1} = k) G(s)^k = G_{n-1}(G(s)). \quad (39.169)$$

□

Théorème 39.48.

La probabilité d'extinction η est donnée par

$$\eta = P\left(\bigcup_{n \geq 1} (Z_n = 0)\right) = \lim_{n \rightarrow \infty} P(Z_n = 0). \quad (39.170)$$

Ce nombre est la plus petite solution positive de l'équation $G(s) = s$.

De plus la classification des cas est comme suit.

(1) Si $P(\xi = 0) = 0$ alors $\eta = 0$.

(2) Si $P(\xi = 0) \neq 0$ alors

(2a) si $m \leq 1$ alors $\eta = 1$,

(2b) si $m > 1$ alors $\eta \in]0, 1[$.

Le cas $m < 1$ est dit **sous-critique**, le cas $m = 1$ est dit **critique**. Le cas $m > 1$ est dit **sur-critique**.

Démonstration. Commençons par prouver que G est une fonction continue. En utilisant la théorie de transfert comme pour l'équation (39.157) nous trouvons que

$$G(s) = E(s^\xi) = \sum_{k=0}^{\infty} p_k s^k \quad (39.171)$$

où nous avons noté $p_k = P(\xi = k)$. Si $r < 1$, alors la suite $p_k r^k$ est bornée, donc le critère d'Abel (16.10) nous indique que la série (39.171) converge absolument et la théorie générale des séries entières conclut que la fonction G est en particulier dérivable terme à terme pour tout $s \in]-1, 1[$.

Le probabilité d'extinction est un point fixe de G En utilisant la continuité de G en 0 nous passons à la limite dans $G_{n+1}(0) = G(G_n(0))$ et nous obtenons

$$\eta = G(\eta), \quad (39.172)$$

ce qui signifie que la probabilité d'extinction est un point fixe de G .

η est le plus petit point fixe de G Nous démontrons maintenant que η est plus précisément le plus petit point fixe de G sur $[0, 1]$. Nous allons effectuer cette partie en décomposant selon les valeurs de p_0 et de p_1 .

Au vu de l'écriture (39.171), si $p_1 = 1$ alors $G(s) = s$ pour tout $s \in [0, 1]$. Mais dans ce cas nous savons par ailleurs que l'extinction est impossible. Zéro est bien la plus petite solution de $G(s) = s$.

Supposons maintenant que $p_1 < 1$ et $p_0 + p_1 = 1$. Alors $G(s) = p_0 + p_1 s$ et $s = 1$ est l'unique solution. Mais vu que nous savons que η est solution, c'est que $\eta = 1$ et l'extinction est certaine.

Nous passons au cas général : $p_0 + p_1 < 1$. D'abord nous remarquons que $s = 1$ est solution parce que

$$G(1) = p_0 + p_1 + \dots = 1. \quad (39.173)$$

Remarquons aussi que dans ce cas $s = 0$ n'est plus solution.

La fonction G est strictement convexe sur $[0, 1]$ (parce que $G'' > 0$). Cela se voit en effectuant deux dérivations termes à termes (le rayon de convergence de la dérivée est le même que celui de la fonction). Cette stricte convexité entraîne que l'équation $G(s) = s$ a au maximum une autre solution que $s = 1$. Nous nommons s_0 la plus petite solution dans $[0, 1]$. Étant donné que G est croissante on a

$$G(0) \leq G(s_0) = s_0. \quad (39.174)$$

En appliquant G à cette équation nous obtenons $G(G(s_0)) \leq G(s_0) = s_0$ et en appliquant n fois,

$$G_n(0) \leq s_0. \quad (39.175)$$

En passant à la limite, $\eta \leq s_0$ mais η étant solution, nous avons $\eta = s_0$. Nous avons donc prouvé que la probabilité d'extinction η est la plus petite solution de $G(s) = s$.

Classification des cas Nous devons encore discuter les cas. Si $P(\xi = 0) = 0$, alors $p_0 = 0$ et $G(0) = 0$, ce qui signifie que $s_0 = \eta = 0$ et l'extinction est impossible.

Nous passons au cas $p_0 \neq 0$. Si $p_0 + p_1 = 1$, alors $m = p_1 < 1$ et nous avons déjà vu que dans le cas $p_0 + p_1 = 1$, la probabilité d'extinction est $\eta = 1$.

Il nous reste à traiter le cas $p_0 + p_1 < 1$. Encore une fois, la courbe G est strictement convexe sur $[0, 1]$ et elle est en particulier plus grande que sa tangente en $s = 1$, c'est-à-dire

$$G(s) > G'(1)(s - 1) + G(1). \quad (39.176)$$

Nous savons que $G(1) = 1$. En ce qui concerne $G'(1)$, nous dérivons encore terme à termes :

$$G'(s) = \sum_{k=1}^{\infty} k p_k s^{k-1}, \quad (39.177)$$

donc

$$G'(1) = \sum_{k=1}^{\infty} kp_k = E(\xi) = m. \quad (39.178)$$

Ce que nous avons donc est

$$G(s) > 1 + m(s - 1). \quad (39.179)$$

Nous nous particularisons au cas sous-critique ($m \leq 1$). En nous rappelant que $s - 1 < 0$,

$$G(s) > 1 + (s - 1) = s, \quad (39.180)$$

donc $s = 1$ est la plus petite solution et effectivement nous avons déjà vu que $\eta = 1$ dans ce cas.

Si $m > 1$, alors on a

$$G(s) > 1 + m(s - 1). \quad (39.181)$$

Mais dire $m > 1$ revient à dire $G'(1) > 1$ et donc dans un voisinage de $s = 1$ on a

$$\frac{G(s) - G(1)}{s - 1} > 1, \quad (39.182)$$

ce qui implique que

$$G(s) < s - 1 + G(1) = s. \quad (39.183)$$

Nous avons donc $G(s) < s$ dans un voisinage de 1. Mais $G(0) - 0 = p_0 > 0$, donc la fonction $f(s) = G(s) - s$ est positive en 0 et négative proche de $s = 1$. Le théorème de la valeur intermédiaire nous indique alors qu'il existe un $s \in]0, 1[$ tel que $f(s) = 0$, c'est-à-dire tel que $G(s) = s$.

□

Chapitre 40

Martingales

40.1 Convergence de martingales

Définition 40.1.

Si \mathcal{A} est une tribu, une **filtration** de \mathcal{A} est une suite croissante de sous-tribus $\mathcal{B}_i \subseteq \mathcal{B}_{i+1} \subseteq \mathcal{A}$.

Nous disons qu'une suite de variables aléatoires (X_n) est **adaptée** à une filtration (\mathcal{F}_n) si X_i est \mathcal{F}_i -mesurable pour tout i .

Ces définitions impliquent immédiatement que si (X_n) est adapté à (\mathcal{F}_n) alors X_n est \mathcal{F}_k -mesurable pour $k \geq n$.

Définition 40.2.

Une **martingale** adaptée à la filtration $(\mathcal{B}_n)_{n \in \mathbb{N}}$ est une suite de variables aléatoires M_n telle que

- (1) $M_n \in L^1(\Omega, \mathcal{A}, P)$,
- (2) M_n est \mathcal{B}_n -mesurable,
- (3) $E(M_{n+1} | \mathcal{B}_n) = M_n$.

Le processus M_n est une **sur-martingale** si $E(M_{n+1} | \mathcal{B}_n) \leq M_n$ pour tout n , et c'est une **sous-martingale** si $E(M_{n+1} | \mathcal{B}_n) \geq M_n$.

Exemple 40.3

Si $M \in L^1(\Omega, \mathcal{A}, P)$ et si $(\mathcal{B}_n)_{n \in \mathbb{N}}$ est une filtration, nous pouvons considérer la martingale $M_n = E(M | \mathcal{B}_n)$. △

Exemple 40.4

Soit $(X_i)_{i \geq 1}$ une suite de variables aléatoires indépendantes et centrées. On pose

$$S_n = X_1 + \cdots + X_n \tag{40.1}$$

et la filtration $\mathcal{B}_n = \sigma(X_1, \dots, X_n)$. Pour montrer que cela est une martingale, nous commençons par remarquer que

$$E(X_{n+1} | \mathcal{B}_n) = E(X_{n+1}) = 0 \tag{40.2}$$

par indépendance des tribus \mathcal{B}_n et $\sigma(X_{n+1})$. Ici c'est le lemme 37.38 qui joue.

Ensuite nous argumentons que $E(X_1 + \cdots + X_n | \mathcal{B}_n) = X_1 + \cdots + X_n$. En effet d'une part $X_1 + \cdots + X_n$ est \mathcal{B}_n -mesurable et évidemment la condition intégrale de l'espérance conditionnelle est satisfaite.

Plus généralement si X est une variable aléatoire et si $\sigma(X) \subset \mathcal{B}$ alors $E(X | \mathcal{B}) = X$. △

Lemme 40.5.

Soit (M_n) une martingale adaptée à la filtration (\mathcal{F}_n) et $n \geq k$. Alors

$$E(M_n | \mathcal{F}_k) = M_k \quad (40.3a)$$

$$E(M_k | \mathcal{F}_n) = M_k. \quad (40.3b)$$

Démonstration. La seconde relation revient seulement à dire que M_k est \mathcal{F}_n -mesurable, ce qui est évident parce que $\mathcal{F}_k \subset \mathcal{F}_n$.

Nous prouvons la première par récurrence (à l'envers) sur k . D'abord si $k = n$, l'égalité $E(M_n | \mathcal{F}_n) = M_n$. Nous supposons maintenant que $E(M_n | \mathcal{F}_k) = M_k$, et nous prouvons que $E(M_n | \mathcal{F}_{k-1}) = M_{k-1}$. Si $B_{k-1} \in \mathcal{F}_{k-1}$, nous avons

$$\int_{B_{k-1}} M_{k-1} = \int_{B_{k-1}} M_k = \int_{B_{k-1}} M_n. \quad (40.4)$$

La première égalité est la définition d'une martingale, et la seconde est l'hypothèse de récurrence. \square

Théorème 40.6 ([500, 501]).

Soit $(M_n)_{n \geq 0}$ une martingale bornée dans $L^2(\Omega)$, c'est-à-dire telle que

$$\alpha = \sup_{n \geq 0} E(M_n^2) < \infty. \quad (40.5)$$

Alors la suite M_n converge dans $L^2(\Omega)$.

Démonstration. Nous écrivons M_n en somme télescopique

$$M_n = M_0 + \sum_{k=1}^n \Delta_k \quad (40.6)$$

où $\Delta_k = M_k - M_{k-1}$. Nous commençons par montrer que les incréments sont orthogonaux au sens où $E(\Delta_n \Delta_k) = 0$. Pour $n > k$, la variable aléatoire $E(\Delta_n \Delta_k | \mathcal{F}_{n-1})$ est la variable aléatoire \mathcal{F}_{n-1} -mesurable telle que

$$\int_{B_{n-1}} E(\Delta_n \Delta_k | \mathcal{F}_{n-1}) = \int_{B_{n-1}} \Delta_n \Delta_k \quad (40.7)$$

pour tout $B_{n-1} \in \mathcal{F}_{n-1}$. En particulier avec $B_{n-1} = \Omega$ nous trouvons

$$E\left(E(\Delta_n \Delta_k | \mathcal{F}_{n-1})\right) = E(\Delta_n \Delta_k) \quad (40.8)$$

par la définition de l'espérance (37.51). Par conséquent, en utilisant le lemme 40.5 nous avons¹

$$E(\Delta_n \Delta_k) = E\left(E(\Delta_n \Delta_k | \mathcal{F}_{n-1})\right) = E\left(\Delta_k E(\Delta_n | \mathcal{F}_{n-1})\right) = 0 \quad (40.9)$$

parce que $E(\Delta_n | \mathcal{F}_{n-1}) = E(M_n | \mathcal{F}_{n-1}) - E(M_{n-1} | \mathcal{F}_{n-1}) = 0$.

Utilisant l'orthogonalité des incréments, nous avons

$$E(M_n^2) = E(M_0^2) + \sum_{k=1}^n E(\Delta_k^2). \quad (40.10)$$

En prenant le supremum (par rapport à n des deux côtés),

$$E(M_0^2) + \sum_{k=1}^{\infty} E(\Delta_k^2) = \alpha < \infty. \quad (40.11)$$

1. À ce niveau je crois qu'il y a une faute dans [501] qui conditionne par rapport à \mathcal{F}_n .

Cela prouve que la suite $\sum_{k=1}^n \Delta_k$ converge dans $L^2(\Omega)$. Nous en déduisons immédiatement que (M_n) est de Cauchy dans $L^2(\Omega)$ parce que si $k, l > n$, nous avons (en utilisant encore l'orthogonalité des incréments)

$$E(|M_k - M_l|^2) = \sum_{i=k+1}^l E(\Delta_i^2) \leq \sum_{i=k+1}^{\infty} E(\Delta_i^2), \quad (40.12)$$

qui tend vers zéro lorsque $n \rightarrow \infty$. \square

Le théorème suivant complète la conclusion du théorème 40.6.

Théorème 40.7 ([501]).

Soit $(M_n)_{n \in \mathbb{N}}$ une martingale bornée dans L^2 . Alors (M_n) converge dans $L^2(\Omega)$ et presque sûrement vers une même variable aléatoire M_∞ qui vérifie

$$M_n = E(M_\infty | \mathcal{F}_n). \quad (40.13)$$

Notons en particulier que la variable aléatoire M_∞ est presque sûrement finie parce qu'en vertu de (40.13) nous avons

$$\int_{\Omega} M_\infty = \int_{\Omega} M_n < \infty. \quad (40.14)$$

Exemple 40.8

Soient des variables aléatoires indépendantes $V_k \sim \mathcal{E}(2^k \lambda)$ et la variable aléatoire somme

$$S_n = \sum_{k=1}^n V_k. \quad (40.15)$$

Nous allons montrer que $S_n \xrightarrow{p.s.} X$ où X est une variable aléatoire presque sûrement finie. Nous posons

$$M_n = S_n - \sum_{k=1}^n \frac{1}{2^k \lambda} \quad (40.16)$$

Cela est une martingale adaptée à la filtration $\mathcal{F}_n = \sigma(V_1, \dots, V_n)$ en vertu de l'exemple 40.4. Nous montrons à présent qu'elle est bornée dans $L^2(\Omega)$ au sens où $\sum_{n \geq 1} E(M_n^2) < \infty$. Nous avons

$$E(M_n^2) = E\left(\left[S_n - \sum_{k=1}^n \frac{1}{2^k \lambda}\right]^2\right) = E\left(\left[\sum_{k=1}^n \left(V_k - \frac{1}{2^k \lambda}\right)\right]^2\right). \quad (40.17)$$

La variable aléatoire $V_k - 1/2^k \lambda$ est une variable aléatoire centrée de variance $1/(2^k \lambda)^2$ (voir proposition 37.97). Étant donné que M_n est centrée, $\text{Var}(M_n) = E(M_n^2)$ et nous avons

$$E(M_n^2) = \sum_{k=1}^n \text{Var}\left(V_k - \frac{1}{2^k \lambda}\right) = \sum_{k=1}^n \frac{1}{(2^k \lambda)^2}, \quad (40.18)$$

cette dernière somme étant bornée par $l = \sum_{k=1}^{\infty} \frac{1}{(2^k \lambda)^2}$, nous avons

$$E(M_n^2) \leq l \quad (40.19)$$

avec l indépendant de n . C'est pour cela que $(M_n)_{n \in \mathbb{N}}$ est une martingale bornée dans $L^2(\Omega)$. Par le théorème 40.7 nous avons $M_n \rightarrow M_\infty$ et en faisant $n \rightarrow \infty$ dans

$$S_n = M_n + \sum_{k=1}^n \frac{1}{2^k \lambda}, \quad (40.20)$$

nous trouvons

$$S_n \rightarrow M_\infty + \sum_{k=1}^{\infty} \frac{1}{2^k \lambda} = M_\infty + \frac{1}{\lambda} \quad (40.21)$$

qui est presque sûrement finie. \triangle

40.2 Temps d'arrêt et martingale terminée

Définition 40.9.

Soit $(\Omega, \mathcal{F}_n, \mathcal{F}, P)$ un espace de probabilité filtré. Une application $T: \Omega \rightarrow \bar{\mathbb{N}}$ est un **temps d'arrêt** adapté à la filtration (\mathcal{F}_n) si pour tout $n \in \mathbb{N}$ nous avons $\{T \leq n\} \in \mathcal{F}_n$.

Le temps d'arrêt T est **borné** s'il existe $k \in \mathbb{N}$ tel que $T(\omega) \leq k$ pour presque tout $\omega \in \Omega$.

Lemme 40.10.

Si T est un temps d'arrêt presque sûrement fini, alors²

- $T \wedge n \xrightarrow{p.s.} T$,
- $\lim_{n \rightarrow \infty} E(T \wedge n) = E(T)$.

Démonstration. Vu que T est presque sûrement fini, il suffit de prouver que

$$(T \wedge n)(\omega) \rightarrow T(\omega) \quad (40.22)$$

pour tout ω tel que $T(\omega) = k$ pour tout $k \in \mathbb{N}$. Soient donc $\omega \in \Omega$ tel que $T(\omega) = k$ et $n > k$. Nous avons

$$(T \wedge n)(\omega) = T(\omega) \wedge n = k = T(\omega). \quad (40.23)$$

En ce qui concerne la seconde assertion, la suite de variables aléatoires $X_x = T \wedge n$ est croissante et positive, donc le théorème de la convergence monotone 15.160 montre que

$$\lim_{n \rightarrow \infty} E(T \wedge n) = E(T). \quad (40.24)$$

□

Remarque 40.11.

Notons la différence subtile entre $S_T(\omega)$ et $(S_T)(\omega)$. La première est la variable aléatoire

$$\omega' \mapsto S_{T(\omega')}(\omega) \quad (40.25)$$

et la seconde est le nombre $S_{T(\omega)}(\omega)$.

Théorème 40.12 (Théorème d'arrêt borné[501]).

Soit (X_n) une sur-martingale et $S \leq T$, deux temps d'arrêts bornés. Alors

- (1) les variables aléatoires X_S et X_T sont intégrables,
- (2) $E(X_T | \sigma(S)) \leq X_S$ presque sûrement.

Si par contre (X_n) est une martingale alors X_S et X_T sont bornées, et

$$E(X_T | \sigma(S)) = X_S. \quad (40.26)$$

Remarque 40.13.

Un cas particulier intéressant de ce théorème 40.12 est le cas $S = 0$ qui est un temps d'arrêt vérifiant $\mathcal{F}_0 = \{\Omega, \emptyset\}$. Si X est n'importe quelle variable aléatoire, la tribu engendrée $\sigma(X)$ est toujours indépendante de la tribu $\{\Omega, \emptyset\}$, donc le résultat $E(X_T | \mathcal{F}_S) = X_S$ donne

$$E(X_T) = X_0. \quad (40.27)$$

Théorème 40.14 (Premier théorème d'arrêt de Doob[502]).

Soient (X_n) une martingale et T un temps d'arrêt; deux pour la filtration (\mathcal{F}_n) . Nous supposons qu'une des trois propriétés suivantes soit vérifiée :

- (1) T est presque sûrement bornée.

². Dans [43], dans le problème de la ruine du joueur, la seconde assertion est avec une limite sup et non avec une limite normale.

(2) $E(T) < \infty$ et il existe une constante c telle que

$$E(|X_{n+1} - X_n| | \mathcal{F}_n) \leq c \quad (40.28)$$

sur l'événement $\{T \geq n\}$.

(3) Il existe une constante c telle que $|X_{T \wedge n}| \leq c$ presque sûrement³.

Alors X_T est une variable aléatoire presque sûrement bien définie nous avons

$$E(X_T) = E(X_0). \quad (40.29)$$

Si (X_n) est une sur-martingale, alors la conclusion est $E(X_T) \leq E(X_0)$ et si (X_n) est une sous-martingale, la conclusion est $E(X_T) \geq E(X_0)$.

Remarque 40.15.

Sous l'hypothèse (3), il est possible d'avoir $T = \infty$ sur un ensemble de mesure non nulle. Sur cet ensemble, la variable aléatoire X_T doit être définie de façon plus fine.

Problèmes et choses à faire

D'après la [page de discussion](#) de l'article sur Wikipédia, il semblerait que la seconde condition soit mal énoncée. Je n'ai pas vérifié.

Définition 40.16.

Nous disons que la martingale $(M_n)_{n \geq 1}$ est **terminée** s'il existe $M \in L^1(\Omega, \mathcal{A}, P)$ telle que $E(M | \mathcal{A}_n) = M$ pour tout $n > 1$.

Définition 40.17.

Un ensemble $H \subset L^1(\Omega, \mu)$ est **équi-intégrable** si

$$\lim_{a \rightarrow \infty} \left(\sup_{f \in H} \int_{|f| > a} |f(x)| d\mu(x) \right) = 0. \quad (40.30)$$

Notons dans cette définition que vu que $f \in L^1$ nous avons toujours

$$\lim_{a \rightarrow \infty} \int_{|f| > a} |f(x)| d\mu(x) = 0. \quad (40.31)$$

L'équi-intégrabilité donne une sorte d'uniformité en f de cette limite.

Théorème 40.18.

Si (M_n) est une martingale, nous avons équivalence entre

- (1) (M_n) converge dans L^1 ;
- (2) (M_n) est terminée ;
- (3) l'ensemble $\{M_n\}_{n \geq 1}$ est équi-intégrable.

Attention : en vertu de la proposition 28.19 et surtout de l'exemple 28.20, la convergence L^1 n'implique pas la convergence presque partout.

Théorème 40.19 (Théorème de Doob[236]).

À propos de convergence de martingales.

- (1) Toute martingale terminée converge presque sûrement et pour la norme L^1 .
- (2) Toute martingale bornée dans L^2 converge presque sûrement et pour la norme L^2 .

Proposition 40.20 ([503]).

Soient (M_n) une martingale et T un temps d'arrêt (pour la même filtration (\mathcal{B}_n)). Alors le processus $V_n = M_{n \wedge T}$ est une martingale.

3. Il est d'usage assez classique de noter $a \wedge b$ le minimum de a et b .

Démonstration. Nous décomposons V_n de la façon suivante :

$$V_n = M_{n \wedge T} = M_n \mathbb{1}_{T \geq n} + M_T \mathbb{1}_{T < n} = M_n \mathbb{1}_{T \geq n} + \sum_{k < n} M_k \mathbb{1}_{T=k}. \quad (40.32)$$

Nous avons, grâce au lemme 15.3,

$$\{T \geq n\} = \mathbb{C}\{T < n\} = \mathbb{C}\{T \leq n-1\} \in \mathcal{B}_{n-1} \quad (40.33)$$

et, si $k \leq n$,

$$\{T = k\} = \underbrace{\{T \leq k\}}_{\in \mathcal{B}_k} \setminus \underbrace{\{T \leq k-1\}}_{\in \mathcal{B}_{k-1}} \in \mathcal{B}_k \subset \mathcal{B}_n. \quad (40.34)$$

La forme (40.32) donne donc manifestement la \mathcal{B}_n -mesurabilité de V_n .

En ce qui concerne l'espérance nous devons calculer

$$E(V_{n+1} | \mathcal{B}_n) = E(M_{n+1} \mathbb{1}_{T \geq n+1} | \mathcal{B}_n) + \sum_{k < n+1} E(M_k \mathbb{1}_{T=k} | \mathcal{B}_n) \quad (40.35)$$

où nous avons utilisé la proposition 37.25. Étant donné que $\mathbb{1}_{T \geq n+1}$ et $\mathbb{1}_{T=k}$ sont des variables aléatoires \mathcal{B}_n -mesurables nous pouvons utiliser la proposition 37.41 pour les sortir :

$$E(V_{n+1} | \mathcal{B}_n) = \mathbb{1}_{T \geq n+1} M_n + \sum_{k \leq n} \mathbb{1}_{T=k} M_k = M_{T \wedge n} = V_n. \quad (40.36)$$

Pour cette ligne, nous avons aussi utilisé les égalités suivantes :

- $E(M_{n+1} | \mathcal{B}_n) = M_n$ parce que (M_n) est une martingale
- $E(M_k | \mathcal{B}_n) = M_k$ parce que M_k est \mathcal{B}_n -mesurable.

□

Définition 40.21.

Si (X_n) est un processus adapté à la filtration (\mathcal{F}_n) et si T est un temps d'arrêt \mathcal{F}_n -mesurable alors le **processus arrêté** à l'instant T est le processus $Y_n = X_{n \wedge T}$.

Nous avons déjà vu par la proposition 40.20 que si (X_n) est une martingale alors son processus arrêté est encore une martingale.

40.3 Décomposition de martingales

Définition 40.22 (Processus croissant prévisible[501]).

Un processus X_n adapté à la filtration \mathcal{F}_n est un processus **croissant prévisible** si

- (1) $A_0 = 0$
- (2) $A_n \leq A_{n+1}$; c'est cette condition qui correspond à « croissant »,
- (3) A_{n+1} est \mathcal{F}_n -mesurable ; c'est cette condition qui correspond à « prévisible ».

Proposition 40.23 (Décomposition de Doob pour une sous-martingale[501]).

Toute sous-martingale (X_n) s'écrit de façon unique sous la forme

$$X_n = M_n + A_n \quad (40.37)$$

où (M_n) est une martingale et (A_n) est un processus croissant prévisible.

Démonstration. Nous considérons le processus

$$\begin{cases} A_0 = 0 & (40.38a) \\ A_{n+1} = A_n + E(X_{n+1} - X_n | \mathcal{F}_n). & (40.38b) \end{cases}$$

Nous vérifions que cela est un processus croissant prévisible. D'abord

$$E(X_{n+1} - X_n | \mathcal{F}_n) = E(X_{n+1} | \mathcal{F}_n) - E(X_n | \mathcal{F}_n). \quad (40.39)$$

Le second terme est égal à X_n parce que cette variable aléatoire est \mathcal{F}_n -mesurable tandis que (X_n) étant une sous-martingale nous avons $E(X_{n+1} | \mathcal{F}_n) \geq X_n$. Nous avons donc bien $A_{n+1} \geq A_n$ et le processus (A_n) est croissant.

En ce qui concerne la prévisibilité nous devons prouver que A_{n+1} est \mathcal{F}_n -mesurable. D'une part A_n est \mathcal{F}_n -mesurable et d'autre part par définition de l'espérance conditionnelle, la variable aléatoire $E(X_{n+1} - X_n | \mathcal{F}_n)$ est également \mathcal{F}_n -mesurable.

Nous posons alors $M_n = X_n - A_n$ et nous devons prouver que cela est une martingale. Nous avons

$$E(M_{n+1} - M_n | \mathcal{F}_n) = E(X_{n+1} - X_n | \mathcal{F}_n) - E(A_{n+1} - A_n | \mathcal{F}_n). \quad (40.40)$$

Le second terme vaut

$$E(A_{n+1} - A_n | \mathcal{F}_n) = E\left(E(X_{n+1} - X_n | \mathcal{F}_n) | \mathcal{F}_n\right) = E(X_{n+1} - X_n | \mathcal{F}_n) \quad (40.41)$$

par la proposition 37.36. Le processus (M_n) est donc une martingale. La preuve de l'existence d'une décomposition (40.37) est achevée.

Nous passons maintenant à l'unicité en posant $X_n = M_n + A_n = M'_n + A'_n$. Nous avons $A_0 = A'_0 = 0$ et $A'_n = X_n - M'_n$, donc

$$A'_{n+1} - A'_n = X_{n+1} - X_n + M'_{n+1} - M'_n = X_{n+1} - X_n - (M'_{n+1} - M'_n). \quad (40.42)$$

Nous appliquons $E(\cdot | \mathcal{F}_n)$ des deux côtés de cette égalité :

$$\underbrace{E(A'_{n+1} - A'_n | \mathcal{F}_n)}_{=A'_{n+1} - A'_n} = E(X_{n+1} - X_n | \mathcal{F}_n) - \underbrace{E(M'_{n+1} - M'_n | \mathcal{F}_n)}_{=0}. \quad (40.43)$$

Nous avons utilisé le fait que (M_n) étant une martingale, $E(M_{n+1} - M_n | \mathcal{F}_n) = 0$, et idem avec (M'_n) . Donc

$$A'_{n+1} - A'_n = E(X_{n+1} - X_n | \mathcal{F}_n) = E(M_{n+1} - M_n | \mathcal{F}_n) + E(A_{n+1} - A_n | \mathcal{F}_n) = A_{n+1} - A_n. \quad (40.44)$$

Nous avons donc montré que $A_{n+1} - A_n = A'_{n+1} - A'_n$ et donc que $A_n = A'_n$ pour tout n . Nous en déduisons immédiatement que $M_n = M'_n$ pour tout n et l'unicité de la décomposition. \square

Lemme 40.24.

Si (X_n) est une martingale de carré intégrable adaptée à la filtration (\mathcal{F}_n) alors

- (1) Le processus (X_n^2) est une sous-martingale.
- (2) Si $X_n^2 = M_n + A_n$ est la décomposition de Doob, alors

$$A_n = \sum_{i=1}^n \left(E(X_i^2 | \mathcal{A}_{i-1}) - X_{i-1}^2 \right) = \sum_{i=1}^n E\left((X_i - X_{i-1})^2 | \mathcal{A}_{i-1} \right). \quad (40.45)$$

Démonstration. Pour la première assertion, nous utilisons l'inégalité de Jensen 37.52 :

$$E(X_n^2 | \mathcal{F}_{n-1}) \geq \left(E(X_n | \mathcal{F}_{n-1}) \right)^2 = X_{n-1}^2 \quad (40.46)$$

parce que $E(X_n | \mathcal{F}_{n-1}) = X_{n-1}$ du fait que (X_n) soit une martingale.

En ce qui concerne la seconde assertion, nous nous souvenons que le processus prévisible de la décomposition de Doob d'une sous-martingale est donné par la récurrence (40.38) que nous récrivons ici :

$$\begin{cases} A_0 = 0 & (40.47a) \\ A_{n+1} = A_n + E(X_{n+1}^2 - X_n^2 | \mathcal{F}_n) & (40.47b) \end{cases}$$

Vu que X_n^2 est \mathcal{F}_n -mesurable, il peut sortir de l'espérance :

$$A_{n+1} = A_n + E(X_{n+1}^2 | \mathcal{F}_n) - X_n^2 \quad (40.48)$$

et donc

$$A_n = \sum_{i=1}^n \left(E(X_i^2 | \mathcal{F}_{i-1}) - X_{i-1}^2 \right). \quad (40.49)$$

Pour obtenir la dernière partie de (40.45) nous travaillons un peu :

$$E((X_i - X_{i-1})^2 | \mathcal{F}_{i-1}) = E(X_i^2 + X_{i-1}^2 - 2X_i X_{i-1} | \mathcal{F}_{i-1}) \quad (40.50a)$$

$$= E(X_i^2 | \mathcal{F}_{i-1}) + X_{i-1}^2 - 2E(X_i X_{i-1} | \mathcal{F}_{i-1}) \quad (40.50b)$$

$$= E(X_i^2 | \mathcal{F}_{i-1}) + X_{i-1}^2 - 2X_{i-1} E(X_i | \mathcal{F}_{i-1}) \quad (40.50c)$$

$$= E(X_i^2 | \mathcal{F}_{i-1}) + X_{i-1}^2 - 2X_{i-1} X_i \quad (40.50d)$$

où nous avons utilisé la proposition 37.41 pour obtenir (40.50c). \square

40.4 Problème de la ruine du joueur

Nous considérons un joueur compulsif qui joue à un jeu très simple⁴ : il joue à pile ou face contre la banque avec une pièce truquée. Si pile sort, la banque donne 1 au joueur et si c'est face, c'est le joueur qui donne 1 à la banque. Nous nommons a la fortune initiale du joueur, b celle de la banque et p la probabilité d'obtenir pile.

Nous supposons que le jeu se poursuit jusqu'à la ruine du joueur ou de la banque. La modélisation est comme suit : nous considérons (Y_n) une suite de variables aléatoires indépendantes et identiquement distribuées de loi

$$Y_n \sim p\delta_1 + (1-p)\delta_{-1}. \quad (40.51)$$

C'est le résultat financier pour le joueur du n^e lancé. La fortune du joueur au bout de n lancers est la variable aléatoire

$$S_n = a + \sum_{j=1}^n Y_j. \quad (40.52)$$

Nous notons $Y_0 = a$.

Nous considérons la filtration

$$\mathcal{A}_n = \sigma(S_i \text{ tel que } 0 \leq i \leq n) = \sigma(Y_i \text{ tel que } 0 \leq i \leq n), \quad (40.53)$$

et le temps d'arrêt du jeu :

$$T = \inf\{n \geq 1 \text{ tel que } S_n \in \{0, a+b\}\}; \quad (40.54)$$

c'est le temps qu'il faut pour que tout l'argent appartienne soit au joueur soit à la banque.

Nous voulons étudier les paramètres suivants :

- (1) $\rho = P(S_T = a+b)$, c'est-à-dire la probabilité que ce soit le joueur qui gagne contre la banque.
- (2) $P(T < \infty)$, c'est-à-dire la probabilité que le jeu se finisse.
- (3) $E(T)$, la durée moyenne du jeu.

Lemme 40.25.

Le processus S_n du problème de la ruine du joueur est vérifié

$$E(S_n | \mathcal{A}_{n-1}) = S_{n-1} + p - q. \quad (40.55)$$

De plus le processus S_n est

4. Le gros des choses dites à propos de la ruine du joueur provient de [43].

- (1) une martingale si $p = q = \frac{1}{2}$,
 (2) une sous-martingale si $p > q$.

Démonstration. Pour $n \geq 1$ nous avons

$$E(S_n | \mathcal{A}_{n-1}) = a + \sum_{j=1}^n E(Y_j | \mathcal{A}_{n-1}) = a + \sum_{j=1}^{n-1} E(Y_j | \mathcal{A}_{n-1}) + E(Y_n | \mathcal{A}_{n-1}). \quad (40.56)$$

Si $j \leq n-1$ alors $Y_j \in m(\mathcal{A}_{n-1})$. Mais nous savons que si X est \mathcal{F} -mesurable, alors $E(X | \mathcal{F}) = X$ (c'est la définition de l'espérance conditionnelle), donc $\sum_{j=1}^{n-1} E(Y_j | \mathcal{A}_{n-1}) = \sum_{j=1}^{n-1} Y_j$.

En ce qui concerne le terme $j = n$ nous utilisons le fait que $\sigma(Y_n)$ soit une tribu indépendante de \mathcal{A}_{n-1} ; nous avons donc au final pour tout j que $E(Y_j | \mathcal{A}_{n-1}) = E(Y_j) = p - q$. Nous avons donc

$$E(S_n | \mathcal{A}_{n-1}) = S_{n-1} + p - q. \quad (40.57)$$

Si $p = q = \frac{1}{2}$ alors c'est une martingale, et si $p > q$ c'est une sous-martingale. \square

Lemme 40.26.

La variable aléatoire T est un temps d'arrêt.

Démonstration. Par définition $T = \inf\{n \geq 1 \text{ tel que } S_n \in \{0, a+b\}\}$. Vu que les variables aléatoires S_i avec $i \leq n$ sont \mathcal{F}_n -mesurables, les ensembles $\{S_k \notin \{0, a+b\}\}$ avec $k \leq n$ sont \mathcal{F}_n -mesurables. Donc les ensembles

$$\{T = n\} = \bigcap_{k \leq n} \{S_k \notin \{0, a+b\}\} \cap \{S_n \in \{0, a+b\}\} \quad (40.58)$$

sont \mathcal{F}_n -mesurables. Nous en concluons que l'ensemble $\{T \leq n\}$ est également mesurable. \square

40.4.1 Le cas où la pièce est truquée

Nous supposons être dans le cas $p > q$.

40.4.1.1 Introduction d'une martingale

Considérons le processus

$$\begin{cases} A_0 = 0 & (40.59a) \\ A_n = A_{n-1} + E(S_n - S_{n-1} | \mathcal{A}_{n-1}). & (40.59b) \end{cases}$$

Vu que $E(S_n | \mathcal{A}_{n-1}) = S_{n-1} + p - q$ (lemme 40.25) et que $E(S_{n-1} | \mathcal{A}_{n-1}) = S_{n-1}$ (parce que S_{n-1} est dans la tribu de \mathcal{A}_{n-1}), nous avons $A_n = A_{n-1} + (p - q)$ et donc

$$A_n = n(p - q). \quad (40.60)$$

Ce processus (A_n) est croissant et prévisible. Nous introduisons le processus

$$M_n = S_n - A_n \quad (40.61)$$

et nous montrons que c'est une martingale⁵. Nous conditionnons la définition (40.61) par rapport à \mathcal{A}_{n-1} :

$$E(M_n | \mathcal{A}_{n-1}) = E(S_n | \mathcal{A}_{n-1}) - \underbrace{E(A_n | \mathcal{A}_{n-1})}_{=A_n} \quad (40.62a)$$

$$= A_n - A_{n-1} + E(S_{n-1} | \mathcal{A}_{n-1}) - A_n \quad (40.62b)$$

$$= E(S_{n-1} | \mathcal{A}_{n-1}) - A_{n-1}. \quad (40.62c)$$

Mais S_{n-1} est \mathcal{A}_{n-1} -mesurable, donc $E(S_{n-1} | \mathcal{A}_{n-1}) = S_{n-1}$ et

$$E(M_n | \mathcal{A}_{n-1}) = S_{n-1} - A_{n-1} = M_{n-1}, \quad (40.63)$$

ce qui signifie que (M_n) est une martingale.

5. Ceci est un peu le contraire de la décomposition de Doob.

40.4.1.2 Finitude du temps d'arrêt

Nous montrons maintenant, en étudiant $M_{T \wedge n}$ que T est intégrable et nous prouvons que $P(T = \infty) = 0$.

Proposition 40.27.

La variable aléatoire T est intégrable, et $P(T = \infty) = 0$, c'est-à-dire que le jeu se termine presque certainement après un temps fini.

Démonstration. Nous rappelons que le lemme 40.26 nous indique que T est un temps d'arrêt. Le temps d'arrêt $T \wedge n$ est borné (par n évidemment) et nous pouvons donc lui appliquer le théorème d'arrêt 40.14 pour dire que

$$E(M_{T \wedge n}) = E(M_0). \quad (40.64)$$

Le membre de droite est simple parce que $M_0 = S_0 - A_0 = S_0 = a$ parce que c'est l'argent de départ du joueur. Pour l'autre :

$$E(M_{T \wedge n}) = E(S_{T \wedge n}) - E(A_{T \wedge n}). \quad (40.65)$$

D'une part, $E(A_{T \wedge n}) = E((T \wedge n)(p - q))$ et d'autre part, $E(S_{T \wedge n}) \leq a + b$ parce que S_T vaut zéro ou $a + b$ (avec des probabilités encore inconnues⁶). En combinant avec ce qui était dit juste au dessus et remarquant que $(p - q)E(T \wedge n) \geq 0$ nous pouvons écrire

$$0 \leq (p - q)E(T \wedge n) \leq b. \quad (40.66)$$

La suite de variables aléatoires $T \wedge n$ est donc croissante, positive et intégrable⁷ et donc nous avons du travail pour le théorème de la convergence monotone 15.160. La variable aléatoire T est alors mesurable et⁸

$$\lim_{n \rightarrow \infty} E(T \wedge n) = E(T). \quad (40.67)$$

Notons que nous n'avons pas encore prouvé que $E(T) < \infty$, mais en passant à la limite dans (40.66) nous écrivons

$$0 \leq (p - q)E(T) \leq b. \quad (40.68)$$

Maintenant nous avons prouvé que T est intégrable et même L^1 . Par conséquent

$$P(T = \infty) = 0. \quad (40.69)$$

Le jeu se termine donc presque certainement après un temps fini. \square

40.4.1.3 Temps moyen de jeu

Le lemme 40.10 nous indique que $S_{T \wedge n} \xrightarrow{p.s.} S_T$.

