

# Analyse numérique

Thomas Cluzeau

École Nationale Supérieure d'Ingénieurs de Limoges

16 rue d'atlantis, Parc ester technopole

87068 Limoges CEDEX

cluzeau@ensil.unilim.fr

[http://www.unilim.fr/pages\\_perso/thomas.cluzeau](http://www.unilim.fr/pages_perso/thomas.cluzeau)





# Table des matières

<b>1</b>	<b>Arithmétique des ordinateurs et analyse d'erreurs</b>	<b>7</b>
1.1	L'arithmétique flottante . . . . .	7
1.1.1	Le système des nombres à virgule flottante . . . . .	7
1.1.2	Représentation effective des réels et sa formalisation . . . . .	9
1.1.3	Unité d'erreur d'arrondi, estimation d'erreurs . . . . .	10
1.1.4	Modèle de l'arithmétique flottante . . . . .	10
1.2	L'analyse d'erreurs . . . . .	10
1.2.1	Non-associativité . . . . .	10
1.2.2	Erreurs d'arrondi sur une somme . . . . .	11
1.2.3	Erreurs d'arrondi sur un produit . . . . .	11
1.2.4	Phénomènes de compensation . . . . .	11
1.2.5	Phénomènes d'instabilité numérique . . . . .	12
1.2.6	Erreur amont et erreur aval . . . . .	13
1.2.7	Outils théoriques de l'analyse d'erreurs . . . . .	13
<b>2</b>	<b>Résolution d'un système d'équations linéaires (Partie 1) : méthodes directes</b>	<b>15</b>
2.1	Introduction et motivation . . . . .	15
2.1.1	Objet . . . . .	15
2.1.2	Motivation . . . . .	16
2.1.3	Résolution d'un système triangulaire . . . . .	17
2.1.4	Les méthodes directes étudiées . . . . .	18
2.2	Méthode de Gauss et factorisation LU . . . . .	19
2.2.1	Description de la méthode . . . . .	19
2.2.2	Point de vue numérique : stratégies de choix du pivot . . . . .	21
2.2.3	Lien avec la factorisation LU d'une matrice . . . . .	23
2.2.4	Coût de l'algorithme . . . . .	26
2.3	Méthode de Cholesky . . . . .	27
2.4	Méthode de Householder et factorisation QR . . . . .	29
2.4.1	Transformation (élémentaire) de Householder . . . . .	29
2.4.2	Principe de la méthode de Householder . . . . .	30
2.4.3	Exemple de résolution d'un système linéaire par la méthode de Householder . . . . .	30

2.4.4	Factorisation QR d'une matrice . . . . .	31
<b>3</b>	<b>Conditionnement d'une matrice pour la résolution d'un système linéaire</b>	<b>33</b>
3.1	Normes matricielles . . . . .	33
3.1.1	Normes vectorielles . . . . .	33
3.1.2	Normes matricielles et normes subordonnées . . . . .	33
3.2	Conditionnement d'une matrice . . . . .	34
3.2.1	Exemple classique . . . . .	34
3.2.2	Définition du conditionnement . . . . .	35
3.2.3	Estimation théorique de l'erreur a priori . . . . .	37
3.2.4	Estimation théorique de l'erreur a posteriori . . . . .	38
<b>4</b>	<b>Résolution d'un système d'équations linéaires (Partie 2) : méthodes itératives</b>	<b>39</b>
4.1	Motivation . . . . .	39
4.2	Notions générales . . . . .	41
4.2.1	Modèle général d'un schéma itératif . . . . .	41
4.2.2	Convergence . . . . .	42
4.2.3	Vitesse de convergence . . . . .	43
4.3	Les méthodes itératives classiques . . . . .	44
4.3.1	Principe . . . . .	44
4.3.2	Méthode de Jacobi . . . . .	45
4.3.3	Méthode de Gauss-Seidel . . . . .	46
4.3.4	Méthode de relaxation . . . . .	47
4.3.5	Résultats de convergence dans des cas particuliers . . . . .	47
4.4	Méthode du gradient conjugué . . . . .	48
4.4.1	Méthodes du gradient . . . . .	49
4.4.2	Méthode de la plus forte pente . . . . .	51
4.4.3	Gradient conjugué . . . . .	51
4.4.4	Gradient conjugué avec préconditionnement . . . . .	55
<b>5</b>	<b>Interpolation polynomiale</b>	<b>59</b>
5.1	Le problème considéré . . . . .	59
5.2	La méthode d'interpolation de Lagrange . . . . .	60
5.3	Effectivité de l'interpolation : interpolant de Newton . . . . .	63
5.3.1	Base d'interpolation de Newton . . . . .	63
5.3.2	Expression de l'interpolant de Newton . . . . .	63
5.3.3	Algorithme de calcul des différences divisées . . . . .	65
5.4	Erreur d'interpolation . . . . .	65
<b>6</b>	<b>Intégration numérique</b>	<b>67</b>
6.1	Introduction et méthodes classiques . . . . .	67
6.2	Formalisation de l'intégration approchée . . . . .	70

6.3	Formules de Newton-Côtes . . . . .	71
6.4	Stabilité des méthodes d'intégration . . . . .	72
6.5	Formules d'intégration composées . . . . .	72
<b>7</b>	<b>Résolution d'équations et de systèmes d'équations non linéaires</b>	<b>75</b>
7.1	Méthode de dichotomie . . . . .	76
7.2	Méthode du point fixe . . . . .	77
7.3	Méthode de Newton . . . . .	80
7.4	Méthode de la sécante . . . . .	82
7.5	Systèmes d'équations non linéaires . . . . .	83



# Chapitre 1

## Arithmétique des ordinateurs et analyse d'erreurs

Dans tout ce cours, nous manipulerons des nombres réels. L'objet de ce premier chapitre est de décrire comment ces nombres réels sont représentés dans un ordinateur.

### 1.1 L'arithmétique flottante

#### 1.1.1 Le système des nombres à virgule flottante

**Théorème 1.1.** *Soit  $\beta$  un entier strictement supérieur à 1. Tout nombre réel  $x$  non nul peut se représenter sous la forme*

$$x = \text{sgn}(x) \beta^e \sum_{k \geq 1} \frac{d_k}{\beta^k},$$

où  $\text{sgn}(x) \in \{+, -\}$  est le signe de  $x$ , les  $d_k$  sont des entiers tels que  $0 < d_1 \leq \beta - 1$  et  $0 \leq d_k \leq \beta - 1$  pour  $k \geq 2$ , et  $e \in \mathbb{Z}$ . De plus, cette écriture est unique (sauf pour les décimaux :  $2,5 = 2,499999\dots$ ).

D'ordinaire, nous utilisons le système décimal, *i.e.*,  $\beta = 10$  et les chiffres 0, 1, 2, 3, 4, 5, 6, 7, 8, 9 pour les  $d_k$ . Nous avons par exemple :

- $0,0038 = 0,38 \cdot 10^{-2} = +10^{-2} \left( \frac{3}{10} + \frac{8}{10^2} \right)$ .  
Remarque : en MATLAB, on peut écrire  $0.38e-2$  au lieu de  $0.38 * 10^{(-2)}$ .
- $\frac{1}{7} = 0,142857\dots = +10^0 \left( \frac{1}{10} + \frac{4}{10^2} + \frac{2}{10^3} + \frac{8}{10^4} + \dots \right)$ . Notons que le développement décimal d'un nombre rationnel est périodique (ici,  $\frac{1}{7} = 0, \mathbf{142857}142857142857\dots$ ).
- $-\sqrt{2} = -1,4142\dots = -10^1 \left( \frac{1}{10} + \frac{4}{10^2} + \frac{1}{10^3} + \frac{4}{10^4} + \dots \right)$ .
- $\pi = 3,14159\dots = +10^1 \left( \frac{3}{10} + \frac{1}{10^2} + \frac{4}{10^3} + \frac{1}{10^4} + \dots \right)$ .

Historiquement, le choix  $\beta = 10$  est lié à une particularité anatomique de la race humaine (nous avons 10 doigts). Les ordinateurs utilisent quant à eux  $\beta = 2$  (numération binaire),  $\beta = 8$  (numération octale), ou encore  $\beta = 16$  (numération hexadécimale).

Remarquons que l'on perd l'unicité si on autorise  $d_1 = 0$  : en effet, on a par exemple :

$$\begin{aligned} 0,0038 &= 0,38 \cdot 10^{-2} = +10^{-2} \left( \frac{3}{10} + \frac{8}{10^2} \right) \\ &= 0,038 \cdot 10^{-3} = +10^{-1} \left( \frac{0}{10} + \frac{3}{10^2} + \frac{8}{10^3} \right). \end{aligned}$$

On définit l'ensemble  $F \subset \mathbb{R}$  par :

$$F = \left\{ y \in \mathbb{R} \mid y = \pm \beta^e \left( \frac{d_1}{\beta} + \frac{d_2}{\beta^2} + \dots + \frac{d_t}{\beta^t} \right), e_{\min} \leq e \leq e_{\max} \right\},$$

ou encore

$$F = \left\{ y \in \mathbb{R} \mid y = \pm m \beta^{e-t}, e_{\min} \leq e \leq e_{\max} \right\}.$$

Ceci correspond aux deux écritures  $0,0038 = +10^{-2} \left( \frac{3}{10} + \frac{8}{10^2} \right) = +38 \cdot 10^{-4}$  avec  $e = -2$ ,  $t = 2$ ,  $e - t = -4$ . Le nombre  $m$  s'appelle *la mantisse* et on utilise la notation  $m = \overline{d_1 d_2 \dots d_t}^\beta$ .

Notons que  $0 \notin F$ .

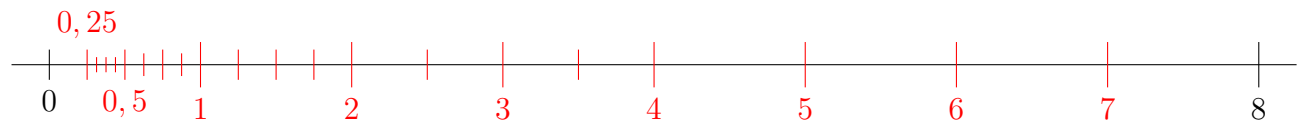
Pour  $y \neq 0$ , on a  $m \beta^{e-t} = \beta^e \left( \frac{d_1}{\beta} + \frac{d_2}{\beta^2} + \dots + \frac{d_t}{\beta^t} \right) \geq \beta^e \frac{1}{\beta}$  car  $d_1 \geq 1$ . D'où  $m \geq \beta^{t-1}$ . D'autre part,  $m = \overline{d_1 d_2 \dots d_t}^\beta = d_1 \beta^{t-1} + \dots + d_{t-k} \beta^k + \dots + d_{t-1} \beta + d_t < \beta^t$ . On a donc montré que

$$\beta^{t-1} \leq m < \beta^t.$$

$F$  est un *système de nombres à virgule flottante* (floating point number system) noté  $F(\beta, t, e_{\min}, e_{\max})$ . Il dépend de quatre paramètres :

1. *la base*  $\beta$  (chiffres utilisés  $0, 1, \dots, \beta - 1$ ),
2. *la précision*  $t$  (nombre de chiffres utilisés pour représenter la mantisse),
3.  $e_{\min}$  et  $e_{\max}$  qui définissent *le domaine des exposants*.

Par exemple, pour  $F(2, 3, -1, 3)$ , on obtient les nombres représentés en rouge sur la figure ci-dessous :



On constate que l'écart entre deux nombres consécutifs est multiplié par 2 à chaque puissance de 2.

Dans le standard IEEE 754 utilisé par MATLAB, on a  $\beta = 2$  et :



- en simple précision :  $t = 24$ ,  $e_{\min} = -125$ ,  $e_{\max} = 128$ ,
- en double précision :  $t = 53$ ,  $e_{\min} = -1021$ ,  $e_{\max} = 1024$ .

**Définition 1.2.** On appelle epsilon machine et on note  $\epsilon_M$  la distance de 1 au nombre flottant suivant.

Par exemple, pour  $F(2, 3, -1, 3)$ , on a  $\epsilon_M = 0, 25$ . Dans MATLAB, c'est `eps`.

**Proposition 1.3.** Pour  $F(\beta, t, e_{\min}, e_{\max})$ , on a  $\epsilon_M = \beta^{1-t}$ .

*Démonstration.* On a  $1 = \beta \frac{1}{\beta} = \beta^{1-t} \overline{10 \dots 0}^\beta$ . Le nombre suivant dans le système de nombres à virgule flottante  $F(\beta, t, e_{\min}, e_{\max})$  est alors  $\beta^{1-t} \overline{10 \dots 1}^\beta = \beta^{1-t} (\beta^{t-1} + 1) = 1 + \beta^{1-t}$ .  $\square$

**Lemme 1.4.** Dans le système de nombres à virgule flottante  $F(\beta, t, e_{\min}, e_{\max})$ , l'écart  $|y - x|$  entre un nombre flottant  $x$  (non nul) et un nombre flottant  $y$  (non nul) adjacent vérifie  $\beta^{-1} \epsilon_M |x| \leq |y - x| \leq \epsilon_M |x|$ .

*Démonstration.* On a  $x = m \beta^{e-t}$ , donc  $|y - x| = 1 \beta^{e-t}$ . Or  $\beta^{t-1} \leq m < \beta^t$  donc  $m \beta^{-t} < 1$  et  $m \beta^{1-t} \geq 1$ . Il vient donc  $m \beta^{-t} \beta^{e-t} < 1 \beta^{e-t} \leq m \beta^{1-t} \beta^{e-t}$  d'où le résultat puisque  $\epsilon_M = \beta^{1-t}$ .  $\square$

## 1.1.2 Représentation effective des réels et sa formalisation

- *Représentation « physique »* : par exemple, en simple précision 32 bits (bit = binary digit), 8 bits sont réservés à l'exposant et 24 bits (dont 1 pour le signe) à la mantisse. En double précision 64 bits, 11 bits sont réservés à l'exposant et 53 bits (dont 1 pour le signe) à la mantisse.