Nous avons les bornes $0 \leq S_{T \wedge n} \leq a + b$ et comme $a + b$ est intégrable, $S_{T \wedge n}$ l'est aussi et nous pouvons parler de $E(S_{T \wedge n})$. Repartons de (40.65) :

$$a = E(M_0) = E(M_{T \wedge n}) = E(S_{T \wedge n}) - E(A_{T \wedge n}) = E(S_{T \wedge n}) - (p - q)E(T \wedge n). \quad (40.70)$$

La variable aléatoire $S_{T \wedge n}$ est majorée par $a + b$ indépendamment de n ; donc le théorème de la convergence dominée 15.184 donne $\lim_{n \rightarrow \infty} E(S_{T \wedge n}) = E(S_T)$. En ce qui concerne le second terme, la convergence dominée ne fonctionne pas parce que $T \wedge n$ n'est pas a priori majoré par quelque chose d'indépendant de n , mais le théorème de la convergence monotone donne $\lim_{n \rightarrow \infty} E(T \wedge n) = E(T)$. Au final en passant à la limite dans (40.70) nous avons

$$a = E(S_T) - (p - q)E(T). \quad (40.71)$$

6. Mais on y travaille.

7. Je rappelle que les constantes sont des fonctions intégrables sur Ω . Oui, je sais, quand on est habitué à faire de l'analyse sur \mathbb{R}^n c'est un truc qu'on perd toujours un peu de vue.

8. Dans [43], l'équation (40.67) vient avec une \limsup et non une limite normale. Je ne comprends pas pourquoi.

Étant donné que $T > 0$ et $p - q > 0$ nous pouvons récrire cela sous la forme

$$0 \leq (p - q)E(T) = E(S_T) - a. \quad (40.72)$$

Par définition de T nous avons aussi

$$E(S_T) = (a + b)P(S_T = a + b) + 0 \cdot P(S_T = 0) = \rho(a + b). \quad (40.73)$$

Nous déduisons

$$E(T) = \frac{(a + b)\rho - a}{p - q}. \quad (40.74)$$

Ne crions pas victoire trop vite : nous n'avons pas encore d'expression de $\rho = P(S_T = a + b)$. Le temps moyen de jeu n'est donc pas encore tout à fait connu.

40.4.1.4 Probabilité de victoire du joueur

Nous avons besoin d'exprimer ρ en termes de a , b et p . Pour cela nous introduisons la variable aléatoire⁹

$$U_n = \left(\frac{p}{q}\right)^{S_n}. \quad (40.75)$$

Nous commençons par prouver que c'est une martingale en calculant

$$E(U_n | \mathcal{A}_{n-1}) = E\left(\left(\frac{q}{p}\right)^{S_{n-1}} \left(\frac{q}{p}\right)^{Y_n} \mid \mathcal{A}_{n-1}\right) \quad (40.76)$$

Nous utilisons la proposition 37.41. Dans notre cas, S_{n-1} et Y_n sont des variables aléatoires \mathcal{A}_n -mesurables ; la variable aléatoire Y_n est même \mathcal{A}_{n-1} -mesurable et sort donc du conditionnement ; nous avons donc

$$E(U_n | \mathcal{A}_{n-1}) = \left(\frac{q}{p}\right)^{S_{n-1}} E\left(\left(\frac{q}{p}\right)^{Y_n}\right) \quad (40.77)$$

Nous allons utiliser le théorème de transfert 37.60 :

$$E(s^{Y_n}) = \int_{\Omega} s^{Y_n(\omega)} dP(\omega) = \int_{Y_n=1} s dP(\omega) + \int_{Y_n=-1} \frac{1}{s} dP(\omega). \quad (40.78)$$

Mais nous savons que $P(Y_n = 1) = p$ et $P(Y_n = -1) = 1 - p = q$, donc

$$E(s^{Y_n}) = ps + \frac{1-p}{s} \quad (40.79)$$

et

$$E\left(\left(\frac{p}{q}\right)^{Y_n}\right) = p + q = 1. \quad (40.80)$$

Donc

$$E(U_n | \mathcal{A}_{n-1}) = \left(\frac{q}{p}\right)^{S_{n-1}} = U_{n-1}, \quad (40.81)$$

ce qui prouve que (U_n) est une martingale.

Par définition nous avons toujours $S_n \geq 0$ tant que $n \leq T$ ¹⁰, donc $U_{T \wedge n} \in [0, 1]$. Il est donc évident que si $a \geq 1$ nous avons

$$\int_{|U_{T \wedge n}| > a} |U_{T \wedge n}| dP = 0 \quad (40.82)$$

9. Nous dirons un mot sur ce choix dans le « petit complément » plus bas

10. Pour $n > T$ le jeu est terminé, donc on ne se pose pas la question.

parce que le domaine d'intégration est vide. Donc les variables aléatoires $V_n = U_{T \wedge n}$ sont équi-intégrables¹¹ et le théorème 40.18 montre que la martingale (V_n) est terminée; par ricochet¹² le théorème de Doob 40.19 montre qu'il existe une variable aléatoire X telle que $V_n \xrightarrow{p.s.} X$. Nous allons prouver que $X = U_T$ presque partout. Nous savions déjà (voir l'équation (40.22) et ses alentours) que

$$S_{n \wedge T} \xrightarrow{p.s.} S_T. \quad (40.83)$$

Nous avons alors (au sens du presque sûrement) :

$$\lim_{n \rightarrow \infty} V_n = \lim_{n \rightarrow \infty} U_{T \wedge n} = \lim_{n \rightarrow \infty} \left(\frac{q}{p}\right)^{S_{T \wedge n}} = \left(\frac{q}{p}\right)^{S_T} = U_T. \quad (40.84)$$

Donc par unicité de la limite presque partout nous avons $X = U_T$ presque partout. Par le théorème de transfert 37.60 nous évaluons

$$E(U_T) = \left(\frac{q}{p}\right)^0 P(S_T = 0) + \left(\frac{q}{p}\right)^{a+b} P(S_T = a+b) = (1 - \rho) + \left(\frac{q}{p}\right)^{a+b} \rho. \quad (40.85)$$

La remarque 40.13 nous permet de dire que

$$E(U_{T \wedge n}) = U_0. \quad (40.86)$$

Mais par définition

$$U_0 = \left(\frac{q}{p}\right)^{S_0} = \left(\frac{q}{p}\right)^a, \quad (40.87)$$

donc nous avons

$$E(U_{T \wedge n}) = \left(\frac{q}{p}\right)^a. \quad (40.88)$$

Nous voudrions passer à la limite $n \rightarrow \infty$ dans cette équation. Pour permuter la limite et l'espérance, il faut utiliser le théorème de la convergence dominée 15.184. Vu que nous avons choisi $q > p$, nous avons $q/p > 1$ et donc $U_{T \wedge n} \leq (q/p)^{a+b}$, ce qui montre que la fonction $\omega \mapsto (U_{T \wedge n})(\omega)$ est majorée par une constante (qui est une fonction intégrable). Nous pouvons donc permuter la limite et l'espérance :

$$\lim_{n \rightarrow \infty} E(U_{T \wedge n}) = E\left(\lim_{n \rightarrow \infty} U_{T \wedge n}\right). \quad (40.89)$$

Mais nous avons déjà montré que $U_{T \wedge n} \xrightarrow{p.s.} U_T$. Donc

$$E(U_T) = \left(\frac{q}{p}\right)^a. \quad (40.90)$$

En égalisant avec l'expression (40.85) de $E(U_T)$ nous trouvons

$$\rho = \frac{\left(\frac{q}{p}\right)^a - 1}{\left(\frac{q}{p}\right)^{a+b} - 1} \quad (40.91)$$

et ensuite nous trouvons $E(T)$ en remettant ce ρ dans l'expression (40.74) donnée plus haut.

40.4.2 Le cas où la pièce est non truquée

Maintenant $p = q = 1/2$.

11. Définition 40.17.

12. Nous rappelons que la convergence L^1 n'implique pas la convergence presque partout.

40.4.2.1 Probabilité de gagner

Le lemme 40.25 nous indique alors que (S_n) est une martingale et le lemme 40.24 nous permet de dire que (S_n^2) est alors une sous-martingale. Le processus croissant prévisible de (S_n^2) est donné par (40.38) qui en adaptant les notations est

$$\begin{cases} B_0 = 0 \\ B_n = B_{n-1} + E\left((S_n - S_{n-1})^2 | \mathcal{A}_{n-1}\right). \end{cases} \quad (40.92a) \quad (40.92b)$$

Nous avons toujours $S_n - S_{n-1} = \pm 1$ parce que soit le joueur gagne soit le joueur perd, mais de toutes façons sa fortune varie de 1 à chaque étape du jeu. Donc (40.92b) nous donne $B_n = B_{n-1} + 1$ et

$$B_n = n. \quad (40.93)$$

Cela nous dit que la variable aléatoire

$$S_n^2 - B_n = S_n^2 - n \quad (40.94)$$

est une martingale (une sur-martingale moins son processus prévisible croissant). Nous lui appliquons le théorème d'arrêt 40.12 avec les temps d'arrêt 0 et $T \wedge n$:

$$E(S_{T \wedge n}^2 - T \wedge n | \mathcal{F}_0) = S_0^2 - 0 \quad (40.95)$$

où \mathcal{F}_0 est la tribu engendrée par la variable aléatoire 0, c'est-à-dire $\{\Omega, \emptyset\}$. Cette tribu est indépendante de toute autre tribu et nous pouvons donc supprimer le conditionnement dans (40.95). Nous avons aussi $S_0 = a$ par définition. Avec tout ça nous avons la majoration

$$E(T \wedge n) = E(S_{T \wedge n}^2) - a^2 \leq (a+b)^2 - a^2 \quad (40.96)$$

parce que S_k est toujours positif et entre 0 et $a+b$. En utilisant le lemme 40.10 et en passant à la limite,

$$E(T) \leq (a+b)^2 - a^2. \quad (40.97)$$

En particulier, $T \in L^1(\Omega)$ et $P(T < \infty) = 1$.

En suivant exactement les mêmes étapes que dans le lemme 40.10 nous avons aussi

$$\lim_{n \rightarrow \infty} S_{T \wedge n} = S_T \quad (40.98)$$

presque partout. De plus nous savons que

$$0 \leq S_{T \wedge n}^2 \leq (a+b)^2, \quad (40.99)$$

et nous pouvons donc utiliser le théorème de la convergence dominée 15.184 pour dire que

$$\lim_{n \rightarrow \infty} E(S_{T \wedge n}^2) = E(S_T^2). \quad (40.100)$$

Nous montrons à présent que $S_{T \wedge n} \xrightarrow{L^2} S_T$. Pour cela nous devons évaluer la limite

$$\lim_{n \rightarrow \infty} \int_{\Omega} |S_{T \wedge n} - S_T|^2. \quad (40.101)$$

La fonction $|S_{T \wedge n} - S_T|^2$ est majorée par $(a+b)^2$ et nous pouvons à nouveau appliquer la convergence dominée :

$$\lim_{n \rightarrow \infty} E(|S_{T \wedge n} - S_T|^2) = \lim_{n \rightarrow \infty} \int_{\Omega} |S_{T(\omega) \wedge n}(\omega) - S_T(\omega)|^2 dP(\omega) = \int_{\Omega} \lim_{n \rightarrow \infty} |S_{T \wedge n} - S_T|^2 = 0. \quad (40.102)$$

La même chose en n'écrivant pas les carrés montre que l'on a aussi $S_{T \wedge n} \xrightarrow{L^1} S_T$.

Il n'y a pas que $n \mapsto S_n^2 - n$ qui est une martingale. Il y a aussi (S_n) lui-même (lemme 40.25). Nous pouvons lui appliquer le théorème d'arrêt 40.12 pour les temps d'arrêts $T \wedge n$ et 0 :

$$E(S_{T \wedge n}) = E(S_0) = a. \quad (40.103)$$

En passant à la limite, $E(S_T) = a$. L'espérance $E(S_T)$ peut par ailleurs être calculée comme

$$E(S_T) = 0 \cdot P(S_T = 0) + (a + b)P(S_T = a + b). \quad (40.104)$$

En égalisant les valeurs (40.103) et (40.104) de $E(S_T)$ nous trouvons

$$\rho = \frac{a}{a + b}. \quad (40.105)$$

Cette formule est assez logique : la probabilité que le joueur gagne est égale à la proportion d'argent en jeu qu'il a amené.

40.4.2.2 Temps moyen de jeu

Nous calculons maintenant l'espérance $E(T)$ du temps de jeu (sans compter les pauses ni les jours de fermeture du casino¹³).

Nous recopions la première égalité de (40.96) sous la forme

$$a^2 = E(S_{T \wedge n}^2 - T \wedge n) \quad (40.106)$$

et nous passons à la limite¹⁴ en sachant que $E(S_T^2) = \rho(a + b)^2$:

$$a^2 = \rho(a + b)^2 - E(T). \quad (40.107)$$

En reprenant la valeur (40.105) de ρ ,

$$E(T) = ab. \quad (40.108)$$

Et là, on voit que si le joueur amène 1000 euros contre une banque qui en a un million, et si ils jouent toutes les secondes, on en a pour 32 ans de jeu en moyenne.

Voilà. C'est fini pour la ruine du joueur.

40.4.3 Un petit complément

Nous avons introduit lors de l'équation (40.75) la variable aléatoire $U_n = (p/q)^{S_n}$. Sans aller jusqu'à motiver complètement ce choix, nous nous proposons maintenant de voir que parmi les variables aléatoires $U_n = s^{S_n}$, le choix $s = p/q$ est le seul qui donne une martingale.

Soit donc $U_n = s^{S_n}$ et exprimons le fait que ce soit une martingale. Nous avons

$$E(U_n | \mathcal{A}_{n-1}) = E(s^{S_{n-1} Y_n} | \mathcal{A}_{n-1}) \quad (40.109a)$$

$$= s^{S_{n-1}} E(s^{Y_n} | \mathcal{A}_{n-1}) \quad (40.109b)$$

$$= s^{S_{n-1}} E(s^{Y_n}). \quad (40.109c)$$

Le passage à (40.109b) se justifie en disant que $s^{S_{n-1}}$ est une variable aléatoire bornée et \mathcal{A}_{n-1} -mesurable, et en invoquant proposition 37.41. La variable aléatoire Y_n vaut 1 avec probabilité p et -1 avec probabilité q ; donc l'espérance est vite vue :

$$E(s^{Y_n}) = ps + q \frac{1}{s} \quad (40.110)$$

et nous avons

$$E(U_n | \mathcal{A}_{n-1}) = \left(ps + q \frac{1}{s} \right) s^{S_{n-1}} = \left(ps + \frac{q}{s} \right) U_{n-1}. \quad (40.111)$$

13. Le joueur est un *vrai* joueur compulsif.

14. Comme il est dit dans La Grande Illusion, à quoi sert un n ? À passer à la limite.

Pour que (U_n) soit une martingale il faut (et il suffit) que

$$ps + \frac{q}{s} = 1. \quad (40.112)$$

Les solutions de cette équation sont $s \in \{1, \frac{p}{q}\}$. C'est évidemment $s = p/q$ qui donne une martingale non triviale. Attention pour être complet, il faut se demander ce qu'il se passe si $s = 0$ séparément parce que manifestement l'équation (40.112) ne traite pas ce cas. Encore une fois, en repartant du début, $s = 0$ ne se révèle pas être une martingale très excitante.

Bref, nous devons poser

$$U_n = \left(\frac{p}{q}\right)^{S_n} \quad (40.113)$$

pour avoir une martingale.

Chapitre 41

Processus de Poisson

41.1 Processus de Poisson

Définition 41.1.

Une famille de variables aléatoires $(N_t)_{t \geq 0}$ est une **processus de Poisson** d'intensité λ s'il existe une suite de variables aléatoires indépendantes et identiquement distribuées $(T_k)_{k \in \mathbb{N}}$ de loi $\mathcal{E}(\lambda)$ telles que

$$N_t = \sup\{n \geq 0 \text{ tel que } \sum_{k=1}^n T_k \leq t\}. \quad (41.1)$$

Si nous posons $S_n = \sum_{k=1}^n T_k$, alors nous avons une expression plus pratique pour N_t :

$$N_t = \sum_{n=1}^{\infty} \mathbb{1}_{\{S_n \leq t\}}. \quad (41.2)$$

Nous avons par la proposition 37.103 vu que $N_t \sim \mathcal{P}(\lambda t)$.

Pour chaque $\omega \in \Omega$, la fonction $t \mapsto N_t(\omega)$ est une fonction (pas du tout strictement) croissante à valeurs dans \mathbb{N} . Cette fonction part de 0 et fait un saut de taille 1 après des intervalles de temps $T_1(\omega), T_2(\omega), \dots$. Elle est continue à droite.

Nous avons les égalités d'événements suivantes qui sont pratiques :

$$\{s < S_n \leq t\} = \{N_t \geq n > N_s\} \quad (41.3a)$$

$$\{N_t = n\} = \{S_n \leq t \leq S_{n+1}\}. \quad (41.3b)$$

Théorème 41.2.

Les variables aléatoires $(N_t)_{t \geq 0}$ est un processus de Poisson d'intensité λ si et seulement si elles vérifient les trois propriétés suivantes.

Accroissements indépendants Pour tout choix $0 < t_0 < t_1 < \dots < t_n$, les variables aléatoires $N_{t_{i+1}} - N_{t_i}$ sont indépendantes.

Accroissements stationnaires Si $0 < s < t$ et $h > 0$ alors

$$N_{t+h} - N_{s+h} \stackrel{\mathcal{L}}{=} N_t - N_s, \quad (41.4)$$

c'est-à-dire que les accroissements décalés suivent les mêmes lois.

Poisson Pour tout t nous avons $N_t \sim \mathcal{P}(\lambda t)$.

Une conséquence des accroissements stationnaires est que $N_t - N_s \stackrel{\mathcal{L}}{=} N_{t-s} - N_0 = N_{t-s}$ parce que $N_0 = 0$.

Proposition 41.3.

Si (N_t) est un processus de Poisson d'intensité λ , alors

$$\lim_{t \rightarrow \infty} N_t = +\infty \quad (41.5)$$

presque surement. De plus

$$\lim_{t \rightarrow \infty} \frac{N_t}{t} = \lambda \quad (41.6)$$

presque surement.

La relation (41.6) est appelée **loi des grands nombres**.

Démonstration. Par définition nous savons que

$$N_t = \sup\{n \geq 0 \text{ tel que } S_n \leq t\}. \quad (41.7)$$

Évidemment la fonction $t \mapsto N_t$ est croissante, donc la limite

$$\lim_{t \rightarrow \infty} N_t(\omega) \quad (41.8)$$

existe dans $[0, \infty]$. Nous pouvons nous restreindre à $t \in \mathbb{N}$ et considérer $L(\omega) = \lim_{n \rightarrow \infty} N_n(\omega)$. Par somme télescopique avec $N_0 = 0$,

$$\frac{N_n}{n} = \frac{\sum_{k=1}^n (N_k - N_{k-1})}{n}. \quad (41.9)$$

Étant donné que le processus est de Poisson, les variables aléatoires $(N_k - N_{k-1})_{k=1, \dots, n}$ sont indépendantes et suivent toutes la loi de $N_1 - N_0$, c'est-à-dire la loi de N_1 . Encore par le fait que N_t soit de Poisson nous savons que $N_1 \sim \mathcal{P}(\lambda)$. La loi des grands nombres (37.83) appliquée aux variables aléatoires $N_k - N_{k-1}$ nous dit que

$$\frac{N_n}{n} \xrightarrow{p.s.} E(N_1) = \lambda > 0. \quad (41.10)$$

Du coup $N_n \rightarrow \infty$ et $L(\omega) = \infty$.

Nous démontrons maintenant la loi des grands nombres pour les processus de Poisson. Étant donné que pour les entiers $N_n/n \rightarrow \lambda$, pour les réels, si la limite existe, ça ne peut pas être autre chose. Si nous notons \bar{t} la partie entière de $t \in \mathbb{R}^+$,

$$\frac{N_t}{t} = \frac{N_t - N_{\bar{t}}}{t} + \frac{N_{\bar{t}}}{t}. \quad (41.11)$$

Le second terme est relativement simple à traiter :

$$\frac{N_{\bar{t}}}{t} = \underbrace{\frac{N_{\bar{t}}}{\bar{t}}}_{\rightarrow \lambda} \cdot \underbrace{\frac{\bar{t}}{t}}_{\rightarrow 1}. \quad (41.12)$$

où nous avons utilisé le premier point, \bar{t} étant entier. Pour le premier terme nous savons que $t \mapsto N_t$ est croissante et donc que

$$\frac{N_t - N_{\bar{t}}}{t} \leq \frac{N_{\bar{t}+1} - N_{\bar{t}}}{t} = \frac{N_{\bar{t}+1} - N_{\bar{t}}}{\bar{t} + 1} \frac{\bar{t} + 1}{t}. \quad (41.13)$$

Le second facteur tend vers 1 lorsque $t \rightarrow \infty$. Le premier s'écrit

$$\frac{N_n - N_{n-1}}{n} \quad (41.14)$$

et tend vers zéro en tant que terme général de la série (41.9) qui converge. \square

Proposition 41.4.

La variable aléatoire N_t/t est un estimateur sans biais de λ . De plus il converge vers λ en moyenne quadratique.

Démonstration. Vu que $N_t/t \rightarrow \lambda$ presque sûrement, la variable aléatoire N_t/t est un estimateur de λ . Le fait qu'il soit sans biais a été fait dans l'exemple 38.25.

D'autre part nous avons (voir théorème 37.96)

$$\text{Var}\left(\frac{N_t}{t}\right) = \frac{1}{t^2} \text{Var}(N_t) = \frac{\lambda}{t}. \quad (41.15)$$

En appliquant la formule $\text{Var}(X) = E(X^2) - E(X)^2$ à $X = N_t/t$ nous trouvons

$$E\left(\frac{N_t^2}{t^2}\right) = \frac{\lambda}{t} + \lambda^2. \quad (41.16)$$

Cela montre que $\frac{N_t}{t} \xrightarrow{L^2} \lambda$. □

Pour le théorème central limite d'un processus de Poisson, nous visons un résultat du style de

$$\frac{\frac{1}{n} \sum_i X_i - mn}{\sigma\sqrt{n}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1). \quad (41.17)$$

Nous écrivons le théorème central limite pour le nombre de sauts que le processus de Poisson a connu en un temps t . Le rôle de la moyenne empirique est joué par N_t . Nous considérons avoir fait *une seule expérience* qui a duré un temps t . Donc le rôle de n est joué par 1 (et non t comme on pourrait le croire). Pour le reste, le nombre de succès en un temps t d'une variable aléatoire exponentielle de paramètre λ est une variable aléatoire de Poisson de paramètre λt , en vertu de ce qui est raconté au point 37.5.8. C'est cela qui motive l'énoncé suivant.

Théorème 41.5 (Théorème central limite pour les processus de Poisson).

Si $(N_t)_{t>0}$ est un processus de Poisson de paramètre λ , alors nous avons

$$\frac{N_t - \lambda t}{\sqrt{\lambda t}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1). \quad (41.18)$$

Remarque 41.6.

Avant de nous lancer dans la démonstration, remarquons que si nous nous limitons à $t \in \mathbb{N}$, alors nous avons

$$\frac{N_n - \lambda n}{\sqrt{\lambda n}} = \frac{\sum_{k=1}^n (N_k - N_{k-1}) - \lambda n}{\sqrt{\lambda n}} \quad (41.19)$$

or par définition nous avons les égalités de lois

$$N_k - N_{k-1} \sim N_1 \sim \mathcal{P}(\lambda), \quad (41.20)$$

donc

$$\frac{S_n - \lambda n}{\sqrt{\lambda n}} = \frac{\frac{1}{n} S_n - \lambda}{\frac{\sqrt{\lambda n}}{n}} = \frac{\frac{1}{n} S_n - \lambda}{\sqrt{\lambda}/\sqrt{n}}, \quad (41.21)$$

ce qui est exactement le théorème central limite pour une suite de lois de Poisson¹.

Démonstration. Nous écrivons \bar{t} la partie entière de \bar{t} et nous décomposons :

$$\frac{N_t - \lambda t}{\sqrt{\lambda t}} = \underbrace{\frac{N_t - N_{\bar{t}}}{\sqrt{\lambda t}}}_A + \underbrace{\frac{N_{\bar{t}} - \lambda \bar{t}}{\sqrt{\lambda t}}}_B + \underbrace{\frac{\lambda \bar{t} - \lambda t}{\sqrt{\lambda t}}}_C. \quad (41.22)$$

En ce qui concerne le terme B , nous avons

$$B = \sqrt{\frac{\bar{t}}{t}} \frac{N_{\bar{t}} - \lambda \bar{t}}{\sqrt{\lambda \bar{t}}} \rightarrow \mathcal{N}(0, 1). \quad (41.23)$$

1. Au fait près que nous devrions encore montrer que S_n est de carré intégrable.

Notons que nous utilisons le fait que si $a_n \rightarrow 1$ (en tant que suite de nombre) et si $X_n \rightarrow \mathcal{N}(0, 1)$ (limite en loi), alors $a_n X_n \rightarrow \mathcal{N}(0, 1)$ en loi.

Le terme C est également facile parce que $\lambda\bar{t} - \lambda t$ est majoré en norme par λ . Du coup

$$-\frac{\lambda}{\sqrt{\lambda t}} \leq C \leq \frac{\lambda}{\sqrt{\lambda t}}. \quad (41.24)$$

Donc $\lim_{t \rightarrow \infty} C = 0$.

Reste à travailler sur A . Vu que $t \mapsto N_t$ est croissante, la différence $N_t - N_{\bar{t}}$ est positive. Soit $\eta > 0$, nous avons

$$P(|A| > \eta) = P(N_t - N_{\bar{t}} > \sqrt{\lambda t} \eta) \leq P(N_{\bar{t}+1} - N_{\bar{t}} \geq \sqrt{\lambda t} \eta) = P(N_1 \geq \eta \sqrt{\lambda t}) \quad (41.25)$$

parce que nous savons que $N_{\bar{t}+1} - N_{\bar{t}} \sim N_1 \sim \mathcal{P}(\lambda)$. En vertu des propriétés de la loi de Poisson,

$$\lim_{t \rightarrow \infty} P(N_1 \geq \eta \sqrt{\lambda t}) = 0. \quad (41.26)$$

En effet si Z est une variable aléatoire de Poisson de paramètre λ nous avons

$$P(Z > l) = \sum_{k=l}^{\infty} P(Z = k) = e^{-\lambda} \sum_{k=l}^{\infty} \frac{\lambda^k}{k!}. \quad (41.27)$$

Nous reconnaissons la queue de série de e^λ , qui tend donc vers zéro lorsque $l \rightarrow \infty$. Nous avons donc prouvé que

$$\lim_{t \rightarrow \infty} P(|A| > \eta) = 0, \quad (41.28)$$

c'est-à-dire la convergence en probabilité de A vers zéro.

Nous avons montré que

$$B + C \xrightarrow{\mathcal{L}} U \sim \mathcal{N}(0, 1) \quad (41.29a)$$

$$A \xrightarrow{P} 0. \quad (41.29b)$$

Le lemme de Slutsky (37.74) nous avons une convergence du couple

$$(A, B + C) \xrightarrow{\mathcal{L}} (0, U). \quad (41.30)$$

Utilisant le corollaire 37.76, nous trouvons la convergence en loi

$$A + (B + C) \xrightarrow{\mathcal{L}} 0 + U, \quad (41.31)$$

ce qu'il fallait. □

41.2 Quelques trucs sur la simulation

Le théorème ergodique dit que

$$\pi(x) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \mathbb{1}_{X_k = x}. \quad (41.32)$$

C'est avec cela qu'on calcule $\pi(x)$ à partir d'une simulation de chaîne de Markov.

41.2.1 Le théorème central limite pour Markov

Théorème 41.7 (Version allégée).

Si (X_n) est irréductible et positive récurrente, alors pour toute fonction f ,

$$\frac{1}{\sqrt{N}} \left[\sum_{k=1}^N -N \int f d\pi \right] \xrightarrow{\mathcal{L}} \mathcal{N}(0, \sigma^2) \quad (41.33)$$

où σ^2 dépend de la fonction f et de la chaîne de Markov.

Ici, $\int f d\pi = \sum_{x \in E} f(x)\pi(x)$.

Nous allons simuler la variable aléatoire

$$Z = \frac{1}{\sqrt{N}} \left[\sum f(X_k) - N \sum_{x \in E} f(x)\pi(x) \right] \quad (41.34)$$

et puis on va mettre sa réalisation dans un histogramme. Dans le cas où on prend $f(i) = \mathbb{1}_{i=i_0}$, il y a de la simplification dans l'intégrale qui devient

$$Z = \frac{1}{\sqrt{N}} \left[\sum_{i=1}^N \mathbb{1}_{X_k=i_0} - N\pi(i_0) \right]. \quad (41.35)$$

41.2.2 Feuille 5

On pose

$$D_n = \sqrt{n} \sup_{x \in \mathbb{R}} |F_n(x) - F(x)|. \quad (41.36)$$

On en génère un milliers de fois D_n , on note $D_n^{(k)}$ ces réalisations, et on regarde ce que vaut

$$\frac{1}{1000} \sum_{k=1}^{1000} \mathbb{1}_{D_n^{(k)} \geq c}. \quad (41.37)$$

Cela nous donne une approximation de

$$P(\sqrt{n} \sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \geq c). \quad (41.38)$$

Note que chacun des $D_n^{(k)}$ demande de créer un nouveau vecteur Y_i de lois qu'on veut regarder. Par exemple de loi exponentielle.

41.2.3 Feuille 6

Pour créer une fonction qui renvoie i avec probabilité p_i pour $i = 1, 2, 3$, on peut faire

$$U \sim \mathcal{U}[0, 1] \quad (41.39)$$

et puis on a

$$P(U < p_0) = p_0 \quad (41.40a)$$

$$P(p_0 < U < p_0 + p_1) = p_1 \quad (41.40b)$$

$$P(p_0 + p_1 < U < p_2) = p_2. \quad (41.40c)$$

Une façon de faire une loi uniforme $[0, 1]$ est de faire `rand`

41.2.4 Feuille 7

L'échantillon est (Y_1, \dots, Y_n) et nous écrivons le vecteur

$$Y = X\beta + \epsilon \quad (41.41)$$

où $Y \sim \mathcal{N}(X\beta, \sigma^2 \text{Id})$ et $\epsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$. Nous utilisons le principe de maximum de vraisemblance. Soit (y_1, \dots, y_n) un échantillon et

$$P_\theta(y_1, \dots, y_n) = \prod_i \frac{1}{\sigma\sqrt{2\pi}} \exp \left[-\frac{1}{2} \left(\frac{y_i - X_i^t \beta}{\sigma} \right)^2 \right]. \quad (41.42)$$

L'astuce est de faire que $y_i - X_i^t \beta$ est la i ème composante du vecteur $Y - X\beta$ et donc la somme qui est dans l'exponentielle devient la norme de $Y - X\beta$:

$$f_\theta(y_1, \dots, y_n) = \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^n \exp \left[-\frac{1}{2} \|Y - X\beta\|^2 \right]. \quad (41.43)$$

On passe au logarithme et on dérive par rapport à σ^2 . Attention : la variable est σ^2 , donc la dérivée de σ^2 est 1 et non 2σ . Bref, on trouve

$$\sigma^2 = \frac{1}{2n} \|U + X\beta\|. \quad (41.44)$$

41.2.5 Simuler des lois conditionnelles

Nous voulons générer des couples (X, Y) tels que Y prend les valeurs 0 ou 1 et tels que

$$\begin{cases} P(X|Y = 0) \sim \mathcal{E}(\lambda_0) & (41.45a) \\ P(X|Y = 1) \sim \mathcal{E}(\lambda_1). & (41.45b) \end{cases}$$

Le plus simple est de générer une liste

$$(X_1, 0) \qquad (X_4, 1) \qquad (41.46a)$$

$$(X_2, 0) \qquad (X_5, 1) \qquad (41.46b)$$

$$(X_3, 0) \qquad (X_6, 1) \qquad (41.46c)$$

avec $X_1, X_2, X_3 \sim \mathcal{E}(\lambda_0)$ et $X_4, X_5, X_6 \sim \mathcal{E}(\lambda_1)$.

Avec cette méthode cependant la liste est triée et en plus on a autant de 1 que de 0. On peut faire un peu plus technologique pour corriger cela. Pour créer un couple, on commence par $Y \sim \mathcal{B}(p)$ et puis suivant que $Y = 0$ ou $y = 1$, on génère $X \sim \mathcal{E}(\lambda_0)$ ou $X \sim \mathcal{E}(\lambda_1)$.

Chapitre 42

Langages

42.1 Langages

42.1.1 Alphabets et mots

Définition 42.1.

Un **alphabet** est un ensemble fini de symboles appelés **lettres**.

On utilise aussi parfois le terme **vocabulaire** pour désigner un alphabet.

Définition 42.2.

Un **mot** sur l'alphabet Σ est une suite finie et ordonnée, éventuellement vide, de lettres de Σ . Le **mot vide** est toujours noté ε .

Définition 42.3.

La **longueur d'un mot** w , noté $|w|$, est le nombre de lettres constituant le mot w . Le mot vide a une longueur de 0.

Soit w un mot de longueur k , on peut désormais noter $w = w_1 \cdots w_k$, où chacun des w_i , $1 \leq i \leq k$ représente une lettre de w . Par convention, si $k = 0$, alors le mot w est le mot vide.

Définition 42.4.

Soient w un mot sur l'alphabet Σ et $a \in \Sigma$ une lettre, le **nombre d'occurrences** de la lettre a dans le mot w , noté $|w|_a$, est le nombre de fois où apparaît la lettre a dans le mot w , c'est-à-dire le cardinal de l'ensemble $\{i \mid w_i = a, 1 \leq i \leq |w|\}$.

Définition 42.5.

Soit Σ un alphabet, l'**ensemble des mots non-vides** sur l'alphabet Σ , noté Σ^+ , est l'ensemble :

$$\Sigma^+ = \{w = w_1 \dots w_n, n > 0\} \quad (42.1)$$

Définition 42.6.

Soit Σ un alphabet, l'**ensemble des mots** sur l'alphabet Σ , noté Σ^* , est l'ensemble :

$$\Sigma^* = \{w = w_1 \dots w_n, n \geq 0\} \quad (42.2)$$

Des deux définitions précédentes, on tire l'égalité suivante :

$$\Sigma^* = \Sigma^+ \cup \{\varepsilon\} \quad (42.3)$$

Définition 42.7.

Soient Σ un alphabet et $x, y \in \Sigma^*$ deux mots sur l'alphabet Σ de longueur respective n et m , le **produit** w de x et y , noté $x \cdot y$ est défini par $w = x_1 \dots x_n y_1 \dots y_m$.

Le produit est également appelé **concaténation**.

Proposition 42.8 (Longueur du produit de deux mots).

La longueur du produit de deux mots x et y est la somme des longueurs des mots x et y .

$$|x \cdot y| = |x| + |y| \quad (42.4)$$

Proposition 42.9 (Monoïde $(\Sigma^*, \cdot, \varepsilon)$).

L'ensemble Σ^* munie de l'opération produit d'élément neutre ε est un monoïde¹.

Démonstration. Soient $x, y, z \in \Sigma^*$, avec les définitions précédentes, on peut vérifier facilement que :

- le produit est une loi interne : $x \cdot y \in \Sigma^*$;
- le produit est associatif : $x \cdot (y \cdot z) = (x \cdot y) \cdot z$;
- ε est l'élément neutre du produit : $x \cdot \varepsilon = \varepsilon \cdot x = x$.

□

Le produit n'est pas commutatif.

Définition 42.10.

Soient Σ un alphabet et $w \in \Sigma^*$, la **puissance** n^{e} d'un mot w , notée w^n , est définie par :

$$w^n = \begin{cases} \varepsilon & \text{si } n = 0 \\ w \cdot w^{n-1} & \text{si } n > 0 \end{cases}$$

42.1.2 Langage

Définition 42.11.

Un **langage** sur un alphabet Σ est un sous-ensemble de Σ^* . C'est un ensemble de mots sur l'alphabet Σ .

Un langage étant défini comme un ensemble, on peut appliquer toutes les notions de la théorie des ensembles aux langages.

Définition 42.12.

Le **langage vide**, noté \emptyset est le langage qui ne contient aucun mot.

Définition 42.13.

Le **langage unité** est le langage qui contient uniquement le mot vide : $\{\varepsilon\}$.

Définition 42.14.

Soient Σ un alphabet et $L_1, L_2 \subseteq \Sigma^*$ deux langages sur l'alphabet Σ , on définit le **produit** L de L_1 et L_2 , noté $L_1.L_2$ par :

$$L = L_1.L_2 = \{u_1 \cdot u_2, u_1 \in L_1, u_2 \in L_2\} \quad (42.5)$$

Le produit de langages est également appelé **concaténation**. Il ne faut pas confondre le produit de langage avec le produit cartésien de deux ensembles. Le langage unité est l'élément neutre du produit de langages.

Proposition 42.15 (Distributivité du produit de langage par rapport à l'union).

Le produit de langage est distributif par rapport à l'union. Soient Σ un alphabet et $L_1, L_2, L_3 \subseteq \Sigma^*$, alors :

$$L_1.(L_2 \cup L_3) = (L_1.L_2) \cup (L_1.L_3) \text{ et } (L_1 \cup L_2).L_3 = (L_1.L_3) \cup (L_2.L_3) \quad (42.6)$$

Démonstration. Soit $w \in L_1.(L_2 \cup L_3)$, montrons que $w \in (L_1.L_2) \cup (L_1.L_3)$. $\exists w_1 \in L_1, w' \in L_2 \cup L_3, w = w_1 \cdot w'$. Donc $w' \in L_2$ ou $w' \in L_3$. Si $w' \in L_2$, alors $w = w_1 \cdot w' \in L_1.L_2$. Si $w' \in L_3$, alors $w = w_1 \cdot w' \in L_1.L_3$. Donc, $w \in (L_1.L_2) \cup (L_1.L_3)$. Donc $L_1.(L_2 \cup L_3) \subseteq (L_1.L_2) \cup (L_1.L_3)$.

1. Définition 6.19.

Soit $w \in (L_1.L_2) \cup (L_1.L_3)$, montrons que $w \in L_1.(L_2 \cup L_3)$. $w \in L_1.L_2$ ou $w \in L_1.L_3$. Si $w \in L_1.L_2$, i avec $i \in \{2, 3\}$ alors $\exists w_1 \in L_1, w_i \in L_i, w = w_1 \cdot w_i$. Donc $w \in L_1.(L_2 \cup L_3)$. Donc, $(L_1.L_2) \cup (L_1.L_3) \subseteq L_1.(L_2 \cup L_3)$

Donc $(L_1.L_2) \cup (L_1.L_3) = L_1.(L_2 \cup L_3)$

L'autre partie de la proposition se montre de manière analogue. \square

Définition 42.16.

Soient Σ un alphabet et $L \subseteq \Sigma^*$, la **puissance** n^e du langage L , notée L^n est définie par :

$$L^n = \begin{cases} \{\varepsilon\} & \text{si } n = 0 \\ L.L^{n-1} & \text{si } n > 0 \end{cases}$$

Définition 42.17 ([504]).

L'**étoile de Kleene** est un opérateur unaire noté $*$. L'**itéré** d'un langage L , noté L^* , est l'application de l'étoile de Kleene à un langage L et est défini par :

$$L^* = \bigcup_{i \geq 0} L^i \quad (42.7)$$

En particulier, on remarque que le mot vide fait toujours partie de l'itéré d'un langage, y compris quand ce même langage ne contient pas le mot vide.

Définition 42.18.

L'**itéré strict** d'un langage L , noté L^+ , est défini par :

$$L^+ = \bigcup_{i > 0} L^i \quad (42.8)$$

Proposition 42.19 (Relations entre itéré et itéré strict).

Soit L un langage, alors on a :

$$L^* = L^+ \cup \{\varepsilon\} \quad (42.9)$$

$$L^+ = L.L^+ = L^+.L \quad (42.10)$$

Chapitre 43

Utilisation dans les autres sciences

Dans ce chapitre nous donnons des applications de divers théorèmes dans les autres sciences que la mathématique.

43.1 Démystification du MRUA

43.1.1 Preuve de la formule

Nous sommes maintenant en mesure de donner une démonstration complète de la formule du MRUA :

$$x(t) = \frac{at^2}{2} + v_0t + x_0. \quad (43.1)$$

Au niveau de la physique, nous considérons un mobile qui se déplace avec une accélération constante a . Nous notons par v_0 sa vitesse initiale et par x_0 sa position initiale.

Nous savons que, pour tout mouvement, si $x(t)$ est la position en fonction du temps, et si $v(t)$ et $a(t)$ représentent la vitesse et l'accélération en fonction du temps, alors

$$v(t) = x'(t) \quad \text{et} \quad a(t) = v'(t) = x''(t). \quad (43.2)$$

Afin de trouver $x(t)$ en connaissant $a(t)$, il « suffit » donc de prendre deux fois la primitive. Essayons ça dans le cas facile du MRUA où $a(t) = a$ est constante.

La vitesse $v(t)$ doit être une primitive de la constante a . Il est facile de voir que $v(t) = at$ est une primitive de a . Par le corollaire 13.136(bis),

$$v(t) = at + C_1 \quad (43.3)$$

pour une certaine constante C_1 . Afin de fixer C_1 , il faut faire appel à la physique : d'après la formule (43.3), la vitesse initiale est $v(0) = C_1$. Donc il faut identifier C_1 à la vitesse initiale : $C_1 = v_0$. Nous avons donc déjà obtenu que

$$v(t) = at + v_0. \quad (43.4)$$

Afin de trouver $x(t)$, il faut trouver une primitive de $v(t)$. Il n'est pas très difficile de voir que $at^2/2 + v_0t$ fonctionne, donc il existe une constante C_2 telle que

$$x(t) = \frac{at^2}{2} + v_0t + C_2. \quad (43.5)$$

Encore une fois, regardons la condition initiale : la formule donne comme position initiale $x(0) = C_2$, et donc nous devons identifier C_2 avec la position initiale x_0 . En définitive, nous avons bien

$$x(t) = \frac{at^2}{2} + v_0t + x_0. \quad (43.6)$$

Cette formule est donc maintenant *démontrée* à partir de la seule définition de la vitesse comme dérivée de la position et de l'accélération comme dérivée de la vitesse.

Remarquons cependant que la preuve complète fut *très* longue. En effet, nous avons utilisé les règles de dérivation de la proposition 13.114, pour la démonstration desquels, les résultats 13.19 et 13.18 ont été utiles. Mais nous avons surtout utilisé le corollaire 13.136(bis) qui repose sur le théorème de Rolle 13.126, qui lui-même demande le théorème de Borel-Lebesgue 8.6 dans lequel la notion d'ensemble compact a été cruciale.

43.1.2 Interprétation graphique

La distance parcourue $x(t)$ en un temps t est la primitive de la vitesse. Nous avons, par ailleurs, que l'opération inverse de la dérivée donnait la surface. Pour reprendre les mêmes notations, nous notons $S_v(t)$ la surface contenue en dessous de la fonction v entre 0 et x . Nous ne serions donc pas étonné que

$$S_v(t) = \frac{at^2}{2} + v_0t + x_0 \quad (43.7)$$

soit la surface en dessous de la fonction $v(t) = at + v_0$. Nous voyons que la surface totale sous la fonction $v(t) = at + v_0$ est exactement

$$S_v(t) = \frac{at^2}{2} + v_0t. \quad (43.8)$$

Cela est un bon début, mais hélas nous ne retrouvons pas le terme « $+x_0$ » de la formule (43.7). Cela n'est pas tout à fait étonnant parce que nous savons que la surface sous une fonction était *une* primitive de la fonction, mais nous n'avons pas dit *laquelle*. D'après le fameux corollaire 13.136(bis), la primitive n'est définie qu'à une constante près. Ici, c'est la constante x_0 qu'on a perdue en chemin.

Nous parlerons plus en détail du lien entre les surfaces et les primitives dans la section dédiée à l'intégration.

43.2 Relativité en mécanique newtonienne

43.2.1 Relativité du mouvement

Prenons quelqu'un qui cours le cent mètres en onze secondes. Par rapport à un spectateur dans les gradins, il se sera déplacé de cent mètres. Mais si je cours à côté de lui de telle façon à avoir parcouru 80 mètres le temps qu'il en fasse cent, alors par rapport à moi l'athlète ne se sera déplacé que de 20 mètres. Par contre, par rapport à mon chronomètre, il aura également mit onze secondes : ce n'est pas parce que je cours que mon chronomètre s'affole !

Entre moi et les spectateurs, on a donc une loi de transformation

$$x' = x - vt \quad t' = t. \quad (43.9)$$

C'est-à-dire que la distance x' qu'aura parcouru l'athlète par rapport à moi vaut la distance x parcourue par le spectateur moins la vitesse que j'ai courue moi-même, c'est-à-dire moins vt .

43.2.2 Bob et Alice

Formalisons le concept de changement de repères. Pour cela, prenons deux amoureux, Bob et Alice¹. Mettons que Bob reste assis sur un banc pendant qu'Alice cours en ligne droite à une vitesse v . Tout deux déclenchent leur chronomètre quand Alice passe devant Bob. À tout moment, Bob et Alice ont leur repères de temps et d'espace. Par exemple si après un temps t , Alice voir une peau de banane à 1 mètre devant elle, elle va dire « Il y a une peau de banane à un mètre. », tandis que Bob va dire « Il y a une peau de banane à $(1 + vt)$ mètres ».

Plus généralement, s'il se passe quelque chose à la position x au temps t pour Bob, ce quelque chose se passera au temps $t' = t$ à l'endroit $x' = x - vt$ pour Alice parce qu'en un temps t , elle aura déjà avancé d'une distance vt .

Ça c'est ce dont tout le monde était persuadé depuis Galilée jusqu'au début du vingtième siècle.

1. C'est plus poétique que dire « soient A et B deux observateurs ».

43.3 Invariance de la vitesse de la lumière

43.3.1 Champ de gravitation et électrique

Nous savons que que la force de gravitation s'écrit :

$$F_{grav} = G \frac{mm'}{r^2},$$

tandis que la force électrique entre deux charges q et q' est donnée par

$$F_{elec} = k \frac{qq'}{r^2}. \quad (43.10)$$

Nous avons aussi fait remarquer que dans le cas de la gravitation, la force a l'air d'être instantanée, et que cela posait quelques problèmes conceptuels. La force électrique a apparemment le même problème. Une différence entre les deux est qu'une charge électrique c'est tout petit et qu'on peut expérimenter à souhait, tandis que pour avoir une masse dont on peut mesurer le champ de gravitation correctement, il faut quelque chose grand comme la Terre².

43.3.1.1 Finitude de la vitesse de propagation de la force électrique

Si un micro est placé juste à côté de ton oreille, et qu'il commence à faire biiiiip, tu l'entends directement. Quand il s'arrête, tu ne l'entends plus. Si le micro est placé à 600 m de toi, tu ne commenceras à l'entendre que deux secondes après le commencement du son, et tu continueras à l'entendre deux secondes après qu'il ait fini.

Eh bien, pour la force électrique, on a pu mesurer que c'est la même chose (sauf que ça va beaucoup plus vite). Si on place une charge quelque part, on ne ressent la force (43.10) qu'après qu'elle ait eut le temps d'arriver. Si on déplace la charge électrique, on continue à ressentir la même force pendant un certain temps : il faut que la modification du champ électrique ait le temps d'arriver. Exactement comme quand on fait des remous quelque part dans un étang : il faut du temps que les remous arrivent plus loin.

On a pu faire des dizaines d'expériences de ce type avec l'électricité, le magnétisme et la lumière ; et les résultats sont clairs : il faut du temps pour que ça se déplace. Tout cela provoque des ondes électromagnétiques qui se déplacent à une vitesse finie. On peut produire de telles ondes avec n'importe quel courant électrique alternatif.

43.3.1.2 Pourquoi pas la gravitation ?

La gravitation telle que donnée par Newton pose le même problème de vitesse de propagation que l'électricité. Est-ce qu'en réalité la gravitation se propage également à une vitesse finie ?

Avec la gravitation c'est beaucoup plus compliqué parce qu'elle est beaucoup plus faible, et donc c'est beaucoup plus difficile à détecter. D'après la théorie d'Einstein de la gravitation, la gravitation devrait également produire des ondes gravitationnelles. Seulement, si un simple courant électrique suffit pour mesurer une onde électromagnétique, afin de mesurer une onde gravitationnelle, il faudrait un déplacement de masse de l'ampleur d'une étoile qui explose. Or ça, on ne sait pas produire dans un laboratoire. Les physiciens sont donc pour l'instant (2009) en train d'attendre qu'une étoile explose pas trop loin d'ici afin d'être capable de mesurer une onde gravitationnelle.

L'existence de ces ondes de gravitation ne fait aucun doute dans la tête d'aucun physicien parce qu'elles sont une conséquence logique (et mathématique) de la théorie de la relativité générale, laquelle a déjà eut beaucoup de confirmations expérimentales. Mais comme on est dans le cadre d'une science expérimentale, il faut être patient et attendre d'en avoir effectivement observée une avant de dire avec certitude que ça existe.

2. Une autre différence fondamentale est qu'il existe des charges électriques négatives, mais pas de masses négatives ; de ce fait on ne peut pas construire d'isolant gravitationnel, contrairement aux isolants électriques qui existent. Cela augmente encore la difficulté de faire des expériences avec la gravitation.

43.3.2 Support du champ : pas d'éther

Nous avons dit qu'une onde électromagnétique se propage comme une onde sonore (quoique beaucoup plus vite). Une question se pose alors. En effet, une onde sonore est matérialisée par de l'air qui vibre. Qu'est-ce qui vibre pour une onde électromagnétique ?

Étant donné que les ondes électromagnétiques se propagent dans le vide (c'est pour ça que la radio fonctionne dans l'espace), la question est problématique. Les physiciens ont donc supposé que tout l'univers était rempli d'un fluide invisible appelé **l'éther**. L'électromagnétisme consiste en une vibration de l'éther, exactement comme l'acoustique consiste en une vibration de l'air.

En fait, vérifier cette hypothèse n'est pas très compliqué. En effet il n'y a aucune raison que l'éther suive la Terre dans son mouvement. Or la Terre se déplace à environ 30 km/s autour du Soleil. Donc les ondes électromagnétiques doivent se propager plus vite dans le sens du mouvement de la Terre que dans le sens perpendiculaire. Tout comme le son se propage plus vite dans le sens du vent.

La célèbre expérience de **Michelson-Morley** a mesuré cet effet ... et ce fut la consternation : il n'y a aucun effet ! Or, la lumière se déplace à 300.000 km/s ; une variation de 30 km/s devrait être détectable !

Mais rien ! On a recommencé les expériences dans tout les sens, à tous les mois de l'année, à tous les endroits de la Terre. On n'a pas observé un poil de variation de la vitesse de la lumière. Et ça, ça pose un gros problème à la physique.

43.3.3 Le problème

Si je joue au football dans un train qui avance à 100 km/h et que je lance une balle à 20 km/h, quelqu'un au sol mesura la vitesse de la balle soit à 120 km/h soit à 80 km/h d'après que l'on ait shooté vers l'avant ou l'arrière du train. Cela paraît logique. Mais ce qu'on vient de voir c'est que ça ne marche pas avec la lumière.

Si un train avance à 100.000 km/s et qu'on y allume une lampe de poche, la lumière avancera à 300.000 km/s par rapport au train et 400.000 km/s par rapport au sol. Non ! Justement pas ! La lumière avancera quand même à 300.000 km/s par rapport au sol.

Là encore, on a fait des dizaines d'expériences partout, sur Terre, dans des avions, dans l'espace avec des atomes, des lampes de poche et des horloges atomique, et dans tous les sens, le sens de déplacement de la Terre, le sens inverse, le sens perpendiculaire, vers le haut, vers le bas : rien ! Personne n'a jamais observé un rayon de lumière se déplacer à une autre vitesse que 300.000 km/s.

Le problème est que le principe d'addition des vitesses est faux pour la lumière. Puisque l'expérience nous force, nous devons faire avec.

Loi numéro 1.

La réalité est que la vitesse de la lumière est la même dans tous les référentiels. On note c cette vitesse. C'est une constante fondamentale de la Nature.

Étant donné que c'est une loi expérimentale, nous n'en pouvons rien. C'est la nature qui est comme ça. En particulier tu ne peux pas en vouloir à ton prof de physique d'avoir inventé une théorie compliquée. Ce n'est pas lui qui l'a inventée et ce n'est pas de sa faute.

43.4 Conséquences

C'est maintenant que les choses vraiment graves commencent (cela soit dit sans vouloir te faire peur). Afin d'un peu simplifier les choses, nous n'allons étudier que les mouvements en une dimension, c'est-à-dire sur une droite.

43.4.1 Ligne d'univers

Un événement a une coordonnée (t, x) . Si je pose un objet juste à mes pieds (disons en $x = 0$), ses coordonnées seront à tout moment $(t, 0)$. Il est bon de voir cette coordonnée comme l'équation

paramétrique d'une droite horizontale dans le plan des coordonnées t et x . Plus généralement quand un mobile effectue un mouvement $x(t)$, on appelle la **ligne d'univers** du mobile la ligne (pas forcément droite) $(t, x(t))$. Dans le premier exemple, on avait $x(t) = 0$ pour tout t .

Le cas d'un mobile se déplaçant à vitesse constante v donne comme ligne d'univers la droite³ $(t, x_0 + vt)$, et un objet qui se déplace selon un MRUA a comme ligne d'univers

$$\left(t, x_0 + v_0t + \frac{at^2}{2}\right).$$

43.4.2 Transformations de Lorentz

Reprenons les amours scientifiques de Bob et Alice, mais cette-fois, analysons celles-ci en tenant compte du fait que la vitesse de la lumière soit invariante. Maintenant, si Bob voit se passer quelque chose au temps t à l'endroit x , on va dire qu'Alice voit cette chose au temps t' à la position x' , et on va chercher (t', x') en fonction de (t, x) .

Posé en termes mathématiques, le problème s'énonce ainsi : trouver les fonctions f et g telles que les formules

$$t' = f(t, x) \tag{43.11a}$$

$$x' = g(t, x) \tag{43.11b}$$

donnent les coordonnées vues par Alice pour un événement vu par Bob à l'instant t au point x . Une première étape importante est franchie par la proposition suivante⁴.

Proposition 43.1.

Les fonctions f et g contenues dans les transformations (43.11) sont nécessairement linéaires (affines), c'est-à-dire qu'elles doivent s'écrire sous la forme

$$t' = \alpha t + \beta x + p$$

$$x' = \gamma t + \delta x + q$$

pour certaines fonctions $\alpha, \beta, \gamma, \delta, p$ et q de la vitesse d'Alice relativement à Bob.

Démonstration. Pendant qu'Alice court et que Bob la regarde, Ève tente de lancer une pierre sur Alice (Ève est jalouse). Bob et Alice regardent deux événements. Le premier est la pierre qui quitte la main de Ève, et le second est la pierre qui percute le sol. Pour Bob, le jet s'est passée au temps t_0 au point x_0 , et la pierre touche le sol un petit peu plus tard, au temps $t_0 + \Delta t$ et un peu plus loin, au point $x_0 + \Delta x$. Bob écrit donc ceci sur sa feuille de papier :

$$E_1 = (t_0, x_0)$$

$$E_2 = (t_0 + \Delta t, x_0 + \Delta x),$$

tandis qu'Alice, en observant les mêmes deux événements, aura noté

$$E'_1 = (f(t_0, x_0), g(t_0, x_0))$$

$$E'_2 = (f(t_0 + \Delta t, x_0 + \Delta x), g(t_0 + \Delta t, x_0 + \Delta x)).$$

Bob et Alice se demandent combien de temps la pierre est restée en l'air et quelle distance elle a parcourue. Par le principe général d'homogénéité, les deux seules quantités pertinentes (qui ont un sens physique) pour Bob sont $(t_0 + \Delta t) - t_0$ et $(x_0 + \Delta x) - x_0$, c'est-à-dire Δt et Δx . En effet, si Bob avait choisi de s'asseoir autre part et si Alice avait commencé à courir un peu plus tard, ça n'aurait rien changé à la longueur du jet de Ève.

D'une façon ou d'une autre, il doit exister une façon de déduire les mesures de Alice en connaissant celles de Bob ; je ne connais pas avec quelles formules, mais ces formules ne peuvent contenir que Δt , Δx et v parce que ce sont les seules quantités qui définissent tous les événements.

3. bon exercice de révision de ton cours de math de vérifier que c'est une droite.

4. dont je te suggère fortement de ne pas lire la preuve si tu ne veux pas que ton cerveau éclate.

Cela dit, Alice va caractériser le mouvement de la pierre avec la différence des coordonnées entre le jet et la chute sur le sol mesurées par elle-même. En d'autres termes, pour Alice ce qui compte c'est la différence entre E'_1 et E'_2 , soit

$$(f(t_0 + \Delta t, x_0 + \Delta x), g(t_0 + \Delta t, x_0 + \Delta x)) - (f(t_0, x_0), g(t_0, x_0)). \quad (43.12)$$

Mais nous venons de signaler que ce qu'Alice mesurait devait pouvoir être exprimé en termes de Δt et Δx . Nous concluons que la différence (43.12) ne dépend en fait pas de x et t mais seulement de Δt et Δx .

Prenons maintenant une notation plus compacte et notons $X = (t, x)$, $\Delta X = (\Delta t, \Delta x)$ puis $F = (f, g)$. Avec ça, l'expression (43.12) se note $F(X + \Delta X) - F(X)$. Comme mentionné, cette expression ne dépend que de Δx . En particulier, elle ne dépend pas de X .

Maintenant tu vas comprendre pourquoi on apprend les dérivées dans ton cours de math. Comme $F(X + \Delta X) - F(X)$ ne dépend pas de X , le rapport $(F(X + \Delta X) - F(X))/\Delta X$ non plus. La limite de ce rapport quand ΔX tend vers zéro non plus :

$$\lim_{\Delta X \rightarrow 0} \frac{F(X + \Delta X) - F(X)}{\Delta X} \quad (43.13)$$

ne dépend pas de X . Tu reconnais là la dérivée de F au point X . En d'autres termes, $F'(X)$ est constante, elle ne dépend pas de X . Disons donc que $F'(X) = a$. Tu connais beaucoup de fonctions dont la dérivée est constante? Non? En effet, il n'y en a pas beaucoup. Les fonctions qui vérifient $F'(X) = a$ signifie sont toutes de la forme

$$F(X) = aX + b.$$

À ce niveau, il convient de re-déballer les notations compactes : si $a = \begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix}$ et $b = (p, q)$ on trouve

$$f(t, x) = \alpha t + \beta x + p \quad (43.14a)$$

$$g(t, x) = \gamma t + \delta x + q, \quad (43.14b)$$

comme annoncé. □

Nous savons que lorsque $(t, x) = (0, 0)$, alors $(t', x') = (0, 0)$. En effet, Bob et Alice ont lancés leurs chronos en même temps au moment où ils étaient au même endroit. En mettant $(t, x) = (0, 0)$ dans les équations (43.14), on trouve $(t', x') = (p, q)$, et donc $p = q = 0$. Ça fait une chose de réglée ; on se retrouve avec

$$\begin{cases} t' = \alpha t + \beta x \\ x' = \gamma t + \delta x. \end{cases} \quad (43.15a)$$

$$\quad (43.15b)$$

Quelles sont les contraintes à vérifier pour que ces transformations décrivent correctement la physique que l'on cherche à écrire ?

- (1) Il faut que les transformations décrivent correctement que Alice avance à une vitesse v par rapport à Bob,
- (2) dans le même ordre d'idée, il faut que l'on trouve que Bob avance à la vitesse $-v$ par rapport à Alice,
- (3) il faut que si Alice et Bob observent un rayon lumineux, ce rayon aille à la vitesse c par rapport à Alice et à la même vitesse c par rapport à Bob,
- (4) enfin, il faut avoir le principe de relativité, c'est-à-dire que comme les équations (43.15) disent ce que Alice voit en fonction de ce que Bob voit, on demande que les équations qui disent ce que Bob voit en fonction de ce que Alice voit soient les mêmes. En d'autres termes, il faut que les transformations et les transformations inverses soient les mêmes au changement près du signe de v .

Étudions une à une ce que chacune de ses contraintes impose. Rappelons que (t, x) et (t', x') sont les coordonnées que Bob et Alice mettent sur le même événement. Par exemple sur l'événement qui consiste à ce que Ève, par jalousie envers Bob, jette une peau de banane sous les pieds d'Alice. Cet événement a lieu à un certain moment, à un certain endroit. C'est ce moment et cet endroit qui sont notés (t, x) et (t', x') .

- (1) Les coordonnées (t, x) et (t', x') peuvent décrire n'importe quoi. Regardons les coordonnées de Alice qui cours. Pour Alice, cela correspond à $(t', x') = (t', 0)$ parce que si x' désigne la position de Alice par rapport à Alice, alors x' est toujours nul. Pour Bob par contre, Alice ne reste pas en place, mais se déplace à une vitesse v . C'est-à-dire que si (t, x) sont les coordonnées de Alice pour Bob, alors $x/t = v$. Écrivons les équations (43.15) en tenant compte de tout ça : avec $x' = 0$, la seconde équation donne

$$0 = \gamma t + \delta x, \quad (43.16)$$

d'où on déduit que $x/t = -\gamma/\delta$. En imposant que cela soit v , on trouve $\gamma = -v\delta$, et on ré-écrit les transformations en tenant compte de ça :

$$\begin{cases} t' = \alpha t + \beta x \\ x' = -v\delta t + \delta x. \end{cases} \quad (43.17a)$$

$$(43.17b)$$

Nous voilà débarrassé d'un paramètre.

- (2) Maintenant, on regarde ce qu'il se passe quand (t, x) et (t', x') décrivent les positions de Bob. On a $(t, x) = (t, 0)$ parce que selon Bob, Bob est au repos. Les équations deviennent :

$$t' = \alpha t \qquad x' = -v\delta t. \quad (43.18)$$

La vitesse de Bob par rapport à Alice est $-v$, donc on exige que $x'/t' = -v$, c'est-à-dire que

$$\frac{-v\delta t}{\alpha t} = -v,$$

ce qui implique que $\delta = \alpha$. On avance encore un peu. Écrivons à nouveau les lois de transformation en en tenant compte :

$$\begin{cases} t' = \alpha t + \beta x \\ x' = -v\alpha t + \alpha x. \end{cases} \quad (43.19a)$$

$$(43.19b)$$

- (3) Si maintenant Bob et Alice regardent un même rayon de lumière (comme c'est romanesque!), alors (t, x) et (t', x') expriment les coordonnées d'un rayon lumineux expriment les coordonnées d'un rayon lumineux. Le fait que Bob regarde un rayon lumineux fait que $x = ct$, et donc que les coordonnées du rayon lumineux, observé par Alice sont :

$$t' = \alpha t + \beta ct \qquad x' = -\alpha vt + \alpha ct. \quad (43.20)$$

L'invariance de la vitesse de la lumière exige que Alice mesure une vitesse c pour le rayon de lumière, c'est-à-dire $x' = ct'$. On exige donc que

$$-\alpha vt + \alpha ct = \alpha ct + \beta c^2 t,$$

ce qui implique que

$$\beta = -\frac{\alpha v}{c^2}.$$

Une fois de plus, l'avant-dernière, on ré-écrit les lois de transformations en tenant compte de ce fait ; mais cette fois, on fait l'effort d'écrire aussi les transformations inverses :

$$t' = \alpha t - \frac{\alpha v}{c^2} x \qquad t = \frac{1}{\Delta} \left(\alpha t' + \frac{\alpha v}{c^2} x' \right) \quad (43.21)$$

$$x' = -\alpha vt + \alpha x \qquad x = \frac{1}{\Delta} (\alpha vt' + \alpha x') \quad (43.22)$$

où $\Delta = \alpha^2 - \frac{\alpha^2 v^2}{c^2}$ que tu noteras au passage être toujours positif, et nul uniquement quand $v = c$.

- (4) Maintenant il reste à imposer le principe de relativité. Les transformations (43.21) montrent comment Alice voit le monde (c'est-à-dire (t', x')) en fonction de la façon dont Bob voit le monde (c'est-à-dire (t, x)). On se demande donc quelle seraient, pour Bob, les coordonnées (t, x) d'un point vu en (t', x') par Alice. Cela signifie que l'on impose que les deux systèmes (43.21) soient en réalité les mêmes, à un changement de signe près.

Attention : il *a priori* faux de dire qu'en changeant le signe de v dans $\alpha v/c^2$, j'obtiens $-\alpha v/c^2$ parce que α est une fonction de v . En réalité, il faut noter $\alpha(v)v/c^2$ et donc le changement de signe de v donne $-\alpha(-v)v/c^2$. Ceci étant clair, on peut un petit peu calculer.

Commençons par égaliser le coefficient de x dans t' à celui de x' dans t , en changeant le signe de v :

$$\frac{\alpha(-v)v}{c^2} = \frac{\alpha(v)v}{c^2},$$

et donc $\alpha(v) = \alpha(-v)$. Ça c'est une bonne nouvelle. Égalisons maintenant le coefficient de t dans t' à celui de t' dans t en changeant le signe de v :

$$\alpha(-v) = \frac{\alpha(v)}{\Delta(v)} = \frac{\alpha(v)}{\alpha(v)^2(1 - \frac{v^2}{c^2})}.$$

Comme $\alpha(-v) = \alpha(v)$, on en déduit que

$$\alpha(v) = \frac{1}{\sqrt{1 - \frac{v^2}{c^2}}}. \quad (43.23)$$

Maintenant qu'on a tout, on peut écrire les transformations de Lorentz. On met donc l'expression (43.23) dans les lois de transformations (43.21) :

$$\begin{aligned} t' &= \frac{1}{\sqrt{1 - \frac{v^2}{c^2}}} \left(t - \frac{v}{c^2} x \right) & t &= \frac{1}{\sqrt{1 - \frac{v^2}{c^2}}} \left(t' + \frac{v}{c^2} x' \right) \\ x' &= \frac{1}{\sqrt{1 - \frac{v^2}{c^2}}} (x - vt) & x &= \frac{1}{\sqrt{1 - \frac{v^2}{c^2}}} (vt' + x'). \end{aligned} \quad (43.24)$$

Tu remarqueras que $\Delta = 1$; si tu ne sais pas ce qu'est le déterminant d'une application linéaire, ça n'a pas d'importance. Mais si tu sais ce qu'est le déterminant d'une application linéaire, alors ce $\Delta = 1$ est crucial !

Afin d'avoir des équations un peu plus courtes, à partir de maintenant nous allons noter

$$\gamma(v) = \sqrt{1 - \frac{v^2}{c^2}}.$$

43.4.3 Conditions d'existence

Comme tu vois une racine carrée et un dénominateur dans ces formules, tu dois te demander quelles sont les conditions d'existence. Étant donné que $v < c$, on a $v^2/c^2 \leq 1$ et en particulier, $v^2/c^2 = 1$ si et seulement si $v = c$.

Ce qui se trouve dans la racine carrée ne pose donc jamais de problèmes parce que ce n'est jamais négatif.

Le dénominateur est par contre plus problématique : quand $v = c$ il n'y a plus rien qui fonctionne. Quelle est la physique de ce problème ? Pour le comprendre, il faut se souvenir ce que représente v . Nous avons dit que v est la vitesse à laquelle Alice court. Ce que la condition d'existence nous enseigne, c'est que personne ne peut courir à la vitesse de la lumière. C'est une vitesse que l'on ne peut pas atteindre.

Dit en termes plus savants, on ne peut pas choisir un repère qui se déplace à la vitesse de la lumière.

La question qui se pose alors est « ah bon, on ne peut pas atteindre la vitesse de la lumière ! Et la lumière, comment elle fait ? ». Bonne question, merci de l'avoir posée. Hélas la réponse sort du cadre de ce cours.

Loi numéro 2.

Aucun objet ne peut atteindre la vitesse de la lumière.

Loi numéro 3.

Tu ne dois pas te demander pourquoi la lumière elle-même se déplace à la vitesse de la lumière malgré la loi numéro 2.

43.4.4 La notion d'intervalle

Un **événement** est quelque chose qui se passe à un endroit à un certain moment. C'est donc caractérisé par le moment et le lieu. Comme on travaille à une dimension, c'est un couple de réels (t, x) .

Regardons un rayon de lumière. Un événement est le fait d'allumer une lampe de poche, et un autre est le fait que la lumière arrive sur l'objet qu'on éclaire. Appelons-les (t_1, x_1) et (t_2, x_2) . Comme d'habitude, on note $\Delta t = t_2 - t_1$ et $\Delta x = x_2 - x_1$. Comme le rayon de lumière va à la vitesse c , on a $c = \Delta x / \Delta t$, ou encore

$$c^2 \Delta t^2 - \Delta x^2 = 0.$$

Pour cette raison, on va dire que l'**intervalle** entre deux événements (t_1, x_1) et (t_2, x_2) vaut en général

$$s^2 = c^2(t_2 - t_1)^2 - (x_2 - x_1)^2. \quad (43.25)$$

Par invariance de la vitesse de la lumière, si un intervalle est nul pour un observateur, il sera nul pour tous les observateurs.

43.4.4.1 En mécanique newtonienne

Afin de voir un peu mieux l'enjeu de l'invariance de l'intervalle, regardons un exemple chiffré. Si par exemple je me déplace de 10 m en 5 s, mon intervalle mesuré par une personne extérieure est

$$c^2 \Delta t^2 - \Delta x^2 = (300.000.000)^2 \cdot (5)^2 - (10)^2 = 2,25 \cdot 10^{18} \text{ m}.$$

Si je fais le calcul pour moi, j'ai que $\Delta x' = 0$ parce que je ne me déplace pas, et $\Delta t' = 5$ parce que je me suis déplacé en 5 secondes. Le truc est que à côté de $(300.000.000)^2$, l'intervalle spatial Δx ne pèse pas grand chose. Ça ne change presque rien qu'il soit de 5 mètres ou de zéro. Ça ne change pas grand chose, mais ça change quand même ! Entre moi qui calcule ou une personne extérieure, l'intervalle change de 100 sur un nombre de la grandeur de 200.000.000.000.000.000.000 !

Reprenons plus clairement le raisonnement. D'après la mécanique classique, l'intervalle mesuré par deux personnes est différent, mais très peu différent. Inutile de dire que du temps de Newton, on n'avait pas les moyens techniques de mesurer si cet intervalle est effectivement différent ou bien s'il est en réalité égal. C'est un peu comme si on te mettait un spot dans les yeux et qu'on te demandait si c'est un spot de 1000 W ou de 1001 W. Bonne chance pour le dire !

43.4.4.2 En mécanique relativiste

Maintenant qu'on a des moyens techniques nettement plus poussés que Newton, on a pu mesurer que l'intervalle est égal. L'intervalle est un invariant. Cela n'est pas un nouveau principe physique parce qu'il découle des transformations de Lorentz.

43.4.5 Le cône de lumière d'un point

Il est intéressant de dessiner dans le plan (t, x) l'ensemble des points atteints par le rayon lumineux. Le point (t, x) est atteint si $c^2 t^2 - x^2 = 0$, ou encore si $x = \pm ct$. Cela forme deux droites dans le plan tracé par les coordonnées t et x . Ces deux droites forment ce qu'on appelle le **cône de lumière** du point $(0, 0)$.

43.4.6 Contraction des longueurs

Bob prend un morceau de bois qu'il mesure de longueur l et le dépose devant lui. À l'instant t (de Bob), les deux extrémités sont aux coordonnées $e_1 = (t, 0)$ et $e_2 = (t, l)$.

Afin de savoir quelle est la longueur de ce même morceau de bois pour Alice, il faut qu'elle mesure les deux extrémités en même temps (pour elle), et qu'elle fasse la différence. Comme Bob et Alice déclenchent leurs chronomètres en même temps, le plus simple est de faire la mesure à cet instant.

Pour Bob, c'est clair : les coordonnées des deux extrémités sont $e_1 = (0, 0)$ et $e_2 = (0, l)$. La longueur du bois est l . Pour savoir quelle est la longueur mesurée par Alice, on utilise les transformations de Lorentz qui donnent les coordonnées e'_1 et e'_2 relatives à Alice. On trouve $e'_1 = (0, 0)$ et

$$e'_2 = \left(\frac{-vl/c^2}{\gamma(v)}, \frac{l}{\gamma(v)} \right). \quad (43.26)$$

En d'autres termes, on a $x_1 = 0$ et $x_2 = l/\gamma(v)$, ce qui fait que la longueur observée par Alice est $l' = x_2 - x_1 = l/\gamma(v)$.

Eh bien ce résultat est faux. Si tu vois pourquoi sans lire la suite, tu es très fort.

Pour mesurer la longueur d'un objet, il faut mesurer la position des deux bouts *en même temps* puis faire la différence entre les deux. Effectivement, e_1 et e_2 sont en même temps pour Bob, et donc Bob peut mesurer la longueur de son bout de bois en faisant la différence $x_2 - x_1$. Mais comme le montre les coordonnées (43.26), les événements e'_1 et e'_2 ne se passent pas en même temps pour Alice ! Eh oui : $t'_1 = 0$ et $t'_2 = -vl/c^2\gamma(v)$; c'est pas la même chose.

Il faut donc trouver un événement qui pour Alice correspond à l'extrémité du bout de bois au temps $t' = 0$. Comme l'événement général qui correspond au bout du bois pour Bob est (t, l) , l'événement général est pour Alice

$$t' = \frac{t - \frac{v}{c^2}l}{\gamma(v)} \quad x' = \frac{l - vt}{\gamma(v)}. \quad (43.27)$$

Afin d'avoir $t' = 0$, il faut $t = vl/c^2$. En mettant cette valeur de t dans x' , on trouve

$$x' = \frac{l - v \left(\frac{vl}{c^2} \right)}{\sqrt{1 - \frac{v^2}{c^2}}} = \frac{l \left(1 - \frac{v^2}{c^2} \right)}{\sqrt{1 - \frac{v^2}{c^2}}} = l\gamma(v).$$

Et là, c'est la bonne formule. Si un objet a une longueur l dans le référentiel où il est au repos, il aura une longueur

$$l' = l\sqrt{1 - \frac{v^2}{c^2}} \quad (43.28)$$

dans un référentiel qui se déplace à la vitesse v par rapport à l'objet.

43.4.7 Dilatation des intervalles de temps

Encore un petit effort et promis, je te donne une application concrète que tu connais des bizarreries de la relativité. Mais en attendant, regarde bien ta montre, tu ne va pas en croire tes yeux !

Reprenons Bob et Alice. On se rappelle que Bob et Alice avaient déclenchés leurs chronomètres en même temps quand Alice était passée devant Bob. Un peu plus tard, Alice regarde sa montre qui indique un temps t . Et elle se demande si Bob a aussi à ce moment une montre qui indique un temps t .

Ce serait dingue que non hein ! ? ! En effet, si je synchronise ma montre avec quelqu'un et que je pars faire un tour, ma montre ne sera pas tout d'un coup désynchronisée. Oui, mais Alice, elle cours presque à la vitesse de la lumière ... et à ces vitesses-là, on a déjà vu des choses incroyables. Calculons donc pour en avoir le cœur net.

Le fait qu'Alice regarde sa montre est un événement qui se passe pour Alice aux coordonnées $(t', 0)$ (le zéro c'est parce que par rapport à elle-même, Alice est toujours au repos). À quelles coordonnées (t, x) pour Bob correspond cet événement ?

L'équation de t en fonction de t' et x' dans les transformations de Lorentz (43.24) prise avec $x' = 0$ donnent

$$t = \frac{t'}{\gamma(v)}.$$

Et si, juste pour le plaisir, on faisait l'inverse ? Bob regarde sa montre, il voit un temps t et sa coordonnée spatiale est $x = 0$. À quel temps d'Alice cela correspond ? Mettons $x = 0$ dans la transformation de Lorentz de t' en fonction de t et x . Ce qu'on obtient c'est

$$t' = \frac{t}{\gamma(v)}.$$

N'est-ce pas génial ? C'est la même ! Évidemment, ça ne pouvait pas être autre chose : le principe de relativité demande qu'on ne puisse pas faire la différence entre Alice qui cours vers la droite avec Bob assis et Alice assise avec Bob qui cours vers la gauche. C'est exactement pour ça que dans une gare, quand le train d'à côté démarre, il t'arrive de croire que c'est ton train qui démarre : tu ne peux pas faire la différence, c'est un principe physique.

43.4.8 Invariance de l'intervalle

Dans deux secondes, je vais te montrer comment une utilisation intelligente des exponentielles permet de trouver un résultat très fort en relativité. Quoi ? Les exponentielles, les mêmes qu'au cours de math ? Eh oui : la même exponentielle que celle qu'on t'a introduit avec des populations de bactéries qui se multiplient, cette même exponentielle qui a la miraculeuse propriété d'être égale à sa propre dérivée.

Mais n'anticipons pas.

Nous avons déjà signalé que si la quantité $\Delta s^2 = c^2 \Delta t^2 - \Delta x^2$ était nulle pour un observateur, alors elle était nulle pour tous les observateurs. Supposons deux événements A et B observés par Alice et Bob. Bob les note aux coordonnées (t_a, x_a) , et (t_b, x_b) tandis qu'Alice les note en (t'_a, x'_a) et (t'_b, x'_b) .

L'intervalle entre les deux événements mesuré par Bob sera

$$s^2 = c^2(t_b - t_a)^2 - (x_b - x_a)^2,$$

tandis que ce même intervalle mesuré par Alice sera

$$s'^2 = c^2(t'_b - t'_a)^2 - (x'_b - x'_a)^2.$$

On peut bien entendu remplacer dans la première équation les t_a , t_b , x_a et x_b par leurs valeurs en termes de t'_a , t'_b , x'_a et x'_b données par les transformations de Lorentz. Tu paries que les trois quart des termes dans le calcul se simplifient et qu'il restera exactement s'^2 ? Je te dis que oui, et je te conseille de me croire sur parole, sinon tu vas devoir lire le calcul suivant :

$$\begin{aligned} s^2 = c^2(t_b - t_a)^2 - (x_b - x_a)^2 &= c^2 \left(\frac{1}{\gamma(v)}(t'_b + \frac{v}{c^2}x'_b) - \frac{1}{\gamma(v)}(t'_a + \frac{v}{c^2}x'_a) \right)^2 \\ &\quad - \left(\frac{1}{\gamma(v)}(vt'_b + x'_b) - \frac{1}{\gamma(v)}(vt'_a + x'_a) \right)^2. \end{aligned}$$

Jusqu'ici, on n'a fait que remplacer les choses par leurs valeurs données par les transformations de Lorentz. Maintenant on regroupe à l'intérieur de chaque parenthèse les termes de façon à faire

apparaître $\Delta x'$ et $\Delta t'$:

$$\begin{aligned} s^2 &= \frac{c^2}{\gamma(v)^2} \left((t'_b - t'_a) + \frac{v}{c^2} (x'_b - x'_a) \right)^2 \\ &\quad - \frac{1}{\gamma(v)^2} \left((x'_b - x'_a) + v(t'_b - t'_a) \right)^2 \\ &= \frac{c^2}{\gamma(v)^2} \left((\Delta t')^2 + 2\frac{v}{c^2} \Delta t' \Delta x' + \frac{v^2}{c^4} (\Delta x')^2 \right) \\ &\quad - \frac{1}{\gamma(v)^2} \left((\Delta x')^2 + 2v \Delta x' \Delta t' + v^2 (\Delta t')^2 \right). \end{aligned}$$

Là, on a utilisé le produit remarquable $(a+b)^2 = a^2 + 2ab + b^2$, et on a systématiquement renommé tous les intervalles avec la notation Δ pour être plus compact. Maintenant, on va regrouper tous les termes contenant $(\Delta t')^2$ ensemble, tous ceux qui contiennent $\Delta t' \Delta x'$ ensemble et ceux qui contiennent $(\Delta x')^2$ ensemble. Autre manière de le dire, on met les Δ en évidence comme on peut. On trouve ceci :

$$\begin{aligned} (\Delta t')^2 \left(\frac{c^2}{\gamma(v)^2} - \frac{v^2}{\gamma(v)^2} \right) &+ \Delta t' \Delta x' \left(\frac{2vc^2}{\gamma(v)^2 c^2} - \frac{2v}{\gamma(v)^2} \right) \\ &+ (\Delta x')^2 \left(\frac{c^2 v^2}{c^4 \gamma(v)^2} - \frac{1}{\gamma(v)^2} \right). \end{aligned}$$

À partir de là, je te laisse vérifier (en utilisant le fait que $\gamma(v)^2 = 1 - v^2/c^2$) que les coefficients se simplifient beaucoup et valent finalement respectivement c^2 , 0 et -1 comme il se doit. Avec tout ça, nous avons montré le résultat très important suivant :

L'intervalle entre deux événements est invariant sous les changements de repères d'inertie, c'est-à-dire que la valeur mesurée par n'importe qui qui se déplace en MRU sera toujours la même.