- Arrondi :

1. par troncature : par exemple avec 3 chiffres, 0,8573... devient 0,857.
2. au plus près : 0,8573... devient 0,857.
3. au représentant le plus proche dont la dernière décimale est paire (rounding to even) : 0,8573... devient 0,858.

- Formalisation :

**Définition 1.5.** Soit  $G = G(\beta, t) = \{y \in \mathbb{R} \mid y = \pm m \beta^{e-t}\}$  sans conditions sur l'exposant  $e$ . L'application  $\text{fl} : \mathbb{R} \rightarrow G$ ,  $x \mapsto \text{fl}(x)$  est appelée opération d'arrondi.

Étant donné un domaine  $F(\beta, t, e_{\min}, e_{\max})$ , il y a alors dépassement de capacité si :

1.  $|\text{fl}(x)| > \max\{|y| \mid y \in F\}$ . On parle d'« overflow ».
2.  $|\text{fl}(x)| < \min\{|y| \mid y \in F\}$ . On parle d'« underflow ».

Sinon,  $x$  est dans le domaine de  $F$ .

### 1.1.3 Unité d'erreur d'arrondi, estimation d'erreurs

**Définition 1.6.** Soit  $x$  un réel et  $\bar{x}$  une valeur approchée de  $x$ . L'erreur absolue  $e$  est défini par  $e = |x - \bar{x}|$ . L'erreur relative est  $|\frac{e}{x}|$ . Le pourcentage d'erreur est l'erreur relative multipliée par 100.

En pratique, on ne connaît en général pas la valeur exacte  $x$  (c'est le cas dans la plupart des mesures physiques) mais on peut souvent avoir une idée de l'erreur maximale  $e$  que l'on a pu commettre : dans ce cas, on majore la quantité  $|\frac{e}{x}|$ .

**Théorème 1.7** (Estimation de l'erreur d'arrondi). Soit  $x$  un réel. Si  $x$  est dans le domaine  $F(\beta, t, e_{\min}, e_{\max})$ , alors il existe  $\delta \in \mathbb{R}$  avec  $|\delta| < u = \frac{1}{2}\beta^{1-t} = \frac{1}{2}\epsilon_M$  tel que  $\text{fl}(x) = x(1 + \delta)$ .

*Démonstration.* Admis pour ce cours. □

L'erreur relative sur l'arrondi est égale à  $|\delta| < u$  : le nombre  $u$  s'appelle *unité d'erreur d'arrondi*.

Par exemple, dans le standard IEEE 754 utilisé par MATLAB, on a  $u = 2^{-24} \approx 5,96 \cdot 10^{-8}$  en simple précision et  $u = 2^{-53} \approx 1,11 \cdot 10^{-16}$  en double précision.

### 1.1.4 Modèle de l'arithmétique flottante

Le modèle suivant de l'arithmétique flottante est celui utilisé par le standard IEEE.

**Modèle Standard** : Soit  $x, y \in F(\beta, t, e_{\min}, e_{\max})$ . Pour  $\text{op} \in \{+, -, \times, \div, \sqrt{\cdot}\}$ , on définit  $x \boxed{\text{op}} y = \text{fl}(x \text{ op } y) = (x \text{ op } y)(1 + \delta)$  avec  $|\delta| < u = \frac{1}{2}\beta^{1-t} = \frac{1}{2}\epsilon_M$ .

Nous allons maintenant nous intéresser aux erreurs faites par  $\boxed{\text{op}}$ .

## 1.2 L'analyse d'erreurs

### 1.2.1 Non-associativité

En général, contrairement à  $\text{op}$ , l'opération  $\boxed{\text{op}}$  n'est pas associative. Ceci est dû aux erreurs d'arrondi. Par exemple, supposons que les réels soient calculés avec 3 chiffres significatifs et arrondis à la décimale la plus proche et cherchons à calculer la somme  $x \boxed{+} y \boxed{+} z$  avec  $x = 8,22$ ,  $y = 0,00317$  et  $z = 0,00432$ .

- $x \boxed{+} y = 8,22$  donc  $(x \boxed{+} y) \boxed{+} z = 8,22$ ,
- $y \boxed{+} z = 0,01$  donc  $x \boxed{+} (y \boxed{+} z) = 8,23$ .

### 1.2.2 Erreurs d'arrondi sur une somme

Supposons que l'on souhaite calculer une somme  $S = u_1 + u_2 + \dots + u_n$  de  $n$  réels positifs dans  $F(\beta, t, e_{\min}, e_{\max})$ . On calcule alors les sommes partielles  $S_i$  par la récurrence  $S_0 = 0$ ,  $S_i = S_{i-1} + u_i$ . Si l'on suppose les  $u_i$  connus exactement, alors les erreurs d'arrondi  $\Delta S_i$  commises sur le calcul des sommes partielles  $S_i$  vérifient  $\Delta S_i \leq \Delta S_{i-1} + \delta(S_{i-1} + u_i) = \Delta S_{i-1} + \delta S_i$  où  $|\delta| < u$ . L'erreur globale sur  $S = S_n$  vérifie donc

$$\Delta S \leq \delta(S_2 + \dots + S_n),$$

ou encore

$$\Delta S \leq \delta(u_n + 2u_{n-1} + 3u_{n-2} + \dots + (n-1)u_2 + (n-1)u_1).$$

On voit donc que pour minimiser cette erreur on a tout intérêt à sommer d'abord les termes les plus petits (Cf exemple de la sous-section précédente).

### 1.2.3 Erreurs d'arrondi sur un produit

Supposons que l'on souhaite calculer un produit  $P = u_1 u_2 \dots u_n$  de  $n$  réels positifs dans  $F(\beta, t, e_{\min}, e_{\max})$ . On calcule alors les produits  $P_i$  par la récurrence  $P_0 = 1$ ,  $P_i = P_{i-1} u_i$ . Si l'on suppose les  $u_i$  connus exactement, alors les erreurs d'arrondi  $\Delta P_i$  commises sur le calcul des produits  $P_i$  vérifient  $\Delta P_i \leq \Delta P_{i-1} u_i + \delta(S_{i-1} u_i) = \Delta P_{i-1} u_i + \delta P_i$  où  $|\delta| < u$ . L'erreur globale sur  $P = P_n$  vérifie donc

$$\Delta P \leq (k-1)\delta P_n.$$

On voit donc que contrairement au cas de l'addition, la majoration de l'erreur ne dépend pas de l'ordre des facteurs.

### 1.2.4 Phénomènes de compensation

Les phénomènes de compensation sont ceux qui se produisent lorsque l'on tente de soustraire des nombres très proches. Nous illustrons ces phénomènes sur deux exemples et donnons des astuces pour les contourner.

**Exemple 1 :** On considère l'expression  $E = \sqrt{x+1} - \sqrt{x}$  avec  $x > 0$ . Sous MATLAB, en calculant  $E$  pour  $x = 10^9$ , on va obtenir  $1,5811.10^{-5}$  mais pour  $x = 10^{16}$ , on va obtenir 0 ! Si l'on remarque que  $E = \frac{1}{\sqrt{x+1} + \sqrt{x}}$ , alors, en utilisant cette nouvelle formule, on trouvera  $E = 1,5811.10^{-5}$  pour  $x = 10^9$  et  $E = 5,000.10^{-9}$  pour  $x = 10^{16}$ .

**Exemple 2 :** On considère l'équation du second degré  $x^2 - 1634x + 2 = 0$ . Supposons que les calculs soient effectués avec 10 chiffres significatifs. Les formules habituelles donnent  $\Delta' = (\frac{1634}{2})^2 - 2 = 667487$ ,  $\sqrt{\Delta'} = 816,9987760$ , d'où les solutions

$$x_1 = \frac{1634}{2} + \sqrt{\Delta'} = 817 + 816,9987760 = 1633,998776,$$

$$x_2 = \frac{1634}{2} - \sqrt{\Delta'} = 817 - 816,9987760 = 0,0012240.$$

On voit donc qu'en procédant ainsi on a une perte de 5 chiffres significatifs sur  $x_2$ . Pour y remédier, on peut utiliser la relation  $x_1 x_2 = 2$  et calculer

$$x_2 = \frac{2}{x_1} = \frac{2}{1633,998776} = 0,001223991125.$$

### 1.2.5 Phénomènes d'instabilité numérique

Les phénomènes d'instabilité numérique sont des phénomènes d'amplification d'erreur d'arrondi. Ils se produisent en général pour des calculs récurrents ou itératifs. Nous illustrons ces phénomènes sur deux exemples.

**Exemple 1 :** On considère l'intégrale

$$I_n = \int_0^1 \frac{x^n}{10+x} dx, \quad n \in \mathbb{N},$$

que l'on cherche à évaluer numériquement. Un calcul direct montre que  $I_0 = \ln(\frac{11}{10})$ . De plus, on a

$$I_n = \int_0^1 \frac{x}{10+x} x^{n-1} dx = \int_0^1 \left(1 - \frac{10}{10+x}\right) x^{n-1} dx = \frac{1}{n} - 10 I_{n-1}.$$

On peut donc calculer successivement les valeurs de  $I_n$  en utilisant la récurrence  $I_0 = \ln(\frac{11}{10})$ ,  $I_n = \frac{1}{n} - 10 I_{n-1}$ . Numériquement, cela conduit à des résultats très mauvais. Cela provient du fait que l'erreur d'arrondi  $\Delta I_n$  sur le calcul de  $I_n$  vérifie  $\Delta I_n \approx 10 \Delta I_{n-1}$ , si l'on néglige l'erreur d'arrondi sur  $\frac{1}{n}$ . On voit donc que l'erreur croît exponentiellement : l'erreur sur  $I_0$  est multipliée par  $10^n$  sur  $I_n$ . Par conséquent cette formule de récurrence ne peut pas nous permettre de calculer la valeur de  $I_{36}$  par exemple. Pour remédier à ce problème, on peut *renverser la récurrence* c'est-à-dire considérer la formule :

$$I_{n-1} = \frac{1}{10} \left( \frac{1}{n} - I_n \right).$$

Toujours en négligeant l'erreur d'arrondi sur  $\frac{1}{n}$ , on obtient alors  $\Delta I_{n-1} \approx \frac{1}{10} \Delta I_n$ . En utilisant l'encadrement  $10 \leq 10+x \leq 11$  pour  $x \in [0, 1]$ , on montre que

$$\frac{1}{11(n+1)} \leq I_n \leq \frac{1}{10(n+1)}.$$

L'approximation  $I_n \approx \frac{1}{11(n+1)}$  nous permet alors de calculer une valeur de départ pour notre récurrence renversée. Par exemple, si l'on part de  $I_{46} \approx \frac{1}{11(46+1)}$ , on obtiendra pour  $I_{36}$  une erreur relative meilleure que  $10^{-10}$ .

On constate ici l'importance du coefficient d'amplification d'erreur pour ce genre de calcul.

**Exemple 2 :** On considère la suite définie par :

$$\begin{cases} u_0 = 2, \\ u_1 = -4, \\ u_n = 111 - \frac{1130}{u_{n-1}} + \frac{3000}{u_{n-1} u_{n-2}}, \end{cases}$$

introduite par J.-M. Muller. On peut alors montrer que la limite de cette suite est égale à 6 et malgré cela, quelque soit le système et la précision utilisés, cette suite semblera tendre vers 100. L'explication de ce phénomène étrange est assez simple : la solution générale de la récurrence  $u_n = 111 - \frac{1130}{u_{n-1}} + \frac{3000}{u_{n-1} u_{n-2}}$  est donnée par :

$$u_n = \frac{\alpha 100^{n+1} + \beta 6^{n+1} + \gamma 5^{n+1}}{\alpha 100^n + \beta 6^n + \gamma 5^n},$$

où  $\alpha$ ,  $\beta$  et  $\gamma$  dépendent des valeurs initiales  $u_0$  et  $u_1$ . Par conséquent, si  $\alpha \neq 0$ , la suite converge vers 100 et sinon (si  $\beta \neq 0$ ) la suite converge vers 6. Dans notre exemple, les valeurs initiales  $u_0 = 2$  et  $u_1 = -4$  correspondent à  $\alpha = 0$ ,  $\beta = -3$  et  $\gamma = 4$ . Par conséquent, la limite exacte de la suite est 6. Mais à cause des erreurs d'arrondi, même les premiers termes calculés seront différents des termes exacts et donc la valeur de  $\alpha$  correspondant à ces termes calculés sera très petite mais non-nulle ce qui suffira à faire en sorte que la suite converge vers 100 au lieu de 6.

### 1.2.6 Erreur amont et erreur aval

Considérons un problème que l'on résout à l'aide d'un algorithme numérique et appelons  $f$  la fonction qui fait correspondre à l'entrée  $x$  de l'algorithme la solution algébrique  $y = f(x)$ . En pratique, compte tenu des erreurs d'arrondis, étant donnée une entrée  $x$ , nous allons obtenir une sortie  $\bar{y}$  qui sera distincte de la solution algébrique  $y = f(x)$ . L'*erreur aval* est alors la différence entre le résultat  $\bar{y}$  obtenu et la solution algébrique  $y$ . L'*erreur amont* ou *erreur inverse* est le plus petit  $\delta x$  tel que la solution algébrique  $f(x + \delta x)$  correspondant à l'entrée  $x + \delta x$  soit égale à  $\bar{y}$ . Ces deux erreurs sont liées par le *conditionnement* (voir Chapitre 3) : l'erreur aval étant du même ordre que l'erreur amont multipliée par le conditionnement. L'erreur amont est en général plus intéressante que l'erreur aval car elle nous renseigne sur le problème qui est réellement résolu par l'algorithme numérique. De plus, en pratique, nous ne connaissons en général qu'une valeur approchée de l'entrée (par exemple obtenue par une mesure).

### 1.2.7 Outils théoriques de l'analyse d'erreurs

Considérons la formule  $(x \times y) + z$  pour trois réels  $x$ ,  $y$  et  $z$  d'un domaine  $F(\beta, t, e_{\min}, e_{\max})$ .

On a alors :

$$\begin{aligned} \text{fl}((x \times y) + z) &= [\text{fl}(x \times y) + z] (1 + \delta_1) \\ &= [(x \times y) (1 + \delta_2) + z] (1 + \delta_1) \\ &= (x \times y) (1 + \delta_2) (1 + \delta_1) + z (1 + \delta_1), \end{aligned}$$

d'après le théorème 1.7, avec de plus  $|\delta_i| < u$ , pour  $i = 1, 2$ .

**Lemme 1.8.** *Si pour tout  $i = 1, \dots, k$ , on a  $|\delta_i| < u$  et si  $ku < 1$ , alors il existe  $\theta_k$  tel que  $|\theta_k| \leq \frac{ku}{1-ku}$  et  $\prod_{i=1}^k (1 + \delta_i) \leq 1 + \theta_k$ .*

*Démonstration.* La démonstration se fait par récurrence. □

En utilisant la notation  $\langle k \rangle \ll = \gg \prod_{i=1}^k (1 + \delta_i)$  avec  $\langle j \rangle \cdot \langle k \rangle = \langle j + k \rangle$ , on a alors

$$\begin{aligned} \text{fl}((x \times y) + z) &= (x \times y) \langle 2 \rangle + z \langle 1 \rangle \\ &\leq (x \times y) \left(1 + \frac{2u}{1-2u}\right) + z \left(1 + \frac{u}{1-u}\right). \end{aligned}$$

# Chapitre 2

## Résolution d'un système d'équations linéaires (Partie 1) : méthodes directes

Beaucoup de problèmes se réduisent à la résolution numérique d'un système d'équations linéaires. Il existe deux grandes classes de méthodes pour résoudre ce type de systèmes :

1. les méthodes *directes* qui déterminent explicitement la solution après un nombre fini d'opérations arithmétiques,
2. les méthodes *itératives* qui consistent à générer une suite qui converge vers la solution du système.

Remarquons que les méthodes itératives ne s'appliquent que dans le cas de systèmes à coefficients dans  $\mathbb{R}$  ou  $\mathbb{C}$  mais pas dans le cas des corps finis  $\mathbb{F}_p$ . Notons qu'il existe aussi des méthodes intermédiaires telles que les méthodes de *Splitting* ou de décomposition incomplètes, et des méthodes probabilistes comme celle de Monte-Carlo qui ne seront pas abordées dans ce cours.

Dans ce deuxième chapitre nous nous intéressons aux méthodes directes. Les méthodes itératives seront abordées au chapitre 4.

### 2.1 Introduction et motivation

#### 2.1.1 Objet

On considère un système  $(S)$  de  $n$  équations à  $n$  inconnues de la forme

$$(S) \begin{cases} a_{1,1} x_1 + a_{1,2} x_2 + \cdots + a_{1,n} x_n = b_1, \\ a_{2,1} x_1 + a_{2,2} x_2 + \cdots + a_{2,n} x_n = b_2, \\ \vdots \\ a_{n,1} x_1 + a_{n,2} x_2 + \cdots + a_{n,n} x_n = b_n. \end{cases}$$

Les données sont les coefficients  $a_{i,j}$  du système qui appartiennent à un corps  $\mathbb{K}$  avec  $\mathbb{K} = \mathbb{R}$  ou  $\mathbb{C}$  ainsi que les coefficients du second membre  $b_1, \dots, b_n$ . Les inconnues sont  $x_1, \dots, x_n$  qui

appartiennent à  $\mathbb{K}$ . Un système est dit *homogène* si son second membre est nul c'est-à-dire si tous les  $b_i$  sont nuls. Nous écrirons ce système sous forme matricielle

$$(S) \quad Ax = b,$$

$$\text{avec } A = \begin{pmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ a_{n,1} & \cdots & \cdots & a_{n,n} \end{pmatrix} \in \mathbb{M}_{n \times n}(\mathbb{K}), \quad x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{K}^n \text{ et } b = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} \in \mathbb{K}^n.$$

Dans tout ce chapitre, on supposera que la matrice  $A$  est **inversible** !

## 2.1.2 Motivation

Le lecteur peut se demander pourquoi le problème de la résolution d'un tel système se pose alors que les formules de Cramer nous donnent la solution d'un tel système (S) :

$$\forall i \in \{1, \dots, n\}, \quad x_i = \frac{\begin{vmatrix} a_{1,1} & \cdots & a_{1,(i-1)} & b_1 & a_{1,(i+1)} & \cdots & a_{1,n} \\ \vdots & & & \vdots & & & \vdots \\ a_{n,1} & \cdots & a_{n,(i-1)} & b_n & a_{n,(i+1)} & \cdots & a_{n,n} \end{vmatrix}}{\det(A)}.$$

Pour comprendre le problème, essayons de compter le nombre d'opérations nécessaires pour calculer la solution en utilisant ces formules de Cramer. Nous devons calculer  $n + 1$  déterminants de taille  $n$ . On sait que le déterminant de  $A$  est donné par la formule

$$\det(A) = \sum_{\sigma \in \mathcal{S}_n} \epsilon_\sigma a_{\sigma(1),1} a_{\sigma(2),2} \cdots a_{\sigma(n),n},$$

où  $\mathcal{S}_n$  est l'ensemble des permutations de  $\{1, \dots, n\}$  et  $\epsilon_\sigma$  est la signature de la permutation  $\sigma$  ( $\epsilon_\sigma \in \{-1, 1\}$ ). L'ensemble  $\mathcal{S}_n$  contient  $n!$  éléments donc le calcul de  $\det(A)$  se ramène à  $n! - 1$  additions et  $n!(n - 1)$  multiplications soit au total  $n! - 1 + n!(n - 1) = nn! - 1$  opérations à virgule flottante. Comme nous avons  $n + 1$  déterminants à calculer, le nombre d'opérations nécessaires pour résoudre le système à l'aide des formules de Cramer est de  $(n + 1)(nn! - 1)$  opérations à virgule flottante qui est équivalent quand la dimension  $n$  tend vers l'infini à  $n(n + 1)!$ . Essayons d'évaluer cette quantité lorsque  $n = 100$ . Pour ceci on peut utiliser la formule de Stirling  $n! \sim n^{n+\frac{1}{2}} e^{-n} \sqrt{2\pi}$  qui est très précise pour évaluer la factorielle :

$$\begin{aligned} 100.101! &= 100.101.100! \\ &\simeq 100.101.100^{100,5} \cdot e^{-100} \cdot \sqrt{2\pi} \\ &\simeq 100.101.100^{100,5} \cdot 10^{-43,43} \cdot \sqrt{2\pi} \quad (\log_{10}(e) \simeq 0,4343) \\ &\simeq 10^{205-44} \cdot 10^{0,57} \cdot \sqrt{2\pi} \\ &\simeq 9,4 \cdot 10^{161}. \end{aligned}$$





De façon analogue, lorsque  $A$  est triangulaire inférieure, on obtient *l'algorithme de substitution progressive (ou substitution avant)*.

**Proposition 2.1.** *La résolution d'un système d'équations linéaires triangulaire se fait en  $n^2$  opérations à virgule flottante.*

*Démonstration.* Calculer  $x_n$  nécessite 1 opération (division), calculer  $x_{n-1}$  nécessite 3 opérations (une multiplication, une soustraction et une division), et calculer  $x_i$  nécessite  $2(n-i)+1$  opérations ( $n-i$  multiplications,  $(n-i-1)$  additions, 1 soustraction et 1 division). Au total, le nombre d'opérations est donc

$$\sum_{i=1}^n (2(n-i)+1) = 2 \sum_{i=1}^n (n-i) + n = 2 \sum_{j=0}^{n-1} j + n = 2 \frac{(n-1)n}{2} + n = n^2.$$

□

**Proposition 2.2.** *Soient  $A, B \in \mathbb{M}_{n \times n}(\mathbb{K})$  deux matrices triangulaires supérieures. On a alors les résultats suivants :*

1.  $AB$  est triangulaire supérieur ;
2. Si  $A$  et  $B$  sont à diagonale unité (i.e., n'ont que des 1 sur la diagonale), alors  $AB$  est à diagonale unité ;
3. Si  $A$  est inversible, alors  $A^{-1}$  est aussi triangulaire supérieure ;
4. Si  $A$  est inversible et à diagonale unité, alors  $A^{-1}$  est aussi à diagonale unité.

*Démonstration.* Exercice.

□

## 2.1.4 Les méthodes directes étudiées

Dans la suite de ce chapitre, nous allons considérer les trois méthodes directes suivantes pour la résolution de  $(S) : Ax = b$  :

1. Méthode de Gauss : le principe est de réduire le système à  $(MA)x = Mb$  avec  $MA$  triangulaire supérieure sans calculer explicitement  $M$ . On se ramène donc à la résolution d'un système triangulaire supérieur. Cette méthode est associée à la factorisation  $A = LU$  de la matrice  $A$  avec  $L$  triangulaire inférieure (**L**ower) et  $U$  triangulaire supérieure (**U**pper). Étant donnée une telle factorisation  $A = LU$ , on peut résoudre le système  $Ax = b$  en résolvant successivement les deux systèmes triangulaires  $Ly = b$  puis  $Ux = y$ .

En MATLAB, on peut appliquer la méthode de Gauss pour résoudre  $(S) : Ax = b$  en tapant  $\mathbf{x} = \mathbf{A} \setminus \mathbf{b}$  (attention, le symbole entre la matrice  $A$  et le vecteur colonne  $b$  est « backslash » et non un « slash » !) et calculer la factorisation LU d'une matrice en utilisant  $[\mathbf{L}, \mathbf{U}] = \text{lu}(\mathbf{A})$ .

2. Méthode de Cholesky associée à la factorisation de Cholesky  $A = R^T R$  avec  $R$  triangulaire supérieure. Cette méthode est valable pour une matrice  $A$  symétrique et définie positive (voir Définitions 2.14 et 2.16). On résout alors  $Ax = b$  en résolvant successivement les systèmes triangulaires  $R^T y = b$  puis  $Rx = y$ .

En MATLAB, on peut calculer la factorisation de Cholesky d'une matrice  $A$  en tapant `R=chol(A)`. On obtient un message d'erreur lorsque  $A$  n'est pas symétrique définie positive.

3. Méthode de Householder associée à la factorisation  $A = QR$  avec  $R$  triangulaire supérieure et  $Q$  orthogonale (voir Définition 2.20). La matrice  $Q$  est un produit de  $n - 1$  « matrices de Householder »  $H_i$  (voir Définition 2.19). Le système  $Ax = b$  s'écrit alors  $H_{n-1} \cdots H_2 H_1 Ax = H_{n-1} \cdots H_2 H_1 b$  que l'on résout facilement grâce au fait que le produit  $H_{n-1} \cdots H_2 H_1 A$  est une matrice triangulaire supérieure.

En MATLAB, on calcule la factorisation QR d'une matrice  $A$  en tapant `[Q,R]=qr(A)`.

## 2.2 Méthode de Gauss et factorisation LU

### 2.2.1 Description de la méthode

On considère le système linéaire  $(S) : Ax = b$  en supposant toujours que  $A$  est inversible et on pose  $b^{(1)} = b$  et  $A^{(1)} = A = (a_{i,j}^{(1)})_{1 \leq i,j \leq n}$ . Le système  $(S)$  s'écrit alors  $A^{(1)}x = b^{(1)}$  que l'on note  $S^{(1)}$ .