Pourquoi cela est tellement important ? À cause de Pythagore et d'une petite démonstration à coups d'exponentielles⁵.

43.4.8.1 Rappel de trigonométrie hyperbolique

Les fonctions de trigonométrie hyperboliques sont :

$$\cosh(x) = \frac{e^x + e^{-x}}{2} \qquad \sinh(x) = \frac{e^x - e^{-x}}{2}. \qquad (43.29)$$

Elles ont pas mal de propriétés en commun avec les sinus cosinus et normaux. D'abord, leurs dérivées sont faciles à calculer :

$$\begin{aligned} \cosh'(x) &= \sinh(x) \\ \sinh'(x) &= \cosh(x) \end{aligned}$$

où tu noteras qu'il n'y a pas de signe moins qui apparaît, contrairement au cas de la trigonométrie normale. Une autre propriété qui ressemble fort à une propriété de la trigonométrie est :

Proposition 43.2.

Pour tout $x \in \mathbb{R}$,

$$\cosh^2(x) - \sinh^2(x) = 1 \qquad (43.30)$$

avec un signe moins comme différence avec la trigonométrie.

5. oui oui tout ton cours de math va finir par y passer.

Démonstration. La preuve revient simplement à calculer en utilisant le produit remarquable de $(a + b)^2$. D'abord, on a :

$$\cosh^2(x) = \frac{1}{4}(e^x + e^{-x})^2 = \frac{1}{4}(e^{2x} + 2e^x e^{-x} + e^{-2x}) = \frac{1}{4}(e^{2x} + 2 + e^{-2x})$$

où l'on a utilisé le fait que $(e^x)^2 = e^{2x}$ et que $e^x e^{-x} = 1$. Il te reste à faire la même chose pour $\sinh^2(x)$, la réponse est :

$$\sinh^2(x) = \frac{1}{4}(e^{2x} - 2 + e^{-2x}).$$

En faisant la différence entre les deux, il reste 1. □

Une propriété qui est par contre très différente entre la trigonométrie plane et la trigonométrie hyperbolique, c'est la périodicité. Les fonctions usuelles \cos et \sin sont périodiques. Pas les fonctions hyperboliques.

Proposition 43.3.

La fonction $\sinh: \mathbb{R} \rightarrow \mathbb{R}$ est bijective.

Démonstration. Il faut démontrer que sinus hyperbolique est injective et surjective. Calculons d'abord les limites. Comme tu sais que $\lim_{x \rightarrow \infty} e^x = \infty$ et $\lim_{x \rightarrow -\infty} e^x = 0$, tu vois facilement que

$$\lim_{x \rightarrow -\infty} \sinh(x) = -\infty \qquad \lim_{x \rightarrow \infty} \sinh(x) = \infty. \qquad (43.31)$$

Par ailleurs, la fonction sinus hyperbolique est continue et respecte donc le théorème de la valeur intermédiaire⁶ 13.50. Soit $y \in \mathbb{R}$. Voyons s'il existe un $x \in \mathbb{R}$ tel que $\sinh(x) = y$. Les deux limites indiquent qu'il existe $x_1 \in \mathbb{R}$ tel que $\sinh(x_1) < y$ et $x_2 \in \mathbb{R}$ tel que $\sinh(x_2) > y$. Le théorème de la valeur intermédiaire conclut qu'il existe un x entre x_1 et x_2 tel que $\sinh(x) = y$. Cela prouve la surjectivité.

Pour l'injectivité, on va utiliser le **théorème de Rolle** et une petite preuve par l'absurde. Supposons que $\sinh(x_1) = \sinh(x_2)$ avec $x_1 \neq x_2$. Dans ce cas, le théorème de Rolle nous dit qu'il existe un x entre x_1 et x_2 tel que $\sinh'(x) = 0$. La dérivée de sinus hyperbolique étant cosinus hyperbolique, il faut se demander il existe un x tel que $\cosh(x) = 0$. Étant donné que $e^x > 0$ pour tout x , en fait le cosinus hyperbolique ne s'annule jamais. □

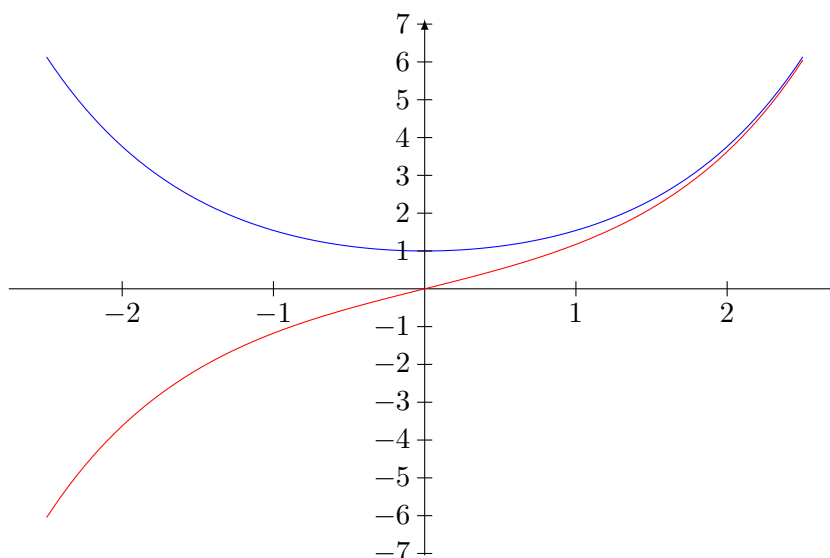


FIGURE 43.1 – En rouge, la fonction $x \mapsto \sinh(x)$ et en bleu, la fonction $x \mapsto \cosh(x)$.

6. Je t'avais dit que tout tons cours de math allait y passer hein.

Un très bon exercice serait de faire un étude complète des fonctions cosinus et sinus hyperbolique. Leur graphes sont donnés à la figure 43.1

Un corollaire de la surjectivité de \sinh sur \mathbb{R} est que si je prends n'importe quel deux nombres dont la différence des carrés vaut 1, alors ces carrés sont représentables avec des fonctions hyperboliques :

$$\forall x, y \in \mathbb{R} \text{ tels que } x^2 - y^2 = 1, \exists \xi \in \mathbb{R} \text{ tel que } x^2 = \cosh(\xi) \text{ et } y^2 = \sinh(\xi).$$

La **tangente hyperbolique** est définie par

$$\tanh(x) = \frac{\sinh(x)}{\cosh(x)}. \quad (43.32)$$

Un bon exercice est de prouver les deux relations suivantes :

$$\sinh(x) = \frac{\tanh(x)}{\sqrt{1 - \tanh^2(x)}} \quad \cosh(x) = \frac{1}{\sqrt{1 - \tanh^2(x)}}. \quad (43.33)$$

43.4.8.2 Les transformations de Lorentz (bis)

Nous avons prouvé qu'en relativité, l'intervalle est un invariant. Pour cela, nous avons utilisé les transformations de Lorentz démontrées à partir de l'hypothèse d'invariance de la vitesse de la lumière. Eh bien, oublions un instant que la vitesse de la lumière soit invariante, et posons à la place comme hypothèse que l'intervalle soit invariant. C'est-à-dire que si Bob mesure un événement aux coordonnées (t, x) et Alice en (t', x') , alors $c^2 t^2 - x^2 = c^2 (t')^2 - (x')^2$.

Théorème 43.4.

Les transformations de Lorentz sont les seules qui laissent l'intervalle invariant.

Démonstration. Toute la partie comme quoi les transformations doivent être linéaires reste parce que cette partie ne demandait pas l'invariance de la vitesse de la lumière.

Nous cherchons donc les transformations entre Alice et Bob sous la forme

$$\begin{aligned} t' &= \alpha t + \beta x \\ x' &= \gamma t + \delta x \end{aligned}$$

telles que $c^2 (t')^2 - (x')^2 = c^2 t^2 - x^2$. Lorsque Alice passe devant Bob, ils déclenchent tout deux leurs chronomètre et leurs axes. C'est-à-dire que si à ce moment un événement se trouve à droite pour Alice, il est aussi à droite pour Bob. On doit donc avoir, quand $t = t' = 0$, que $x > 0$ implique $x' > 0$. Cela donne la contrainte que $\delta > 0$. D'autre part, comme leurs chronomètres vont dans le même sens (ils choisissent tout les deux de *compter* le temps et non *décompter*), on a $\alpha > 0$.

En développant l'expression de $(s')^2$ en termes de t et x , on trouve la condition d'invariance de l'intervalle sous la forme :

$$c^2(\alpha^2 t^2 + 2\alpha\beta tx + \beta^2 x^2) - (\gamma^2 t^2 + 2\gamma\delta tx + \delta^2 x^2) = c^2 t^2 - x^2, \quad (43.34)$$

qui doit être valable pour tout t et pour tout x . En $t = 0$ on trouve la condition

$$\delta^2 - c^2 \beta^2 = 1. \quad (43.35)$$

Cela implique qu'il existe un $\xi \in \mathbb{R}$ tel que $\delta^2 = \cosh^2(\xi)$ et $c^2 \beta^2 = \sinh(\xi)$. La première équation donne $\delta = \cosh(\xi)$ (il faut rejeter $\delta = -\cosh(\xi)$ parce qu'on a demandé que $\delta > 0$). Pour la seconde, on trouve $c\beta = \sinh(\xi)$ où l'on peut oublier la possibilité $c\beta = -\sinh(\xi)$ parce que cela revient juste à renommer $\xi \rightarrow -\xi$ (la fonction sinus hyperbolique est impaire). Bref, il existe un ξ tel que

$$\begin{aligned} \delta &= \cosh(\xi) \\ \beta &= \frac{\sinh(\xi)}{c} \end{aligned} \quad (43.36)$$

En mettant maintenant $x = 0$ dans la condition (43.34), on trouve la condition

$$\alpha^2 - \frac{\gamma^2}{c^2} = 1.$$

Pour les mêmes raisons qu'avant, il existe un $\eta \in \mathbb{R}$ tel que

$$\begin{aligned}\alpha &= \cosh(\eta) \\ \gamma &= c \sinh(\eta).\end{aligned}\tag{43.37}$$

Rien qu'en regardant deux cas très particuliers, on a déjà bien avancé, non ? Remettons maintenant les valeurs (43.36) et (43.37) dans la condition (43.34). En utilisant l'identité $\cosh^2(x) - \sinh^2(x) = 1$, et en séparant les termes en t^2 , x^2 et tx pour satisfaire la condition, il faut

$$\cosh(\eta) \sinh(\xi) = \sinh(\eta) \cosh(\xi)\tag{43.38}$$

parce que les termes en t^2 et x^2 donnent exactement $c^2 t^2 - x^2$ et qu'il faut que le terme en tx s'annule. Mettons la condition (43.38) au carré, et substituons $\cosh^2(\eta) = 1 + \sinh^2(\eta)$ et $\cosh^2(\xi) = 1 + \sinh^2(\xi)$, il reste

$$\sinh^2 \xi = \sinh^2 \eta,$$

ce qui signifie $\sinh \xi = \pm \sinh \eta$, ou encore $\xi = \pm \eta$. On voit que $\xi = -\eta$ ne fonctionne pas dans (43.38), donc on reste avec $\xi = \eta$ et les transformations prennent la forme

$$\begin{aligned}t' &= \cosh(\xi)t + \frac{\sinh(\xi)}{c}x \\ x' &= c \sinh(\xi)t + \cosh(\xi)x.\end{aligned}\tag{43.39}$$

Ce que nous avons prouvé, c'est qu'il existe un $\xi \in \mathbb{R}$ tel que les transformations entre Alice et Bob aient cette forme. Il faut trouver ce que vaut ξ en fonction de la vitesse v à laquelle Alice court.

Pour ce faire, étudions le mouvement d'Alice. Bob la voit aux coordonnées (t, vt) , ce qui correspond à

$$x' = c \sinh(\xi)t + \cosh(\xi)vt$$

pour Alice. Mais ces coordonnées sont celles de Alice elle-même, donc $x' = 0$, ce qui donne⁷ $vt = -c \sinh(\xi)t / \cosh(\xi)$, ou encore

$$\tanh(\xi) = -\frac{v}{c}\tag{43.40}$$

En utilisant les relations (43.33), on trouve

$$\cosh(\xi) = \frac{1}{\sqrt{1 - \frac{v^2}{c^2}}} \qquad \sinh(\xi) = \frac{-v/c}{\sqrt{1 - \frac{v^2}{c^2}}}.\tag{43.41}$$

En remettant ces valeurs dans les transformations (43.39), on trouve

$$t' = \frac{t - \frac{v}{c^2}x}{\gamma(v)}\tag{43.42}$$

$$x' = \frac{x - vt}{\gamma(v)},\tag{43.43}$$

exactement les transformations de Lorentz !

□

Ce résultat est important pour une raison assez simple : maintenant, la théorie de la relativité est indépendante de toute considérations sur la lumière. En effet, ce que nous venons de prouver, c'est que s'il existe une vitesse c telle que $c^2 t^2 - x^2 = c^2 (t')^2 - (x')^2$, alors (t, x) et (t', x') sont liés par les transformations de Lorentz.

⁷ Conditions d'existence : $\cosh(\xi) \neq 0$; heureusement, nous avons vu que le cosinus hyperbolique ne s'annule jamais.

43.4.9 Vitesse limite

Afin de nous passer de l'hypothèse d'invariance de la vitesse de la lumière, nous avons prouvé que l'hypothèse d'invariance de l'intervalle était suffisante. Mais il faut avouer que cette hypothèse n'est pas très intuitive. Nous allons montrer maintenant que l'existence d'une vitesse limite est une troisième hypothèse qui peut être utilisée comme alternative aux deux premières.

43.5 Applications

Une première application sympa est le logiciel⁸ *lightspeed*. Si tu es sous Ubuntu-Linux, installe juste le paquet nommé *lightspeed*, et régales-toi ! Tu verras c'est marrant. Si tu utilise des fenêtres, je laisse faire l'adage « Windows c'est facile ».

43.5.1 Le GPS

Pour qu'un système **GPS** puisse te localiser, en gros, il t'envoie un signal, tu lui réponds et il mesure le temps qu'il a fallu à la lumière pour faire l'aller-retour. Déjà, tu remarques que cela n'est possible que grâce au fait que la vitesse de la lumière soit finie. Sinon, le GPS ne fonctionnerait pas. Mais il y a mieux.

Comme pour te localiser il faut plusieurs satellites en plus de ton appareil, il faut que les horloges internes de tout ce petit monde soient bien synchronisées, sinon pour mesurer des intervalles de temps et calculer des distances, c'est mal parti. Eh mais un satellite, ça bouge assez vite (surtout que les mesures doivent être très précises), et en plus ça ne fait même pas un MRU, vu que ça tourne en rond. Comme tu vois tout le travail qu'il a fallu faire pour trouver les transformations de Lorentz d'un MRU, tu t'imagines le travail pour un mouvement circulaire ! Eh bien ce travail a été fait, et le résultat est que si on en tient pas compte, les contractions temporelles liées à la relativité sont suffisamment grandes pour complètement dérégler le GPS.

43.5.2 Les ondes électromagnétiques

Tu te souviens qu'au début du chapitre, nous avons dit que le problème qui a amené la relativité était la propagation du champ électrique. Maintenant que nous avons déjà vu une partie des conséquences du problème, il est temps de se rendre compte que les champs électriques et magnétiques sont les objets les plus soumis aux bizarreries relativistes du monde : elles se propagent à la vitesse de la lumière. Regarde un coup autour de toi ; tout ce qui est champ électromagnétique a besoin de la relativité pour être bien compris : GSM, lumière, four à micro-onde, radio, wifi, fibre optique, ...

Si un jour un ingénieur te dit qu'il n'y a pas besoin de connaître la relativité pour inventer la radio (c'est vrai : la radio a été inventée avant la relativité), ni pour construire une fibre optique, dis lui en pensant à moi qu'il utilise tout le temps les équations de Maxwell⁹, et que ces équations sont relativistes.

Bref, soit convaincu que tu vis dans un monde relativiste et que les transformations de Lorentz te suivent à chacun de tes pas.

43.6 Mécanique relativiste

Cela est bien beau, mais la dilatation du temps, et les contractions de longueurs doivent bien avoir des répercussions sur la cinématique et la dynamique des objets. Est-ce que le théorème de l'énergie cinétique est encore valable ? est-ce que les lois de Newton tiennent encore la route ?

8. jeu de mot sur « application » ! ah ah !

9. C'est sous ce nom là qu'on nomme l'ensemble des équations de l'électromagnétisme comme la loi de l'induction.

43.6.1 Des problèmes, toujours des problèmes

Attardons-nous un peu pour faire quelques commentaires sur cette citation du chevalier pégase dans [les chevaliers du zodiaque](#) :

Ses coups vont à la vitesse de la lumière et pourtant je les vois distinctement arriver.

Est-ce possible ? Nous avons vu qu'il y avait des dénominateurs qui s'annulent quand des objets se déplacent plus vite que la lumière ; or pour voir venir un rayon de lumière qui vient vers soi, il faudrait que le rayon émette de la lumière devant elle. Ça semble un peu mal parti pour respecter les lois de la relativité, non ?

Cela pose en tout cas une question qu'il faudra résoudre. On *entend* venir une ambulance parce qu'elle émet du son qui avance plus vite qu'elle. Pas de problèmes avec ça. Mais quid de la *voir* venir ?

On peut voir venir un tram parce qu'il émet de la lumière ; cette lumière allant plus vite que le tram, elle arrive à nos yeux avant le tram lui-même. Cela est très bien. Mettons que le tram avance à 50 km/h ; pour le conducteur, la lumière de son phare avant avance devant lui à la vitesse c . Par conséquent pour un observateur au sol, cette même lumière devrait avancer à la vitesse $c + 50$. Encore une fois, on a un problème d'invariance de la vitesse de la lumière ; mais comme c'est de la lumière, on est habitué à ce que des trucs bizarres arrivent. On ne sera pas étonné que $c + 50 = c$ d'une manière ou d'une autre¹⁰. Pire. Si un vaisseau spatial avance à la vitesse 200000 km/s et qu'il envoie en reconnaissance un vaisseau devant lui à la vitesse de 150000 km/s, le vaisseau de reconnaissance ira à la vitesse 150000 km/s par rapport au vaisseau principal. Et par rapport au sol, il ira à la vitesse $150000 + 200000 = 350000$ km/s, ce qui est impossible. Il faudra trouver quelque chose pour que ça se passe bien.

Un autre problème maintenant.

Prenons une masse m que l'on soumet à une force constante F . Par la loi de Newton, $a = F/m$ est constante et la vitesse après un temps t vaut $v = Ft/m$. Pas de bol, ça devient plus grand que la vitesse de la lumière à partir du temps $t = cm/F$. Ça est un problème hein ? Il faut trouver un truc pour qu'avec une force constante, l'accélération diminue.

43.6.2 Loi d'addition des vitesses

Si Bob observe un objet se déplacer à la vitesse V , alors Alice devrait l'observer bouger à la vitesse $V - v$. Tout comme si une vache voit passer un train à 90 km/h, alors le vélo qui avance à 25 km/h le voit passer à 65 km/h.

Maintenant, tu es habitué à ce que rien ne se passe comme d'habitude, donc tu te doutes bien qu'en réalité la bonne formule ne va pas être $V - v$.

Bob observe l'objet aux coordonnées (t, Vt) , ce qui fait pour Alice :

$$\left(\frac{t - \frac{v}{c^2} Vt}{\gamma(v)}, \frac{Vt - vt}{\gamma(v)} \right).$$

En divisant le x' d'Alice par le t' d'Alice, on trouve la vitesse mesurée par Alice :

$$V' = \frac{(V - v)t}{\gamma(v)} \frac{\gamma(v)}{t(1 - \frac{vV}{c^2})} = \frac{V - v}{1 - \frac{vV}{c^2}}.$$

La loi de transformation des vitesses relativiste est donc

$$V' = \frac{V - v}{1 - \frac{vV}{c^2}}. \quad (43.44)$$

Qu'en est-il de notre $c + 50 = c$? Disons que Bob lance un bisou à Alice pendant qu'elle arrive vers lui. Le bisou arrive à la vitesse de la lumière (càd $V = c$) tandis que Alice s'approche de Bob

10. et je ne te cache pas que c'est ce qui va arriver.

à la vitesse 50 m/s (càd $v = 50$). Donc la vitesse à laquelle Alice devrait voir arriver le bisou est bien $c + 50$. En utilisant la formule d'addition relativiste des vitesses (43.44), nous trouvons

$$V' = \frac{c - 50}{1 - \frac{50c}{c^2}} = \frac{c - 50}{1 - \frac{50}{c}} = \frac{c(c - 50)}{c - 50} = c.$$

Donc effectivement en relativité quand on additionne des vitesses il faut penser à la règle du « $c + 50 = c$ ».

43.6.3 L'action d'une force

L'équation fondamentale de la mécanique classique est

$$F = ma.$$

Or tu n'es pas sans savoir que l'accélération est la dérivée seconde de la position par rapport au temps. Nous noterions donc $F = mx''(t)$. Le problème est évidemment que si F est constante, on trouve $v = Ft/m$ qui dépasse toujours la vitesse c quand t est assez grand. Il faudra donc modifier la loi $F = ma$. Pour cela, posons-nous des questions sur la dérivée $x'(t)$. On dérive par rapport au temps ; oui mais nous avons vu que le temps n'est pas le même pour tout le monde. Introduisons donc la notation

$$v = \frac{dx}{dt} \quad (43.45)$$

qui ne signifie rien d'autre que nous dérivons x par rapport à t et non par rapport au temps t' de quelqu'un d'autre. Dans le cadre de la relativité, ce que signifie l'équation (43.45) est que v est la dérivée de x par rapport à t . Dans le cas où x et t sont les coordonnées de la position d'Alice mesurées par Bob, cela signifie qu'on dérive la position *mesurée par Bob* par rapport au temps *mesuré par Bob*.

Ce que dit la relativité est que cette quantité v ne peut pas varier proportionnellement à la force sous peine de dépasser la vitesse de la lumière. La subtilité est de modifier la loi de Newton en disant que la quantité qui varie sous l'action d'une force n'est plus $dx/dt = v$, mais

$$\frac{dx}{dt'} = \frac{v}{\sqrt{1 - \frac{v^2}{c^2}}},$$

c'est-à-dire la dérivée de la position *mesurée par Bob* par rapport au temps *mesuré par Alice* ! La loi de Newton $v = Ft/m$ devient donc

$$\frac{v}{\sqrt{1 - \frac{v^2}{c^2}}} = \frac{Ft}{m}. \quad (43.46)$$

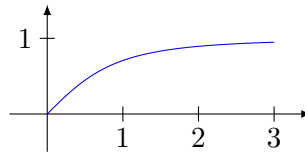
Est-ce que cela résout le problème ? Pour le savoir, regardons la vitesse acquise par le mobile de masse m soumis à la force F pendant un temps t . Il faut résoudre l'équation (43.46) par rapport à v et voir si cela reste bien toujours inférieur à c . On commence par mettre la racine à droite et à élever toute l'équation au carré :

$$\begin{aligned} v^2 &= \frac{F^2 t^2}{m^2} \left(1 - \frac{v^2}{c^2} \right) \\ v^2 \left(1 + \frac{F^2 t^2}{c^2 m^2} \right) &= \frac{F^2 t^2}{m^2} \\ v &= \frac{\sqrt{F^2 t^2 / m^2}}{\sqrt{1 + \frac{F^2 t^2}{c^2 m^2}}}, \end{aligned}$$

et donc finalement

$$v(t) = \frac{Ft}{m\sqrt{1 + \frac{F^2 t^2}{c^2 m^2}}}. \quad (43.47)$$

Tu dois remarquer que si F et t ne sont pas trop grands, l'expression $F^2 t^2 / c^2 m^2$ est minuscule parce que c est énorme. Si on fait l'approximation $F^2 t^2 / c^2 m^2 = 0$ dans cette expression, on retrouve $v = Ft/m$. Cela montre qu'à moins de faire des expériences avec de très grandes forces pendant énormément de temps, on ne peut pas voir la différence entre la mécanique de Newton et la mécanique relativiste.



Sur le graphique suivant, la vitesse en fonction du temps lorsqu'une particule de masse $m = 1$ est soumise à une force constante. Pour les besoins du graphique, nous avons mis à 1 la vitesse c . Tu vois que quand la vitesse n'est pas très grande, le graphique est presque celui d'une droite ; et à partir d'un certain moment, la courbe s'infléchit pour tendre vers 1 sans l'atteindre.

Remarque que si on maintient une accélération constante égale à celle de la gravité terrestre pendant deux heures, on arrive déjà sur la Lune, à une vitesse de 75 km/s, c'est-à-dire encore rien par rapport à la vitesse de la lumière ! Cela pour te dire que la formule (43.47) a l'air d'être très différente de la formule classique $v = Ft/m$, mais en réalité tant qu'on n'atteint pas des forces énormes, elle ressemble très fort.

Vérifions maintenant que la formule (43.47) n'est pas en contradiction avec l'impossibilité de dépasser la vitesse de la lumière. Pour cela, regardons ce qu'il se passe si on applique une force constante F sur un objet de masse m pendant un temps très long. C'est-à-dire : calculons la limite

$$\lim_{t \rightarrow \infty} v(t).$$

Tu vois tout de suite qu'on est sur un cas $\frac{\infty}{\infty}$, ce qui t'oblige à utiliser la règle de l'Hospital. On peut cependant un peu réfléchir et deviner la réponse sans passer par des math trop compliquées.

En effet, quand t est vraiment énorme, l'expression $\frac{F^2 t^2}{m^2 c^2}$ devient très grande, et le 1 qui se trouve à côté ne vaut plus grand chose, on peut le négliger.

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{Ft}{m \sqrt{1 + \frac{F^2 t^2}{c^2 m^2}}} &= \lim_{t \rightarrow \infty} \frac{Ft}{m \sqrt{\frac{F^2 t^2}{m^2 c^2}}} \\ &= \lim_{v \rightarrow \infty} \frac{Ft}{m \frac{Ft}{cm}} \\ &= c. \end{aligned} \tag{43.48}$$

Tout est bien : on arrive au maximum à la vitesse de la lumière, mais il faut un temps infini pour y parvenir. Conclusion : il n'est pas possible d'accélérer un objet jusqu'à atteindre la vitesse de la lumière.

43.6.4 Équivalence entre la masse et l'énergie

Le moment est venu de montrer ce que signifie la fameuse formule $E = mc^2$.

43.7 Principe de correspondance

Nous ne sommes pas parvenu à démontrer la formule (43.46) de la mécanique relativiste qui montre comme un objet accélère sous l'effet d'une force constante. Nous avons juste montré qu'il fallait modifier la loi $v = Ft/m$ et nous avons prit la première modification qui nous soit tombée sous la main, à savoir qu'il faut dériver la position par rapport au temps de l'objet qu'on observe plutôt que par rapport au temps de l'observateur.

En fait, il est possible de prouver rigoureusement ¹¹ la formule

$$\frac{Ft}{m} = \frac{\alpha v}{\sqrt{1 - \frac{v^2}{c^2}}}.$$

Mais il n'y a pas moyen de trouver la valeur de la constante α . Tout ce qu'il y a moyen de trouver avec l'hypothèse de l'invariance de la vitesse de la lumière est l'existence d'une constante telle que cette formule soit vraie.

Afin de fixer la constante α , il faut faire intervenir un principe physique supplémentaire, le **principe de correspondance**

Loi numéro 4.

Lorsque la vitesse d'une particule est faible, les équations doivent être en première approximation les mêmes que celles de la mécanique classique.

Que signifie *en première approximation* ? Tu sais qu'une fonction $x \mapsto f(x)$ peut être approximée (pour des petits x) par la formule

$$f(x) \simeq f(0) + xf'(0).$$

Nous voudrions donc que Ft/m soit en première approximation égal à v . Nous devons étudier la fonction

$$f(v) = \frac{\alpha v}{\sqrt{1 - \frac{v^2}{c^2}}}.$$

Voir ce que vaut cette fonction en première approximation lorsque v est petit est un exercice de dérivation. En utilisant la règle de dérivation des fractions, on trouve que

$$f'(v) = \frac{\alpha}{\sqrt{1 - \frac{v^2}{c^2}}} + \frac{\alpha v^2}{c^2 \left(1 - \frac{v^2}{c^2}\right)^{3/2}},$$

et donc que $f'(0) = \alpha$. Bien entendu, $f(0) = 0$. En première approximation, nous trouvons donc

$$f(v) \simeq \alpha v \tag{43.49}$$

qui doit être égal à la quantité non relativiste v . Nous en déduisons qu'il faut fixer $\alpha = 1$, et on tombe sur la formule relativiste proposée plus haut

$$\frac{Ft}{m} = \frac{v}{\sqrt{1 - \frac{v^2}{c^2}}}$$

L'utilisation cruciale du principe de correspondance a une répercussion énorme sur notre vision de la physique. En effet, la relativité d'Einstein ne parvient pas à *remplacer* la mécanique de Newton. On a besoin d'invoquer la mécanique de Newton pour fixer la théorie. On peut écrire l'axiome suivant :

$$\lim_{v \rightarrow 0} \text{Einstein} = \text{Newton}. \tag{43.50}$$

Cela n'est pas une propriété de la théorie d'Einstein, mais un de ses axiomes !

La relativité ne fait donc pas table rase des principes physiques de la mécanique newtonienne : elle les complète et les contient.

11. Mais il n'existe pas de démonstrations simples à ma connaissance.

Chapitre 44

Exemples avec Sage

Ce chapitre est un fourre-tout de choses que l'on peut faire avec Sage.

44.0.1 Graphiques

Pour afficher le graphe d'une fonction, vous pouvez faire

```
+-----+
| SageMath version 8.1, Release Date: 2017-12-07          |
| Type "notebook()" for the browser-based notebook interface. |
| Type "help()" for help.                                |
+-----+
sage: plot(cos(x),0,5)
Launched png viewer for Graphics object consisting of 1 graphics primitive
sage: f(x)=sin(x)
sage: f.plot(-pi,pi)
Launched png viewer for Graphics object consisting of 1 graphics primitive
```

Un programme externe se lance automatiquement pour afficher le graphique que vous avez demandé.

Il se peut qu'aucun programme ne se lance et vous ayez, au lieu de `Launched png viewer for Graphics object ...` uniquement `Created graphics object ...`. Disons pour faire court que Sage a produit un `png` et qu'il ne sait pas quel programme externe utiliser pour l'afficher.

La solution est à l'adresse <http://doc.sagemath.org/html/en/reference/misc/sage/misc/viewer.html>

44.0.2 Autres

Dans le but d'automatiser certaines tâches, j'ai écrit ce module, nommé `outilsINGE.sage`, dans le cadre d'un cours de première année donné à des ingénieurs. Certaines des fonctions définies ici sont utilisées dans les exemples qui suivent.

```
1 # -*- coding: utf8 -*-
2 from sage.all import *
3
4 """
5 This module provides _pragmatic_ tools for solving exercise for
6 a first year in general mathematics.
7 """
8
9 # TODO : trouver une bonne traduction pour "point de selle."
10
```

```

11 def automatedVar(symbol,n):
12     """ If symbol = "x" and n=4, return the string 'x1,x2,x3,x4' ←
13         """
14     s = ",".join([ symbol+str(i) for i in range(1,n+1)])
15     return s
16
17 class SolveLinearSystem(object):
18     """
19     Solve Ax=v and print it in a nice way
20
21     Example :
22
23     A=matrix([ [1,-2,3,-2,0],[3,-7,-2,4,0],[4,3,5,2,0] ])
24     v=vector((0,0,0,0,0))
25     print SolveLinearSystem(A,v)
26     """
27     def __init__(self,A,v):
28         self.matrix = A
29         self.vector = v
30         self.nvars = A.ncols()
31         s = automatedVar("x",self.nvars)
32         self.xx=var(s)
33     def equations(self):
34         """Return the equations corresponding to the
35             self.matrix and self.vector as a list of equations←
36             """
37         X=matrix( [self.xx[i] for i in range(0,self.nvars) ] ).←
38             transpose()
39         eqs=[]
40         for i in range(0,self.matrix.nrows()):
41             exp = (self.matrix*X)[i][0]==self.vector[i]
42             eqs.append(exp)
43         return eqs
44     def solutions(self):
45         return solve(self.equations(),self.xx)
46     def latex(self):
47         """Return the LaTeX's code of the system."""
48         a=[]
49         a.append(r"""
50
51             \item
52             $
53             \left\{
54             \begin{array}{ll}
55             """ )
56         for eq in self.equations():
57             a.append(" "+str(eq).replace("x","x_").replace("*"←
58                 ,").replace("==","=")+"\\\\ \n")
59         a.append(r"""
60             \end{array}
61             \right.
62             $
63             """ )

```

```

60     return "".join(a)
61 def __str__(self):
62     a = []
63     a.append("The given matrix corresponds to the system")
64     for eq in self.equations():
65         a.append(str(eq))
66     a.append("And the solutions are")
67     a.append(str(self.solutions()))
68     return "\n".join(a)
69
70 def QuadraticMap(A,v):
71     """
72     Return the result of the quadratic form associated
73         with A applied on the vector v, that is the number
74      $A_{ij} v^i v^j$ 
75     using the summation convention.
76     """
77     n = A.nrows()
78     if not A.is_symmetric():
79         print "Warning : Given matrix is not symmetric"
80     if not A.is_square():
81         raise TypeError,"Error : The matrix A is not square"
82     if not v.degree()==n :
83         raise TypeError,"The size do not agree"
84     return sum([ A[i,j]*v[i]*v[j] for i in range(n) for j in range(←
85         n) ] ).simplify_full()
86
87 class SymmetricMatrix(object):
88     """
89     Provide informations about the matrix A assuming it is symmetric←
90     """
91     def __init__(self,A):
92         if not A.is_square():
93             print "Error : A symmetric matrix must be square"
94             raise TypeError
95         self.matrix = A
96         self.degree = A.nrows()
97         self.matrix.set_immutable()
98     def primary_principal_submatrix(self,n):
99         """
100         Return the primary principal submatrix of order n, that is the←
101             matrix obtained
102             by removing the n last lines and columns from self←
103             """
104         taille=self.degree-n
105         v=[]
106         for i in range(0,taille):
107             v.append(self.matrix[i][0:taille])
108         return matrix(v)
109     def principal_minors(self):
110         """

```



```

109     Return the list of principal minors. The principal minor of ←
        order k is
110         the determinant of the primary principal matrix of ←
            order k.
111     """
112     a=[]
113     for i in range(self.degree):
114         a.append(self.primary_principal_submatrix(i).determinant())
115     return a
116 def genre_list(self):
117     """
118     Return the genus of the matrix as a list of booleans in the ←
        order
119     positive defined, negative defined;
120         semidefinite positive, semidefinite negative, ←
            indefinite.
121
122     """
123     defpos = True
124     defneg = True
125     semidefpos = True
126     semidefneg = True
127     indefinie=True
128     mineurs = self.principal_minors()
129     for i in range(len(mineurs)):
130         m = mineurs[i]
131         if m == 0:
132             defneg=False
133             defpos=False
134         if m < 0:
135             defpos=False
136             semidefpos=False
137             if i%2==0:
138                 defneg=False
139         if m > 0:
140             semidefneg=False
141             if i%2==1:
142                 defneg=False
143     if 0 not in mineurs:
144         semidefneg=False
145         semidefpos=False
146     if (defpos==True) or (defneg==True) or (semidefpos==True) or (←
        semidefneg==True): indefinie=False
147     return [defpos,defneg,semidefpos,semidefneg,indefinie]
148 def __str__(self):
149     return str(self.matrix)
150
151 class QuadraticForm(SymmetricMatrix):
152     """
153     From a symmetric matrix A, provide informations concerning the ←
        associated quadratic form.
154     """
155     def __init__(self,A):

```

```

156     SymmetricMatrix.__init__(self,A)
157     if not A.is_symmetric():
158         print "Warning : matrix is not symmetric"
159 def evaluate(self,v):
160     """
161     Return the value of the quadratic form on the vector v.
162     """
163     return QuadraticMap(self.matrix,v)
164 def diagonalizing_martrix(self):
165     """
166     Return the matrix B such that  $B^tAB$  is diagonal.
167     """
168     # The transposition is because, in the matrix B, the ←
169     # eigenvectors have
170     # to be read as column while Sage's matrix constructor takes ←
171     # rows.
172     return matrix(self.orthonormal_basis()).transpose()
173 def new_variables(self):
174     """
175     Give the change of variables needed to put the quadratic form ←
176     under its normal form
177     X=BY
178     where X are the "old" variables
179     """
180     variables = var(automatedVar("y",self.degree))
181     Y = vector(variables)
182     return self.diagonalizing_martrix()*Y
183 def eigenmatrix_left(self):
184     return self.matrix.eigenmatrix_left()
185 def eigenvectors(self):
186     """
187     Return a list of eigenvectors of the matrix.
188     """
189     As the matrix is symmetric, that list has to be a basis.
190     """
191     D,P = self.eigenmatrix_left()
192     return [P[i] for i in range(P.nrows())]
193 def eigenvalues(self):
194     """
195     Return a list of eigenvalues of the matrix in the same order ←
196     as the list of eigenvectors given in
197     self.eigenvectors()
198     """
199     D,P = self.eigenmatrix_left()
200     return [ D[i,i] for i in range(D.nrows()) ]
201 def orthonormal_basis(self):
202     """
203     Return a basis of eigenvectors normalised to 1 as a list.
204     """
205     M,mu = matrix(self.eigenvectors()).gram_schmidt()
206     return [ v/v.norm() for v in M ]
207 def verification(self):
208     """

```

```

205     return the value of the quadratic form on the vector ←
           new_variables()
206     """
207     return self.evaluate(self.new_variables())
208
209 def __str__(self):
210     a = []
211     a.append("Hi guy; I'm the quadratic form associated with the ←
           matix")
212     a.append(str(self.matrix))
213     a.append("My eigenvalues and eigenvectors are : ")
214     veps = self.eigenvectors()
215     vaps = self.eigenvalues()
216     for i in range(len(veps)):
217         a.append("%s -> %s"%(str(vaps[i]),str(veps[i])))
218     a.append("I've the following orthonormal basis of eigenvectors←
           :")
219     for v in self.orthonormal_basis():
220         a.append(str(v))
221     a.append("A matrix B such that B^tAB is diagonal is ")
222     a.append(str(self.diagonalizing_martrix()))
223     a.append("I'm quite pretty in the following variables ...")
224     for i in range(self.degree):
225         a.append("x%s = %s"%(str(i+1),str(self.new_variables()[i]))←
           )
226     a.append("Look at me when I wear my cool variables")
227     a.append(str(self.verification()))
228     return "\n".join(a)
229
230 class Extrema(object):
231     """
232     From a function f, provides the informations for the study of ←
           the extrema :
233     partial derivative
234     critical points
235     Hessian matrix at the critical points
236     Genius of the Hessian and conclusion as local min/max
237
238     Dear student : remember that this class does not furnish any ←
           informations
239         concerning *global* extrema. The latter have to be found
240         among the critical points OR on the border of the domain.
241     """
242     def __init__(self,f):
243         var('x,y')
244         self.fun = f
245         self.gx=self.fun.diff(x).full_simplify()
246         self.gy=self.fun.diff(y).full_simplify()
247         self.gxx=self.gx.diff(x).simplify_full()
248         self.gxy=self.gx.diff(y).full_simplify()
249         self.gyy=self.gy.diff(y).full_simplify()
250         self.cp = solve( [self.gx(x,y)==0,self.gy(x,y)==0],[x,y] )
251     def critical_points(self):

```

```

252     """
253     Return the critical points as a list of tuples (x,y)
254     """
255     a = []
256     for pt in self.cp :
257         try :
258             px = SR(pt[0].rhs())
259             py = SR(pt[1].rhs())
260             a.append((px,py))
261         except TypeError :
262             a.append(" I'm not able to solve these equations.")
263     return a
264 def hessienne(self,a,b):
265     return matrix(SR,2,2,[self.gxx(a,b),self.gxy(a,b),self.gxy(a,←
        b),self.gyy(a,b)])
266 def __str__(self):
267     a = []
268     a.append ("The function :")
269     a.append(str(self.fun))
270     a.append ("Derivative x and y :")
271     a.append(str(self.gx))
272     a.append(str(self.gy))
273     a.append ("Hessian matrix :")
274     a.append(str(self.hessienne(x,y)))
275     a.append ("Critical points :")
276     for pt in self.critical_points() :
277         a.append(str(pt))
278     for pt in self.critical_points():
279         try :
280             px = pt[0]
281             py = pt[1]
282             a.append("At (%s,%s), the Hessian is"%(str(px),str(py)))
283         try :
284             Hess = SymmetricMatrix(self.hessienne(px,py))
285             for l in Hess.matrix:
286                 a.append(" "+str(l))
287             a.append(" Primary principal minors are %s"%str(Hess.←
                principal_minors()))
288             l = Hess.genre_list()
289             if l[0]==True:
290                 a.append(" Hessian positive defined")
291                 a.append(" local minimum")
292             if l[1]==True:
293                 a.append(" Hessian negative defined")
294                 a.append(" local maximum")
295             if l[2]==True:
296                 a.append(" Hessian positive semidéfinite")
297                 a.append(" I don't conclude")
298             if l[3]==True:
299                 a.append(" Hessian negative semidefinite")
300                 a.append(" I don't conclude")
301             if l[4]==True:
302                 a.append(" Undefined Hessian")

```

```

303         a.append(" «selle» point")
304     except RuntimeError,data :
305         a.append(" "+str(data))
306     except TypeError :
307         a.append(" I'm not able to solve these equations.")
308     return "\n".join(a)

```

tex/sage/outilsINGE.sage

Exemple 44.1

Calculer la limite

$$\lim_{x \rightarrow \infty} \frac{\sin(x) \cos(x)}{x} \quad (44.1)$$

```

var('x')
f(x)=sin(x)*cos(x)/x
limit(f(x),x=oo)

```

La première ligne déclare que la lettre x désignera une variable. Pour la troisième ligne, notez que l'infini est écrit par deux petits « o ». △

Exemple 44.2

Quelques limites et graphes avec Sage.