**Étape 1 :** Puisque  $A$  est inversible, quitte à échanger la première ligne de  $A^{(1)}$  avec une autre, on peut supposer que  $a_{1,1}^{(1)} \neq 0$ . Le nombre  $a_{1,1}^{(1)}$  est le premier *pivot* de l'élimination de Gauss. Pour  $i = 2, \dots, n$ , on multiplie la première équation de  $(S^{(1)})$  par  $g_{i,1} = \frac{a_{i,1}^{(1)}}{a_{1,1}^{(1)}}$  et on retranche l'équation obtenue à la  $i$ ème équation de  $(S^{(1)})$ . La  $i$ ème ligne  $L_i$  de  $(S^{(1)})$  devient donc  $L_i - g_{i,1} L_1$ . On obtient alors un nouveau système  $(S^{(2)}) : A^{(2)}x = b^{(2)}$  avec :

$$\begin{cases} a_{1,j}^{(2)} = a_{1,j}^{(1)}, & j = 1, \dots, n, \\ a_{i,1}^{(2)} = 0, & i = 2, \dots, n, \\ a_{i,j}^{(2)} = a_{i,j}^{(1)} - g_{i,1} a_{1,j}^{(1)}, & i, j = 2, \dots, n, \\ b_1^{(2)} = b_1^{(1)}, \\ b_i^{(2)} = b_i^{(1)} - g_{i,1} b_1^{(1)}, & i = 2, \dots, n. \end{cases}$$

La matrice  $A^{(2)}$  et le vecteur  $b^{(2)}$  sont donc de la forme :

$$A^{(2)} = \begin{pmatrix} a_{1,1}^{(1)} & a_{1,2}^{(1)} & \cdots & a_{1,n}^{(1)} \\ 0 & a_{2,2}^{(2)} & \cdots & a_{2,n}^{(2)} \\ 0 & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ 0 & a_{n,2}^{(2)} & \cdots & a_{n,n}^{(2)} \end{pmatrix}, \quad b^{(2)} = \begin{pmatrix} b_1^{(1)} \\ b_2^{(2)} \\ \vdots \\ b_n^{(2)} \end{pmatrix}.$$

**Étape  $k$  :** On a ramené le système à  $(S^{(k)}) : A^{(k)} x = b^{(k)}$  avec

$$A^{(k)} = \begin{pmatrix} a_{1,1}^{(1)} & & \cdots & \cdots & a_{1,k}^{(1)} & \cdots & a_{1,n}^{(1)} \\ 0 & a_{2,2}^{(2)} & & & a_{2,k}^{(2)} & \cdots & a_{2,n}^{(2)} \\ 0 & 0 & a_{3,3}^{(3)} & & a_{3,k}^{(3)} & \cdots & a_{3,n}^{(3)} \\ \vdots & \ddots & \ddots & \ddots & \vdots & & \vdots \\ 0 & \cdots & 0 & 0 & a_{k,k}^{(k)} & \cdots & a_{k,n}^{(k)} \\ \vdots & & \vdots & 0 & a_{k+1,k}^{(k)} & \cdots & a_{k+1,n}^{(k)} \\ \vdots & & \vdots & \vdots & \vdots & & \vdots \\ 0 & \cdots & 0 & 0 & a_{n,k}^{(k)} & \cdots & a_{n,n}^{(k)} \end{pmatrix}.$$

On peut alors se ramener au cas  $a_{k,k}^{(k)} \neq 0$  et  $a_{k,k}^{(k)}$  est le  $k$ ème *pivot* de l'élimination de Gauss. Par le même principe qu'à l'étape 1 et en utilisant  $g_{i,k} = \frac{a_{i,k}^{(k)}}{a_{k,k}^{(k)}}$  pour  $i > k$ , on obtient alors  $(S^{(k+1)}) : A^{(k+1)} x = b^{(k+1)}$  avec

$$A^{(k+1)} = \begin{pmatrix} a_{1,1}^{(1)} & & \cdots & \cdots & a_{1,k+1}^{(1)} & \cdots & \cdots & a_{1,n}^{(1)} \\ 0 & a_{2,2}^{(2)} & & & a_{2,k}^{(2)} & \cdots & \cdots & a_{2,n}^{(2)} \\ 0 & 0 & a_{3,3}^{(3)} & & a_{3,k}^{(3)} & \cdots & \cdots & a_{3,n}^{(3)} \\ \vdots & \ddots & \ddots & \ddots & \vdots & & & \vdots \\ 0 & \cdots & 0 & 0 & a_{k,k}^{(k)} & \cdots & \cdots & a_{k,n}^{(k)} \\ \vdots & & \vdots & 0 & 0 & a_{k+1,k+1}^{(k+1)} & \cdots & a_{k+1,n}^{(k+1)} \\ \vdots & & \vdots & \vdots & \vdots & \vdots & & \vdots \\ 0 & \cdots & 0 & 0 & 0 & a_{n,k+1}^{(k+1)} & \cdots & a_{n,n}^{(k+1)} \end{pmatrix}.$$

**Étape  $n - 1$**  : Le système  $(S^{(n)}) : A^{(n)} x = b^{(n)}$  obtenu est triangulaire supérieure avec

$$A^{(n)} = \begin{pmatrix} a_{1,1}^{(1)} & \dots & \dots & a_{1,n}^{(1)} \\ 0 & a_{2,2}^{(2)} & & a_{2,n}^{(2)} \\ 0 & 0 & a_{3,3}^{(3)} & a_{3,n}^{(3)} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & 0 & a_{n,n}^{(n)} \end{pmatrix},$$

et peut donc être résolu par l'algorithme de substitution rétrograde de la sous-section 2.1.3.

**Exemple** : Considérons le système suivant :

$$(S) = (S^{(1)}) \begin{cases} x_1 + 2x_2 + 5x_3 = 1, \\ 3x_1 + 2x_2 - x_3 = \frac{1}{2}, \\ 5x_2 + 3x_3 = 1. \end{cases}$$

Le premier pivot de l'élimination de Gauss est donc  $a_{1,1}^{(1)} = 1$  et on a  $g_{2,1}^{(1)} = 3$ ,  $g_{3,1}^{(1)} = 0$ . La première étape fournit donc

$$(S^{(2)}) \begin{cases} x_1 + 2x_2 + 5x_3 = 1, \\ -4x_2 - 16x_3 = -\frac{5}{2}, \\ 5x_2 + 3x_3 = 1. \end{cases}$$

Le second pivot de l'élimination de Gauss est donc  $a_{2,2}^{(2)} = -4$  et on a  $g_{3,2}^{(2)} = -\frac{5}{4}$ . On obtient donc le système

$$(S^{(3)}) \begin{cases} x_1 + 2x_2 + 5x_3 = 1, \\ -4x_2 - 16x_3 = -\frac{5}{2}, \\ -17x_3 = -\frac{17}{8}. \end{cases}$$

On résout alors ce système triangulaire supérieur par l'algorithme de substitution rétrograde pour trouver  $x_1 = x_2 = x_3 = \frac{1}{8}$  (voir l'exemple de la sous-section 2.1.3).

## 2.2.2 Point de vue numérique : stratégies de choix du pivot

Au cours de l'exécution de l'élimination de Gauss, si on tombe sur un pivot nul, alors on permute la ligne en question avec une ligne en dessous pour se ramener à un pivot non nul (ceci est toujours possible car  $A$  est supposée inversible). Cependant, certains choix de pivots peuvent s'avérer plus judicieux que d'autres.

### Exemple :

Considérons le système  $(S) : Ax = b$  où

$$A = \begin{pmatrix} \alpha & 1 \\ 1 & 1 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 2 \end{pmatrix},$$

avec  $\alpha$  réel non nul. On suppose de plus  $\alpha \neq 1$  de sorte que  $A$  est inversible. La solution de ce système est alors donnée par  $x_1^* = \frac{1}{1-\alpha}$ ,  $x_2^* = \frac{1-2\alpha}{1-\alpha}$ . Supposons maintenant que  $\alpha$  est « très petit » ( $0 \ll \alpha < 1$ ) et appliquons l'élimination de Gauss décrite dans la sous-section précédente. Le premier pivot est alors  $\alpha$  et  $g_{2,1} = \frac{1}{\alpha}$ . La première étape transforme donc le système  $(S)$  en  $(S^{(2)}) : A^{(2)}x = b^{(2)}$  avec

$$A^{(2)} = \begin{pmatrix} \alpha & 1 \\ 0 & 1 - \frac{1}{\alpha} \end{pmatrix}, \quad b^{(2)} = \begin{pmatrix} 1 \\ 2 - \frac{1}{\alpha} \end{pmatrix}.$$

La deuxième équation entraîne  $-\frac{1}{\alpha}x_2 \approx -\frac{1}{\alpha}$  d'où  $x_2 \approx 1$  et la première nous donne alors  $x_1 \approx 0$  ce qui est faux. L'erreur ne provient pas seulement du fait que  $\alpha$  est « très petit » car si on multiplie la première ligne par une puissance de 10 quelconque, on va trouver la même erreur. Notons  $x_2 = x_2^* + \delta x_2$  où  $|\delta x_2|$  est l'erreur absolue sur  $x_2$ . On a alors

$$x_1 = \frac{1 - x_2}{\alpha} = \frac{1 - x_2^*}{\alpha} - \frac{\delta x_2}{\alpha},$$

et on voit donc que l'erreur  $\delta x_1 = \frac{1}{\alpha} \delta x_2$  sur  $x_1$  est très amplifiée par rapport à celle sur  $x_2$ . L'anomalie provient du déséquilibre entre les coefficients de  $x_1$  et  $x_2$  de la ligne du pivot. Pour y remédier, échangeons maintenant les deux lignes de notre système et appliquons l'élimination de Gauss avec 1 comme premier pivot. On obtient alors

$$A^{(2)} = \begin{pmatrix} 1 & 1 \\ 0 & 1 - \alpha \end{pmatrix}, \quad b = \begin{pmatrix} 2 \\ 1 - 2\alpha \end{pmatrix},$$

qui entraîne alors  $x_2 \approx 1$  et  $x_1 \approx 1$  ce qui est correct.

### Élimination de Gauss à pivot partiel

À l'étape  $k$ , on échange les lignes  $k$  et  $k'$  ( $k' \geq k$ ) de  $A^{(k)}$  de telle sorte que  $|a_{k,k}^{(k)}| = \max\{|a_{i,k}^{(k)}|, i \geq k\}$ .

Par exemple, pour le système  $(S)$  de la sous-section précédente, à l'étape 1, on va permuter les lignes 1 et 2 et considérer

$$(S') : \begin{cases} 3x_1 + 2x_2 - x_3 = \frac{1}{2}, \\ x_1 + 2x_2 + 5x_3 = 1, \\ 5x_2 + 3x_3 = 1. \end{cases}$$

## Élimination de Gauss à pivot total

À l'étape  $k$ , on échange à la fois les lignes  $k$  et  $k'$  ( $k' \geq k$ ) et les colonnes  $k$  et  $k''$  ( $k'' \geq k$ ) de telle sorte que :  $|a_{k,k}^{(k)}| = \max\{|a_{i,j}^{(k)}|, i \geq k, j \geq k\}$ .

**Attention !** Si on fait des échanges de colonnes cela modifie l'ordre des composantes du vecteur solution  $x$  donc il faut penser à rétablir le bon ordre des composantes à la fin.

Par exemple, pour le système  $(S)$  de la sous-section précédente, à l'étape 1, on va permuter les colonnes 1 et 3 et considérer

$$(S') : \begin{cases} 5x_3 + 2x_2 + x_1 = 1, \\ -x_3 + 2x_2 + 3x_1 = \frac{1}{2}, \\ 3x_3 + 5x_2 = 1. \end{cases}$$

### 2.2.3 Lien avec la factorisation LU d'une matrice

**Définition 2.3.** Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$ . On appelle factorisation LU de  $A$  une factorisation  $A = LU$  avec  $L \in \mathbb{M}_{n \times n}(\mathbb{K})$  triangulaire inférieure et  $U \in \mathbb{M}_{n \times n}(\mathbb{K})$  triangulaire supérieure.

**Lemme 2.4.** Avec les notations de la sous-section 2.2.1, à l'étape  $k$  de l'élimination de Gauss, on a  $A^{(k+1)} = G_k A^{(k)}$  où

$$G_k = \begin{pmatrix} 1 & (0) & & 0 & \dots & 0 \\ & \ddots & & \vdots & & \vdots \\ & & (0) & 1 & 0 & \dots & 0 \\ 0 & \dots & 0 & -g_{k+1,k} & 1 & & (0) \\ \vdots & & \vdots & \vdots & & \ddots & \\ 0 & \dots & 0 & -g_{n,k} & (0) & & 1 \end{pmatrix}, \quad g_{i,k} = \frac{a_{i,k}^{(k)}}{a_{k,k}^{(k)}}, \quad i = k+1, \dots, n.$$

On a de plus  $b^{(k+1)} = G_k b^{(k)}$ .

*Démonstration.* Il suffit d'effectuer le produit  $G_k A^{(k)}$  et de vérifier que l'on retrouve  $A^{(k+1)}$ . □

**Définition 2.5.** Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$ . Les mineurs fondamentaux  $D_k$ ,  $k = 1, \dots, n$  de  $A$  sont les déterminants des sous-matrices de  $A$  formées par les  $k$  premières lignes et les  $k$  premières colonnes de  $A$  :  $D_k = \det((a_{i,j})_{1 \leq i,j \leq k})$  pour  $k = 1, \dots, n$ .

**Théorème 2.6.** Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  une matrice carrée inversible. Les propriétés suivantes sont équivalentes :

- (i) L'élimination de Gauss s'effectue sans permutation de lignes ;
- (ii) Il existe  $L \in \mathbb{M}_{n \times n}(\mathbb{K})$  triangulaire inférieure inversible et  $U \in \mathbb{M}_{n \times n}(\mathbb{K})$  triangulaire supérieure inversible telles que  $A = LU$  ;
- (iii) Tous les mineurs fondamentaux de  $A$  sont non nuls.

*Démonstration.* Pour la suite de ce cours, la partie de la démonstration qui nous intéresse est l'implication (i)  $\Rightarrow$  (ii) qui nous permettra de calculer la factorisation LU de  $A$ . D'après le lemme 2.4, on a  $A^{(n)} = G_{n-1} G_{n-2} \cdots G_1 A$ . La matrice  $G_{n-1} G_{n-2} \cdots G_1$  étant inversible, on peut donc écrire  $A = (G_{n-1} G_{n-2} \cdots G_1)^{-1} A^{(n)}$ . Les matrices  $G_k$  étant triangulaires inférieures, la proposition 2.2 affirme que  $(G_{n-1} G_{n-2} \cdots G_1)^{-1}$  est aussi triangulaire inférieure. De plus, par construction  $A^{(n)}$  est triangulaire supérieure d'où le résultat en posant  $L = (G_{n-1} G_{n-2} \cdots G_1)^{-1}$  et  $U = A^{(n)}$ .  $\square$

**Corollaire 2.7.** *Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  une matrice carrée inversible. La matrice  $A$  admet une factorisation LU si et seulement si tous ses mineurs fondamentaux sont non nuls.*

Le lemme suivant nous permet d'expliciter la matrice  $L$  de la factorisation LU de  $A$  obtenue à partir de l'élimination de Gauss.