(1) $\lim_{x \rightarrow 0} \frac{\sin(ax)}{\sin(bx)}$.

Pour effectuer cet exercice avec Sage, il faut taper les lignes suivantes :

```

sage: var('x,a,b')
(x, a, b)
sage: f(x)=sin(a*x)/sin(b*x)
sage: limit( f(x),x=0 )
a/b

```

Noter qu'il faut déclarer les variables x , a et b .

(2) $\lim_{x \rightarrow \pm\infty} \frac{\sqrt{x^2+1}-x}{x-2}$

```

sage: f(x)=(sqrt(x**2+1))/(x-2)
sage: limit(f(x),x=oo)
1
sage: limit(f(x),x=-oo)
-1

```

Noter la commande pour la racine carré : `sqrt`. Étant donné que cette fonction diverge en $x = 2$, si nous voulons la tracer, il faut procéder en deux fois :

```

sage: plot(f,(-100,1.9))
Launched png viewer for Graphics object consisting of 1 graphics primitive
sage: plot(f,(2.1,100))
Launched png viewer for Graphics object consisting of 1 graphics primitive

```

La première ligne trace de -100 à 1.9 et la seconde de 2.1 à 100 . Ces graphiques vous permettent déjà de voir les limites. Attention : ils ne sont pas des *preuves* ! Mais ils sont de sérieux indices qui peuvent vous inspirer dans vos calculs.



Exemple 44.3

Calculer les dérivées partielles $\partial_x f$, $\partial_y f$, $\partial_x^2 f$, $\partial_{xy}^2 f$, $\partial_{yx}^2 f$ et $\partial_y^2 f$ des fonctions suivantes.

- | | |
|---------------------------|------------------------|
| (1) $2x^3 + 3x^2y - 2y^2$ | (3) $\tan(x/y)$ |
| (2) $\ln(xy^2)$ | (4) $\frac{xy^2}{x+y}$ |

Le script Sage suivant (exoDV002.sage) résout l'exercice :

```

1 # -*- coding: utf8 -*-
2
3 def LesCalculs(f):
4     print "Pour la fonction %s"%str(f)
5     print "d_x",f.diff(x).simplify_full()
6     print "d_y",f.diff(y).simplify_full()
7     print "d^2_x",f.diff(x).diff(x).simplify_full()
8     print "d_xd_y",f.diff(x).diff(y).simplify_full()
9     print "d_yd_x",f.diff(y).diff(x).simplify_full()
10    print "d^2_y",f.diff(y).diff(y).simplify_full()
11    print ""
12
13 def exercice_DV002():
14     var('x,y')
15     fa(x,y)=2*x**3+3*x**2*y-2*y**2
16     fb(x,y)=ln(x*y**2)
17     fc(x,y)=tan(x/y)
18     fd(x,y)=x*y**2/(x+y)
19     LesCalculs(fa)
20     LesCalculs(fb)
21     LesCalculs(fc)
22     LesCalculs(fd)

```

tex/sage/exoDV002.sage

La sortie est :

```

Pour la fonction (x, y) |--> 2*x^3 + 3*x^2*y - 2*y^2
d_x (x, y) |--> 6*x^2 + 6*x*y
d_y (x, y) |--> 3*x^2 - 4*y
d^2_x (x, y) |--> 12*x + 6*y
d_xd_y (x, y) |--> 6*x
d_yd_x (x, y) |--> 6*x
d^2_y (x, y) |--> -4

Pour la fonction (x, y) |--> log(x*y^2)
d_x (x, y) |--> 1/x
d_y (x, y) |--> 2/y
d^2_x (x, y) |--> -1/x^2
d_xd_y (x, y) |--> 0
d_yd_x (x, y) |--> 0
d^2_y (x, y) |--> -2/y^2

```

```

Pour la fonction (x, y) |--> tan(x/y)
d_x (x, y) |--> 1/(y*cos(x/y)^2)
d_y (x, y) |--> -x/(y^2*cos(x/y)^2)
d^2_x (x, y) |--> 2*sin(x/y)/(y^2*cos(x/y)^3)
d_xd_y (x, y) |--> -(2*x*sin(x/y) + y*cos(x/y))/(y^3*cos(x/y)^3)
d_yd_x (x, y) |--> -(2*x*sin(x/y) + y*cos(x/y))/(y^3*cos(x/y)^3)
d^2_y (x, y) |--> 2*(x^2*sin(x/y) + x*y*cos(x/y))/(y^4*cos(x/y)^3)

Pour la fonction (x, y) |--> x*y^2/(x + y)
d_x (x, y) |--> y^3/(x^2 + 2*x*y + y^2)
d_y (x, y) |--> (2*x^2*y + x*y^2)/(x^2 + 2*x*y + y^2)
d^2_x (x, y) |--> -2*y^3/(x^3 + 3*x^2*y + 3*x*y^2 + y^3)
d_xd_y (x, y) |--> (3*x*y^2 + y^3)/(x^3 + 3*x^2*y + 3*x*y^2 + y^3)
d_yd_x (x, y) |--> (3*x*y^2 + y^3)/(x^3 + 3*x^2*y + 3*x*y^2 + y^3)
d^2_y (x, y) |--> 2*x^3/(x^3 + 3*x^2*y + 3*x*y^2 + y^3)

```

△

Exemple 44.4

Résoudre les systèmes suivants.

$$\begin{array}{ll}
 (1) \begin{cases} x_1 - 2x_2 + 3x_3 - 2x_4 = 0 \\ 3x_1 - 7x_2 - 2x_3 + 4x_4 = 0 \\ 4x_1 + 3x_2 + 5x_3 + 2x_4 = 0 \end{cases} & (8) \begin{cases} x_1 - 2x_2 + 3x_3 - 2x_4 = 0 \\ 3x_1 - 7x_2 - 2x_3 + 4x_4 = 0 \\ 4x_1 + 3x_2 + 5x_3 + 2x_4 = 0 \end{cases} \\
 (2) \begin{cases} 2x_1 + x_2 - 2x_3 + 3x_4 = 0 \\ 3x_1 + 2x_2 - x_3 + 3x_4 = 4 \\ 3x_1 + 3x_2 + 3x_3 - 3x_4 = 9 \end{cases} & (9) \begin{cases} 2x_1 + x_2 - 2x_3 + 3x_4 = 0 \\ 3x_1 + 2x_2 - x_3 + 3x_4 = 4 \\ 3x_1 + 3x_2 + 3x_3 - 3x_4 = 9 \end{cases} \\
 (3) \begin{cases} x_1 + 2x_2 - 3x_3 = 0 \\ 2x_1 + 5x_2 + 2x_3 = 0 \\ 3x_1 - x_2 - 4x_3 = 0 \end{cases} & (10) \begin{cases} x_1 + 2x_2 - 3x_3 = 0 \\ 2x_1 + 5x_2 + 2x_3 = 0 \\ 3x_1 - x_2 - 4x_3 = 0 \end{cases} \\
 (4) \begin{cases} x_1 + 2x_2 - x_3 = 0 \\ 2x_1 + 5x_2 + 2x_3 = 0 \\ x_1 + 4x_2 + 7x_3 = 0 \\ x_1 + 3x_2 + 3x_3 = 0 \end{cases} & (11) \begin{cases} x_1 + 2x_2 - x_3 = 0 \\ 2x_1 + 5x_2 + 2x_3 = 0 \\ x_1 + 4x_2 + 7x_3 = 0 \\ x_1 + 3x_2 + 3x_3 = 0 \end{cases} \\
 (5) \begin{cases} x_1 + x_2 + x_3 + x_4 = 0 \\ x_1 + x_2 + x_3 - x_4 = 4 \\ x_1 + x_2 - x_3 + x_4 = -4 \\ x_1 - x_2 + x_3 + x_4 = 2 \end{cases} & (12) \begin{cases} x_1 + x_2 + x_3 + x_4 = 0 \\ x_1 + x_2 + x_3 - x_4 = 4 \\ x_1 + x_2 - x_3 + x_4 = -4 \\ x_1 - x_2 + x_3 + x_4 = 2 \end{cases} \\
 (6) \begin{cases} x_1 + 3x_2 + 3x_3 = 1 \\ x_1 + 3x_2 + 4x_3 = 0 \\ x_1 + 4x_2 + 3x_3 = 3 \end{cases} & (13) \begin{cases} x_1 + 3x_2 + 3x_3 = 1 \\ x_1 + 3x_2 + 4x_3 = 0 \\ x_1 + 4x_2 + 3x_3 = 3 \end{cases} \\
 (7) \begin{cases} x_1 - 3x_2 + 2x_3 = -6 \\ -3x_1 + 3x_2 - x_3 = 17 \\ 2x_1 - x_2 = 3 \end{cases} & (14) \begin{cases} x_1 - 3x_2 + 2x_3 = -6 \\ -3x_1 + 3x_2 - x_3 = 17 \\ 2x_1 - x_2 = 3 \end{cases}
 \end{array}$$

Nous résolvons les systèmes en utilisant Sage avec le script suivant.

```

1 # -*- coding: utf8 -*-
2 """

```

```

3 Ce script Sage résout un certain nombre
4 de systèmes d'équations linéaires du cours INGE1121
5 """
6
7 import outilsINGE
8
9 def exercice_1_1_bcdefhi():
10     # Exercice 1.1.b (INGE1121)
11     A=matrix([ [1,-2,3,-2,0],[3,-7,-2,4,0],[4,3,5,2,0] ])
12     v=vector((0,0,0,0,0))
13     print outilsINGE.SolveLinearSystem(A,v)
14     # Exercice 1.1.c (INGE1121)
15     A=matrix([ [2,1,-2,3],[3,2,-1,3],[3,3,3,-3] ])
16     v=vector((0,4,9))
17     print outilsINGE.SolveLinearSystem(A,v)
18     # Exercice 1.1.d (INGE1121)
19     A=matrix([ [1,2,-3],[2,5,2],[3,-1,-4] ])
20     v=vector((0,0,0))
21     print outilsINGE.SolveLinearSystem(A,v)
22     # Exercice 1.1.e (INGE1121)
23     A=matrix([ [1,2,-1],[2,5,2],[1,4,7],[1,3,3] ])
24     v=vector((0,0,0,0))
25     print outilsINGE.SolveLinearSystem(A,v)
26     # Exercice 1.1.f (INGE1121)
27     A=matrix([ [1,1,1,1],[1,1,1,-1],[1,1,-1,1],[1,-1,1,1] ])
28     v=vector((0,4,-4,2))
29     print outilsINGE.SolveLinearSystem(A,v)
30     # Exercice 1.1.h (INGE1121)
31     A=matrix([ [1,3,3],[1,3,4],[1,4,3] ])
32     v=vector((1,0,3))
33     print outilsINGE.SolveLinearSystem(A,v)
34     # Exercice 1.1.i (INGE1121)
35     A=matrix([ [1,-3,2],[-3,3,-1],[2,-1,0] ])
36     v=vector((-6,17,3))
37     print outilsINGE.SolveLinearSystem(A,v)

```

tex/sage/exo11.sage

Le résultat est le suivant :

The given matrix corresponds to the system

$$x_1 - 2x_2 + 3x_3 - 2x_4 == 0$$

$$3x_1 - 7x_2 - 2x_3 + 4x_4 == 0$$

$$4x_1 + 3x_2 + 5x_3 + 2x_4 == 0$$

And the solutions are

[

$$[x_1 == -23/16*r_{19}, x_2 == -5/16*r_{19}, x_3 == 15/16*r_{19}, x_4 == r_{19}, x_5 == r_{18}]$$

]

The given matrix corresponds to the system

$$2x_1 + x_2 - 2x_3 + 3x_4 == 0$$

$$3x_1 + 2x_2 - x_3 + 3x_4 == 4$$

$$3x_1 + 3x_2 + 3x_3 - 3x_4 == 9$$

And the solutions are

[

$$[x_1 == 3*r_{20} - 7, x_2 == -4*r_{20} + 11, x_3 == r_{20}, x_4 == 1]$$


```

]
The given matrix corresponds to the system
x1 + 2*x2 - 3*x3 == 0
2*x1 + 5*x2 + 2*x3 == 0
3*x1 - x2 - 4*x3 == 0
And the solutions are
[
[x1 == 0, x2 == 0, x3 == 0]
]
The given matrix corresponds to the system
x1 + 2*x2 - x3 == 0
2*x1 + 5*x2 + 2*x3 == 0
x1 + 4*x2 + 7*x3 == 0
x1 + 3*x2 + 3*x3 == 0
And the solutions are
[
[x1 == 9*r21, x2 == -4*r21, x3 == r21]
]
The given matrix corresponds to the system
x1 + x2 + x3 + x4 == 0
x1 + x2 + x3 - x4 == 4
x1 + x2 - x3 + x4 == -4
x1 - x2 + x3 + x4 == 2
And the solutions are
[
[x1 == 1, x2 == -1, x3 == 2, x4 == -2]
]
The given matrix corresponds to the system
x1 + 3*x2 + 3*x3 == 1
x1 + 3*x2 + 4*x3 == 0
x1 + 4*x2 + 3*x3 == 3
And the solutions are
[
[x1 == -2, x2 == 2, x3 == -1]
]
The given matrix corresponds to the system
x1 - 3*x2 + 2*x3 == -6
-3*x1 + 3*x2 - x3 == 17
2*x1 - x2 == 3
And the solutions are
[
[x1 == 37, x2 == 71, x3 == 85]
]

```

△

Exemple 44.5

Pour chacun des systèmes suivants $A \cdot X = B$,

- (1) Résoudre le système par échelonnement,
- (2) Calculer A^{-1} ,
- (3) Vérifier votre réponse en calculant $A^{-1}B$. Qu'êtes-vous censé obtenir ?

Les énoncés sont

(1)

$$A = \begin{pmatrix} 2 & 1 & -2 \\ 3 & 2 & 2 \\ 5 & 4 & 3 \end{pmatrix}, \quad B = \begin{pmatrix} 10 \\ 1 \\ 4 \end{pmatrix} \quad (44.2)$$

Nous utilisons Sage pour fournir la réponse. Le code suivant résout le système et donne l'inverse de la matrice :

```

1 # -*- coding: utf8 -*-
2
3 import outilsINGE
4
5 def exercise_1_3():
6     A=matrix([[2,1,-2],[3,2,2],[5,4,3]])
7     v=vector((10,1,4))
8     print outilsINGE.SolveLinearSystem(A,v)
9     print "Matrice inverse :"
10    print A.inverse()

```

tex/sage/exo13.sage

La sortie est ici :

The given matrix corresponds to the system

$$2x_1 + x_2 - 2x_3 == 10$$

$$3x_1 + 2x_2 + 2x_3 == 1$$

$$5x_1 + 4x_2 + 3x_3 == 4$$

And the solutions are

```
[
[x1 == 1, x2 == 2, x3 == -3]
]
```

Matrice inverse :

$$\begin{bmatrix} 2/7 & 11/7 & -6/7 \end{bmatrix}$$

$$\begin{bmatrix} -1/7 & -16/7 & 10/7 \end{bmatrix}$$

$$\begin{bmatrix} -2/7 & 3/7 & -1/7 \end{bmatrix}$$

△

Exemple 44.6

Sachant que $(-1, 0, 1, 0)$ est un vecteur propre de la matrice

$$A = \begin{pmatrix} 2 & 1 & -1 & 1 \\ 1 & 0 & 1 & 1 \\ -1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix} \quad (44.3)$$

- (1) Diagonaliser A au moyen d'une matrice orthogonale
- (2) Écrire la forme quadratique $X^t A X$ sous forme d'une somme pondérée de carrés.

Calculons Av afin de savoir la valeur propre associée au vecteur donné :

$$\begin{pmatrix} 2 & 1 & -1 & 1 \\ 1 & 0 & 1 & 1 \\ -1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} -1 \\ 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} -3 \\ 0 \\ 3 \\ 0 \end{pmatrix}. \quad (44.4)$$

La valeur propre est donc 3. Nous savons donc que $(\lambda - 3)$ pourra être factorisé dans le polynôme caractéristique.

Pour le reste de l'exercice c'est standard et c'est résolu de la façon suivante :

```

1 # -*- coding: utf8 -*-
2
3 import outilsINGE
4
5 def exercise_6_5():
6     A=matrix(QQ,4,4,[2,1,-1,1,1,0,1,1,-1,1,2,1,1,1,1,0])
7     x=outilsINGE.QuadraticForm(A)
8     print x

```

tex/sage/exo65.sage

qui retourne

Hi guy; I'm the quadratic form associated with the matix

```

[ 2  1 -1  1]
[ 1  0  1  1]
[-1  1  2  1]
[ 1  1  1  0]

```

My eigenvalues and eigenvectors are :

```

3 -> (1, 0, -1, 0)
3 -> (0, 1, 2, 1)
-1 -> (1, 0, 1, -2)
-1 -> (0, 1, 0, -1)

```

I've the following orthonormal basis of eigenvectors :

```

(1/2*sqrt(2), 0, -1/2*sqrt(2), 0)
(1/2, 1/2, 1/2, 1/2)
(1/6*sqrt(6), 0, 1/6*sqrt(6), -1/3*sqrt(6))
(-1/2*sqrt(1/3), 3/2*sqrt(1/3), -1/2*sqrt(1/3), -1/2*sqrt(1/3))

```

A matrix B such that B^tAB is diagonal is

```

[  1/2*sqrt(2)      1/2    1/6*sqrt(6) -1/2*sqrt(1/3)]
[          0         1/2          0  3/2*sqrt(1/3)]
[ -1/2*sqrt(2)      1/2    1/6*sqrt(6) -1/2*sqrt(1/3)]
[          0         1/2   -1/3*sqrt(6) -1/2*sqrt(1/3)]

```

I'm quite pretty in the following variables ...

```

x1 = -1/2*sqrt(1/3)*y4 + 1/2*sqrt(2)*y1 + 1/6*sqrt(6)*y3 + 1/2*y2
x2 = 3/2*sqrt(1/3)*y4 + 1/2*y2
x3 = -1/2*sqrt(1/3)*y4 - 1/2*sqrt(2)*y1 + 1/6*sqrt(6)*y3 + 1/2*y2
x4 = -1/2*sqrt(1/3)*y4 - 1/3*sqrt(6)*y3 + 1/2*y2

```

Look at me when I wear my cool variables

```

3*y1^2 + 3*y2^2 - y3^2 - y4^2

```

△

Exemple 44.7

Rechercher les extrema des fonctions suivantes

(1) $f(x, y) = 2 - \sqrt{x^2 + y^2}$

(2) $f(x, y) = x^3 + 3xy^2 - 15x - 12y$

$$(3) f(x, y) = \frac{x^3}{3} + \frac{4y^3}{3} - x^2 - 3x - 4y - 3$$

Les corrigés sont créés par le script Sage `exo101.sage`

```
# -*- coding: utf8 -*-
```

```
import outilsINGE
```

```
def exercice_10_1_A():
    var('x,y')
    f(x,y)=2-sqrt(x**2+y**2)
    print outilsINGE.Extrema(f)
def exercice_10_1_B():
    var('x,y')
    f(x,y)=x**3+3*x*y**2-15*x-12*y
    print outilsINGE.Extrema(f)
def exercice_10_1_C():
    var('x,y')
    f(x,y)=x**3/3+4*y**3/3-x**2-3*x-4*y-3
    print outilsINGE.Extrema(f)
```

Des réponses :

(1) The function :

$$(x, y) \mapsto -\sqrt{x^2 + y^2} + 2$$

Derivative x and y :

$$(x, y) \mapsto -x/\sqrt{x^2 + y^2}$$

$$(x, y) \mapsto -y/\sqrt{x^2 + y^2}$$

Hessian matrix :

$$\begin{bmatrix} -\sqrt{x^2 + y^2} * y^2 / (x^4 + 2 * x^2 * y^2 + y^4) & x * y / (x^2 + y^2)^{(3/2)} \\ x * y / (x^2 + y^2)^{(3/2)} & -\sqrt{x^2 + y^2} * x^2 / (x^4 + 2 * x^2 * y^2 + y^4) \end{bmatrix}$$

Critical points :

$$(0, 0)$$

At (0,0), the Hessian is

`power::eval(): division by zero`

Ici nous voyons que Sage a du mal à calculer la matrice hessienne en (0,0). En effet, nous tombons sur une division par zéro. Pour résoudre l'exercice, il faut se rendre compte que la fonction $(x, y) \mapsto \sqrt{x^2 + y^2}$ est toujours positive et est nulle seulement au point (0,0). Donc f est toujours plus petite ou égale à deux tandis que $f(0,0) = 2$. Le point est donc un maximum global.

(2) The function :

$$(x, y) \mapsto x^3 + 3 * x * y^2 - 15 * x - 12 * y$$

Derivative x and y :

$$(x, y) \mapsto 3 * x^2 + 3 * y^2 - 15$$

$$(x, y) \mapsto 6 * x * y - 12$$

Hessian matrix :

$$\begin{bmatrix} 6 * x & 6 * y \\ 6 * y & 6 * x \end{bmatrix}$$

$$\begin{bmatrix} 6 * y & 6 * x \end{bmatrix}$$

Critical points :

$$(2, 1)$$

$$(1, 2)$$

$$(-1, -2)$$

$$(-2, -1)$$

At (2,1), the Hessian is

$$(12, 6)$$

(6, 12)
 Primary principal minors are [108, 12]
 Hessian positive defined
 local minimum
 At (1,2), the Hessian is
 (6, 12)
 (12, 6)
 Primary principal minors are [-108, 6]
 Undefinite Hessian
 «selle» point
 At (-1,-2), the Hessian is
 (-6, -12)
 (-12, -6)
 Primary principal minors are [-108, -6]
 Undefinite Hessian
 «selle» point
 At (-2,-1), the Hessian is
 (-12, -6)
 (-6, -12)
 Primary principal minors are [108, -12]
 Hessian negative defined
 local maximum

(3) The function :

$(x, y) \mapsto \frac{1}{3}x^3 - x^2 + \frac{4}{3}y^3 - 3x - 4y - 3$

Derivative x and y :

$(x, y) \mapsto x^2 - 2x - 3$

$(x, y) \mapsto 4y^2 - 4$

Hessian matrix :

$[2x - 2 \quad 0]$

$[\quad 0 \quad 8y]$

Critical points :

(3, 1)

(-1, 1)

(3, -1)

(-1, -1)

At (3,1), the Hessian is

(4, 0)

(0, 8)

Primary principal minors are [32, 4]

Hessian positive defined

local minimum

At (-1,1), the Hessian is

(-4, 0)

(0, 8)

Primary principal minors are [-32, -4]

Undefinite Hessian

«selle» point

At (3,-1), the Hessian is

(4, 0)

(0, -8)

Primary principal minors are [-32, 4]

Undefinite Hessian

«selle» point

At (-1,-1), the Hessian is
 (-4, 0)
 (0, -8)
 Primary principal minors are [32, -4]
 Hessian negative defined
 local maximum

△

Exemple 44.8

Déterminer les valeurs extrêmes et les points de selle des fonctions suivantes.

- (1) $f(x, y) = x^2 + 4x + y^2 - 2y$. (3) $f(x, y) = e^x \sin(y)$.
 (2) $f(x, y) = e^{x^2+xy}$.

Certains corrigés de cet exercice ont été réalisés par Sage. Le script utilisé est `exo103.sage`

```

1 # -*- coding: utf8 -*-
2
3 import outilsINGE
4
5 def exercise_10_3_A():
6     var('x,y')
7     f(x,y)=x**2+4*x+y**2-2*y
8     print outilsINGE.Extrema(f)
9
10 def exercise_10_3_H():
11     var('x,y')
12     f(x,y)=exp(x**2+x*y)
13     print outilsINGE.Extrema(f)
14
15 def exercise_10_3_Q():
16     var('x,y')
17     f(x,y)=exp(x)*sin(y)
18     print outilsINGE.Extrema(f)

```

tex/sage/exo103.sage

Des réponses :

- (1) The function :
 (x, y) |--> x² + y² + 4*x - 2*y
 Derivative x and y :
 (x, y) |--> 2*x + 4
 (x, y) |--> 2*y - 2
 Hessian matrix :
 [2 0]
 [0 2]
 Critical points :
 (-2, 1)
 At (-2,1), the Hessian is
 (2, 0)
 (0, 2)

Primary principal minors are [4, 2]
 Hessian positive defined
 local minimum

(2) The function :

(x, y) |--> e^(x^2 + x*y)

Derivative x and y :

(x, y) |--> (2*x*e^(x^2) + y*e^(x^2))*e^(x*y)

(x, y) |--> x*e^(x^2 + x*y)

Hessian matrix :

[(4*x*y*e^(x^2) + y^2*e^(x^2) + 2*(2*x^2 + 1)*e^(x^2))*e^(x*y)

[(x*y*e^(x^2) + (2*x^2 + 1)*e^(x^2))*e^(x*y)

Critical points :

(0, 0)

At (0,0), the Hessian is

(2, 1)

(1, 0)

Primary principal minors are [-1, 2]

Undefinite Hessian

«selle» point

(3) The function :

(x, y) |--> e^x*sin(y)

Derivative x and y :

(x, y) |--> e^x*sin(y)

(x, y) |--> e^x*cos(y)

Hessian matrix :

[e^x*sin(y) e^x*cos(y)]

[e^x*cos(y) -e^x*sin(y)]

Critical points :

I'm not able to solve these equations.

I'm not able to solve these equations.

At (,I), the Hessian is

I'm not able to solve these equations.

At (,I), the Hessian is

I'm not able to solve these equations.

Ici, Sage n'est pas capable de résoudre les équations qui annulent le jacobien. Les équations à résoudre sont pourtant faciles :

$$\begin{cases} e^x \cos(y) = 0 & (44.5a) \\ e^x \sin(y) = 0 & (44.5b) \end{cases}$$

Étant donné que l'exponentielle ne s'annule jamais, il faudrait avoir en même temps $\cos(y) = 0$ et $\sin(y) = 0$, ce qui est impossible. La fonction n'a donc aucun extrema local.

△

Exemple 44.9

Considérons la fonction

$$f(x, y) = xy^2 e^{-(x^2+y^2)/4}. \quad (44.6)$$

(1) Montrer qu'il y a une infinité de points critiques.

(2) Déterminer leur nature.

Voici la fonction Sage qui fournit les informations :

```

1 # -*- coding: utf8 -*-
2
3 import outilsINGE
4
5 def exercise_10_4():
6     var('x,y')
7     f(x,y)=x*y**2*exp(-(x**2+y**2)/4)
8     print outilsINGE.Extrema(f)

```

tex/sage/exo104.sage

La sortie est

The function :

$(x, y) \mapsto x*y^2*e^{-(1/4*x^2 - 1/4*y^2)}$

Derivative x and y :

$(x, y) \mapsto -1/2*(x^2 - 2)*y^2*e^{-(1/4*x^2 - 1/4*y^2)}$

$(x, y) \mapsto -1/2*(x*y^3 - 4*x*y)*e^{-(1/4*x^2 - 1/4*y^2)}$

Hessian matrix :

[
 $1/4*(x^3 - 6*x)*y^2*e^{-(1/4*x^2 - 1/4*y^2)}$ $1/4*((x^2 - 2)*y^3 - 4*(x^2 - 2)*y)*e^{-(1/4*x^2 - 1/4*y^2)}$
 $1/4*((x^2 - 2)*y^3 - 4*(x^2 - 2)*y)*e^{-(1/4*x^2 - 1/4*y^2)}$ $1/4*(x*y^4 - 10*x*y^2 + 4*y^2)*e^{-(1/4*x^2 - 1/4*y^2)}$

Critical points :

$(r17, 0)$

$(-\sqrt{2}, -2)$

$(\sqrt{2}, -2)$

$(-\sqrt{2}, 2)$

$(\sqrt{2}, 2)$

At $(r17,0)$, the Hessian is

$(0, 0)$

$(0, 2*r17*e^{-(1/4*r17^2)})$

Primary principal minors are $[0, 0]$

Hessian positive semidéfinitive

I don't conclude

Hessian negative semidefinite

I don't conclude

At $(-\sqrt{2},-2)$, the Hessian is

$(4*\sqrt{2}*e^{-3/2}, 0)$

$(0, 4*\sqrt{2}*e^{-3/2})$

Primary principal minors are $[32*e^{-3}, 4*\sqrt{2}*e^{-3/2}]$

Hessian positive defined

local minimum

At $(\sqrt{2},-2)$, the Hessian is

$(-4*\sqrt{2}*e^{-3/2}, 0)$

$(0, -4*\sqrt{2}*e^{-3/2})$

Primary principal minors are $[32*e^{-3}, -4*\sqrt{2}*e^{-3/2}]$

Hessian negative defined

local maximum

At $(-\sqrt{2},2)$, the Hessian is

$(4*\sqrt{2}*e^{-3/2}, 0)$

$(0, 4*\sqrt{2}*e^{-3/2})$

Primary principal minors are $[32*e^{-3}, 4*\sqrt{2}*e^{-3/2}]$

Hessian positive defined


```

local minimum
At (sqrt(2),2), the Hessian is
(-4*sqrt(2)*e^(-3/2), 0)
(0, -4*sqrt(2)*e^(-3/2))
Primary principal minors are [32*e^(-3), -4*sqrt(2)*e^(-3/2)]
Hessian negative defined
local maximum

```

Notez la présence de `r1` comme paramètres dans les solutions. Tous les points avec $y = 0$ sont des points critiques. Cependant, Sage¹ ne parvient pas à conclure la nature de ces points $(x, 0)$.

Notons que le nombre $f(x, y)$ a toujours le signe de x parce que y^2 et l'exponentielle sont positives. Toujours? En tout cas lorsque $x \neq 0$. Prenons un point $(a, 0)$ avec $a > 0$. Dans un voisinage de ce point, nous avons $f(x, y) > 0$ parce que si $a > 0$, alors $x > 0$ dans un voisinage de a . Le point $(a, 0)$ est un minimum local parce que $0 = f(a, 0) \leq f(x, y)$ pour tout (x, y) dans un voisinage de $(a, 0)$.

De la même façon, les points $(a, 0)$ avec $a < 0$ sont des maxima locaux parce que dans un voisinage, la fonction est négative.

Le point $(0, 0)$ n'est ni maximum ni minimum local. C'est un point de selle.

△

Exemple 44.10

Dériver les fonctions suivantes.

(1) $\sin(\ln(x))$

(2) $\frac{\sin x}{x}$;

(3) e^{x^2}

(4) $\cos(x)^{\sin(x)}$

Le programme suivant par Sage résout l'exercice :

```

1  #! /usr/bin/sage -python
2  # -*- coding: utf8 -*-
3
4  from sage.all import *
5
6  var('x')
7  f=sin(ln(x))
8  print f.diff(x)
9  f=sin(x)/x
10 print f.diff(x)
11 f=exp(x**2)
12 print f.diff(x)
13 f=cos(x)**(sin(x))
14 print f.diff(x)

```

tex/sage/corrDerive_0002.sage

Le résultat est :

```

cos(log(x))/x
cos(x)/x - sin(x)/x^2

```

1. ou, plus précisément, le programme que j'ai écrit avec Sage.

```
2*x*e^(x^2)
(log(cos(x))*cos(x) - sin(x)^2/cos(x))*cos(x)^sin(x)
```

△

Exemple 44.11

Donner une approximation de $\ln(1.0001)$.

```
-----
| Sage Version 4.5.3, Release Date: 2010-09-04           |
| Type notebook() for the GUI, and license() for information. |
-----
```

```
sage: numerical_approx(ln(1.0001))
0.0000999950003332973
```

△

Chapitre 45

Épilogue : la constante de Weiner

Nous voici à la fin du Frido. Nous avons étudié beaucoup de math, et beaucoup reste à voir. En guise de conclusion, je voudrais vous parler de la constante de Weiner, introduite dans [505]. Il s'agit d'une constante, qui comme π ou e intervient dans à peu près tous les domaines de la mathématique.

Comme toujours, il existe énormément de définitions équivalentes différentes ; nous choisissons celle-ci, motivée par le le théorème de Weinersmith 28.42.

Définition 45.1.

La **constante de Weiner** W_c est l'unique entier p tel que l'espace L^p soit un espace de Hilbert.

Cette constante intervient de façon centrale dans de nombreux résultats dans tous les domaines ; nous en citons quelques-uns.

- (1) La moyenne de tout couple de réels peut être calculée en divisant leur somme par la constante de Weiner¹.
- (2) La constante de Weiner donne l'indice du groupe alterné dans le groupe symétrique pour tous les ordres, théorème 5.30.
- (3) La constante de Weiner donne une borne inférieure optimale pour l'ensemble des nombres premiers.
- (4) L'unique point fixe non trivial de la fonction factorielle est la constante de Weiner.
- (5) Tout automorphisme d'anneau a un polynôme minimal dont le degré est donné par la constante de Weiner.
- (6) L'égalité $ab = 0$ dans un anneau n'implique pas spécialement $a = 0$ ou $b = 0$ lorsque la caractéristique de l'anneau est égale à la constante de Weiner, et seulement dans ce cas.
- (7) Pour l'anecdote, la constante de Weiner donne le rapport τ/π ; elle est aussi la partie entière de e .

Il reste encore de nombreuses conjectures mettant en valeur la constante de Weiner :

- (1) La fameuse droite critique de la conjecture de Riemann est donnée par l'inverse de la constante de Weiner.
- (2) Soit \mathcal{P} l'ensemble de nombres premiers. Est-ce que l'ensemble

$$\{(p, q) \in \mathcal{P} \times \mathcal{P} \text{ tel que } 0 < |p - q| = W_c\} \quad (45.1)$$

est fini ?

D'aucuns pourraient objecter que tout cela n'est que fantaisie et trivialité. Il n'en est rien. La preuve que la constante de Weiner est centrale en mathématique est précisément qu'elle avait déjà un nom et un symbole réservé bien avant le début de l'histoire des mathématiques.

Le fait est que toutes les mathématiques que vous connaissez se basent sur les nombres entiers ; cela n'est pas du tout une trivialité.

1. C'est historiquement la première propriété énoncée de la constante de Weiner ; elle suggère également une notion de constante de Weiner généralisée pour moyennner un nombre arbitraire de nombres. La construction des nombres de Weiner généralisés est en projet dans la section 1.2.

Chapitre 46

GNU Free Documentation License

Version 1.3, 3 November 2008

Copyright © 2000, 2001, 2002, 2007, 2008 Free Software Foundation, Inc.

<http://fsf.org/>

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Preamble

The purpose of this License is to make a manual, textbook, or other functional and useful document “free” in the sense of freedom : to assure everyone the effective freedom to copy and redistribute it, with or without modifying it, either commercially or noncommercially. Secondly, this License preserves for the author and publisher a way to get credit for their work, while not being considered responsible for modifications made by others.

This License is a kind of “copyleft”, which means that derivative works of the document must themselves be free in the same sense. It complements the GNU General Public License, which is a copyleft license designed for free software.

We have designed this License in order to use it for manuals for free software, because free software needs free documentation : a free program should come with manuals providing the same freedoms that the software does. But this License is not limited to software manuals ; it can be used for any textual work, regardless of subject matter or whether it is published as a printed book. We recommend this License principally for works whose purpose is instruction or reference.

APPLICABILITY AND DEFINITIONS

This License applies to any manual or other work, in any medium, that contains a notice placed by the copyright holder saying it can be distributed under the terms of this License. Such a notice grants a world-wide, royalty-free license, unlimited in duration, to use that work under the conditions stated herein. The “**Document**”, below, refers to any such manual or work. Any member of the public is a licensee, and is addressed as “**you**”. You accept the license if you copy, modify or distribute the work in a way requiring permission under copyright law.

A “**Modified Version**” of the Document means any work containing the Document or a portion of it, either copied verbatim, or with modifications and/or translated into another language.

A “**Secondary Section**” is a named appendix or a front-matter section of the Document that deals exclusively with the relationship of the publishers or authors of the Document to the Document’s overall subject (or to related matters) and contains nothing that could fall directly within that overall subject. (Thus, if the Document is in part a textbook of mathematics, a Secondary Section may not explain any mathematics.) The relationship could be a matter of historical

connection with the subject or with related matters, or of legal, commercial, philosophical, ethical or political position regarding them.

The “**Invariant Sections**” are certain Secondary Sections whose titles are designated, as being those of Invariant Sections, in the notice that says that the Document is released under this License. If a section does not fit the above definition of Secondary then it is not allowed to be designated as Invariant. The Document may contain zero Invariant Sections. If the Document does not identify any Invariant Sections then there are none.

The “**Cover Texts**” are certain short passages of text that are listed, as Front-Cover Texts or Back-Cover Texts, in the notice that says that the Document is released under this License. A Front-Cover Text may be at most 5 words, and a Back-Cover Text may be at most 25 words.

A “**Transparent**” copy of the Document means a machine-readable copy, represented in a format whose specification is available to the general public, that is suitable for revising the document straightforwardly with generic text editors or (for images composed of pixels) generic paint programs or (for drawings) some widely available drawing editor, and that is suitable for input to text formatters or for automatic translation to a variety of formats suitable for input to text formatters. A copy made in an otherwise Transparent file format whose markup, or absence of markup, has been arranged to thwart or discourage subsequent modification by readers is not Transparent. An image format is not Transparent if used for any substantial amount of text. A copy that is not “Transparent” is called “**Opaque**”.

Examples of suitable formats for Transparent copies include plain ASCII without markup, Texinfo input format, LaTeX input format, SGML or XML using a publicly available DTD, and standard-conforming simple HTML, PostScript or PDF designed for human modification. Examples of transparent image formats include PNG, XCF and JPG. Opaque formats include proprietary formats that can be read and edited only by proprietary word processors, SGML or XML for which the DTD and/or processing tools are not generally available, and the machine-generated HTML, PostScript or PDF produced by some word processors for output purposes only.

The “**Title Page**” means, for a printed book, the title page itself, plus such following pages as are needed to hold, legibly, the material this License requires to appear in the title page. For works in formats which do not have any title page as such, “Title Page” means the text near the most prominent appearance of the work’s title, preceding the beginning of the body of the text.

The “**publisher**” means any person or entity that distributes copies of the Document to the public.

A section “**Entitled XYZ**” means a named subunit of the Document whose title either is precisely XYZ or contains XYZ in parentheses following text that translates XYZ in another language. (Here XYZ stands for a specific section name mentioned below, such as “**Acknowledgements**”, “**Dedications**”, “**Endorsements**”, or “**History**”.) To “**Preserve the Title**” of such a section when you modify the Document means that it remains a section “Entitled XYZ” according to this definition.

The Document may include Warranty Disclaimers next to the notice which states that this License applies to the Document. These Warranty Disclaimers are considered to be included by reference in this License, but only as regards disclaiming warranties : any other implication that these Warranty Disclaimers may have is void and has no effect on the meaning of this License.

VERBATIM COPYING

You may copy and distribute the Document in any medium, either commercially or noncommercially, provided that this License, the copyright notices, and the license notice saying this License applies to the Document are reproduced in all copies, and that you add no other conditions whatsoever to those of this License. You may not use technical measures to obstruct or control the reading or further copying of the copies you make or distribute. However, you may accept compensation in exchange for copies. If you distribute a large enough number of copies you must also follow the conditions in section 3.

You may also lend copies, under the same conditions stated above, and you may publicly display copies.

COPYING IN QUANTITY

If you publish printed copies (or copies in media that commonly have printed covers) of the Document, numbering more than 100, and the Document's license notice requires Cover Texts, you must enclose the copies in covers that carry, clearly and legibly, all these Cover Texts : Front-Cover Texts on the front cover, and Back-Cover Texts on the back cover. Both covers must also clearly and legibly identify you as the publisher of these copies. The front cover must present the full title with all words of the title equally prominent and visible. You may add other material on the covers in addition. Copying with changes limited to the covers, as long as they preserve the title of the Document and satisfy these conditions, can be treated as verbatim copying in other respects.

If the required texts for either cover are too voluminous to fit legibly, you should put the first ones listed (as many as fit reasonably) on the actual cover, and continue the rest onto adjacent pages.

If you publish or distribute Opaque copies of the Document numbering more than 100, you must either include a machine-readable Transparent copy along with each Opaque copy, or state in or with each Opaque copy a computer-network location from which the general network-using public has access to download using public-standard network protocols a complete Transparent copy of the Document, free of added material. If you use the latter option, you must take reasonably prudent steps, when you begin distribution of Opaque copies in quantity, to ensure that this Transparent copy will remain thus accessible at the stated location until at least one year after the last time you distribute an Opaque copy (directly or through your agents or retailers) of that edition to the public.

It is requested, but not required, that you contact the authors of the Document well before redistributing any large number of copies, to give them a chance to provide you with an updated version of the Document.

MODIFICATIONS

You may copy and distribute a Modified Version of the Document under the conditions of sections 2 and 3 above, provided that you release the Modified Version under precisely this License, with the Modified Version filling the role of the Document, thus licensing distribution and modification of the Modified Version to whoever possesses a copy of it. In addition, you must do these things in the Modified Version :

- A. Use in the Title Page (and on the covers, if any) a title distinct from that of the Document, and from those of previous versions (which should, if there were any, be listed in the History section of the Document). You may use the same title as a previous version if the original publisher of that version gives permission.
- B. List on the Title Page, as authors, one or more persons or entities responsible for authorship of the modifications in the Modified Version, together with at least five of the principal authors of the Document (all of its principal authors, if it has fewer than five), unless they release you from this requirement.
- C. State on the Title page the name of the publisher of the Modified Version, as the publisher.
- D. Preserve all the copyright notices of the Document.
- E. Add an appropriate copyright notice for your modifications adjacent to the other copyright notices.
- F. Include, immediately after the copyright notices, a license notice giving the public permission to use the Modified Version under the terms of this License, in the form shown in the Addendum below.

- G. Preserve in that license notice the full lists of Invariant Sections and required Cover Texts given in the Document's license notice.
- H. Include an unaltered copy of this License.
- I. Preserve the section Entitled "History", Preserve its Title, and add to it an item stating at least the title, year, new authors, and publisher of the Modified Version as given on the Title Page. If there is no section Entitled "History" in the Document, create one stating the title, year, authors, and publisher of the Document as given on its Title Page, then add an item describing the Modified Version as stated in the previous sentence.
- J. Preserve the network location, if any, given in the Document for public access to a Transparent copy of the Document, and likewise the network locations given in the Document for previous versions it was based on. These may be placed in the "History" section. You may omit a network location for a work that was published at least four years before the Document itself, or if the original publisher of the version it refers to gives permission.
- K. For any section Entitled "Acknowledgements" or "Dedications", Preserve the Title of the section, and preserve in the section all the substance and tone of each of the contributor acknowledgements and/or dedications given therein.
- L. Preserve all the Invariant Sections of the Document, unaltered in their text and in their titles. Section numbers or the equivalent are not considered part of the section titles.
- M. Delete any section Entitled "Endorsements". Such a section may not be included in the Modified Version.
- N. Do not retitle any existing section to be Entitled "Endorsements" or to conflict in title with any Invariant Section.
- O. Preserve any Warranty Disclaimers.

If the Modified Version includes new front-matter sections or appendices that qualify as Secondary Sections and contain no material copied from the Document, you may at your option designate some or all of these sections as invariant. To do this, add their titles to the list of Invariant Sections in the Modified Version's license notice. These titles must be distinct from any other section titles.

You may add a section Entitled "Endorsements", provided it contains nothing but endorsements of your Modified Version by various parties—for example, statements of peer review or that the text has been approved by an organization as the authoritative definition of a standard.

You may add a passage of up to five words as a Front-Cover Text, and a passage of up to 25 words as a Back-Cover Text, to the end of the list of Cover Texts in the Modified Version. Only one passage of Front-Cover Text and one of Back-Cover Text may be added by (or through arrangements made by) any one entity. If the Document already includes a cover text for the same cover, previously added by you or by arrangement made by the same entity you are acting on behalf of, you may not add another; but you may replace the old one, on explicit permission from the previous publisher that added the old one.

The author(s) and publisher(s) of the Document do not by this License give permission to use their names for publicity for or to assert or imply endorsement of any Modified Version.

COMBINING DOCUMENTS

You may combine the Document with other documents released under this License, under the terms defined in section 4 above for modified versions, provided that you include in the combination all of the Invariant Sections of all of the original documents, unmodified, and list them all as Invariant Sections of your combined work in its license notice, and that you preserve all their Warranty Disclaimers.

The combined work need only contain one copy of this License, and multiple identical Invariant Sections may be replaced with a single copy. If there are multiple Invariant Sections with the same name but different contents, make the title of each such section unique by adding at the end of

it, in parentheses, the name of the original author or publisher of that section if known, or else a unique number. Make the same adjustment to the section titles in the list of Invariant Sections in the license notice of the combined work.

In the combination, you must combine any sections Entitled “History” in the various original documents, forming one section Entitled “History”; likewise combine any sections Entitled “Acknowledgements”, and any sections Entitled “Dedications”. You must delete all sections Entitled “Endorsements”.

COLLECTIONS OF DOCUMENTS

You may make a collection consisting of the Document and other documents released under this License, and replace the individual copies of this License in the various documents with a single copy that is included in the collection, provided that you follow the rules of this License for verbatim copying of each of the documents in all other respects.

You may extract a single document from such a collection, and distribute it individually under this License, provided you insert a copy of this License into the extracted document, and follow this License in all other respects regarding verbatim copying of that document.

AGGREGATION WITH INDEPENDENT WORKS

A compilation of the Document or its derivatives with other separate and independent documents or works, in or on a volume of a storage or distribution medium, is called an “aggregate” if the copyright resulting from the compilation is not used to limit the legal rights of the compilation’s users beyond what the individual works permit. When the Document is included in an aggregate, this License does not apply to the other works in the aggregate which are not themselves derivative works of the Document.

If the Cover Text requirement of section 3 is applicable to these copies of the Document, then if the Document is less than one half of the entire aggregate, the Document’s Cover Texts may be placed on covers that bracket the Document within the aggregate, or the electronic equivalent of covers if the Document is in electronic form. Otherwise they must appear on printed covers that bracket the whole aggregate.

TRANSLATION

Translation is considered a kind of modification, so you may distribute translations of the Document under the terms of section 4. Replacing Invariant Sections with translations requires special permission from their copyright holders, but you may include translations of some or all Invariant Sections in addition to the original versions of these Invariant Sections. You may include a translation of this License, and all the license notices in the Document, and any Warranty Disclaimers, provided that you also include the original English version of this License and the original versions of those notices and disclaimers. In case of a disagreement between the translation and the original version of this License or a notice or disclaimer, the original version will prevail.

If a section in the Document is Entitled “Acknowledgements”, “Dedications”, or “History”, the requirement (section 4) to Preserve its Title (section 1) will typically require changing the actual title.

TERMINATION

You may not copy, modify, sublicense, or distribute the Document except as expressly provided under this License. Any attempt otherwise to copy, modify, sublicense, or distribute it is void, and will automatically terminate your rights under this License.

However, if you cease all violation of this License, then your license from a particular copyright holder is reinstated (a) provisionally, unless and until the copyright holder explicitly and finally terminates your license, and (b) permanently, if the copyright holder fails to notify you of the violation by some reasonable means prior to 60 days after the cessation.

Moreover, your license from a particular copyright holder is reinstated permanently if the copyright holder notifies you of the violation by some reasonable means, this is the first time you have received notice of violation of this License (for any work) from that copyright holder, and you cure the violation prior to 30 days after your receipt of the notice.

Termination of your rights under this section does not terminate the licenses of parties who have received copies or rights from you under this License. If your rights have been terminated and not permanently reinstated, receipt of a copy of some or all of the same material does not give you any rights to use it.

FUTURE REVISIONS OF THIS LICENSE

The Free Software Foundation may publish new, revised versions of the GNU Free Documentation License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns. See <http://www.gnu.org/copyleft/>.

Each version of the License is given a distinguishing version number. If the Document specifies that a particular numbered version of this License “or any later version” applies to it, you have the option of following the terms and conditions either of that specified version or of any later version that has been published (not as a draft) by the Free Software Foundation. If the Document does not specify a version number of this License, you may choose any version ever published (not as a draft) by the Free Software Foundation. If the Document specifies that a proxy can decide which future versions of this License can be used, that proxy’s public statement of acceptance of a version permanently authorizes you to choose that version for the Document.

RELICENSING

“Massive Multiauthor Collaboration Site” (or “MMC Site”) means any World Wide Web server that publishes copyrightable works and also provides prominent facilities for anybody to edit those works. A public wiki that anybody can edit is an example of such a server. A “Massive Multiauthor Collaboration” (or “MMC”) contained in the site means any set of copyrightable works thus published on the MMC site.

“CC-BY-SA” means the Creative Commons Attribution-Share Alike 3.0 license published by Creative Commons Corporation, a not-for-profit corporation with a principal place of business in San Francisco, California, as well as future copyleft versions of that license published by that same organization.

“Incorporate” means to publish or republish a Document, in whole or in part, as part of another Document.

An MMC is “eligible for relicensing” if it is licensed under this License, and if all works that were first published under this License somewhere other than this MMC, and subsequently incorporated in whole or in part into the MMC, (1) had no cover texts or invariant sections, and (2) were thus incorporated prior to November 1, 2008.

The operator of an MMC Site may republish an MMC contained in the site under CC-BY-SA on the same site at any time before August 1, 2009, provided the MMC is eligible for relicensing.

ADDENDUM : How to use this License for your documents

To use this License in a document you have written, include a copy of the License in the document and put the following copyright and license notices just after the title page :

Copyright © YEAR YOUR NAME. Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.3 or any later version published by the Free Software Foundation ; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license is included in the section entitled “GNU Free Documentation License”.

If you have Invariant Sections, Front-Cover Texts and Back-Cover Texts, replace the “with ... Texts.” line with this :

with the Invariant Sections being LIST THEIR TITLES, with the Front-Cover Texts being LIST, and with the Back-Cover Texts being LIST.

If you have Invariant Sections without Cover Texts, or some other combination of the three, merge those two alternatives to suit the situation.

If your document contains nontrivial examples of program code, we recommend releasing these examples in parallel under your choice of free software license, such as the GNU General Public License, to permit their use in free software.

Liste des notations

$N \triangleleft G$ Le sous-groupe N est normal dans G , page 141

Algèbre

$[\mathbb{L} : \mathbb{K}]$ degré d'une extension de corps, page 304

$\mathcal{L}(E, F)$ Ensemble des applications linéaires de E dans F , page 223

$\mathbb{K}(A)$ corps contenant \mathbb{K} et A , page 317

$\mathbb{K}[A]$ anneau contenant \mathbb{K} et A , page 317

\mathbb{N}_0 les naturels non nuls : $\mathbb{N}_0 = \mathbb{N} \setminus \{0\}$, page 109

$\mathcal{M}_{n \times m}$ l'ensemble des matrices $n \times m$, page 223

∇f gradient de la fonction f , page 734

proj_V projection de $V \times W$ sur V , page 575

$\text{Span}(A)$ l'ensemble des combinaisons linéaires finies d'éléments de A , page 215

$C^1(U, \mathbb{R}^n)$ Les applications une fois continument dérivables, page 747

$df_a(u)$ Application de la différentielle de f sur le vecteur u , page 719

$f^{(n)}$ La n -ième dérivée de la fonction f , page 799

$o(x)$ fonction tendant rapidement vers zéro, page 801

(p) idéal engendré par p , page 179

\mathbb{F}_p lorsque p est premier, page 289

\mathbb{F}_{p^n} corps fini à p^n éléments, page 1308

$\text{Frac}(\mathbb{A})$ Le corps des fractions de l'anneau \mathbb{A} , page 118

$\text{Fun}(X, Y)$ les applications de X vers Y , page 112

$S(E)$ Les opérateurs autoadjoints de E , page 472

U_n Le groupe des racines n^{e} de l'unité., page 1291

$\text{res}(P, Q)$ résultat des polynômes P et Q , page 464

\sqrt{A} racine d'une matrice hermitienne positive, page 835

$\theta_\alpha(P)$ la multiplicité de α par rapport à P , page 211

$A[X]$ tous les polynômes de degré fini à coefficients dans A , page 207

$A_n[X]$ les polynômes à coefficients dans A et de degré inférieur à n , page 208

$C(P)$ matrice compagnon, page 528

$D \mid P$ D divise P , page 209

$E_\lambda(u)$ Espace propre de u , page 483

$mat_{\mathcal{B}}(q)$ matrice de q dans la base \mathcal{B} , page 1142

$U(A)$ ensemble des inversibles, page 114

Ensembles de matrices

$S^+(n, \mathbb{R})$ matrices symétriques définies positives, page 509

$S^{++}(n, \mathbb{R})$ matrices symétriques strictement définies positives, page 509

$\text{Aut}(E)$ automorphisme de l'espace vectoriel E , page 224

$\text{End}(E)$ les endomorphismes de E , page 224

$L(E, F)$ applications linéaires bornées (continues), page 569

$O(n, \mathbb{R})$ le groupe des matrices orthogonales, page 472

$\Omega(E)$ formes quadratiques non dégénérées, page 1138

$Q(E)$ formes quadratiques réelles sur E , page 511

$Q^+(E)$ formes quadratiques positives, page 1138

$Q^{++}(E)$ formes quadratiques strictement définies positives, page 1138

$S_n^{p,q}(\mathbb{R})$ matrices symétriques réelles de signature (p, q) , page 1139

$\beta(s)$ Vecteur unitaire de la binormale, page 1474

$\gamma \sim g$ Équivalence d'arcs paramétrés, page 1465

$\nu(s)$ Vecteur unitaire de la normale principale, page 1474

$c(s)$ rayon de courbure, page 1474

$t(s)$ Torsion, page 1475

Géométrie

$(x_0 : \dots : x_n)$ coordonnées homogènes dans un espace projectif, page 1526

$\text{Conv}(A)$ enveloppe convexe, page 427

$\mathbb{C}[G]$ combinaisons d'éléments de G à coefficients dans \mathbb{C} , page 1049

$\text{PGL}(E)$ groupe projectif, page 1514

B^o orthogonal dans le dual, page 255

$P(E)$ l'espace projectif de E , page 1507

$P_1(\mathbb{C})$ sphère de Riemann, page 1529

Chaînes de Markov

$\pi(x)$ lié au temps de retour, page 2129

Probabilités et statistique

$\sigma(X)$ La tribu engendrée par la variable aléatoire X , page 2001

$a \wedge b$ $\min(a, b)$, page 2147

K_X matrice de covariance d'un vecteur gaussien, page 2052

$m(\mathcal{A})$ Ensemble des fonctions \mathcal{A} -mesurables, page 861

Théorie des groupes

$(G/H)_g$ classes à gauche, page 153

$[a]_p$ ensemble des $a + kp$, page 177

- $[G, G]$ groupe dérivé, page 144
- $[g, h]$ commutateur dans un groupe, page 144
- $\text{Aff}(\mathbb{R}^n)$ Le groupe des applications affines bijectives de \mathbb{R}^n ., page 437
- gr groupe engendré, page 141
- \hat{G} groupe des caractères de G , page 1045
- σ_x réflexion par rapport à x , page 1159
- A_n groupe alterné, page 275
- $D(G)$ groupe dérivé, page 144
- D_n groupe diédral, page 1202
- G^{ab} groupe abélianisé de G , page 145
- $N \times_{\phi} H$ produit semi-direct, page 161
- S_n le groupe symétrique, page 157

Topologie et théorie des ensembles

- $\text{Adh}(A)$ adhérence de A , page 344
- $\complement A$ Le complémentaire de l'ensemble A , page 107
- $\text{Diam}(A)$ Diamètre de la partie A , page 667
- $\text{Int}(A)$ intérieur de A , page 344
- ∂A La frontière de l'ensemble A , page 480
- $A \Delta B$ différence symétrique, page 108
- A^c complémentaire de A , page 107

Analyse

- $\text{Isom}(X)$ Le groupe des isométries de X , page 388
- $\mu \ll \nu$ La mesure μ est absolument continue par rapport à la mesure ν ., page 930
- $C^{\infty}(\mathbb{R}, \mathcal{S}'(\mathbb{R}^d))$ Fonctions à valeurs dans les distributions., page 1859
- $C^{\infty}(I, \mathcal{D}'(\mathbb{R}^d))$ fonctions à valeurs dans les distributions, page 1789
- $(S, \hat{\mathcal{F}}, \hat{\mu})$ complété de l'espace mesuré $(S, \hat{\mathcal{F}}, \hat{\mu})$, page 867
- $\mathcal{L}(E, F)$ Les applications linéaires de E vers F , page 612
- $\mathcal{L}^{(n)}(V, W)$ L'espace des applications n -linéaires $V^n \rightarrow W$, page 755
- $\arg(z)$ La valeur principale de l'argument de $z \in \mathbb{C}$, page 1633
- L Les applications linéaires continues de E vers F , page 612
- \mathbb{D}_b l'ensemble de écritures décimales en base b , page 589
- \mathbb{R} l'ensemble des réels, page 127
- \mathbb{R}^+ les réels positifs ou nuls, page 130
- \exp exponentielle, page 1004
- \mathcal{H}' dual, page 1579
- $\liminf a_n$ limite inférieure, page 375
- $\limsup a_n$ limite supérieure, page 375

- \mathcal{L}^p espace de Lebesgue, sans les classes, page 1645
- $\mu \perp \nu$ mesures perpendiculaires, page 930
- μ^* La mesure extérieure associée à la mesure μ , page 859
- $\partial_z, \partial_{\bar{z}}$ dérivées partielles d'une fonction complexe, page 1601
- $\text{proj}_K(x)$ projection orthogonale de x sur y , page 1576
- $\sigma(\mathcal{A})$ tribu engendrée par \mathcal{D} , page 848
- $\mathcal{D}(\Omega)$ Les fonctions C^∞ à support compact sur Ω , page 1775
- $\mathcal{S}'(\mathbb{R}^d)$ espace des distributions tempérées, page 1781
- $A^2(\Omega)$ espace de Bergman, page 1728
- A^\perp orthogonal d'une partie., page 1578
- $C_c(I)$ fonctions continues à support compact dans I , page 1665
- $f \sim g$ fonctions ayant des limites équivalentes, page 707
- $H^1(\Omega)$ espace de Sobolev sur Ω , page 1804
- $H^1(I)$ espace de Sobolev, page 1799
- $H^m(M)$ espace de Sobolev, page 1806
- $L^1_{loc}(I)$ fonctions intégrables sur les compacts de I , page 1799
- L^p espace de Lebesgue avec les classes, page 1647
- $M_i\varphi$ La fonction $x \mapsto x_i\varphi(x)$, page 1723
- $S_n f$ somme partielle de série de Fourier, page 1693

Bibliographie

- [1] L'auteur de ces lignes. Invention personnelle de l'auteur de ces lignes, lire avec prudence et merci de me dire si c'est correct.
- [2] A. Casamayou, N Cohen, G. Connan, T. Dumont, L. Fousse, F. Maltey, M. Meulien, M. Mezzarobba, C. Pernet, M. N. Thiéry, and Zimmermann P. Calcul mathématique avec Sage. 2013. URL <http://sagebook.gforge.inria.fr/>.
- [3] Wikipedia contributors. Nth root algorithm — Wikipedia, the free encyclopedia. https://en.wikipedia.org/w/index.php?title=Nth_root_algorithm&oldid=873962836, 2018. [Online; accessed 31-January-2019].
- [4] Sylvie Benzoni. Distributions tempérées. 13 novembre 2013. URL <http://math.univ-lyon1.fr/~benzoni/Agreg/Distributions.pdf>.
- [5] Wikipédia. Axiome du choix — wikipédia, l'encyclopédie libre, 2015. URL http://fr.wikipedia.org/w/index.php?title=Axiome_du_choix&oldid=116151439. [En ligne; Page disponible le 9-septembre-2015].
- [6] O. Teytaud, C. Antonini, J.-B. Bardet, J.-F. Quint, M. De Crisenoy, P. Borgnat, J. Bérard, E. Lebeau, E. Souche, and A. Chateau. Les mathématiques pour l'Agrégation. 21 mai 2002. <http://les.mathematiques.free.fr/pdf/agreg.zip>.
- [7] Pedro Tamaroff. Prove that if A is an infinite set then $A \times 2$ is equipotent to A . 2014. URL <https://math.stackexchange.com/questions/1041731/prove-that-if-a-is-an-infinite-set-then-a-times-2-is-equipotent-to-a>.
- [8] Daniel Daigle. Cardinalité. URL <http://mysite.science.uottawa.ca/asavag2/mat2762/cardinal.pdf>.
- [9] Christine Laurent-Thiébaud. Aximatique des nombres. URL <http://ljk.imag.fr/membres/Bernard.Ycart/mel/ax/ax.pdf>.
- [10] Ilie Grigorescu. Lecture 4 - finite and infinite sets. URL http://www.math.miami.edu/~igrigore/teaching/mth533/lecture4_mth433.pdf.
- [11] Wikipédia. Ensemble dénombrable — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Ensemble_d%C3%A9nombrable&oldid=150759981. [En ligne; Page disponible le 28-juillet-2018].
- [12] Patrice Tauvel. *Cours d'algèbre*. Dunod, 1999. ISBN 2-10-005490-3.
- [13] Wikipédia. Morphisme d'anneaux — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Morphisme_d%27anneaux&oldid=129892066. [En ligne; Page disponible le 23-septembre-2016].
- [14] Wikipédia. Diviseur de zéro — wikipédia, l'encyclopédie libre, 2017. URL http://fr.wikipedia.org/w/index.php?title=Diviseur_de_z%C3%A9ro&oldid=143858271. [En ligne; Page disponible le 26-décembre-2017].

- [15] Wikipédia. Corps commutatif — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Corps_commutatif&oldid=152289565. [En ligne; Page disponible le 18-septembre-2018].
- [16] Wikipédia. Corps des fractions — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Corps_des_fractions&oldid=130674526. [En ligne; Page disponible le 14-octobre-2016].
- [17] Frédéric Paulin. Topologie, analyse et calcul différentiel. 2008-2009. http://www.math-question-center.com/publications-pdf/cours_d'analyse+topologie_et_calcul_differntiel_Frederic_Paulin.pdf.
- [18] Wikipédia. Construction des nombres réels — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Construction_des_nombres_r%C3%A9els&oldid=145402164. [En ligne; Page disponible le 12-février-2018].
- [19] Wikipédia. Nombre positif — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Nombre_positif&oldid=146043202. [En ligne; Page disponible le 3-mars-2018].
- [20] Wikipedia contributors. Least-upper-bound property — Wikipedia, the free encyclopedia, 2019. URL https://en.wikipedia.org/w/index.php?title=Least-upper-bound_property&oldid=893001075. [Online; accessed 30-April-2019].
- [21] Peter Hendrikus Kropholler. Groups rings and fields. 2004-2005. <http://www.maths.gla.ac.uk/~phk/ulearning.pdf>.
- [22] Jean-Louis Rouget. Compléments d'algèbre. 2017. URL <https://www.maths-france.fr/MathSpe/Cours/01-structures.pdf>.
- [23] Anonyme. Théorie des groupes. 2008. Lien original mort : <http://theoriesdesgroupes.perso.sfr.fr/cours/theoriePDF.pdf>
autre lien, qui fonctionne pour l'instant : <http://ekldata.com/CKoLuTDRT9S2sx5BPX9DkCDAcdo.pdf>
et <https://web.archive.org/web/20150701233527/http://theoriesdesgroupes.perso.sfr.fr/cours/theoriePDF.pdf>.
- [24] Jean Luc W. groupe-quotient et création de nouveaux sous-groupes. 2007. URL https://fr.wikipedia.org/wiki/Discussion:Groupe_quotient.
- [25] Pierre Lissy. Théorème de Jordan-Hölder. 6 avril 2010. <http://www.ljll.math.upmc.fr/~lissy/Agreg/developpements/JordanHolder.pdf>.
- [26] Muriel Galley. Groupes résolubles. 24 mai 2012. URL <http://matthieu.gendulpe.com/Galley.pdf>.
- [27] François Combes. *Algèbre et géométrie*. Bréal, 2003. ISBN 2-84291-202-0.
- [28] Fabrice Castel. Groupes finis. 2009-2012. http://agreg-maths.univ-rennes1.fr/documentation/docs/Groupes_finis.pdf.
- [29] Wikipédia. Groupe symétrique — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Groupe_sym%C3%A9trique&oldid=146656747. [En ligne; Page disponible le 21-mars-2018].
- [30] Jean Delcourt. Groupes-permutations. URL http://delcourt.u-cergy.fr/StrucAlg/chap_5.pdf.

- [31] Wikipédia. Permutation — wikipédia, l'encyclopédie libre, 2013. URL <http://fr.wikipedia.org/w/index.php?title=Permutation&oldid=89805226>. [En ligne ; Page disponible le 24-juin-2013].
- [32] Patrick Polo. Idéaux premiers et maximaux, localisation, théorème des zéros de Hilbert. 23 octobre 2005. URL <https://webusers.imj-prg.fr/~patrick.polo/M1Galois/ATG05ch3.pdf>.
- [33] cdrcprds. Questions (9), (15) et (21) d'algèbre. 2017. URL <https://github.com/LaurentClaessens/mazhe/issues/60>.
- [34] Wikiversité. Arithmétique/divisibilité et congruences dans \mathbb{Z} — wikiversité,, 2017. URL https://fr.wikiversity.org/w/index.php?title=Arithm%C3%A9tique/Divisibilit%C3%A9_et_congruences_dans_Z&oldid=650617. [En ligne ; accédé le 25-mai-2017].
- [35] Mongi Amorri. Théorème de bezout, théorème de gauss. http://lux.lyceefrancais-brasilia.net/documents/cours/ts/bezout_gauss.pdf.
- [36] Gilles Costantini. pgcd et ppcm dans \mathbb{Z} , théorème de Bezout, applications. . http://gilles.costantini.pagesperso-orange.fr/agreg_fichiers/bezout.pdf.
- [37] Wikipédia. Lemme d'euclide — wikipédia, l'encyclopédie libre, 2015. URL http://fr.wikipedia.org/w/index.php?title=Lemme_d%27Euclide&oldid=112046247. [En ligne ; Page disponible le 28-juin-2015].
- [38] Wikipédia. Théorème fondamental de l'arithmétique — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_fondamental_de_l%27arithm%C3%A9tique&oldid=110036012. [En ligne ; Page disponible le 26-juin-2015].
- [39] Wikipédia. Théorème de cauchy (groupes) — wikipédia, l'encyclopédie libre, 2017. URL [http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Cauchy_\(groupes\)&oldid=140167231](http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Cauchy_(groupes)&oldid=140167231). [En ligne ; Page disponible le 17-septembre-2017].
- [40] Alexei Pantchichkine. Algèbre2. 2004-2005. <http://www-fourier.ujf-grenoble.fr/~panchish/05ens1.pdf>.
- [41] Niels Borne. Fiche numéro 2 : morphismes, sous-groupe distingué, quotient. 2009-2012. http://math.univ-lille1.fr/~borne/Enseignement/TD2_M308.pdf.
- [42] Georges Skandalis. Algèbre générale et algèbre linéaire. 6 mars 2015. URL http://www.math.univ-paris-diderot.fr/_media/formations/prepa/agreginterne/polycopiealgebre.pdf.
- [43] Arnaud Girand. Développements pour l'agrégation externe de mathématiques. 18 août 2012. URL http://perso.univ-rennes1.fr/arnaud.girand/pdf/dvp_agreg/dvp.pdf.
- [44] Wikipédia. Formule du binôme de newton — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Formule_du_bin%C3%B4me_de_Newton&oldid=144161875. [En ligne ; Page disponible le 4-janvier-2018].
- [45] Wikipédia. Caractéristique d'un anneau — wikipédia, l'encyclopédie libre, 2017. URL http://fr.wikipedia.org/w/index.php?title=Caract%C3%A9ristique_d%27un_anneau&oldid=136316021. [En ligne ; Page disponible le 22-mai-2017].
- [46] Wikipédia. Module sur un anneau — wikipédia, l'encyclopédie libre, 2017. URL http://fr.wikipedia.org/w/index.php?title=Module_sur_un_anneau&oldid=141582982. [En ligne ; Page disponible le 16-octobre-2017].

- [47] Wikipédia. Algèbre sur un corps — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Alg%C3%A8bre_sur_un_corps&oldid=98992040. [En ligne ; Page disponible le 6-mars-2014].
- [48] Wikipédia. Primalité dans un anneau — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Primalit%C3%A9_dans_un_anneau&oldid=127919489. [En ligne ; Page disponible le 30-mai-2017].
- [49] Wikipédia. Idéal premier — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Id%C3%A9al_premier&oldid=125964642. [En ligne ; Page disponible le 7-mai-2016].
- [50] Anonyme. sans titre. <http://www.ac-grenoble.fr/champo/IMG/arithmetique.pdf>.
- [51] Wikipédia. Primalité dans un anneau — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Primalit%C3%A9_dans_un_anneau&oldid=145407116. [En ligne ; Page disponible le 12-février-2018].
- [52] Michel Cretin. Anneaux principaux et factoriels. URL <http://math.univ-lyon1.fr/~cretin/OralAlgebre/annprincfact.pdf>.
- [53] Wikipédia. Idéal premier — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Id%C3%A9al_premier&oldid=145977206. [En ligne ; Page disponible le 1-mars-2018].
- [54] Danny-Jack Mercier. Anneaux factoriels. 5 octobre 2003. <http://megamaths.perso.neuf.fr/cours/ari/cann0001.pdf>.
- [55] Wikipédia. Anneau euclidien — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Anneau_euclidien&oldid=126145344. [En ligne ; Page disponible le 1-juin-2017].
- [56] Wikiversité. Anneau (mathématiques)/exercices/exercices — wikiversité,, 2017. URL [https://fr.wikiversity.org/w/index.php?title=Anneau_\(math%C3%A9matiques\)/Exercices/Exercices&oldid=676468](https://fr.wikiversity.org/w/index.php?title=Anneau_(math%C3%A9matiques)/Exercices/Exercices&oldid=676468). [En ligne ; accédé le 26-octobre-2017].
- [57] Sébastien Pellerin. Développements d'algèbre pour l'oral de l'agrégation. . <http://pellerin.xyz/doc/agreg/algebre.pdf>.
- [58] Wikipédia. Triplet pythagorien — wikipédia, l'encyclopédie libre, 2015. URL http://fr.wikipedia.org/w/index.php?title=Triplet_pythagorien&oldid=114628048. [En ligne ; Page disponible le 24-juin-2015].
- [59] Wikipédia. Primalité dans un anneau — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Primalit%C3%A9_dans_un_anneau&oldid=127919489. [En ligne ; Page disponible le 1-juin-2017].
- [60] Wikipédia. Lemme de Gauss (polynômes) — wikipédia, l'encyclopédie libre, 2015. URL [http://fr.wikipedia.org/w/index.php?title=Lemme_de_Gauss_\(polyn%C3%B4mes\)&oldid=112056131](http://fr.wikipedia.org/w/index.php?title=Lemme_de_Gauss_(polyn%C3%B4mes)&oldid=112056131). [En ligne ; Page disponible le 1-juin-2017].
- [61] Wikipédia. Espace vectoriel — wikipédia, l'encyclopédie libre, 2017. URL http://fr.wikipedia.org/w/index.php?title=Espace_vectoriel&oldid=143370323. [En ligne ; Page disponible le 10-décembre-2017].
- [62] Espaces vectoriels de dimension infinie. . URL <http://www.les-mathematiques.net/b/e/e/node8.php>.

- [63] Wikipedia contributors. Tensor product — Wikipedia, the free encyclopedia, 2018. URL https://en.wikipedia.org/w/index.php?title=Tensor_product&oldid=842778935. [Online ; accessed 18-June-2018].
- [64] Marie-Pierre Lebaux. Rappels sur les applications linéaires. URL <https://perso.univ-rennes1.fr/marie-pierre.lebaud/agint/ecrit/algebre-lineaire/applications-lineaires/V-appli-lin.pdf>.
- [65] Jean-Louis Rouget. Dimension d'un espace vectoriel. 2018. URL <https://www.maths-france.fr/MathSup/Cours/22-dimensions.pdf>.
- [66] Daniel Ferrand. Étendre le corps. Juillet 2007. URL <http://webusers.imj-prg.fr/~daniel.ferrand/ExtCorps.pdf>.
- [67] Ycart Bernard and Luc Rozoy. Déterminants. URL <https://ljk.imag.fr/membres/Bernard.Ycart/mel/de/de.pdf>.
- [68] G. Donald Allen. Matrices and linear algebra. 2004. URL http://www.math.tamu.edu/~dallen/m640_03c/lectures/chapter2.pdf.
- [69] Wikipédia. Mineur (algèbre linéaire) — wikipédia, l'encyclopédie libre, 2017. URL [http://fr.wikipedia.org/w/index.php?title=Mineur_\(alg%C3%A8bre_lin%C3%A9aire\)&oldid=133960288](http://fr.wikipedia.org/w/index.php?title=Mineur_(alg%C3%A8bre_lin%C3%A9aire)&oldid=133960288). [En ligne ; Page disponible le 25-janvier-2017].
- [70] Marc Sage. Dualité en dimension finie. 24 octobre 2005. <http://www.normalesup.org/~sage/Enseignement/Cours/DualDimFinie.pdf>.
- [71] Wikipédia. Théorème de Cauchy (groupes) — wikipédia, l'encyclopédie libre, 2016. URL [http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Cauchy_\(groupes\)&oldid=122684298](http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Cauchy_(groupes)&oldid=122684298). [En ligne ; Page disponible le 1-août-2017].
- [72] Wikiversité. Groupe (mathématiques)/exercices/premiers résultats sur les groupes simples — wikiversité,, 2016. URL [https://fr.wikiversity.org/w/index.php?title=Groupe_\(math%C3%A9matiques\)/Exercices/Premiers_r%C3%A9sultats_sur_les_groupes_simples&oldid=597195](https://fr.wikiversity.org/w/index.php?title=Groupe_(math%C3%A9matiques)/Exercices/Premiers_r%C3%A9sultats_sur_les_groupes_simples&oldid=597195). [En ligne ; accédé le 31-mai-2017].
- [73] Wikiversité. Groupe (mathématiques)/automorphismes d'un groupe cyclique — wikiversité,, 2016. URL [https://fr.wikiversity.org/w/index.php?title=Groupe_\(math%C3%A9matiques\)/Automorphismes_d%27un_groupe_cyclique&oldid=595608](https://fr.wikiversity.org/w/index.php?title=Groupe_(math%C3%A9matiques)/Automorphismes_d%27un_groupe_cyclique&oldid=595608). [En ligne ; accédé le 3-juin-2017].
- [74] Structure des groupes d'ordre pq . URL http://agregmaths.free.fr/doc/docs_nicolas/developpementAlgebre/grouped'ordrepq.pdf.
- [75] Olivier Debarre. Td4 : produit semi-direct. 2015-2016. URL <http://www.math.ens.fr/~debarre/TDC4.pdf>.
- [76] Christophe Mourougane. Théorie des groupes et géométrie. 2009-2010. <http://perso.univ-rennes1.fr/christophe.mourougane/enseignements/2009-10/THGG/poly.groupes.pdf>.
- [77] Michel Emsalem and Pierre Dèbes. Théorème de Sylow. <http://exo7.emath.fr/ficpdf/fic00023.pdf>.
- [78] N.G.J. Pagnon. Le groupe symétrique et le groupe alterné. URL http://amatheux.com/IMG/pdf/groupe_symetrique.pdf.
- [79] Wikipédia. Corps ordonné — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Corps_ordonn%C3%A9&oldid=125185412. [En ligne ; Page disponible le 22-mai-2017].