**Lemme 2.8.** *Avec les notations précédentes, on a*

$$(G_{n-1} G_{n-2} \cdots G_1)^{-1} = \begin{pmatrix} 1 & 0 & \cdots & \cdots & 0 \\ g_{2,1} & 1 & \ddots & & \vdots \\ g_{3,1} & g_{3,2} & 1 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ g_{n,1} & g_{n,2} & \cdots & g_{n,n-1} & 1 \end{pmatrix}.$$

*Démonstration.* Faire le calcul.  $\square$

**Corollaire 2.9.** *Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  une matrice carrée inversible. Si tous les mineurs fondamentaux de  $A$  sont non nuls, alors avec les notations précédentes, l'élimination de Gauss fournit la factorisation LU de  $A$  suivante :*

$$A = \begin{pmatrix} 1 & 0 & \cdots & \cdots & 0 \\ g_{2,1} & 1 & \ddots & & \vdots \\ g_{3,1} & g_{3,2} & 1 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ g_{n,1} & g_{n,2} & \cdots & g_{n,n-1} & 1 \end{pmatrix} \begin{pmatrix} a_{1,1}^{(1)} & \cdots & \cdots & a_{1,n}^{(1)} \\ 0 & a_{2,2}^{(2)} & & a_{2,n}^{(2)} \\ 0 & 0 & a_{3,3}^{(3)} & a_{3,n}^{(3)} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & a_{n,n}^{(n)} \end{pmatrix}.$$

On remarque que la matrice  $L$  obtenue est à diagonale unité.

**Exemple :** En reprenant l'exemple de la sous-section 2.2.1, on trouve :

$$\underbrace{\begin{pmatrix} 1 & 2 & 5 \\ 3 & 2 & -1 \\ 0 & 5 & 3 \end{pmatrix}}_A = \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ 0 & -\frac{5}{4} & 1 \end{pmatrix}}_L \underbrace{\begin{pmatrix} 1 & 2 & 5 \\ 0 & -4 & -16 \\ 0 & 0 & -17 \end{pmatrix}}_U.$$



**Proposition 2.10.** Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  une matrice carrée inversible admettant une factorisation LU. Alors il existe une unique factorisation LU de  $A$  avec  $L$  à diagonale unité.

*Démonstration.* Si  $A$  admet une factorisation LU, alors l'existence d'une factorisation LU avec  $L$  à diagonale unité est claire d'après le corollaire 2.9. Supposons maintenant que  $A$  admette deux factorisations LU :  $A = L_1 U_1$  et  $A = L_2 U_2$  avec  $L_1$  et  $L_2$  à diagonale unité. L'égalité  $L_1 U_1 = L_2 U_2$  implique  $U_1 U_2^{-1} = L_1^{-1} L_2$ . D'après la proposition 2.2,  $U_1 U_2^{-1}$  est triangulaire supérieure et  $L_1^{-1} L_2$  est triangulaire inférieure à diagonale unité. La seule possibilité pour que ces deux matrices soient égales est donc  $U_1 U_2^{-1} = L_1^{-1} L_2 = \mathbb{I}_n$  ce qui implique  $L_1 = L_2$  et  $U_1 = U_2$ .  $\square$

Lorsque  $A$  admet une factorisation LU, la résolution du système d'équations linéaires  $(S) : Ax = b$  se ramène à la résolution de deux systèmes linéaires triangulaires. En effet :

$$Ax = b \iff LUx = b \iff \begin{cases} Ly = b, \\ Ux = y. \end{cases}$$

En pratique, on résout donc d'abord  $Ly = b$  puis connaissant  $y$  on résout  $Ux = y$ .

**Définition 2.11.** On appelle matrice de permutation associée à une permutation  $\sigma \in \mathcal{S}_n$ , la matrice  $\mathcal{P}_\sigma = (\delta_{i\sigma(j)})_{1 \leq i, j \leq n}$  où  $\delta_{ij} = 1$  si  $i = j$  et  $\delta_{ij} = 0$  si  $i \neq j$ .

Par exemple si l'on considère la permutation,  $\sigma : (1, 2, 3, 4, 5) \mapsto (3, 2, 5, 1, 4)$ , on obtient la « matrice de permutation élémentaire »

$$\mathcal{P}_\sigma = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}.$$

Multiplier une matrice  $A$  à gauche (resp. à droite) par une matrice de permutation revient alors à permuter les lignes (resp. les colonnes) de la matrice. Par exemple, en multipliant une matrice à gauche par la matrice  $\mathcal{P}_\sigma$  ci-dessus, la troisième ligne devient la première, la seconde ligne reste inchangée, la cinquième ligne devient la troisième . . . . Les matrices de permutation sont orthogonales (voir Définition 2.20) c'est-à-dire que  $\forall \sigma \in \mathcal{S}_n, \mathcal{P}_\sigma^{-1} = \mathcal{P}_\sigma^T$ .

Le corollaire 2.7 nous donne une condition nécessaire et suffisante pour qu'une matrice inversible admette une factorisation LU. Lorsque cette factorisation LU n'existe pas, on peut tout de même utiliser le théorème suivant :

**Théorème 2.12.** Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  une matrice carrée inversible. Il existe une matrice de permutation  $\mathcal{P}$  telle que  $\mathcal{P}A$  admette une factorisation LU.

*Démonstration.* Admis pour ce cours.  $\square$

Notons que dans ce cas, on a :

$$Ax = b \iff \mathcal{P}Ax = \mathcal{P}b \iff L U x = \mathcal{P}b \iff \begin{cases} Ly = \mathcal{P}b, \\ Ux = y. \end{cases}$$

En pratique, on résout donc d'abord  $Ly = \mathcal{P}b$  puis connaissant  $y$  on résout  $Ux = y$ .

## 2.2.4 Coût de l'algorithme

### Résolution d'un système via l'élimination de Gauss

Nous allons nous intéresser au nombre d'opérations à virgule flottante nécessaires à la résolution d'un système  $(S) : Ax = b$  en utilisant l'élimination de Gauss puis en résolvant le système triangulaire  $A^{(n)}x = b^{(n)}$ .

Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  une matrice carrée inversible et supposons que  $A$  admette une factorisation LU. À l'étape  $k$  de l'élimination de Gauss, on doit faire  $n - k$  divisions pour calculer les  $g_{i,k}$  puis nous avons  $(n - k)^2$  coefficients  $a_{i,j}^{(k+1)}$  à calculer et  $(n - k)$  coefficients  $b_i^{(k+1)}$ . Cela nécessite donc  $(n - k)(n - k + 1)$  multiplications puis  $(n - k)(n - k + 1)$  soustractions. Au total, pour toute l'élimination de Gauss, nous aurons donc à faire

$$\begin{aligned} \sum_{k=1}^{n-1} (n - k) + 2 \sum_{k=1}^{n-1} (n - k)(n - k + 1) &= \sum_{k=0}^{n-1} k + 2 \sum_{k=0}^{n-1} k(k + 1) \\ &= \frac{n(n-1)}{2} + 2 \left( \sum_{k=1}^{n-1} k + \sum_{k=1}^{n-1} k^2 \right) \\ &= \frac{n(n-1)}{2} + 2 \left( \frac{n(n-1)}{2} + \frac{n(n-1)(2n-1)}{6} \right) \\ &= n(n-1) \left( \frac{3+6+2(2n-1)}{6} \right) \\ &= \left( \frac{n(n-1)(4n+7)}{6} \right) \\ &= \left( \frac{4n^3+3n^2-7n}{6} \right) \end{aligned}$$

opérations à virgule flottante. D'après la proposition 2.1, la résolution d'un système triangulaire coûte  $n^2$  opérations à virgule flottantes donc au total, la résolution du système  $(S)$  par cette méthode coûtera  $G(n) = n^2 + \left( \frac{4n^3+3n^2-7n}{6} \right) = \frac{4n^3+9n^2-7n}{6}$  opérations à virgule flottante.

Lorsque  $n$  tend vers l'infini, on a  $G(n) \sim \frac{2n^3}{3}$  opérations à virgule flottante. Notons que d'après ce qui précède, cette estimation représente aussi le coût asymptotique du calcul de la factorisation LU d'une matrice. On a donc obtenu :

**Lemme 2.13.** *Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  une matrice carrée inversible. Résoudre un système linéaire  $(S) : Ax = b$  via l'élimination de Gauss nécessite un nombre d'opérations à virgule flottante équivalent à  $\frac{2n^3}{3}$  lorsque  $n$  tend vers l'infini. Ce coût asymptotique est aussi celui du calcul de la factorisation LU de  $A$ .*

Pour comparer avec le résultat obtenu en utilisant les formules de Cramer, si l'on suppose que  $n = 100$ , alors cela donne environ  $6,6 \cdot 10^5$  opérations à virgule flottante et avec un

ordinateur fonctionnant à 100 megaflops, cela prendra moins de 7 millièmes de secondes. Cela ne se fait donc pas à la main mais votre ordinateur personnel pourra le faire en quelques (dizaines de) secondes.

### Faut-il inverser une matrice ?

Nous admettrons qu'étant donnée la factorisation LU de  $A$ , le coût du calcul de l'inverse  $A^{-1}$  de  $A$  lorsque  $n$  tend vers l'infini est de  $\frac{4n^3}{3}$  opérations à virgule flottante. Au total, lorsque  $n$  tend vers l'infini, il faut donc  $2n^3$  opérations à virgule flottante pour calculer l'inverse de  $A$ . D'après ce qui précède cela signifie donc qu'asymptotiquement (*i.e.*, lorsque  $n$  tend vers l'infini), il faut 3 fois plus d'opérations à virgule flottante pour calculer l'inverse de  $A$  que pour résoudre le système linéaire  $Ax = b$  en utilisant l'élimination de Gauss. Cela implique qu'il ne faut pas calculer l'inverse d'une matrice pour résoudre un système linéaire.

### Cas de la résolution de plusieurs systèmes de même matrice $A$

Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  une matrice carrée inversible et supposons que l'on ait à résoudre  $K$  systèmes linéaires avec la même matrice  $A$  et  $N$  seconds membres  $b^{[1]}, \dots, b^{[K]}$ . Si l'on applique l'élimination de Gauss à chacun de ces systèmes pour les résoudre, alors d'après ce qui précède cela nous coûtera  $K \frac{4n^3 + 9n^2 - 7n}{6}$  opérations à virgule flottante. Si maintenant on calcule une fois pour toute une factorisation LU de  $A$  (via l'élimination de Gauss) et que l'on résout ensuite successivement les  $2K$  systèmes triangulaires cela ne nous coutera que  $\left(\frac{4n^3 + 3n^2 - 7n}{6}\right) + 2Kn^2$  opérations à virgule flottante ce qui est asymptotiquement meilleur. Notons que si l'on calcule une fois pour toute l'inverse  $A^{-1}$  de  $A$  et que l'on résout ensuite chaque système en posant  $x^{[i]} = A^{-1}b^{[i]}$  cela coutera  $2n^3$  opérations pour le calcul de l'inverse puis  $n(2n-1)$  opérations pour calculer chacun des  $K$  produits  $A^{-1}b^{[i]}$  soit au total  $2n^3 + 2Kn^2$  ce qui est moins avantageux que d'utiliser la factorisation LU.

## 2.3 Méthode de Cholesky

La méthode de Cholesky est une alternative à l'élimination de Gauss qui s'applique aux matrices symétriques et définies positives.

**Définition 2.14.** Une matrice  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  est dite symétrique si elle est égale à sa transposée, *i.e.*,  $A^T = A$ .

**Définition 2.15.** Soit  $\mathbb{K} = \mathbb{R}$  ou  $\mathbb{C}$ . Le produit scalaire canonique sur  $\mathbb{K}^n$  est défini comme l'application  $\langle \cdot, \cdot \rangle : \mathbb{K}^n \times \mathbb{K}^n \rightarrow \mathbb{K}$ ,  $(u, v) \mapsto \langle u, v \rangle$  qui vérifie :

- Si  $\mathbb{K} = \mathbb{R}$ ,  $\langle u, v \rangle = v^T u = \sum_{i=1}^n u_i v_i$  (produit scalaire euclidien),
- Si  $\mathbb{K} = \mathbb{C}$ ,  $\langle u, v \rangle = \bar{v}^T u = \sum_{i=1}^n u_i \bar{v}_i$  (produit scalaire hermitien).

**Définition 2.16.** Une matrice  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  est dite définie positive, resp. semi définie positive si pour tout  $x \in \mathbb{R}^n$  non nul, on a  $\langle Ax, x \rangle > 0$ , resp.  $\langle Ax, x \rangle \geq 0$ .

On montre facilement les propriétés suivantes :

1. Une matrice définie positive est inversible ;
2. Si  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  est inversible, alors  $A^T A$  est symétrique et définie positive ;
3. Si  $A = (a_{i,j})_{1 \leq i,j \leq n} \in \mathbb{M}_{n \times n}(\mathbb{K})$  est définie positive, alors  $a_{i,i} > 0$  pour tout  $i = 1, \dots, n$ .

**Théorème 2.17** (Caractérisation des matrices symétriques définies positives). *Une matrice réelle  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  est symétrique définie positive si et seulement si il existe une matrice  $L = (l_{i,j})_{1 \leq i,j \leq n} \in \mathbb{M}_{n \times n}(\mathbb{R})$  triangulaire inférieure inversible telle que  $A = L L^T$ . De plus, si pour tout  $i = 1, \dots, n$ ,  $l_{i,i} \geq 0$ , alors  $L$  est unique.*

*Démonstration.* Admis pour ce cours. □

Remarquons que si  $L$  est triangulaire inférieure, alors  $L^T$  est triangulaire supérieure de sorte que la factorisation de Cholesky  $A = L L^T$  peut être vue comme un cas particulier de la factorisation LU d'une matrice étudiée dans la section précédente.

On donne maintenant *l'algorithme de Cholesky* qui étant donnée une matrice carrée réelle  $A$  symétrique et définie positive calcule une matrice  $L$  telle que  $A = L L^T$ . On admettra que cet algorithme est correct.

**Entrée** :  $A = (a_{i,j})_{1 \leq i,j \leq n} \in \mathbb{M}_{n \times n}(\mathbb{R})$  symétrique et définie positive.