- [80] Automorphisme de \mathbb{R} et continuité. 2008. URL <http://forums.futura-sciences.com/mathematiques-superieur/261451-automorphisme-de-r-continuite.html>.
- [81] Wikipédia. Monoïde — wikipédia, l'encyclopédie libre, 2018. URL <http://fr.wikipedia.org/w/index.php?title=Mono%C3%AFde&oldid=148198215>. [En ligne ; Page disponible le 6-mai-2018].
- [82] Sylvain Duchet. Congruence dans \mathbb{Z} , anneaux $\mathbb{Z}/n\mathbb{Z}$, applications. <http://epsilon.2000.free.fr/Csup/congruences.pdf>.
- [83] Wikipédia. Polynôme irréductible — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Polyn%C3%B4me_irr%C3%A9ductible&oldid=132607764. [En ligne ; Page disponible le 12-décembre-2016].
- [84] Patrick Polo. Extension de corps, caractéristique, corps de rupture, corps de décomposition, clôtures algébriques. 2007-2008. URL <https://webusers.imj-prg.fr/~patrick.polo/M1Galois/ATG07chVII.pdf>.
- [85] C. Antonini, JF. Quint, P. Borgnat, J. Bérard, E. Lebeau, E. Souche, A. Chateau, and O. Teytaud. Le théorème de wedderburn. . URL <http://www.les-mathematiques.net/d/a/w/node5.php>.
- [86] cdrcprds. Lemme 5.43 sur les pgcd de polynômes. 2017. URL <https://github.com/LaurentClaessens/mazhe/issues/52#issuecomment-333251728>.
- [87] Wikipédia. Extension de corps — wikipédia, l'encyclopédie libre, 2017. URL http://fr.wikipedia.org/w/index.php?title=Extension_de_corps&oldid=133373097. [En ligne ; Page disponible le 5-janvier-2017].
- [88] Wikipédia. Extension de corps — wikipédia, l'encyclopédie libre, 2017. URL http://fr.wikipedia.org/w/index.php?title=Extension_de_corps&oldid=133373097. [En ligne ; Page disponible le 5-janvier-2017].
- [89] Wikipédia. Extension algébrique — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Extension_alg%C3%A9brique&oldid=148162879. [En ligne ; Page disponible le 5-mai-2018].
- [90] David Harari. Algèbre 1 - notions de théorie des corps. URL <https://www.math.u-psud.fr/~harari/exposes/corps.pdf>.
- [91] Pierre Bernard. Entiers algébriques et racines de l'unité. 7 janvier 2012. URL <http://allken-bernard.org/pierre/weblog/?p=2061>.
- [92] Mortajine Abdellatif. Les extensions de corps. 2012-2013. URL <http://www.iecn.u-nancy.fr/~mortajin/chap3-L3-S5.pdf>.
- [93] Wikipédia. Corps de décomposition — wikipédia, l'encyclopédie libre, 2017. URL http://fr.wikipedia.org/w/index.php?title=Corps_de_d%C3%A9composition&oldid=139881008. [En ligne ; Page disponible le 19-août-2017].
- [94] Patrick Polo. Extensions normales, séparables, galoisiennes. corps fini. 2007-2008. <http://www.math.jussieu.fr/~polo/M1/ATG07chVIII.pdf>.
- [95] Olivier Dodane. Le théorème de zéros de Hilbert. 8 octobre 2007. <http://www.math.ens.fr/~debarre/nullstellensatz.pdf>.
- [96] Patrick Polo. Polynômes symétriques et résolutions d'équations. 13 décembre 2004. <http://www.math.jussieu.fr/~polo/M1/ATGch9.pdf>.

- [97] Vincent Beck, Jérôme Malick, and Gabriel Peyré. Sur c tout est connexe! URL <http://objagr.gforge.inria.fr/documents/files/connexite-polynome.pdf>.
- [98] El Hage. Équation générale de degré n . <http://les.mathematiques.free.fr/pdf/gal10.pdf>.
- [99] Florian Morel Chevillet. Résolubilité par radicaux des équations algébriques. 31 mai 2012. URL <http://matthieu.gendulphe.com/MorelChevillet.pdf>.
- [100] Wikipédia. Espace séparé — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Espace_s%C3%A9par%C3%A9&oldid=97013072. [En ligne ; Page disponible le 27-septembre-2013].
- [101] Wikipédia. Équivalent — wikipédia, l'encyclopédie libre, 2018. URL <http://fr.wikipedia.org/w/index.php?title=%C3%89quivalent&oldid=149344743>. [En ligne ; Page disponible le 8-juin-2018].
- [102] Wikipédia. Base (topologie) — wikipédia, l'encyclopédie libre, 2017. URL [http://fr.wikipedia.org/w/index.php?title=Base_\(topologie\)&oldid=140151583](http://fr.wikipedia.org/w/index.php?title=Base_(topologie)&oldid=140151583). [En ligne ; Page disponible le 29-août-2017].
- [103] C. Antonini, JF. Quint, P. Borgnat, J. Bérard, E. Lebeau, E. Souche, A. Chateau, and O. Teytaud. Généralités. . URL <http://www.les-mathematiques.net/a/a/b/node19.php>.
- [104] Emmanuel Vieillard Baron. Sous-espaces compacts. 2001. URL <http://www.les-mathematiques.net/a/t/c/node5.php>.
- [105] Wikipédia. Compactifié d'alexandrov — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Compactifi%C3%A9_d%27Alexandrov&oldid=146811167. [En ligne ; Page disponible le 26-mars-2018].
- [106] Wikipedia contributors. Filtration (mathematics) — Wikipedia, the free encyclopedia, 2018. URL [https://en.wikipedia.org/w/index.php?title=Filtration_\(mathematics\)](https://en.wikipedia.org/w/index.php?title=Filtration_(mathematics)). [Online ; accessed 14-October-2018].
- [107] C. Antonini, JF. Quint, P. Borgnat, J. Bérard, E. Lebeau, E. Souche, A. Chateau, and O. Teytaud. Espaces métriques compacts. . <http://www.les-mathematiques.net/a/a/b/node22.php>.
- [108] Éric Brunelle. Norme matricielle. <http://www.dms.umontreal.ca/~math1600/6Supplement/Normematricielle-1.pdf>.
- [109] Wikipédia. Théorème de borel-lebesgue — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Borel-Lebesgue&oldid=96560796. [En ligne ; Page disponible le 11-mars-2014].
- [110] Wikipédia. Théorème des suites adjacentes — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_des_suites_adjacentes&oldid=148318442. [En ligne ; Page disponible le 10-mai-2018].
- [111] Mathieu Mansuy. Autour des séries alternées. URL <http://www.mathieu-mansuy.fr/pdf/ECS2-complément1.pdf>.
- [112] Wikipédia. Espace topologique — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Espace_topologique&oldid=152161195. [En ligne ; Page disponible le 13-septembre-2018].
- [113] Wikipédia. Connexité (mathématiques) — wikipédia, l'encyclopédie libre, 2016. URL [http://fr.wikipedia.org/w/index.php?title=Connexit%C3%A9_\(math%C3%A9matiques\)&oldid=123744973](http://fr.wikipedia.org/w/index.php?title=Connexit%C3%A9_(math%C3%A9matiques)&oldid=123744973). [En ligne ; Page disponible le 26-avril-2016].

- [114] C. Antonini, JF. Quint, P. Borgnat, J. Bérard, E. Lebeau, E. Souche, A. Chateau, and O. Teytaud. Connexité. . URL <http://www.les-mathematiques.net/a/a/b/node23.php>.
- [115] Huiqiang Jiang. Functional analysis. URL <http://www.math.pitt.edu/~hqjiang/2303/functional.pdf>.
- [116] Nicolas Bourbaki. General topology 2 chapters 5 - 10. URL <https://books.google.be/books?id=bQwhdmL6IjUC>. Oui je sais c'est une honte de citer Bourbaki en anglais en pointant vers une version partielle disponible sur Googlebooks. Mais c'est surtout une honte que ce livre ne soit pas disponible gratuitement en version électronique.
- [117] Guillaume Carlier. Analyse fonctionnelle. 2009-20010. URL <https://www.ceremade.dauphine.fr/~carlier/poly2010.pdf>.
- [118] Raz Kupferman. Topological vector spaces. URL http://www.ma.huji.ac.il/~razk/iWeb/My_Site/Teaching_files/TVS.pdf.
- [119] Ron Freiwald. Chapter iv completeness and compactness. URL <http://www.math.wustl.edu/~freiwald/ch4.pdf>.
- [120] Wikipédia. Suite de cauchy — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Suite_de_Cauchy&oldid=98046451. [En ligne ; Page disponible le 11-mars-2014].
- [121] Wikipédia. Continuité — wikipédia, l'encyclopédie libre, 2013. URL <http://fr.wikipedia.org/w/index.php?title=Continuit%C3%A9&oldid=96548980>. [En ligne ; Page disponible le 17-septembre-2013].
- [122] Wikipédia. Théorème des fermés emboîtés — Wikipédia, l'encyclopédie libre, 2012. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_des_ferm%C3%A9s_emboit%C3%A9s&oldid=85987870. [En ligne ; Page disponible le 1-avril-2013].
- [123] C. Antonini, JF. Quint, P. Borgnat, J. Bérard, E. Lebeau, E. Souche, A. Chateau, and O. Teytaud. Le théorème de tykhonov. . URL <http://www.les-mathematiques.net/a/a/b/node20.php>.
- [124] Wikipédia. Équicontinuité — wikipédia, l'encyclopédie libre, 2017. URL <http://fr.wikipedia.org/w/index.php?title=%C3%89quicontinuit%C3%A9&oldid=137910699>. [En ligne ; Page disponible le 4-juin-2017].
- [125] Wikipédia. Continuité uniforme — wikipédia, l'encyclopédie libre, 2017. URL http://fr.wikipedia.org/w/index.php?title=Continuit%C3%A9_uniforme&oldid=138238518. [En ligne ; Page disponible le 16-juin-2017].
- [126] Zied Ammari. Analyse fonctionnelle : Pré-requis. <http://perso.univ-rennes1.fr/zied.ammari/other/pdf/chapitre1.pdf>.
- [127] Jean Saint Raymond. Semi-normes. URL <http://www.math.jussieu.fr/~raymond/preprints/seminormes.pdf>.
- [128] Stéphane Mischler. Semi-norme et introduction aux evtcls. Février 2007. URL <https://www.ceremade.dauphine.fr/~mischler/Enseignements/AFAENS/Chap1evtlcs.pdf>.
- [129] Wikipédia. Semi-norme — wikipédia, l'encyclopédie libre, 2013. URL <http://fr.wikipedia.org/w/index.php?title=Semi-norme&oldid=89981857>. [En ligne ; Page disponible le 22-septembre-2013].
- [130] Wikipédia. Théorème de baire — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Baire&oldid=95077621. [En ligne ; Page disponible le 4-août-2013].

- [131] Marie-Claire David, Frédéric Haglund, and Daniel Perrin. Géométrie affine. 8 décembre 2003.
<http://webens.math.u-psud.fr/~mclld/GAEL/Gael/bar.pdf>.
- [132] Marie-Claude David, Frédéric Haglund, and Daniel Perrin. Géométrie affine. Document de travail pour la préparation au CAPES. Deuxième partit : barycentre. 8 décembre 2003. URL <http://omphale.math.u-psud.fr/~mclld/GAEL/Gael/bar.pdf>.
- [133] Marie-Claude David, Frédéric Haglund, and Daniel Perrin. Géométrie affine. Document de travail pour la préparation au CAPES. Troisième partie : convexité. 8 décembre 2003. URL <http://omphale.math.u-psud.fr/~mclld/GAEL/Gael/conv.pdf>.
- [134] Jean-Marc Decauwert. Géométrie affine. 8 novembre 2011. URL <http://ljk.imag.fr/membres/Bernard.Ycart/mel/ga/ga.pdf>.
- [135] Wikiversité. Barycentre/théorème de l'associativité du barycentre — wikiversité,, 2016. URL https://fr.wikiversity.org/w/index.php?title=Barycentre/Th%C3%A9or%C3%A8me_de_l%27associativit%C3%A9_du_barycentre&oldid=582092. [En ligne ; accédé le 16-janvier-2017].
- [136] Wikipédia. Application bilinéaire — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Application_bilin%C3%A9aire&oldid=144184234. [En ligne ; Page disponible le 5-janvier-2018].
- [137] Michel Coste. Formes bilinéaires symétriques, formes quadratiques. URL <https://perso.univ-rennes1.fr/michel.coste/Bil.pdf>.
- [138] Wikipédia. Hermitien — wikipédia, l'encyclopédie libre, 2016. URL <http://fr.wikipedia.org/w/index.php?title=Hermitien&oldid=124476255>. [En ligne ; Page disponible le 18-mars-2016].
- [139] Daniel Li. Espaces de hilbert. .
http://www.editions-ellipses.fr/product_info.php?products_id=11183
http://www.editions-ellipses.fr/product_info.php?products_id=9387.
- [140] user127.0.0.1. Parallelogram law in normed vector space without an inner product. 2014. URL <https://math.stackexchange.com/questions/641077/parallelogram-law-in-normed-vector-space-without-an-inner-product>.
- [141] Frank Jones. Chapter 7 : cross product. 2004.
- [142] Emmanuel Vieillard Baron. Déterminant d'une matrice, d'une application linéaire. 2001.
<http://www.les-mathematiques.net/b/e/d/node5.php>.
- [143] Robert Rolland. Le théorème de Müntz-szász. .
http://megamaths.perso.neuf.fr/rr/fichexo_201.pdf.
- [144] Wikipédia. Résultant — wikipédia, l'encyclopédie libre, Mai 2013. URL <http://fr.wikipedia.org/w/index.php?title=R%C3%A9sultant&oldid=90109864>. [En ligne ; Page disponible le 24-mai-2013].
- [145] Emmanuel Pedon. Cours de géométrie affine et euclidienne pour la licence de mathématiques. 23 mars 2015. URL <http://pedon.perso.math.cnrs.fr/fichiers/enseignement/CoursGeoLicence.pdf>.
- [146] Jean-Etienne Rombaldi. Polynômes d'endomorphismes en dimension finie. applications. .
<http://www-fourier.ujf-grenoble.fr/~rombaldi/AgregInterne/Oral1/110.pdf>.
- [147] Matthieu Romagny. Endomorphismes cycliques. 2012. URL https://perso.univ-rennes1.fr/matthieu.romagny/agreg/dvt/endom_cycliques.pdf.

- [148] Callus. URL <http://math.stackexchange.com/questions/1829332/prove-that-if-v-is-finite-dimensional-then-v-is-even-dimensional>.
- [149] La taverne de l'Irlandais. De la réduction des endomorphismes. 2001. URL <http://tanopah.jo.free.fr/epilogues/reduction.pdf>.
- [150] Wikipédia. Trigonalisation — wikipédia, l'encyclopédie libre, 2014. URL <http://fr.wikipedia.org/w/index.php?title=Trigonalisation&oldid=109085753>. [En ligne; Page disponible le 22-décembre-2014].
- [151] Henry C. King. Unitary diagonalization of matrices. <http://www-users.math.umd.edu/~hck/Normal.pdf>.
- [152] David Delaunay. Nilpotence. 10 juillet 2014. URL <http://mp.cpgedupuydelome.fr/pdf/Réductiondesendomorphismes-Nilpotence.pdf>.
- [153] G. Donald Allen. Lectures on linear algebras. 2004. http://www.math.tamu.edu/~dallen/m640_03c/readings.htm.
- [154] Wikipedia. Spectral theorem — wikipedia, the free encyclopedia, 2013. URL http://en.wikipedia.org/w/index.php?title=Spectral_theorem&oldid=575488135. [Online; accessed 21-October-2013].
- [155] theorem for normal triangular matrices. URL <http://planetmath.org/theoremfornormaltriangularmatrices>.
- [156] Michel Granget. Normes matricielles, conditionnement. URL http://math.univ-angers.fr/~granger/anatum/Chapitre_II.pdf.
- [157] Wikipédia. Forme quadratique — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Forme_quadratique&oldid=103591878. [En ligne; Page disponible le 7-mai-2014].
- [158] Robert Rollan. Produit tensorile d'espaces vectoriels. URL <http://robert.rolland.acrypta.com/telechargements/algebre/tensor.pdf>.
- [159] Marc Sage. Algèbre multilinéaire. . URL <http://www.normalesup.org/~sage/Cours/ProdTens.pdf>.
- [160] Ministère de l'éducation nationale. Rapport de jury de concours – agrégation de mathématiques, concours externe. 2011. <http://agreg.org/Rapports/rapport2011provisoire.pdf>.
- [161] Hervé Carrieu, Maurice Fadel, Etienne Fieux, Patrice Lassère, and Frédéric Rodriguez. Autour des matrices de Frobenius ou compagnon. 2007. <http://www.math.univ-toulouse.fr/~lassere/pdf/vfcomp.pdf>.
- [162] Grégory Vial. Autour du théorème des invariants de similitude. Octobre 2005. <http://w3.bretagne.ens-cachan.fr/math/people/gregory.vial/files/cplts/ivs.pdf>.
- [163] Arnaud Moncet. Invariants de similitude. <http://blogperso.univ-rennes1.fr/arnaud.moncet/public/IVS.pdf>.
- [164] Commutant d'un endomorphisme. . URL <http://minerve.bretagne.ens-cachan.fr/images/Commutant.pdf>.
- [165] Bernard Alken. Dans un commentaire de « dimension du commutant d'une matrice ». 2009. URL <http://www.mathoman.com/index.php/1538-dimension-du-commutant-d-une-matrice>.

- [166] Wikipédia. Forme linéaire — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Forme_lin%C3%A9aire&oldid=154647982. [En ligne; Page disponible le 9-décembre-2018].
- [167] Wikipédia. Théorème de burnsides (problème de 1902) — wikipédia, l'encyclopédie libre, 2013. URL [http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Burnside_\(probl%C3%A8me_de_1902\)&oldid=98510166](http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Burnside_(probl%C3%A8me_de_1902)&oldid=98510166). [En ligne; Page disponible le 2-novembre-2015].
- [168] Pierre Monmarché. Développements. 18 mai 2011. URL <http://laurent.claessens-donadello.eu/pdf/1-total.pdf>.
- [169] Transformation laissant invariante une forme quadratique. URL <http://www.ilemaths.net/forum-sujet-500814.html>.
- [170] Éric Jourdain. Angles. URL <https://perso.univ-rennes1.fr/eric.jourdain/GEEU/Cours/Angles.pdf>.
- [171] Maxime Pouvreau. Pseudo-réduction simultanée. 2013-2014. URL <http://minerve.bretagne.ens-cachan.fr/images/Pseudored.pdf>.
- [172] Christian Squarcini. Dégénérescence. 2005. URL http://christian-squarcini.pagesperso-orange.fr/AgregInterne/Algebrelinaire/2_2.pdf.
- [173] William F. Trench. Introduction to real analysis. 2010. http://ramanujan.math.trinity.edu/wtrench/texts/TRENCH_REAL_ANALYSIS.PDF.
- [174] Wikipédia. Norme d'opérateur — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Norme_d%27op%C3%A9rateur&oldid=125963366. [En ligne; Page disponible le 26-avril-2017].
- [175] Wikipedia. Spectral radius — wikipedia, the free encyclopedia, 2017. URL https://en.wikipedia.org/w/index.php?title=Spectral_radius&oldid=758072091. [Online; accessed 26-April-2017].
- [176] Pierre-Emmanuel Jabin. Analyse numérique, correction du TD 8. 2008-2009. URL <http://math.unice.fr/~jabin/CTD3-8.pdf>.
- [177] Christian Squarcini. Applications linéaires. . URL http://christian-squarcini.pagesperso-orange.fr/AgregInterne/Topo/4_3.pdf.
- [178] Espaces métriques et espaces normés. URL <http://www.les-mathematiques.net/a/a/b/node4.php>.
- [179] Gallouët!Thierry. Norme et conditionnement d'une matrice. URL <https://www.i2m.univ-amu.fr/~gallouet/licence.d/anum.d/anum-tg2.pdf>.
- [180] Gilles Costantini. Espaces vectoriels de dimension infinie, normes usuelles, équivalence de normes. . URL http://gilles.costantini.pagesperso-orange.fr/agreg_fichiers/evn.pdf.
- [181] David Wilkins. Normed vector spaces and functional analysis. 1997-2001. URL <http://www.maths.tcd.ie/~dwilkins/Courses/212/212PtIII.pdf>.
- [182] Gabriel Nagy. Operator theory in Hilbert spaces. URL <https://www.math.ksu.edu/~nagy/real-an/2-07-op-th.pdf>.
- [183] H. Jerome Keisler. Elementaty calculus, an infinitesimal approach. 2010. <http://www.math.wisc.edu/~keisler/keislercalc-810.pdf>.

- [184] Wikipédia. Suite arithmético-géométrique — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Suite_arithm%C3%A9tico-g%C3%A9om%C3%A9trique&oldid=88168841. [En ligne ; Page disponible le 16-juin-2013].
- [185] Prime.mover. Sum of sequence of products of consecutive reciprocals. URL https://proofwiki.org/wiki/Sum_of_Sequence_of_Products_of_Consecutive_Reciprocals.
- [186] Alan D. Sokal. A really simple elementary proof of the uniform boundedness theorem. 2010. URL <https://arxiv.org/pdf/1005.1585.pdf>.
- [187] Wikipédia. Théorème de banach-steinhaus — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Banach-Steinhaus&oldid=92377069. [En ligne ; Page disponible le 31-août-2013].
- [188] Wikipedia contributors. Strong operator topology — Wikipedia, the free encyclopedia, 2018. URL https://en.wikipedia.org/w/index.php?title=Strong_operator_topology&oldid=849989429. [Online ; accessed 16-February-2019].
- [189] Wikipedia contributors. Direct sum of modules — Wikipedia, the free encyclopedia, 2018. URL https://en.wikipedia.org/w/index.php?title=Direct_sum_of_modules&oldid=835636392. [Online ; accessed 22-June-2018].
- [190] John Douglas Moore. Othogonal complement. 2010. URL <http://web.math.ucsb.edu/~moore/orthogonalcomplements.pdf>.
- [191] Wikipedia contributors. Tensor product — Wikipedia, the free encyclopedia, 2018. URL https://en.wikipedia.org/w/index.php?title=Tensor_product&oldid=846586202. [Online ; accessed 21-June-2018].
- [192] Arjeh Cohen. Tensor product. 2004. URL <https://www.win.tue.nl/~amc/ow/lba/lba3.pdf>.
- [193] Jerome E. Marsden and Tudor S. Ratiu. Manifolds, tensor analysis and applications. 2002. URL <http://www.springer.com/gp/book/9780387967905>.
- [194] F. Laudenbach. *Calcul différentiel et intégral*. les Éd. de l'École polytechnique, 2000. ISBN 9782730207249. URL <http://books.google.fr/books?id=Ws9A7ZoRJNEC>.
- [195] Pierre Lairez. Théorème de dérivation d'une limite. Communication personnelle.
- [196] Wikipédia. Droite réelle achevée — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Droite_r%C3%A9elle_achev%C3%A9e&oldid=144800302. [En ligne ; Page disponible le 24-janvier-2018].
- [197] Ud779. Limits of functions and left hand right hand limit. 2015. URL <https://math.stackexchange.com/questions/1418673/limits-of-functions-and-left-hand-right-hand-limit>.
- [198] Wikipédia. Fonction cauchy-continue — wikipédia, l'encyclopédie libre, 2019. URL http://fr.wikipedia.org/w/index.php?title=Fonction_Cauchy-continue&oldid=157229246. [En ligne ; Page disponible le 3-mars-2019].
- [199] Provaticus. Uniforme continuité utilisee sans justification page 656. 2019. URL <https://github.com/LaurentClaessens/mazhe/issues/124>.
- [200] Wikiversité. Topologie générale/complétude — wikiversité,, 2019. URL https://fr.wikiversity.org/w/index.php?title=Topologie_g%C3%A9n%C3%A9rale/Compl%C3%A9tude&oldid=753513. [En ligne ; accédé le 24-juillet-2019].

- [201] Adrien Fontaine. Développement : théorème de Sarkiowski. 8 octobre 2013. URL <http://perso.eleves.bretagne.ens-cachan.fr/~afontain/dvptthmdesarkowski.pdf>.
- [202] Wikipédia. Théorème de heine — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Heine&oldid=145940281. [En ligne ; Page disponible le 28-février-2018].
- [203] J.L. Littlewoor. Every polynomyal has a root. 1941. URL <http://www.math.univ-toulouse.fr/~bauval/Littlewood-95-8.pdf>.
- [204] Wikipédia. Théorème fondamental de l'algèbre — wikipédia, l'encyclopédie libre, 2017. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_fondamental_de_l%27alg%C3%A8bre&oldid=139327371. [En ligne ; Page disponible le 20-août-2017].
- [205] URL https://www.math.hmc.edu/calculus/tutorials/quotient_rule/proof.pdf.
- [206] Wikiversité. Fonctions d'une variable réelle/dérivabilité — wikiversité,, 2018. URL https://fr.wikiversity.org/w/index.php?title=Fonctions_d%27une_variable_r%C3%A9elle/D%C3%A9rivabilit%C3%A9&oldid=718842. [En ligne ; accédé le 24-mai-2018].
- [207] Wikiversité. Fonction dérivée/dérivée d'un quotient — wikiversité,, 2017. URL https://fr.wikiversity.org/w/index.php?title=Fonction_d%C3%A9riv%C3%A9e/D%C3%A9riv%C3%A9e_d%27un_quotient&oldid=675129. [En ligne ; accédé le 24-mai-2018].
- [208] Wikipédia. Opérations sur les dérivées — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Op%C3%A9rations_sur_les_d%C3%A9riv%C3%A9es&oldid=99099598. [En ligne ; Page disponible le 28-juillet-2014].
- [209] Wikiversité. Fonctions d'une variable réelle/dérivabilité — wikiversité,, 2019. URL https://fr.wikiversity.org/w/index.php?title=Fonctions_d%27une_variable_r%C3%A9elle/D%C3%A9rivabilit%C3%A9&oldid=753518. [En ligne ; accédé le 19-juillet-2019].
- [210] Wikipédia. Théorème de Rolle — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Rolle&oldid=132483699. [En ligne ; Page disponible le 28-décembre-2016].
- [211] Guillaume Connan. Dérivation. URL <http://gconnan.free.fr/les%20pdf/Deriv.pdf>.
- [212] Keith Konrad. The remainder in Taylor series. URL <http://www.math.uconn.edu/~kconrad/blurbs/analysis/TaylorRemainder.pdf>.
- [213] Wikipédia. Règle de l'hôpital — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=R%C3%A8gle_de_L%27H%C3%B4pital&oldid=131491717. [En ligne ; Page disponible le 29-décembre-2016].
- [214] Wikipedia contributors. Squeeze theorem — Wikipedia, the free encyclopedia, 2019. URL https://en.wikipedia.org/w/index.php?title=Squeeze_theorem&oldid=903141752. [Online ; accessed 20-July-2019].
- [215] Livio Flaminio. Éléments de géométrie différentielle. 29 octobre 2009. URL http://www-gat.univ-lille1.fr/~flaminio/M403/2009-2010/cours_geo_top.pdf.
- [216] Wikipédia. Convergence uniforme — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Convergence_uniforme&oldid=95493463. [En ligne ; Page disponible le 15-octobre-2013].
- [217] Suites et séries de fonctions. . URL <http://blog.psi945.fr/public/maths-psi/cours-psi-suites-series-fonctions.pdf>.

- [218] F. Poupaud. Analyse fonctionnelle pour la licence. URL <http://math.unice.fr/~rascle/pdf/files/coursanapp/ana-fonc.pdf>.
- [219] Michael Gechele. Théorème de stone-weierstrass. http://michael.gechele.perso.neuf.fr/Agregation/Theoreme_Stone_Weierstrass.pdf.
- [220] Primitives et intégrales. URL http://math.univ-lyon1.fr/capes/IMG/pdf/new_primitive.pdf.
- [221] Nicole Bopp. Un complément à la leçon sur l'équation fonctionnelle de la fonction exponentielle. 2008. URL <http://irma.math.unistra.fr/~bopp/CAPES/cours/equation-felle-exp.pdf>.
- [222] John K. Hunter. Chapter 13 : Metric, normed and topological spaces. URL https://www.math.ucdavis.edu/~hunter/intro_analysis_pdf/ch13.pdf.
- [223] G.B. Folland. Taylor's formula. URL <https://sites.math.washington.edu/~folland/Math425/taylor.pdf>.
- [224] Wikipedia. Développement de taylor, 2016. URL http://fr.wikipedia.org/wiki/D%C3%A9veloppement_de_Taylor.
- [225] Wikiversité. Calcul différentiel/théorèmes utiles — wikiversité,, 2017. URL https://fr.wikiversity.org/w/index.php?title=Calcul_diff%C3%A9rentiel/Th%C3%A9or%C3%A8mes_utiles&oldid=657239. [En ligne; accédé le 10-juillet-2017].
- [226] Rached Mneimé and Frédéric Testard. *Groupes de Lie classiques*. Hermann, 1986. ISBN 2-7056-6040-2.
- [227] Yoann Gelineau. Générateurs du groupe linéaire. . URL http://math.univ-lyon1.fr/~gelineau/devagreg/Generateurs_Groupe_Lineaire.pdf.
- [228] Sandrine Caruso. Générateurs de $GL_n(K)$ et $SL_n(K)$. URL <http://sandrine.toonywood.org/pageperso/agreg/geneSL.pdf>.
- [229] Wikipédia. Théorème de cayley-hamilton — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Cayley-Hamilton&oldid=104470625. [En ligne; Page disponible le 10-juin-2014].
- [230] David Monniaux. 9 juin 2014. URL <http://david.monniaux.free.fr/dotclear/index.php/post/2014/06/09/M%C3%A9moire-%C3%A0-effa%C3%A7age-rapide>. David Monniaux esquisse la preuve dans un commentaire à cette note.
- [231] Pierre Lissy. Décomposition polaire. 30 avril 2010. URL <http://www.ljll.math.upmc.fr/~lissy/Agreg/developpements/DecPol.pdf>.
- [232] Wikipédia. Décomposition polaire — wikipédia, l'encyclopédie libre, avril 2013. URL http://fr.wikipedia.org/w/index.php?title=D%C3%A9composition_polaire&oldid=90240040. [En ligne; Page disponible le 24-avril-2013].
- [233] Jean-Etienne Rombaldi. Endomorphismes symétriques d'un espace vectoriel. applications. . <http://www-fourier.ujf-grenoble.fr/~rombaldi/AgregInterne/Oral1/120.pdf>.
- [234] Brice Loustau. Développements d'algèbre pour l'agrégation. 2008. URL <http://myismail.net/docs/divers/agreg/Sortie/LessonsAlgLoustau.pdf>.
- [235] Michel Coste. Sous-groupes compacts du groupe linéaire. 29 septembre 2006. URL <http://agreg-maths.univ-rennes1.fr/documentation/docs/ssgrpecompact.pdf>.

- [236] Daniel Li. Notions fondamentales de la théorie des probabilités. .
<http://labomathlens.free.fr/Liens/ProbaM1/PROBA01.pdf>.
- [237] Daniel Saada. Tribu de borel et tribu de baire d'un espace topologique. 2009. URL http://www.daniel-saada.eu/fichiers/16-Tribus_de_Baire.pdf.
- [238] Wikipédia. Tribu-trace — wikipédia, l'encyclopédie libre, 2018. URL <http://fr.wikipedia.org/w/index.php?title=Tribu-trace&oldid=155238311>. [En ligne ; Page disponible le 28-décembre-2018].
- [239] Daniel Li. Construction de la mesure de lebesgue. 28 janvier 2008.
http://www.editions-ellipses.fr/product_info.php?products_id=11183
http://www.editions-ellipses.fr/product_info.php?products_id=9387.
- [240] Wikipédia. Mesure (mathématiques) — wikipédia, l'encyclopédie libre, 2018. URL [http://fr.wikipedia.org/w/index.php?title=Mesure_\(math%C3%A9matiques\)&oldid=145203435](http://fr.wikipedia.org/w/index.php?title=Mesure_(math%C3%A9matiques)&oldid=145203435). [En ligne ; Page disponible le 5-février-2018].
- [241] ProofWiki. Outer measure of limit of increasing sequence of sets — proofwiki,, 2012. URL http://www.proofwiki.org/w/index.php?title=Outer_Measure_of_Limit_of_Increasing_Sequence_of_Sets&oldid=85952. [Online ; accessed 2-February-2014].
- [242] C. Antonini, JF. Quint, P. Borgnat, J. Bérard, E. Lebeau, E. Souche, A. Chateau, and O. Teytaud. Quelques résultats d'unicité. . URL <http://www.les-mathematiques.net/a/i/d/node1.php#T42>.
- [243] Daniel Li. Tribus et mesures. 24 mars 2011.
http://www.editions-ellipses.fr/product_info.php?products_id=11183
http://www.editions-ellipses.fr/product_info.php?products_id=9387.
- [244] C. Antonini, JF. Quint, P. Borgnat, J. Bérard, E. Lebeau, E. Souche, A. Chateau, and O. Teytaud. Ensembles négligeables et complétion de tribus. . URL <http://www.les-mathematiques.net/a/i/c/node1.php>.
- [245] Wikipédia. Complétion d'une mesure — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Compl%C3%A9tion_d%27une_mesure&oldid=89819650. [En ligne ; Page disponible le 27-mars-2014].
- [246] Gerald Teschl. Topics in real and functional analysis. 2013. URL <http://www.mat.univie.ac.at/~gerald/ftp/book-fa/fa.pdf>.
- [247] Emmanuel Vieillard Baron. Les fonctions mesurables. 2001. URL <http://www.les-mathematiques.net/a/i/b/node1.php>.
- [248] Wikipédia. Fonction étagée — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Fonction_%C3%A9tag%C3%A9e&oldid=104836104. [En ligne ; Page disponible le 6-septembre-2015].
- [249] B. Mauray. Intégration jusqu'au théorème de Lebesgue. 2006. URL <https://webusers.imj-prg.fr/~bernard.maurey/agreg/Textes/Lebesgue.pdf>.
- [250] John K. Hunter. Measure theory. 2011. URL https://www.math.ucdavis.edu/~hunter/measure_theory/measure_notes.pdf.
- [251] Daniel Li. Intégration sur un espace produit. 29 mars 2011.
http://www.editions-ellipses.fr/product_info.php?products_id=11183
http://www.editions-ellipses.fr/product_info.php?products_id=9387.
- [252] Rémi Peyre and Frédéric Simon. Travaux dirigés de probabilités. 2009. URL <http://www.normalesup.org/~rpeyre/pro/enseignement/td09pbas-e.pdf>.

- [253] Wikipédia. Ensemble de vitali — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Ensemble_de_Vitali&oldid=100925109. [En ligne; Page disponible le 9-septembre-2015].
- [254] christophe c. Re : mesure non nulle et existence d'un ouvert. URL <http://www.les-mathematiques.net/phorum/read.php?4,608327,608430#msg-608430>.
- [255] Daniel Choi. Mathématiques pour la mécanique. 2008. <http://www.meca.unicaen.fr/Enseignement/Document/CoursChoi/mathmeca.pdf>.
- [256] Gilles Dubois. URL http://dubois.gilles.pagesperso-orange.fr/analyse_reelle/intlimites.html.
- [257] Daniel Li. Construction de l'intégrale de lebesgue. 10 février 2011. http://www.editions-ellipses.fr/product_info.php?products_id=11183
http://www.editions-ellipses.fr/product_info.php?products_id=9387.
- [258] David Madore. Sur la rédaction des maths et la recherche de l'inambigüié. 24 octobre 2018. URL <http://www.madore.org/~david/weblog/d.2018-10-24.2562.html>.
- [259] Yao Jianfeng. Statistique et logiciels. 2010. <http://perso.univ-rennes1.fr/jian-feng.yao/pedago/statlog/>.
- [260] Daniel Li. Théorème de radon-nikodým et applications. . http://www.editions-ellipses.fr/product_info.php?products_id=11183
http://www.editions-ellipses.fr/product_info.php?products_id=9387.
- [261] Bruce Driver. Analysis tools with applications. April 10, 2003. URL http://www.math.ucsd.edu/~bdriver/240-01-02/Lecture_Notes/anal.pdf.
- [262] B. Mauray. Préparation à l'agrégation, analyse. 2011. <http://www.math.jussieu.fr/~maurey/agreg/index.html>.
- [263] Thierry Gallouët and Raphaèle Herbin. Mesure, intégration, probabilités. 2012. <http://www.cmi.univ-mrs.fr/~gallouet/licence.d/mes-int-pro.pdf>.
- [264] Bruno Demange. Théorème de changement de variables. 2012. URL <https://www-fourier.ujf-grenoble.fr/~demange/integration/2012/integration-chap5.pdf>.
- [265] Daniel Li. Changement de variables dans les intégrales sur un ouvert de \mathbb{R}^n . 8 juin 2010. URL <http://li.perso.math.cnrs.fr/textes/Integration/change-v.pdf>.
- [266] Wikipédia. Intégrale impropre — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Int%C3%A9grale_impropre&oldid=103902526. [En ligne; Page disponible le 5-août-2014].
- [267] romain Boillaud. Séries de fonctions. 2010-2011. URL <http://www.rblld.fr/cours/chap6.pdf>.
- [268] Gérard Eguether. Ex-séries entières. URL <http://iecl.univ-lorraine.fr/~Gerard.Eguether/zARTICLE/EX.pdf>.
- [269] Georges Skandalis. Analyse, résumés et exercices. 6 mars 2015. URL https://www.math.univ-paris-diderot.fr/_media/formations/prepa/agreginterne/polycopieanalyse.pdf.
- [270] Rahul Krishna. Approximations by taylor polynomials. URL <http://math.columbia.edu/~krishna/08092011.pdf>.

- [271] Vogel Thomas. Analytic functions. 2012. URL <http://www.math.tamu.edu/~tvogel/410/sect74a.pdf>.
- [272] Dan Klain. The matrix exponential (with exercices). 2018. URL <http://faculty.uml.edu/dklain/exponential.pdf>.
- [273] ENS Cachan. Théorème de stabilité de Lyapunov. URL <http://minerve.bretagne.ens-cachan.fr/images/Lyapunov.pdf>.
- [274] Rached Mneimé. *Réduction des endomorphismes*. Calvage et Mounet, 2006. ISBN 2-916352-01-5.
- [275] Matthieu Romagny. Représentations linéaires des groupes finis. 2010-2011. http://www.math.jussieu.fr/~romagny/agreg/theme/rlgf_paysage.pdf.
- [276] Sébastien Pellerin. Développements d'analyse pour l'oral de l'agrégation. . <http://pellerin.xyz/doc/agreg/analyse.pdf>.
- [277] Kevin Quirin. Théorème de Müntz. http://kevin.quirin.free.fr/Trucs/agreg/Dvt_Muntz.pdf.
- [278] Pierre Lissy. Déterminant de cauchy et application au théorème de Müntz. 5 mai 2010. <http://www.ljll.math.upmc.fr/~lissy/Agreg/developpements/CauchyMuntz.pdf>.
- [279] Benjamin Dadoun. Théorème de Müntz. URL <http://benjamin.dadoun.free.fr/muntz.pdf>.
- [280] Wikipédia. Lemme de hadamard — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Lemme_de_Hadamard&oldid=90659396. [En ligne ; Page disponible le 14-juillet-2014].
- [281] Jean-Pierre Demailly. Analyse numérique et équations différentielles. 1991.
- [282] Clémence Minazzo and Kelsey Rider. Théorèmes du point fixe et applications aux équations différentielles. 2006-2007. <http://math.unice.fr/~eaubry/Enseignement/M1/memoire.pdf>.
- [283] Ivan Nourdin. *Leçons d'analyse, probabilités algèbre et géométrie*. Dunod, 2001. ISBN 2-10-005668-9.
- [284] Sandrine Caruso. Théorème de Cauchy-Lipschitz. 2009. <http://boumbo.toonywood.org/sandrine/pageperso/agreg/cauchy-lipschitz.pdf>
Voir aussi la page toute pleine de développement
<http://boumbo.toonywood.org/sandrine/pageperso/agreg.html>.
- [285] Franck Boyer. Équations différentielles ordinaires. 2012. URL http://www.latp.univ-mrs.fr/~fboyer/Enseignement/Agreg/cours_EDO_Agreg_FBoyer.pdf.
- [286] Théo Pierron. Théorème de Cauchy-Lipschitz global. 5 janvier 2014. URL http://perso.eleves.ens-rennes.fr/~tpier758/agreg/dvpt/maths/cauchy_lipschitz.pdf.
- [287] Thierry Audibert. Le problème de Cauchy, résultats fondamentaux. URL <http://www.univ-nice.fr/Documents/PolycopiesAgreg/ED01.pdf>.
- [288] Thomas Budzinski. Théorème de Cauchy-Lipschitz. 2014. URL https://www.eleves.ens.fr/home/budzinsk/polys/Nonolympique/2014_cauchylipschitz.pdf.
- [289] Wikipédia. Théorème d'inversion locale — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_d%27inversion_locale&oldid=100975433. [En ligne ; Page disponible le 19-février-2014].

- [290] Antoine Chambert-Loir. Introduction aux groupes et algèbres de Lie. 2004-2005. URL <http://www.math.u-psud.fr/~chambert/enseignement/2004-05/lie/lie.pdf>.
- [291] E.P. van den Ban. Lie groups. 2003. URL <http://www.math.uu.nl/people/ban/lecnot.html>.
- [292] Joachim Stubbe. URL <http://mathaa.epfl.ch/prob/enseignement/analyse2/series/AnalysisE.pdf>.
- [293] Laurent Guillopé. Optimisation sous contrainte. 2015-2016. URL <https://www.math.sciences.univ-nantes.fr/~guillope/l3-osc/osc.pdf>.
- [294] Bruno Galerne. Optimisation, algorithmique. 2016-2017. URL http://www.math-info.univ-paris5.fr/~bgalerie/m1_opti_algo/poly_opti_algo.pdf.
- [295] Wikipédia. Fonction convexe — wikipédia, l'encyclopédie libre, 2019. URL http://fr.wikipedia.org/w/index.php?title=Fonction_convexe&oldid=157832047. [En ligne; Page disponible le 24-mars-2019].
- [296] Wikiversité. Fonctions convexes/définition et premières propriétés — wikiversité,, 2013. URL http://fr.wikiversity.org/w/index.php?title=Fonctions_convexes/D%C3%A9finition_et_premi%C3%A8res_propri%C3%A9t%C3%A9s&oldid=381775. [En ligne; accédé le 9-février-2014].
- [297] Wikiversité. Fonctions convexes/fonctions convexes dérivables — wikiversité,, 2013. URL http://fr.wikiversity.org/w/index.php?title=Fonctions_convexes/Fonctions_convexes_d%C3%A9rivables&oldid=379779. [En ligne; accédé le 9-février-2014].
- [298] Stephen Boyd and Lieven Vandenbergh. Convex optimization. 2004. URL http://www.stanford.edu/~boyd/cvxbook/bv_cvxbook.pdf.
- [299] Document 33 : fonctions convexes. . URL http://math.univ-lyon1.fr/capes/IMG/pdf/new_convexe.pdf.
- [300] Wikipédia. Inégalité arithmético-géométrique — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=In%C3%A9galit%C3%A9_arithm%C3%A9tico-g%C3%A9om%C3%A9trique&oldid=99682278. [En ligne; Page disponible le 1-juillet-2014].
- [301] Wikipédia. Inégalité de kantorovitch — wikipédia, l'encyclopédie libre, juillet 2014. URL http://fr.wikipedia.org/w/index.php?title=In%C3%A9galit%C3%A9_de_Kantorovitch&oldid=102417007. [En ligne; Page disponible le 1-juillet-2014].
- [302] Jean-François Burnol. Normes lp. URL <http://jf.burnol.free.fr/agregnormeslp.pdf>.
- [303] Georges Comte. Lemme de Morse. 2009-2010. URL <http://gc83.perso.sfr.fr/Agregation/LemmedeMorse.pdf>.
- [304] Wikipédia. Calcul du volume de l'hypersphère — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Calcul_du_volume_de_l%27hypersph%C3%A8re&oldid=147651804. [En ligne; Page disponible le 18-avril-2018].
- [305] Florian Bouguet. Ellipsoïde de John-Lowner. . URL <http://florian.bouguet.free.fr/doc/developpements/john-loewner.pdf>.
- [306] Martin Henk. Löwner-john ellipsoids. URL https://www.math.uni-bielefeld.de/documenta/vol-ismp/24_henk-martin.pdf.
- [307] Shuzhong Zhang. Supplement : the löwner-john ellipsoids. URL http://www.isye.umn.edu/courses/ie8534/pdf/Loewner-John_ellipsoid.pdf.

- [308] Michel Rascle. Analyse fonctionnelle pour la licence. <http://math.unice.fr/~rascle/pdf/files/coursanapp/ana-fonc.pdf>.
- [309] B. Mauray. Analyse fonctionnelle et théorie spectrale (version longue). 2001–2002. <http://www.math.jussieu.fr/~maurey/ts012/poly/lths.pdf>.
- [310] Mourad Besbes. Chapitre 5 : espaces métriques complets, espaces de banach. <http://besbes.mourad.free.fr/Enseignement/Topologie/Chapitre5.pdf>.
- [311] C. Antonini, JF. Quint, P. Borgnat, J. Bérard, E. Lebeau, E. Souche, A. Chateau, and O. Teytaud. Complété d'un espace métrique. . URL <http://www.les-mathematiques.net/a/a/b/node26.php>.
- [312] Wikipédia. Espace complet — Wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Espace_complet&oldid=91125353. [En ligne ; Page disponible le 1-avril-2013].
- [313] David Cimasoni. Cours de géométrie I, semestre de printemps. Printemps 2014. URL <http://www.unige.ch/math/folks/cimasoni/GeometrieI.pdf>.
- [314] Christophe Mourougane. Théorie des groupes et géométrie. 2008-2009. <http://perso.univ-rennes1.fr/christophe.mourougane/enseignements/2008-9/THGG/poly.groupes.pdf>.
- [315] Romain Boillaud. Séries entières. <http://www.rblld.fr/cours/chap8.pdf>.
- [316] Philippe Picart. Définition de cos et sin par les séries entières. URL http://trucsmaths.free.fr/telech/exp_cos_sin.pdf.
- [317] Alfred Gray. Modern differential geometry of curves and surfaces with mathematica. 1998. URL <http://webmath2.unito.it/paginepersonali/sergio.console/CurveSuperfici/>.
- [318] Wikipédia. Pi — wikipédia, l'encyclopédie libre, 2017. URL <http://fr.wikipedia.org/w/index.php?title=Pi&oldid=133197907>. [En ligne ; Page disponible le 3-janvier-2017].
- [319] Sylvain. Générateurs de $\simeq (e)$. URL <https://agreg-maths.fr/developpements/147>.
- [320] Groupes diédraux. <http://theoriesdesgroupes.perso.sfr.fr/cours/Diedraux.pdf>.
- [321] Pierre Lissy. 6 avril 2010. URL <http://www.ljll.math.upmc.fr/~lissy/Agreg/developpements/Died.pdf>.
- [322] Pierre Renfer. Roulettes et colliers. <http://les.mathematiques.free.fr/pdf/collier.pdf>.
- [323] Daniel Perrin. Isométries du plan. URL <https://www.math.u-psud.fr/~perrin/CAPES/geometrie/isometries11-12.pdf>.
- [324] Wikipédia. Angle — wikipédia, l'encyclopédie libre, 2016. URL <http://fr.wikipedia.org/w/index.php?title=Angle&oldid=130370952>. [En ligne ; Page disponible le 15-novembre-2016].
- [325] B. Guesmi. Théorème de l'angle inscrit. URL <http://s999cdfd874dd6835.jimcontent.com/download/version/1427398317/module/5509391568/name/theoremedel'arccapable.pdf1.pdf>.
- [326] Nicole Bopp. Un complément à la leçon sur les isométries d'un polygone régulier. Octobre 2008. URL <http://irma.math.unistra.fr/~bopp/CAPES/cours/groupe-fini-isometrie.pdf>.

- [327] Théo Pierron. Groupes de pavage du plan. 4 juin 2014. URL http://perso.eleves.ens-rennes.fr/~tpier758/agreg/dvpt/useless/gpe_pavage.pdf.
- [328] Frédéric Hélein. Géométrie euclidienne. 2005-2006. URL <https://webusers.imj-prg.fr/~frederic.helein/cours/euclide.pdf>.
- [329] Laurent Dietrich. Les sous-groupes finis de $so(3)$. 16 mai 2011. URL <http://www.math.cmu.edu/~ldietric/doc/ssgrso3.pdf>.
- [330] Sous groupes finis de $so(3,r)$. 2013-2014. URL http://minerve.bretagne.ens-cachan.fr/images/Sous-groupes_finis_S03.pdf.
- [331] Marc Sage. Sons, fréquences, harmoniques, tons : le compromis du piano. . URL <http://www.normalesup.org/~sage/Musique/Harmonie.pdf>.
- [332] Wikipédia. Chiffrement rsa — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=Chiffrement_RSA&oldid=146774926. [En ligne; Page disponible le 25-mars-2018].
- [333] URL <https://www.torproject.org/about/torusers.html.en>.
- [334] URL <https://www.torproject.org/projects/torbrowser.html.en>.
- [335] Wikipédia. Polynôme cyclotomique — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Polyn%C3%B4me_cyclotomique&oldid=93751495. [En ligne; Page disponible le 25-juin-2013].
- [336] Adam Parusinski. Théorème de Dirichlet (version faible). 2010-2011. URL <http://math.unice.fr/~parus/AGREG/Dev/dirichlet.pdf>.
- [337] Michel Merle. corps finis. 2005-2006. URL http://www-math.unice.fr/~merle/Complements/corps_finis.pdf.
- [338] Gabriel Peyré. Corps finis. 2007. URL <https://objectifagregation.github.io/>.
- [339] François Rodier. Corps finis. URL <http://iml.univ-mrs.fr/~rodier/Cours/RappelCorpsfinis.pdf>.
- [340] Stephane Vento. La loi de réciprocité quadratique. 2005. URL <http://www.proba.jussieu.fr/pageperso/nourdin/LeSiteDeLAgregatif/vento2.pdf>.
- [341] Wikipedia contributors. Chevalley–warning theorem — Wikipedia, the free encyclopedia, 2017. URL https://en.wikipedia.org/w/index.php?title=Chevalley%E2%80%9393Warning_theorem&oldid=792119700. [Online; accessed 21-May-2018].
- [342] Christian Squarcini. Anneaux factoriels. . URL http://christian-squarcini.pagesperso-orange.fr/AgregInterne/Anneauxcorps/3_2.pdf.
- [343] Wikipédia. Fonction de möbius — wikipédia, l'encyclopédie libre, mai 2013. URL http://fr.wikipedia.org/w/index.php?title=Fonction_de_M%C3%B6bius&oldid=91236321. [En ligne; Page disponible le 1-mai-2013].
- [344] Wikipédia. Corps fini — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Corps_fini&oldid=92489466. [En ligne; Page disponible le 1-mai-2013].
- [345] Wikipédia. Nombre constructible — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Nombre_constructible&oldid=103130793. [En ligne; Page disponible le 13-mai-2014].

- [346] Tracés usuels « à la règle et au compas ». 2012-2013. URL http://dokeos3.u-cergy.fr/dokeos20/courses/EEMEM2CRPE/scorm/2012_2013SCORMEEMEUE3EC2M2P1_SCORM/page_04.htm.
- [347] Wikipédia. Théorème de wantzel — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Wantzel&oldid=103220511. [En ligne ; Page disponible le 13-mai-2014].
- [348] Wikipédia. Nombre de fermat — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Nombre_de_Fermat&oldid=103080861. [En ligne ; Page disponible le 18-juin-2014].
- [349] Philippe Spindel. Éléments de géométrie différentielle pour la mécanique analytique et le gravitation. Novembre 2010. https://portail.umons.ac.be/EN2/universite/facultes/fs/services/institut_physique/mecanique_et_gravitation/Documents/Textespédagogiques/geodiff.pdf.
- [350] William K Allard. The Brouwer fixed point theorem. 2004. <http://www.math.duke.edu/~wka/math204/fixe.pdf>
see also <http://www.math.duke.edu/~wka/math204/>.
- [351] Kenneth Kuttler. Topics in analysis. Janvier 2012. <http://www.math.byu.edu/~klkuttler/SobolevSpacesB.pdf>
Juste un indice : ce document fait 2086 pages ... j'y arriverais, courage !
- [352] Hervé Le Dret. Notes de cours M2 Équations aux dérivées partielles élliptiques. 4 mars 2010. URL <https://www.ljll.math.upmc.fr/~ledret/M2Elliptique/chapitre4.pdf>.
- [353] Bachir Bekka. Théorème du point fixe de Markov-Kakutani : applications à l'existence de mesures de haar et aux chaînes de Markov. 2005. <http://agreg-maths.univ-rennes1.fr/documentation/docs/PointFixe.pdf>.
- [354] Olivier Casterá. Différentielles totales exactes. http://o.castera.free.fr/pdf/Differentielle_totale_exacte.pdf.
- [355] Intégration de fractions rationnelles. URL <http://blascheck.franck.free.fr/IMG/pdf/IntegFracRat.pdf>.
- [356] Ophélie Rouby. Théorème de Rothstein-Trager. URL http://perso.eleves.bretagne.ens-cachan.fr/~oroub842/dvlpts/Rothstein_Trager.pdf.
- [357] Frédéric Chyzak. Intégration symbolique des fractions rationnelles. Décembre 2004. URL <http://www.enseignement.polytechnique.fr/informatique/profs/Frederic.Chyzak/CF-2004/notes-IntRat.pdf>.
- [358] Sylvain Golénia. Approximation de $\ln(2)$. Novembre 2018. URL <http://mathaoutils.blogspot.com/2018/11/approximation-de-ln2.html>.
- [359] Giuseppe De Marco. *Analisi due*. 1999.
- [360] Louis Fauchier-Magnan. Introduction à la théorie de Morse. URL <http://cqfd.epfl.ch/webdav/site/cqfd/shared/projets/igat/IntroductionàlathéoriedeMorse-LouisFauchier-Magnan.pdf>.
- [361] Reduction formula for integral of power of sine. . URL https://proofwiki.org/wiki/Reduction_Formula_for_Integral_of_Power_of_Sine.
- [362] Michel Stainer. Compléments sur les séries numériques. URL <http://michel.stainer.pagesperso-orange.fr/PSIx/Cours/Ch03-ComplSeriesNum.pdf>.

- [363] Wallis's product. . URL https://proofwiki.org/wiki/Wallis%27s_Product.
- [364] Jacek Cichoń. Stirling approximation formula. URL <http://cs.pwr.edu.pl/cichon/Math/StirlingApp.pdf>.
- [365] Wikipédia. Formule de stirling — wikipédia, l'encyclopédie libre, Mai 2013. URL http://fr.wikipedia.org/w/index.php?title=Formule_de_Stirling&oldid=93224272. [En ligne; Page disponible le 20-mai-2013].
- [366] Jean-Louis Rouget. Formule de stirling. 2008. URL <https://www.maths-france.fr/MathSpe/GrandsClassiquesDeConcours/SeriesNumeriques/FormuleDeStirling.pdf>.
- [367] Stirling's formula. . URL https://proofwiki.org/wiki/Stirling%27s_Formula.
- [368] gb. Cercle circonscrit à une courbe. URL <http://www.les-mathematiques.net/phorum/read.php?8,750556,750586>.
- [369] Mohammad Ghomi. Lecture note on differential geometry. 2007. URL <http://people.math.gatech.edu/~ghomi/LectureNotes/LectureNotes5U.pdf>.
- [370] Wikipedia. Convex curve — wikipedia, the free encyclopedia, 2015. URL https://en.wikipedia.org/w/index.php?title=Convex_curve&oldid=685739052. [Online; accessed 13-April-2016].
- [371] Wikipédia. Théorème de jordan — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Jordan&oldid=124990679. [En ligne; Page disponible le 6-avril-2016].
- [372] Les quadrilatères. URL <http://ddata.over-blog.com/xxxyyy/4/39/08/01/cinquieme/5eme.parallelogramme.jeu-set-et-maths.mathematiques.pdf>.
- [373] Robert Rolland. Géométrie projective. . <http://robert.rolland.acrypta.com/telechargements/geometrie/pr0.pdf>.
- [374] Daniel Bertrand. Géométrie projective. 2009-2010. http://www.math.jussieu.fr/~bertrand/Enseignement/epdf/Algeo09_chap2.pdf.
- [375] Nicolas Jacon. Géométrie projective. URL http://njacon.perso.math.cnrs.fr/jacon_Geometrieproj.pdf.
- [376] François Labourie. Géométrie affine et projective. 9 février 2019. URL <http://math.unice.fr/~labourie/preprints/pdf/geomproj.pdf>.
- [377] Matthieu Romagny. Droites projectives, homographies. 2005-2006. URL https://perso.univ-rennes1.fr/matthieu.romagny/agreg/droite_projective_Halberstadt.pdf.
- [378] Arnaud Bodin. L'inversion. 2012. URL http://math.univ-lille1.fr/~bodin/geometrie/ch_inversion.pdf.
- [379] Laurent Dietrich. Le groupe circulaire. 21 mai 2011. URL <http://www.math.cmu.edu/~ldietric/doc/groupecirc.pdf>.
- [380] Florian Bouget. Le groupe circulaire. URL http://florian.bouguet.free.fr/doc/developpements/Groupe_circulaire.pdf.
- [381] Florian Bouguet. Le groupe circulaire. . URL http://minerve.bretagne.ens-cachan.fr/images/Groupe_circulaire.pdf.
- [382] 182 : utilisation des nombres complexes en géométrie. Homographies. 2015-2016. URL http://minerve.bretagne.ens-cachan.fr/images/182_2015-2016.pdf.

- [383] Sébastien Pellerin. Action du groupe modulaire sur le demi-plan de Poincaré. .
http://dynamaths.free.fr/docs/lecons/developpement_algebre_6.pdf.
- [384] Alexandre Alessandri. Action du groupe modulaire sur le demi-plan de Poincaré. 2011-2012.
http://minerve.bretagne.ens-cachan.fr/images/Dvt_psl2.pdf.
- [385] Aline Kurtzmann. Espérance conditionnelle.
<http://www.iecn.u-nancy.fr/~kurtzman/cours/agreg/esperance-conditionnelle.pdf>.
- [386] Wikipédia. Théorème de projection sur un convexe fermé — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_projection_sur_un_convexe_ferm%C3%A9&oldid=90246160. [En ligne; Page disponible le 3-mai-2013].
- [387] Michel Merle. Espaces de Hilbert. 2004-2005. URL http://www-math.unice.fr/~merle/Algebre_et_geometrie/hilbert.pdf.
- [388] Douglas N. Arnold. Functional analysis. 1997. URL <http://www.math.umn.edu/~arnold/502.s97/functional.pdf>.
- [389] Wikipédia. Base de hilbert. . URL https://fr.wikipedia.org/wiki/Base_de_Hilbert?oldid=128533376.
- [390] Christiane Schomblond. *Théorie quantique des champs, QED, QCD*. Notes de cours, en français.
<http://homepages.ulb.ac.be/~cschomb/intfonc+QED+QCD.pdf>.
- [391] Remsirems. Re : Réponses à quelques questions. 2016. URL <http://linuxfr.org/nodes/110155/comments/1675813>.
- [392] Jonathan Bain. Kochen-Specker and measurement. 2011.
<http://faculty.poly.edu/~jbain/philqm/lectureslides/08.KS&Measurement.pdf>.
- [393] Michel Kern. Introduction à la méthode des éléments finie. 2004-2005. URL http://mms2.ensmp.fr/ef_paris/formulation/polycop/f_coursEF.pdf.
- [394] Antoine Bensalah. Communication électronique privée. 2017.
- [395] B. Maurey. Holomorphie. 7 février 2005.
<http://www.math.jussieu.fr/~maurey/agreg/Textes/Holomorphe045.pdf>.
- [396] Pascal Dingoyan. Chapitre 5. URL <http://www.math.jussieu.fr/~dingoyan/lescours/lm367-pageweb/lm367-chapitre5.pdf>.
- [397] Wikipédia. Anneau principal — wikipédia, l'encyclopédie libre, 2017. URL http://fr.wikipedia.org/w/index.php?title=Anneau_principal&oldid=141994549. [En ligne; Page disponible le 27-octobre-2017].
- [398] Corentin. Re : formule de cauchy et dérivation. URL <http://www.les-mathematiques.net/phorum/read.php?4,385272,385289>.
- [399] Philippe Charpentier. Analyse complexe. Septembre 2010.
http://www.math.u-bordeaux1.fr/~pcharpen/enseignement/fichiers-master1/Analyse_Complexe.pdf.
- [400] Christian Kkein. Séries de laurent, résidus. URL http://math.u-bourgogne.fr/IMB/klein/Complement_dAnalyse_files/chap607.pdf.

- [401] Wikipédia. Intégrale de fresnel — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Int%C3%A9grale_de_Fresnel&oldid=132170602. [En ligne ; Page disponible le 26-novembre-2016].
- [402] David Sirajuddin. Fresnel integrals. 2008. http://itcanbeshown.com/integrals/FresnelIntegrals/fresnel_integrals.pdf.
- [403] Noach Dana-Picard. The logarithm of a complex number. URL <http://ndp.jct.ac.il/tutorials/complex/node26.html>.
- [404] Guillaume Carlier. Analyse complexe. 2012-2013. URL <https://www.ceremade.dauphine.fr/~carlier/analysecomplexe.pdf>.
- [405] Wikiversité. Fonctions d'une variable complexe/le logarithme complexe — wikiversité,, 2016. URL https://fr.wikiversity.org/w/index.php?title=Fonctions_d%27une_variable_complexe/Le_logarithme_complexe&oldid=617412. [En ligne ; accédé le 16-décembre-2016].
- [406] Jean-Marie Lion. Fonctions holomorphes d'une variable. 5 mai 2004. <http://perso.univ-rennes1.fr/jean-marie.lion/coursholo.pdf>.
- [407] Wikipédia. Théorème d'Ascoli — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_d%27Ascoli&oldid=92958118. [En ligne ; Page disponible le 11-mai-2013].
- [408] Keith Conrad. l^p spaces for $0 < p < 1$. URL <https://www.math.uconn.edu/~kconrad/blurbs/analysis/lpspace.pdf>.
- [409] B. Maurey. 9 mars 2011. <http://www.math.jussieu.fr/~maurey/IntegProba/Cours/CoursIP101-14.pdf>.
- [410] Danarmk. 2016. URL <http://linuxfr.org/nodes/110155/comments/1675589>.
- [411] B. Mauray. Convolution, inégalités, approximations et régularisation. 2007. <https://webusers.imj-prg.fr/~bernard.maurey/agreg/Textes/Convolution.pdf>
<https://webusers.imj-prg.fr/~bernard.maurey/agreg/ag045/Convolution045.pdf>.
- [412] Schilling. Minkowsky's inequality for integrals. URL <http://math.ucsd.edu/~lni/math240/suppl-1.pdf>.
- [413] Charles Suquet. Intégration, analyse de fourier, probabilités. 2003-2004. <http://math.univ-lille1.fr/~suquet/ens/IFP/Cours/cours04/CoursIFP04.html>.
- [414] Wikipédia. Théorème de Riesz-Fischer — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Riesz-Fischer&oldid=90074115. [En ligne ; Page disponible le 3-mai-2013].
- [415] JanTinbergen1991. L_p -space is a hilbert space if and only if $p=2$. Novembre 2018. URL <https://math.stackexchange.com/questions/3017814/lp-space-is-a-hilbert-space-if-and-only-if-p-2>.
- [416] .
- [417] Gilles Leborgne. Introduction à la théorie des distributions. 16 mai 2013. URL <http://www.isima.fr/~leborgne/IsimathDistributions/distrib.pdf>.
- [418] Hart Smith. Lecture 2 : convolution. 2013. URL <https://sites.math.washington.edu/~hart/m526/Lecture2.pdf>.

- [419] R. Fortet. Remarques sur les espaces uniformément convexes. 1941. URL http://www.numdam.org/article/BSMF_1941__69__23_0.pdf.
- [420] Kai-Seng Chou and Tianwen Luo. Chapter 4 : The lebesgue spaces. 2014. URL https://www.math.cuhk.edu.hk/course_builder/1415/math5011/MATH5011_Chapter_4.2014.pdf.
- [421] Olof Hanner. On the uniform convexity of l^p and l^p . 1955. URL https://projecteuclid.org/download/pdf_1/euclid.afm.1485893271.pdf.
- [422] Assaf Naor. Proof of the uniform convexity lemma. February 2004. URL <https://web.math.princeton.edu/~naor/homepagefiles/inequality.pdf>.
- [423] David Giraudo. 2011. URL <https://math.stackexchange.com/questions/95307/showing-left-fracab2-right-fracab2-right-leq-frac12/95312#95312>.
- [424] gerw. How to prove clarkson's inequality? 2017. URL <https://math.stackexchange.com/questions/1607683/how-to-prove-clarksons-inequality>.
- [425] Hamed F.A. Usman, M .A and M.O Olayiwola.
- [426] Jean-François Burnol. Le théorème de dualité des espaces l^p . 2006. URL <http://jf.burnol.free.fr/0506L312annexeDualiteLp.pdf>.
- [427] Wikipédia. Espace L_p — wikipédia, l'encyclopédie libre, 2013. URL http://fr.wikipedia.org/w/index.php?title=Espace_Lp&oldid=92618226. [En ligne; Page disponible le 10-mai-2013].
- [428] Didier Smets. Base d'analyse fonctionnelle. 2011-2012. URL https://www.ljll.math.upmc.fr/snets/MM005/MM005_Chapitre_6.pdf.
- [429] Haïm Brézis. *Analyse fonctionnelle*. Masson, Paris, 1983. ISBN 2-225-77198-7. Théorie et applications.
- [430] MoebiusCorzer Giraudo, Davide. Show smooth functions of compact support are dense in the Schwartz space. 2012. URL <http://math.stackexchange.com/questions/229751/show-smooth-functions-of-compact-support-are-dense-in-the-schwartz-space>.
- [431] Richard Gomez. Séries de Fourier. 10 octobre 2010. URL <http://megamaths.perso.neuf.fr/rg/fourier.pdf>.
- [432] David Delaunay. Séries de Fourier. 8 juin 2012. <http://mp.cpgedupuydelome.fr/pdf/SériesdeFourier.pdf>.
- [433] Wikipédia. Théorème de jordan — wikipédia, l'encyclopédie libre, Mai 2013. URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Jordan&oldid=92525205. [En ligne; Page disponible le 25-mai-2013].
- [434] B. Maurey. Notes de cours : analyse hilbertienne et de Fourier. 2008-2009. <http://www.math.jussieu.fr/~maurey/FourHilb/PolyFH.pdf>.
- [435] Bachir Bekka. Le théorème d'inversion de Fourier du point de vue des distributions. URL <http://agreg-maths.univ-rennes1.fr/documentation/docs/Theoreme-Inversion-Fourier.pdf>.
- [436] Erdős László. Fourier transform. 2008. URL <https://www.mathematik.uni-muenchen.de/~lerdos/WS08/QM/four08.pdf>.
- [437] Richard Zekri. Analyse de fourier. 15 septembre 2009. URL <http://lumimath.univ-mrs.fr/infoetudiant/ANALYSEDEFOURIERMASTER2009.pdf>.

- [438] Michel Crouzeix. Compléments d'analyse. URL <https://perso.univ-rennes1.fr/michel.crouzeix/publis/ane.pdf>.
- [439] Isabelle Gallagher. Chapitre b. distributions. URL <https://www.math.ens.fr/~gallagher/ChapitreB2017-2018.pdf>.
- [440] Daniel Li. Notions sur les distributions. . URL <http://labomathlens.free.fr/Liens/AF/distrib.pdf>.
- [441] Claude Zuily. *Éléments de distributions et d'équations aux dérivées partielles*. Dunod, 2002.
- [442] Sabin Lessard. Cours de processus stochastiques. 2013. URL <http://www.dms.umontreal.ca/~lessards/ProcessusStochastiquesLessard2014.pdf>.
- [443] Peter Young. Singular Fourier transforms and the integral representation of the Dirac delta function. November, 2009.
- [444] Daniel Choi. Méthode des éléments-finis par l'exemple. Avril 2010. URL <http://meca.unicaen.fr/~choi/pdf/cours-mef.pdf>.
- [445] Michael Reiter and Arthur Schuste. Fourier transform and Sobolev spaces. 2008. URL http://www.mat.univie.ac.at/~stein/lehre/SoSem08/sobolev_fourier.pdf.
- [446] Norbert Heuer. Ltcc module computation methods. URL <http://www.ltcc.ac.uk/courses/AppliedComputationalMethods/ACM2007Chap1.pdf\etautresfichiersdelamemeserie>.
- [447] Arkhnor. Densité des fonctions Cinf à support compact. 2014. URL <http://www.ilemaths.net/sujet-densite-des-fonctions-c-infty-a-support-compact-587937.html>.
- [448] Gilles Leborgne. Complément : espaces de Sobolev fractionnaires. 2007. URL http://ws.isima.fr/~leborgne/IsimathDistributions/sobolev_frac.pdf.
- [449] Richard B. Melrose. Lecture notes for 18.155, fall 2002. 2002. URL <https://ocw.mit.edu/courses/mathematics/18-155-differential-analysis-fall-2004/lecture-notes/section10.pdf>.
- [450] Julien Vovelle. Équations différentielles, cours no 2. URL <http://math.univ-lyon1.fr/~vovelle/2Cours.pdf>.
- [451] A. Munnier. Théorie des équations différentielles ordinaires. 2006-2007. URL http://www.iecn.u-nancy.fr/~munnier/files/cours_edo.pdf.
- [452] Olivier Debarre. Équations différentielles. URL <http://www.math.ens.fr/~debarre/EquaDiff.pdf>.
- [453] Sylvie Benzoni-Gavage. Propagation d'ondes. URL <http://math.univ-lyon1.fr/~benzoni/Ondes.pdf>.
- [454] Sonia Fliss. L'équation de transport à coefficients variables. URL http://perso.ensta-paristech.fr/~fliss/Sonia_Fliss_web_page/Enseignement_files/Amphi3_2014.pdf.
- [455] Sylvie Benzoni-Gavage. Équations différentielles ordinaires. 11 mai 2007. URL <http://math.univ-lyon1.fr/~benzoni/ED0-M1.pdf>.
- [456] Wikipédia. Flot (mathématiques) — wikipédia, l'encyclopédie libre, 2017. URL [http://fr.wikipedia.org/w/index.php?title=Flot_\(math%C3%A9matiques\)&oldid=134130419](http://fr.wikipedia.org/w/index.php?title=Flot_(math%C3%A9matiques)&oldid=134130419). [En ligne ; Page disponible le 1-février-2017].

- [457] Alexander Grigor'yan. Ordinary differential equation. 2009. URL <https://www.math.uni-bielefeld.de/~grigor/odelec2009.pdf>.
- [458] Christopher P. Grant. Theory of ordinary differential equation. 2007. URL http://www.math.pitt.edu/~bard/bardware/classes/2920/Grant_4july2007.pdf.
- [459] John B. Picard-lindelöf theorem with parameter : what regularity ? 2017. URL <http://math.stackexchange.com/questions/2126539/picard-lindelöf-theorem-with-parameter-what-regularity>.
- [460] Kevin Quirin. Développements. 2011-2012. URL <http://kevin.quirin.free.fr/Trucs/agreg/developpements.pdf>.
- [461] Wikipédia. Équations de lotka-volterra — wikipédia, l'encyclopédie libre, 2015. URL http://fr.wikipedia.org/w/index.php?title=%C3%89quations_de_Lotka-Volterra&oldid=113756354. [En ligne ; Page disponible le 5-juillet-2015].
- [462] Predrag Cvitanović. Partial differential equations. 2012. URL <http://www.cns.gatech.edu/~predrag/courses/PHYS-6124-12/StGoChap6.pdf>.
- [463] Ralph Chill. Introduction aux équations aux dérivées partielles. 2010-2011. URL <http://www.math.univ-metz.fr/~chill/edp.pdf>.
- [464] Gianluca Bontempi. Un cours d'analyse numérique et de matlab. 2009. <http://www.ulb.ac.be/di/map/gbonte/calcul/>
<https://web.archive.org/web/20090203224115/http://www.ulb.ac.be/di/map/gbonte/calcul/>.
- [465] Wikipedia. Two's complement — wikipedia, the free encyclopedia, 2016. URL https://en.wikipedia.org/w/index.php?title=Two%27s_complement&oldid=706335033. [Online ; accessed 6-April-2016].
- [466] Wikipédia. Ieee 754 — wikipédia, l'encyclopédie libre, 2015. URL http://fr.wikipedia.org/w/index.php?title=IEEE_754&oldid=120234746. [En ligne ; Page disponible le 12-avril-2016].
- [467] Comprendre les nombres à virgules flottantes. URL <http://blog.netinfluence.ch/2009/09/24/comprendre-les-nombres-a-virgule-flottante/>.
- [468] Harald Schmidt. IEEE754 converter. URL <http://www.h-schmidt.net/FloatConverter/IEEE754.html>.
- [469] David Goldberg. What every computer scientist should know about floating-point arithmetic. 1991. URL http://docs.oracle.com/cd/E19957-01/806-3568/ngc_goldberg.html#11655.
- [470] Emmanuel Frénod. Calcul scientifique. 2011-2012. URL <http://web.univ-ubs.fr/lmam/frenod/IMG/DocEtudiant/MTH1504/calculscientifique.pdf>.
- [471] Franck Boyer. Agrégation externe de mathématiques. Analyse numérique. 14 octobre 2014. URL http://www.math.univ-toulouse.fr/~fboyer/_media/enseignements/agreg/cours_an_agreg_fboyer_2014.pdf.
- [472] Antoine Chambert-Loir. Autour de la méthode de newton. <http://perso.univ-rennes1.fr/antoine.chambert-loir/2005-06/agreg/newton.pdf>.
- [473] Wikipédia. Vitesse de convergence — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=Vitesse_de_convergence&oldid=110252913. [En ligne ; Page disponible le 22-mai-2016].

- [474] Martin Pinto Campos. Approximations polynomiales de fonctions et et des intégrales. 21 décembre 2006. D'après le cours de Laurence Halpern
<http://documents.lamacs.fr/cours/mac1/cours-macs1-pinto1.pdf>.
- [475] Ming Yang. Matrix decomposition. URL http://users.eecs.northwestern.edu/~mya671/files/Matrix_YM_.pdf.
- [476] Benjamin Ambrosio. Chapitre 3 : Méthodes directes de résolution du système $Ax = b$. URL <http://lmah.univ-lehavre.fr/~ambrosio/CoursMN/Chapitre3.pdf>.
- [477] Emmanuel Frénod. Méthodes directes de résolution des systèmes linéaires. URL <http://web.univ-ubs.fr/lmam/frenod/IMG/DocEtudiant/MTH1504/methodes-directes.pdf>.
- [478] Christian Lécot. Analyse numérique matricielle. 2009-2010. URL https://www.lama.univ-savoie.fr/~lecot/data/M1_ANM.pdf.
- [479] Wikipédia. Matrice à diagonale dominante — wikipédia, l'encyclopédie libre, 2016. URL http://fr.wikipedia.org/w/index.php?title=Matrice_%C3%A0_diagonale_dominante&oldid=132268176. [En ligne ; Page disponible le 4-mai-2017].
- [480] Habib Joulak. Méthodes des différences finies en élasticité. 1974. URL <https://ori-nuxeo.univ-lille1.fr/nuxeo/site/esupversions/0e05f3d2-fe12-4d67-b74f-f622869539f2>.
- [481] Philippe Briand. Probabilités de base. 2005-2006.
http://www.lama.univ-savoie.fr/~briand/proba/g12_cours.pdf.
- [482] Vincent Bansaye. Variables aléatoires, espérance, indépendance.
<http://www.cmap.polytechnique.fr/~bansaye/CoursTD2.pdf>.
- [483] Alfio Marazzi. There are 22 chapters in 22 files.
<http://www.iunsp.ch/Unites/us/Alfio/polybiostat/ch06.pdf>.
- [484] Cédric Boutillier. 2016. URL <https://github.com/LaurentClaessens/mazhe/issues/16#issue-180279120>.
- [485] Automaths. Paradoxe des deux enfants – episode 2! 25 août 2018. URL <https://automaths.blog/2018/08/25/paradoxe-des-deux-enfants-episode-2/>.
- [486] Wikipédia. Inégalité de jensen — wikipédia, l'encyclopédie libre, août 2013. URL http://fr.wikipedia.org/w/index.php?title=In%C3%A9galit%C3%A9_de_Jensen&oldid=90235295. [En ligne ; Page disponible le 9-août-2013].
- [487] Gérard Letac. Calcul des probabilités, Deug 2ième année. juin 2001.
<http://les.mathematiques.free.fr/pdf/proba.zip>.
- [488] Daniel Saada. Note sur le théorème de transfert pour les variables aléatoires réelles. 18 octobre 2012.
http://www.les-mathematiques.net/phorum/file.php?12,file=17435,filename=th_du_tranfert.pdf.
- [489] Christophe Sabot. Différents types de convergence de suites de v.a. URL <http://math.univ-lyon1.fr/~sabot/Convergence.pdf>.
- [490] Antonin Macé and Thomas Pradeau. Exposé de maîtrise. Estimations de densités de probabilité. 2007. Encadré par Gilles Stoltz.
<http://www.eleves.ens.fr/home/amace/travaux/exposemaîtrise.pdf>.
- [491] Wikipedia. Continuous mapping theorem — wikipedia, the free encyclopedia, 2013. URL http://en.wikipedia.org/w/index.php?title=Continuous_mapping_theorem&oldid=545523937. [Online ; accessed 15-June-2013].

- [492] Wikipédia. Inégalité de markov — wikipédia, l'encyclopédie libre, 2014. URL http://fr.wikipedia.org/w/index.php?title=In%C3%A9galit%C3%A9_de_Markov&oldid=100208179. [En ligne ; Page disponible le 11-février-2014].
- [493] Charles Suquet. Simulation. 2005-2006. <http://math.univ-lille1.fr/~suquet/ens/Agr/simul06.pdf>.
- [494] Yoann Gelineau. Vecteurs gaussiens. . http://math.univ-lyon1.fr/~gelineau/files/vecteurs_gaussiens.pdf.
- [495] Wikipédia. Nombre normal — wikipédia, l'encyclopédie libre, avril 2013. URL http://fr.wikipedia.org/w/index.php?title=Nombre_normal&oldid=89848922. [En ligne ; Page disponible le 28-avril-2013].
- [496] Wikipédia. Théorème de glivenko-cantelli — wikipédia, l'encyclopédie libre, . URL http://fr.wikipedia.org/w/index.php?title=Th%C3%A9or%C3%A8me_de_Glivenko-Cantelli&oldid=90841113. [En ligne ; Page disponible le 15-juin-2013].
- [497] Vincent Rivoirard and Gilles Stoltz. Estimation de densité de probabilités. URL <http://www.math.ens.fr/statenaction/PDF/Densite-Principal.pdf>.
- [498] Nils Berglund. Chaînes de Markov. Décembre 2007. URL <http://www.univ-orleans.fr/mapmo/membres/berglund/markov.pdf>.
- [499] Arnaud Guyader. Espérance conditionnelle et chaînes de Markov. 200-2008. <http://www.sites.univ-rennes2.fr/laboratoire-statistique/AGUYADER/doc/proba/poly.pdf>.
- [500] Massimiliano Gubinelli. Comportement asymptotique des martingales. 2011-2012. <http://www.ceremade.dauphine.fr/~mgubi/e1112/pd3.pdf>.
- [501] Jean Lacroix and Pierre Priouret. Probabilités approfondies. 2005-2006. <http://www.proba.jussieu.fr/cours/processus.pdf>.
- [502] Wikipedia. Optional stopping theorem — wikipedia, the free encyclopedia, 2013. URL http://en.wikipedia.org/w/index.php?title=Optional_stopping_theorem&oldid=547547174. [Online ; accessed 14-August-2013].
- [503] Nizar Touzi. Martingales en temps discret et chaîne de Markov. Septembre 2009. URL <http://www.cmap.polytechnique.fr/~touzi/MAP432-Poly.pdf>.
- [504] Wikipédia. Étoile de kleene — wikipédia, l'encyclopédie libre, 2018. URL http://fr.wikipedia.org/w/index.php?title=%C3%89toile_de_Kleene&oldid=149113546. [En ligne ; Page disponible le 1-juin-2018].
- [505] Zach Weinersmith. 2008. URL <http://www.smbc-comics.com/comic/2008-02-21>.