**Sortie** :  $L = (l_{i,j})_{1 \leq i,j \leq n} \in \mathbb{M}_{n \times n}(\mathbb{R})$  tel que  $A = L L^T$ .

1.  $l_{1,1} = \sqrt{a_{1,1}}$  ;
2. Pour  $i$  de 2 à  $n$  par pas de 1, faire :
  - $l_{i,1} = \frac{a_{i,1}}{l_{1,1}}$  ;
3. Pour  $j$  de 2 à  $n$  par pas de 1, faire :
  - Pour  $i$  de 1 à  $j - 1$  par pas de 1, faire :
 
$$l_{i,j} = 0 ;$$
  - $l_{j,j} = \sqrt{a_{j,j} - \sum_{k=1}^{j-1} l_{j,k}^2}$  ;
  - Pour  $i$  de  $j + 1$  à  $n$  par pas de 1, faire :
 
$$l_{i,j} = \frac{a_{i,j} - \sum_{k=1}^{j-1} l_{i,k} l_{j,k}}{l_{j,j}} ;$$
4. Retourner  $L = (l_{i,j})_{1 \leq i,j \leq n} \in \mathbb{M}_{n \times n}(\mathbb{R})$ .

Notons que dans ce cas, comme pour la factorisation LU, on a :

$$A x = b \iff L L^T x = b \iff \begin{cases} L y = b, \\ L^T x = y. \end{cases}$$

En pratique, on résout donc d'abord  $L y = b$  puis connaissant  $y$  on résout  $L^T x = y$ .

**Proposition 2.18.** *L'algorithme de Cholesky décrit ci-dessus nécessite  $n$  extractions de racines carrées et un nombre d'opérations à virgule flottante équivalent à  $\frac{n^3}{3}$  lorsque  $n$  tend vers l'infini.*

*Démonstration.* Admis pour ce cours. □

On voit donc que cet algorithme requiert asymptotiquement « presque » (cela dépendra de la méthode utilisée pour le calcul des racines carrées qui peut-être plus ou moins rapide) deux fois moins d'opérations à virgule flottante que celui calculant la factorisation LU via l'élimination de Gauss. Par conséquent, il est conseillé de l'utiliser lorsque  $A$  est réelle symétrique et définie positive.

## 2.4 Méthode de Householder et factorisation QR

Dans cette section, on suppose que  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  est une matrice réelle inversible.

### 2.4.1 Transformation (élémentaire) de Householder

**Définition 2.19.** *On appelle matrice (élémentaire) de Householder une matrice  $H$  de la forme  $H_u = \mathbb{I}_n - 2uu^T$ , où  $u \in \mathbb{R}^n$  est un vecteur unitaire c'est-à-dire de norme 1 pour la norme associée au produit scalaire canonique sur  $\mathbb{R}^n$  définie par  $\|u\| = \sqrt{\langle u, u \rangle}$ .*

Par exemple, pour  $n = 3$ , on peut considérer le vecteur  $u = \frac{1}{\sqrt{6}} \begin{pmatrix} -1 & 1 & 2 \end{pmatrix}^T$  qui vérifie bien  $\|u\| = 1$ . On obtient alors la matrice de Householder  $H_u = \frac{1}{3} \begin{pmatrix} 2 & 1 & 2 \\ 1 & 2 & -2 \\ 2 & -2 & -1 \end{pmatrix}$ .

**Définition 2.20.** *Une matrice  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  est dite orthogonale si elle est réelle, i.e.,  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  et si  $AA^T = A^T A = \mathbb{I}_n$ .*

Par exemple que les matrices de permutation (voir Définition 2.11) sont orthogonales.

**Proposition 2.21.** *Toute matrice de Householder  $H$  est symétrique et orthogonale.*

*Démonstration.* Facile à vérifier. □

**Proposition 2.22.** *Pour tout vecteur  $u \in \mathbb{R}^n$  tel que  $\|u\| = 1$ , on a  $H_u u = -u$ . De plus, si  $v \in \mathbb{R}^n$  est orthogonal à  $u$ , i.e.,  $\langle u, v \rangle = 0$ , alors  $H_u v = v$ .*

*Démonstration.* Facile à vérifier. □

Géométriquement, la matrice  $H_u$  représente donc la symétrie orthogonale par rapport au plan  $u^T$  orthogonal à  $u$ .

**Lemme 2.23.** *Soit  $x$  et  $y$  deux vecteurs de  $\mathbb{R}^n$  tels que  $x \neq y$  et  $\|x\| = \|y\|$ . Alors il existe un vecteur unitaire  $u \in \mathbb{R}^n$  tel que  $H_u x = y$ .*

*Démonstration.* Prendre  $u = \frac{x-y}{\|x-y\|}$ . □

## 2.4.2 Principe de la méthode de Householder

La méthode de Householder pour la résolution d'un système linéaire  $Ax = b$  est basée sur les deux propositions suivantes que nous admettrons.

**Proposition 2.24.** *Soit  $v$  un vecteur non nul de  $\mathbb{R}^n$ . Alors il existe une matrice de Householder  $H$  et un réel  $\alpha$  tels que  $Hv = \alpha e_1$ , où  $e_1 = (1, 0, \dots, 0)^T$  est le premier vecteur de la base canonique de  $\mathbb{R}^n$ .*

**Proposition 2.25.** *Soit  $u = (u_i)$  un vecteur unitaire de  $\mathbb{R}^n$  tel que  $u_1 = \dots = u_p = 0$  pour  $p < n$ . On décompose alors  $u$  en deux blocs :  $u = (0 \ z)^T$  avec  $z \in \mathbb{R}^{n-p}$ . La matrice de Householder  $H_u$  se décompose alors par blocs de la manière suivante :  $H_u = \begin{pmatrix} \mathbb{I}_p & 0 \\ 0 & H_z \end{pmatrix}$ .*

## 2.4.3 Exemple de résolution d'un système linéaire par la méthode de Householder

Considérons le système linéaire suivant :

$$(S) : \begin{cases} 2x_1 + x_2 + 2x_3 = 1, \\ x_1 + x_2 + 2x_3 = 1, \\ 2x_1 + x_2 + x_3 = 1. \end{cases}$$

**Étape 1 :** On considère le vecteur obtenu à partir de la première colonne de la matrice  $A$  de  $(S)$ , i.e.,  $a_1 = (2 \ 1 \ 2)^T$ . Le but de cette étape est de trouver la matrice  $H$  apparaissant dans la proposition 2.24 et correspondant à  $v = a_1$  de sorte qu'en multipliant la matrice  $A$  du système  $(S)$  à gauche par  $H$  on obtienne des zéros sous le premier coefficient. Pour ceci, en concordance avec la preuve du lemme 2.23, on introduit le vecteur  $v_1 = \frac{a_1}{\|a_1\|} - e_1 = \frac{1}{3}(-1 \ 1 \ 2)^T$ , le vecteur  $u_1 = \frac{v_1}{\|v_1\|} = \frac{1}{\sqrt{6}}(-1 \ 1 \ 2)^T$  et on considère la matrice élémentaire de Householder  $H_{u_1} = \frac{1}{3} \begin{pmatrix} 2 & 1 & 2 \\ 1 & 2 & -2 \\ 2 & -2 & -1 \end{pmatrix}$ . On a alors :

$$(S) : Ax = b \iff H_{u_1} Ax = H_{u_1} b \iff \begin{pmatrix} 9 & 5 & 8 \\ 0 & 1 & 4 \\ 0 & -1 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 5 \\ 1 \\ -1 \end{pmatrix}.$$

**Étape 2 :** On se place maintenant dans  $\mathbb{R}^2$  et on considère le vecteur  $a_2 = (1 \ -1)^T$  auquel nous allons appliquer la même technique qu'au vecteur  $a_1$  à l'étape précédente. On utilisera ensuite la proposition 2.25 pour obtenir une matrice de Householder de la bonne taille. On pose alors  $z'_2 = \frac{a_2}{\|a_2\|} - e'_1$  où  $e'_1 = (1, 0)$  est le premier vecteur de la base canonique de  $\mathbb{R}^2$  et  $z_2 = \frac{z'_2}{\|z'_2\|}$ . On considère ensuite la matrice de Householder  $H_{z_2} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 \\ -1 & -1 \end{pmatrix}$

puis  $u_2 = (0 \ z_2)^T$  et  $H_{u_2} = \begin{pmatrix} 1 & 0 \\ 0 & H_{z_2} \end{pmatrix}$  (Voir Proposition 2.25). On a alors :

$$Ax = b \iff H_{u_2} H_{u_1} Ax = H_{u_2} H_{u_1} b \iff \begin{pmatrix} 9 & 5 & 8 \\ 0 & \sqrt{2} & \frac{5}{\sqrt{2}} \\ 0 & 0 & -\frac{3}{\sqrt{2}} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 5 \\ \sqrt{2} \\ 0 \end{pmatrix},$$

d'où en résolvant ce système triangulaire  $x = (0 \ 1 \ 0)^T$ .

#### 2.4.4 Factorisation QR d'une matrice

**Définition 2.26.** Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  une matrice carrée réelle inversible. On appelle factorisation QR de  $A$  une factorisation de la forme  $A = QR$  avec  $Q \in \mathbb{M}_{n \times n}(\mathbb{R})$  orthogonale et  $R \in \mathbb{M}_{n \times n}(\mathbb{R})$  triangulaire supérieure.

En généralisant la méthode expliquée dans la sous-section précédente pour une matrice carrée réelle inversible quelconque de taille  $n$ , on obtient  $n - 1$  matrices de Householder  $H_1, \dots, H_{n-1}$  telles que le produit  $R = H_{n-1} H_{n-2} \cdots H_1 A$  est triangulaire supérieure. On pose alors  $Q = (H_{n-1} H_{n-2} \cdots H_1)^{-1}$  de sorte que  $A = QR$ . La matrice  $Q$  ainsi défini étant orthogonale, on a obtenu une factorisation QR de  $A$ . On obtient alors le résultat suivant :

**Théorème 2.27.** Pour toute matrice réelle  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$ , il existe une matrice orthogonale  $Q \in \mathbb{M}_{n \times n}(\mathbb{R})$  produit d'au plus  $(n - 1)$  matrices de Householder et une matrice triangulaire supérieure  $R \in \mathbb{M}_{n \times n}(\mathbb{R})$  telles que  $A = QR$ .

**Proposition 2.28.** La méthode de Householder pour résoudre un système linéaire nécessite un nombre d'opérations à virgule flottante équivalent à  $\frac{4n^3}{3}$  lorsque  $n$  tend vers l'infini.

*Démonstration.* Admis pour ce cours. □

Le coût de cette méthode est donc relativement élevé comparé à l'élimination de Gauss ou à la méthode de Cholesky. Cependant un avantage de cette méthode est qu'elle est plus *stable numériquement* que les deux méthodes citées précédemment.

Notons enfin que la factorisation QR d'une matrice existe aussi pour des matrices rectangulaires et qu'elle est utilisée pour des problèmes de *moindres carrés*.

**Factorisation QR et moindres carrés linéaires.** Le problème des moindres carrés linéaires est le suivant : étant donné une matrice  $A \in \mathbb{M}_{m \times n}(\mathbb{R})$  et un vecteur  $b \in \mathbb{R}^m$ , minimiser la quantité  $\|Ax - b\|_2^2$ .

Ce problème peut être résolu en utilisant la factorisation QR de la matrice  $A \in \mathbb{M}_{m \times n}(\mathbb{R})$  en procédant comme suit : la matrice  $A$  étant rectangulaire (on supposera ici  $m > n$ ), on partitionne les matrices  $Q$  et  $R$  de la manière suivante :

$$A = \begin{pmatrix} Q_1 & Q_2 \end{pmatrix} \begin{pmatrix} R_1 \\ 0 \end{pmatrix},$$

où  $Q_1 \in \mathbb{M}_{m \times n}(\mathbb{R})$ ,  $Q_2 \in \mathbb{M}_{m \times (m-n)}(\mathbb{R})$  et  $R_1 \in \mathbb{M}_{n \times n}(\mathbb{R})$ . En remplaçant  $A$  par sa factorisation QR, on a :

$$\|Ax - b\|_2^2 = \|QRx - b\|_2^2.$$

On admettra qu'en appliquant une transformation orthogonale au vecteur des résidus  $Ax - b$  on ne modifie pas la norme euclidienne du vecteur et que le problème de minimisation conduit au même résultat. On multiplie alors le résidu par  $Q^T$  et on est ramené à minimiser

$$\|Q^T QRx - Q^T b\|_2^2 = \|Rx - Q^T b\|_2^2 = \left\| \begin{pmatrix} R_1 \\ 0 \end{pmatrix} x - \begin{pmatrix} Q_1^T \\ Q_2^T \end{pmatrix} b \right\|_2^2,$$

ce qui revient à minimiser

$$\|R_1 x - Q_1^T b\|_2^2 + \|Q_2^T b\|_2^2.$$

On résout alors le système triangulaire  $R_1 x = Q_1^T b$  et la somme des résidus au carré correspond à  $\|Q_2^T b\|_2^2$ .

Dans les problèmes qui apparaissent en économétrie ou statistique, on est aussi intéressé par la matrice des variances et covariances  $\sigma^2 (A^T A)^{-1}$ , où  $\sigma^2 = \frac{\|Ax - b\|_2^2}{m - n}$ . Pour obtenir cette matrice, on a donc besoin d'évaluer le résidu à la solution  $x$  trouvée et de calculer l'inverse de  $A^T A$ . Si on possède la factorisation QR de  $A$ , alors  $A^T A = R^T Q^T Q R = R^T R$  car  $Q$  est orthogonale. On doit donc inverser la matrice  $R^T R$ . De plus, avec les notations précédentes, si on note  $d = Q_2^T b$ , on peut écrire  $\sigma^2 = \frac{d^T d}{m - n}$ .

**Remarque** : cela n'entre pas dans le cadre de ce cours mais il existe une méthode efficace et numériquement stable pour calculer l'inverse de matrices de la forme  $A^T A$ .



# Chapitre 3

## Conditionnement d'une matrice pour la résolution d'un système linéaire

### 3.1 Normes matricielles

#### 3.1.1 Normes vectorielles

Soit  $E$  un espace vectoriel sur  $\mathbb{K} = \mathbb{R}$  ou  $\mathbb{C}$ .

**Définition 3.1.** On appelle norme sur  $E$  une application  $\|\cdot\| : E \rightarrow \mathbb{R}_+$  telle que :

- $\forall x \in E, (\|x\| = 0 \Rightarrow x = 0)$  ;
- $\forall \lambda \in \mathbb{K}, \forall x \in E, \|\lambda x\| = |\lambda| \|x\|$  ;
- $\forall (x, y) \in E^2, \|x + y\| \leq \|x\| + \|y\|$ .

Les exemples classiques de normes sur  $\mathbb{R}^n$  sont les normes  $\|\cdot\|_1, \|\cdot\|_2$  et  $\|\cdot\|_\infty$  définies par :

$$\forall x = (x_i)_{1 \leq i \leq n} \in \mathbb{R}^n, \quad \|x\|_1 = \sum_{i=1}^n |x_i|, \quad \|x\|_2 = \left( \sum_{i=1}^n |x_i|^2 \right)^{\frac{1}{2}} = \langle x, x \rangle^{\frac{1}{2}}, \quad \|x\|_\infty = \max_{1 \leq i \leq n} |x_i|.$$

#### 3.1.2 Normes matricielles et normes subordonnées

**Définition 3.2.** Une norme  $\|\cdot\|$  sur  $\mathbb{M}_{n \times n}(\mathbb{K})$  est une norme matricielle si elle vérifie :  $\forall (A, B) \in \mathbb{M}_{n \times n}(\mathbb{K})^2, \|AB\| \leq \|A\| \|B\|$ .

L'exemple fondamental est celui des normes dites *subordonnées* qui sont associées à une norme vectorielle.

**Théorème et Définition 3.3.** Soit  $\|\cdot\|$  une norme vectorielle sur  $\mathbb{K}^n$ . Pour toute matrice  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$ , on définit  $\|\cdot\|_M : \mathbb{M}_{n \times n}(\mathbb{K}) \rightarrow \mathbb{R}_+$  par  $\|A\|_M = \sup_{x \in \mathbb{K}^n \setminus \{0\}} \frac{\|Ax\|}{\|x\|}$ . Alors  $\|\cdot\|_M$  est une norme matricielle. Elle est dite norme subordonnée à la norme vectorielle  $\|\cdot\|$ .

Dans la suite on notera indifféremment  $\|\cdot\|$  pour une norme vectorielle ou une norme matricielle.

*Démonstration.* On utilise le lemme suivant qui découle de la définition d'une norme subordonnée :

**Lemme 3.4.** *Soit  $\|\cdot\|$  la norme subordonnée à une norme vectorielle  $\|\cdot\|$ . Alors, pour tout vecteur  $x \in \mathbb{K}^n$  et pour toute matrice  $M \in \mathbb{M}_{n \times n}(\mathbb{K})$ , on a  $\|Mx\| \leq \|M\| \|x\|$ .*

Soit  $(A, B) \in \mathbb{M}_{n \times n}(\mathbb{K})^2$  et  $x \neq 0$ . On a :

$$\frac{\|ABx\|}{\|x\|} \leq \frac{\|A\| \|Bx\|}{\|x\|} \leq \frac{\|A\| \|B\| \|x\|}{\|x\|}.$$

d'où le résultat. □

Les normes subordonnées respectivement aux normes vectorielles  $\|\cdot\|_1$ ,  $\|\cdot\|_2$  et  $\|\cdot\|_\infty$  de  $\mathbb{R}^n$  sont données par :  $\forall A = (a_{i,j})_{1 \leq i,j \leq n} \in \mathbb{M}_{n \times n}(\mathbb{K})$  :

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{i,j}|, \quad \|A\|_2 = \sqrt{\rho(AA^*)}, \quad \|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{i,j}|,$$

où  $A^* = \overline{A}^T$  désigne la matrice adjointe de  $A$  et  $\rho(M)$  désigne le rayon spectral d'une matrice  $M$ , i.e., le maximum des modules des valeurs propres de  $M$ .

Notons que ces trois normes matricielles sont équivalentes : pour  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$ , on a :

$$\frac{1}{\sqrt{n}} \|A\|_2 \leq \|A\|_1 \leq \sqrt{n} \|A\|_2, \quad \frac{1}{\sqrt{n}} \|A\|_\infty \leq \|A\|_2 \leq \sqrt{n} \|A\|_\infty, \quad \frac{1}{n} \|A\|_1 \leq \|A\|_\infty \leq n \|A\|_1.$$

## 3.2 Conditionnement d'une matrice

### 3.2.1 Exemple classique

Cet exemple est dû à R. S. Wilson. Considérons le système linéaire (S) :  $Ax = b$  avec

$$A = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix}, \quad b = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}.$$

Remarquons tout d'abord que la matrice  $A$  est symétrique, que son déterminant vaut 1 et que la solution de (S) est donnée par  $x = (1 \ 1 \ 1 \ 1)^T$ .

**Premier cas :  $b$  est perturbé.** Perturbons légèrement le second membre  $b$  et considérons

$$b' = \begin{pmatrix} 32,1 \\ 22,9 \\ 33,1 \\ 30,9 \end{pmatrix}.$$

Si on résout le système ( $S'$ ) :  $Ax' = b'$ , on trouve  $x' = (9,2 \quad -12,6 \quad 4,5 \quad -1,1)^T$ . La « petite » perturbation sur le second membre  $b$  entraîne donc une « forte » perturbation sur la solution du système.

D'une manière générale, on considère les deux systèmes  $Ax = b$  et  $A(x + \delta x) = b + \delta b$ . On a donc  $A\delta x = \delta b$  de sorte que  $\delta x = A^{-1}\delta b$  et on obtient la majoration suivante de l'erreur absolue sur la solution :  $\|\delta x\| \leq \|A^{-1}\| \cdot \|\delta b\|$ . Or  $Ax = b$  donc  $\|b\| \leq \|A\| \cdot \|x\|$ . Finalement, on obtient une majoration de l'erreur relative sur la solution en fonction de l'erreur relative sur la donnée :

$$\frac{\|\delta x\|}{\|x\|} \leq \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\delta b\|}{\|b\|}.$$

Notons que cette majoration est optimale dans le sens où il n'existe pas de borne plus petite qui soit valable pour tout système. En effet, prenons  $A = \begin{pmatrix} 1 & 0 \\ 0 & \frac{1}{2} \end{pmatrix}$ ,  $b = (1 \quad 0)^T$ , et  $\delta b = (0 \quad \frac{1}{2})^T$ . La solution de  $Ax = b$  est alors  $x = (1 \quad 0)^T$  et celle de  $A\delta x = \delta b$  est  $\delta x = (0 \quad 1)^T$ . On a donc  $\frac{\|\delta x\|}{\|x\|} = 1$ ,  $\frac{\|\delta b\|}{\|b\|} = \frac{1}{2}$ . Or  $\|A^{-1}\| \cdot \|A\| = 2$  donc la borne est atteinte.

**Deuxième cas :  $A$  est perturbée.** De façon analogue, si on perturbe légèrement la matrice  $A$  et que l'on considère

$$A'' = \begin{pmatrix} 10 & 7 & 8,1 & 7,2 \\ 7,08 & 5,04 & 6 & 5 \\ 8 & 5,98 & 9,89 & 9 \\ 6,99 & 4,99 & 9 & 9,98 \end{pmatrix},$$

alors la solution du système ( $S''$ ) :  $A''x'' = b$  est  $x'' = (-81 \quad 107 \quad -34 \quad 22)^T$ .

D'une manière générale, on considère les deux systèmes  $Ax = b$  et  $(A + \Delta A)(x + \delta x) = b$ . Il vient donc  $\delta x = A^{-1}\Delta A(x + \delta x)$  d'où

$$\frac{\|\delta x\|}{\|x + \delta x\|} \leq \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\Delta A\|}{\|A\|}.$$

### 3.2.2 Définition du conditionnement

**Définition 3.5.** Soit  $\|\cdot\|$  une norme matricielle subordonnée et  $A$  une matrice inversible. Le nombre  $\text{Cond}(A) = \|A^{-1}\| \cdot \|A\|$  s'appelle le conditionnement de  $A$  relatif à la norme  $\|\cdot\|$ .

Ce nombre mesure la « sensibilité » de la solution par rapport aux données du problème. Une matrice est bien conditionnée si  $\text{Cond}(A) \approx 1$  et mal conditionnée si  $\text{Cond}(A) \gg 1$ . Le fait que les normes matricielles  $\|\cdot\|_1$ ,  $\|\cdot\|_2$  et  $\|\cdot\|_\infty$  soient équivalentes implique la même propriété d'équivalence pour les conditionnements associés : pour toute matrice  $A \in \mathbb{M}_n(\mathbb{K})$ , on a :

$$\frac{1}{n} \text{Cond}_2(A) \leq \text{Cond}_1(A) \leq n \text{Cond}_2(A), \quad \frac{1}{n} \text{Cond}_\infty(A) \leq \text{Cond}_2(A) \leq n \text{Cond}_\infty(A),$$

$$\frac{1}{n^2} \text{Cond}_1(A) \leq \text{Cond}_\infty(A) \leq n^2 \text{Cond}_1(A).$$

Des exemples de matrices bien conditionnées sont les matrices de la forme

$$A = \begin{pmatrix} 4 & 1 & 0 & \dots & 0 \\ 1 & 4 & 1 & \ddots & \vdots \\ 0 & 1 & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & 1 \\ 0 & \dots & 0 & 1 & 4 \end{pmatrix},$$

qui vérifient que, quelque soit la dimension  $n$  de  $A$ ,  $\text{Cond}_\infty(A) \leq 3$ . Inversement, les *matrices de Hilbert*  $H_n$  et les *matrices de Vandermonde*  $V_n$  respectivement définies par

$$H_n = \left( \frac{1}{i+j-1} \right)_{1 \leq i, j \leq n}, \quad V_n = \left( \binom{j}{n}^{i-1} \right)_{1 \leq i, j \leq n},$$

sont très mal conditionnées : leur conditionnement pour la norme  $\|\cdot\|_\infty$  vérifie :

$n$	$\text{Cond}(H_n)$	$\text{Cond}(V_n)$
2	27	8
4	$2,8 \cdot 10^4$	$5,6 \cdot 10^2$
6	$2,9 \cdot 10^7$	$3,7 \cdot 10^4$

**Proposition 3.6.** *Soit  $A$  une matrice réelle et considérons la norme matricielle subordonnée  $\|A\|_2 = \sqrt{\rho(AA^*)}$ . On a*

$$\text{Cond}_2(A) = \|A^{-1}\|_2 \cdot \|A\|_2 = \sqrt{\frac{\rho(AA^T)}{\sigma(AA^T)}},$$

où  $\sigma(M)$  désigne le minimum des modules des valeurs propres de  $M$ . En particulier, si  $A$  est symétrique, alors on obtient  $\text{Cond}_2(A) = \frac{\rho(A)}{\sigma(A)}$ .

*Démonstration.* Admis pour ce cours. □

### 3.2.3 Estimation théorique de l'erreur a priori

**Premier cas :  $b$  est perturbé**

**Théorème 3.7.** Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  inversible et  $b \in \mathbb{R}^n$  tels que  $Ax = b$  et  $A(x + \delta x) = b + \delta b$  avec  $x \neq 0$ . Alors on a :

$$\frac{1}{\text{Cond}(A)} \cdot \frac{\|\delta b\|}{\|b\|} \leq \frac{\|\delta x\|}{\|x\|} \leq \text{Cond}(A) \cdot \frac{\|\delta b\|}{\|b\|},$$

*Démonstration.* La majoration a déjà été prouvée. La minoration s'obtient directement à partir des inégalités  $\|\delta b\| \leq \|A\| \cdot \|\delta x\|$  et  $\|x\| \leq \|A^{-1}\| \cdot \|b\|$ .  $\square$

**Deuxième cas :  $A$  est perturbée**

**Théorème 3.8.** Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  inversible,  $b \in \mathbb{R}^n$  et  $\Delta A \in \mathbb{M}_{n \times n}(\mathbb{R})$  tels que  $\|A^{-1}\| \cdot \|\Delta A\| < 1$ . Alors  $A + \Delta A$  est inversible. De plus si on suppose  $Ax = b$  et  $(A + \Delta A)(x + \delta x) = b$  avec  $x \neq 0$ , alors on a :

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\text{Cond}(A) \cdot \frac{\|\Delta A\|}{\|A\|}}{1 - \text{Cond}(A) \cdot \frac{\|\Delta A\|}{\|A\|}}.$$

*Démonstration.* On a  $(A + \Delta A) = A(\mathbb{I}_n + A^{-1} \Delta A)$  avec par hypothèse  $\|A^{-1} \Delta A\| < 1$ . Le lemme suivant (admis pour ce cours - la preuve n'est pas difficile) implique que  $(A + \Delta A)$  est inversible.

**Lemme 3.9.** Soit  $\|\cdot\|$  une norme matricielle subordonnée et  $B \in \mathbb{M}_{n \times n}(\mathbb{R})$  une matrice telle que  $\|B\| < 1$ . Alors la matrice  $\mathbb{I}_n + B$  est inversible et on a :  $\|(\mathbb{I}_n + B)^{-1}\| \leq \frac{1}{1 - \|B\|}$ .

On a de plus  $(A + \Delta A)^{-1} = (\mathbb{I}_n + A^{-1} \Delta A)^{-1} A^{-1}$ . Maintenant, Si  $Ax = b$  et  $(A + \Delta A)(x + \delta x) = b$ , alors  $(A + \Delta A) \delta x = -\Delta A.x$  d'où  $\delta x = -(\mathbb{I}_n + A^{-1} \Delta A)^{-1} A^{-1} \Delta A.x$ . On a donc

$$\frac{\|\delta x\|}{\|x\|} \leq \|\mathbb{I}_n + A^{-1} \Delta A\|^{-1} \cdot \|A^{-1}\| \cdot \|\Delta A\|,$$

d'où d'après le lemme 3.9 ci-dessus,

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{1}{1 - \|A^{-1} \Delta A\|} \cdot \|A^{-1}\| \cdot \|\Delta A\| \leq \frac{\text{Cond}(A) \cdot \frac{\|\Delta A\|}{\|A\|}}{1 - \text{Cond}(A) \cdot \frac{\|\Delta A\|}{\|A\|}}.$$

$\square$

### Troisième cas : $A$ et $b$ sont perturbés

**Théorème 3.10.** Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  inversible,  $b \in \mathbb{R}^n$  et  $\Delta A \in \mathbb{M}_{n \times n}(\mathbb{R})$  vérifiant  $\|A^{-1}\| \cdot \|\Delta A\| < 1$ . Si l'on suppose que  $Ax = b$  et  $(A + \Delta A)(x + \delta x) = b + \delta b$  avec  $x \neq 0$ , alors on a :

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\text{Cond}(A)}{1 - \text{Cond}(A) \cdot \frac{\|\Delta A\|}{\|A\|}} \left( \frac{\|\delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right).$$

*Démonstration.* Exercice : on utilise les mêmes techniques que pour les preuves précédentes. □

### 3.2.4 Estimation théorique de l'erreur a posteriori

Étant donné un système linéaire  $Ax = b$ , on va maintenant estimer, en fonction du conditionnement, l'erreur commise sur la solution réellement calculée. Soit  $x$  la solution exacte et  $y$  la solution obtenue par la machine. On pose  $r = Ay - b$  ;  $r$  s'appelle le *résidu*. On a alors le résultat suivant :

**Théorème 3.11.** Avec les notations précédentes, si  $x \neq 0$ , alors

$$\|y - x\| \leq \text{Cond}(A) \cdot \frac{\|r\|}{\|b\|} \cdot \|x\|.$$

*Démonstration.* On a  $r = Ay - b = A(y - x)$  d'où  $y - x = A^{-1}r$  et  $\|y - x\| \leq \|A^{-1}\| \cdot \|r\|$  ce qui implique  $\|y - x\| \leq \text{Cond}(A) \cdot \frac{\|r\|}{\|A\|}$ . Or  $Ax = b$  implique  $\|A\| \geq \frac{\|b\|}{\|x\|}$  d'où le résultat. □

On voit donc que si le conditionnement est grand, l'erreur relative peut être grande. Cette majoration n'est pas très facile à utiliser car le conditionnement est en général inconnu. Soit  $C$  une approximation de  $A^{-1}$  (que l'on peut par exemple calculer via la méthode d'élimination de Gauss) et posons  $R = AC - \mathbb{I}_n$ . On a alors le résultat suivant :

**Théorème 3.12.** Avec les notations précédentes, si  $\|R\| < 1$ , alors

$$\|y - x\| \leq \frac{\|r\| \cdot \|C\|}{1 - \|R\|}.$$

*Démonstration.* On a vu dans la preuve précédente que  $\|y - x\| \leq \|A^{-1}\| \cdot \|r\|$ . Il suffit alors de remarquer que  $A^{-1} = C(\mathbb{I}_n + R)^{-1}$  et d'utiliser le lemme 3.9. □

# Chapitre 4

## Résolution d'un système d'équations linéaires (Partie 2) : méthodes itératives

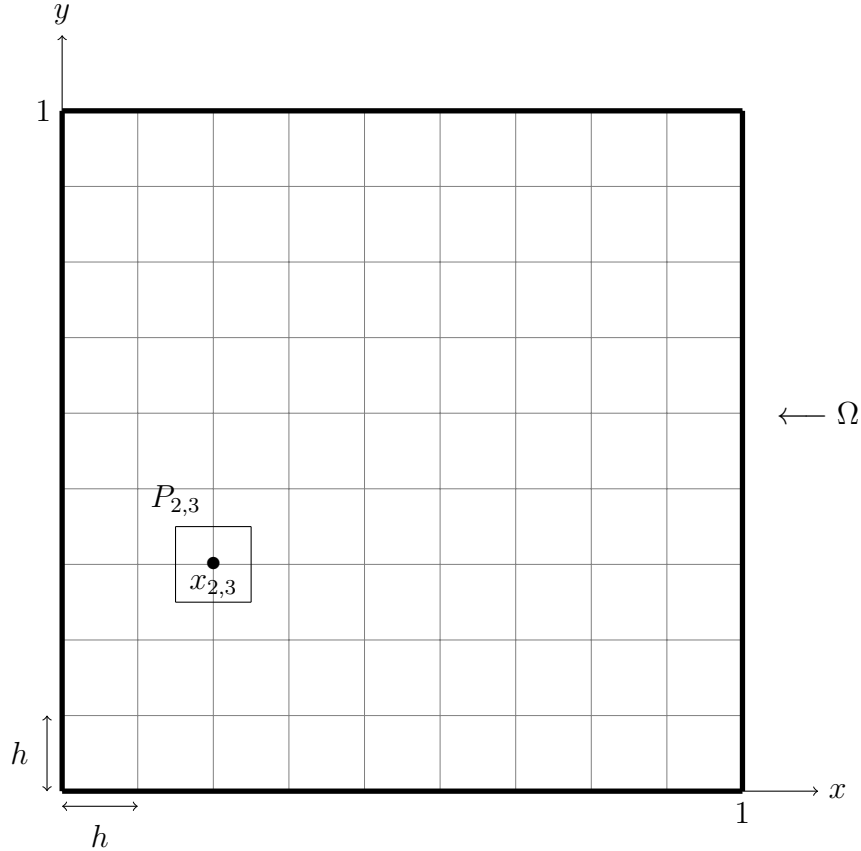
### 4.1 Motivation

Les méthodes itératives deviennent indispensables dès que la taille  $n$  du système est très grande. En effet, les méthodes directes exigent un nombre d'opérations à virgule flottante de l'ordre de  $n^3$  lorsque  $n$  tend vers l'infini ce qui les rend lentes pour de grandes valeurs de  $n$ . De tels systèmes apparaissent par exemple dans les techniques de résolution numérique d'équations aux dérivées partielles. Les matrices des systèmes obtenus sont en général « creuses » (c'est-à-dire qu'elles ont beaucoup de 0) et (semi) définies positives. Voici un exemple classique. Étant donnée une fonction  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , on se propose de trouver une solution approchée  $\tilde{u} : \Omega \subset \mathbb{R}^2 \rightarrow \mathbb{R}$  du problème suivant :

$$\begin{cases} -\Delta \tilde{u} = f, & \forall (x, y) \in \Omega = ]0, 1[ \times ]0, 1[, \\ \tilde{u} = 0, & \forall (x, y) \in \partial\Omega, \end{cases}$$

où  $\Delta \tilde{u} = \frac{\partial^2 \tilde{u}}{\partial x^2} + \frac{\partial^2 \tilde{u}}{\partial y^2}$  désigne le laplacien de la fonction  $\tilde{u}$  et  $\partial\Omega$  la frontière de  $\Omega$ .

Pour ce faire, on se donne un réel  $h > 0$ , on effectue une *discrétisation de  $\Omega = ]0, 1[ \times ]0, 1[$  de pas  $h$*  (i.e., on quadrille  $\Omega$  à l'aide de « petits » pavés d'aire  $h^2$ ) et on cherche une fonction étagée  $u$  (dépendante de  $h$ ) telle que  $u$  tende vers  $\tilde{u}$  lorsque  $h$  tend vers 0. On pose  $h = \frac{1}{n+1}$  et on écrit  $u = \sum_{i,j} u_{i,j} \chi_{i,j}$  où, pour  $1 \leq i, j \leq n$ ,  $\chi_{i,j}$  est la fonction caractéristique du pavé  $P_{i,j} = ](i - \frac{1}{2})h, (i + \frac{1}{2})h[ \times ](j - \frac{1}{2})h, (j + \frac{1}{2})h[$ . On note alors  $x_{i,j} = (ih, jh)$  les nœuds du quadrillage et  $P_{i,j}$  est donc le pavé de centre  $x_{i,j}$ .



En se basant sur la définition de la dérivée d'une fonction, pour  $h$  suffisamment petit, on a les *approximations par différences finies* suivantes :

$$\frac{\partial \tilde{u}}{\partial x} \approx \frac{u(x + \frac{h}{2}, y) - u(x - \frac{h}{2}, y)}{h}, \quad \frac{\partial^2 \tilde{u}}{\partial x^2} \approx \frac{u(x + h, y) - 2u(x, y) + u(x - h, y)}{h^2}.$$

De même,

$$\frac{\partial \tilde{u}}{\partial y} \approx \frac{u(x, y + \frac{h}{2}) - u(x, y - \frac{h}{2})}{h}, \quad \frac{\partial^2 \tilde{u}}{\partial y^2} \approx \frac{u(x, y + h) - 2u(x, y) + u(x, y - h)}{h^2}.$$

En décomposant  $f$  sous la forme  $\sum_{i,j} f_{i,j} \chi_{i,j}$ , notre problème se ramène alors à chercher les  $u_{i,j}$ ,  $1 \leq i, j \leq n$  satisfaisants les équations suivantes :

$$\begin{cases} 4u_{i,j} - u_{i-1,j} - u_{i+1,j} - u_{i,j-1} - u_{i,j+1} = h^2 f_{i,j}, & 1 \leq i, j \leq n, \\ u_{0,j} = u_{n+1,j} = u_{i,0} = u_{i,n+1} = 0, & 1 \leq i, j \leq n. \end{cases}$$

Ce schéma numérique est dit *implicite* par opposition à un schéma *explicite* pour lequel il est possible d'ordonner les inconnues de sorte que chacune d'entre elles puisse être déterminée « explicitement » en fonction des précédentes (système triangulaire). Les schémas numériques implicites ont l'avantage d'être numériquement stables ce qui n'est pas toujours le cas pour



les schémas explicites.

On remarque que ces équations forment un système linéaire que nous allons écrire sous forme matricielle. Pour ceci, on pose

$$M = \begin{pmatrix} 4 & -1 & 0 & \dots & 0 \\ -1 & 4 & -1 & \ddots & \vdots \\ 0 & -1 & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & -1 \\ 0 & \dots & 0 & -1 & 4 \end{pmatrix} \in \mathbb{M}_{n \times n}(\mathbb{R}),$$

$X_j = (u_{1,j} \ u_{2,j} \ \dots \ u_{n,j})^T$  pour  $1 \leq j \leq n$ ,  $X = (X_1^T \ X_2^T \ \dots \ X_n^T)^T$  ainsi que  $F_j = (f_{1,j} \ f_{2,j} \ \dots \ f_{n,j})^T$  pour  $1 \leq j \leq n$ ,  $F = (F_1^T \ F_2^T \ \dots \ F_n^T)^T$ . En définissant de plus  $X_0 = X_{n+1} = 0$ , le système précédent s'écrit

$$-X_{j-1} + M X_j - X_{j+1} = h^2 F_j, \quad 1 \leq j \leq n,$$

ce qui conduit à

$$A X = h^2 F, \quad A = \begin{pmatrix} M & -\mathbb{I}_n & 0 & \dots & 0 \\ -\mathbb{I}_n & M & -\mathbb{I}_n & \ddots & \vdots \\ 0 & -\mathbb{I}_n & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & -\mathbb{I}_n \\ 0 & \dots & 0 & -\mathbb{I}_n & M \end{pmatrix} \in \mathbb{M}_{n^2 \times n^2}(\mathbb{R}).$$

La matrice  $A$  est symétrique réelle et on peut montrer qu'elle est définie positive donc en particulier inversible. En pratique, on a donc à résoudre un système linéaire tridiagonal par blocs de grande taille (notons que faire tendre  $h$  vers 0 équivaut à faire tendre  $n$  vers l'infini) à résoudre. Notons qu'il existe des méthodes efficaces (*e.g.*, l'algorithme de Thomas qui est une simplification de l'algorithme de Gauss dans le cas particulier des systèmes tridiagonaux) pour résoudre les systèmes linéaires tridiagonaux.

## 4.2 Notions générales

On rappelle que les méthodes itératives ne s'appliquent que dans le cas de systèmes à coefficients dans  $\mathbb{R}$  ou  $\mathbb{C}$  mais pas dans le cas des corps finis  $\mathbb{F}_p$ .

### 4.2.1 Modèle général d'un schéma itératif

On considère une matrice  $A \in \mathbb{M}_{n \times n}(\mathbb{K})$  inversible, un vecteur  $b \in \mathbb{K}^n$  et un système linéaire

$$(S) : A x = b.$$

Le principe général d'une méthode itérative pour résoudre (S) est de générer une suite de vecteurs qui converge vers la solution  $A^{-1}b$ . Pour ce faire l'idée est d'écrire le système (S) sous une forme équivalente permettant de voir la solution comme le point fixe d'une certaine fonction, e.g. :

$$(S) \iff Bx + c = x, \quad (4.1)$$

avec  $B \in \mathbb{M}_{n \times n}(\mathbb{K})$  et  $c \in \mathbb{K}^n$  bien choisis c'est-à-dire  $\mathbb{I} - B$  inversible et  $c = (\mathbb{I} - B)A^{-1}b$ . Par exemple, si  $A = M - N$  pour deux matrices  $M, N \in \mathbb{M}_{n \times n}(\mathbb{K})$  avec  $M$  inversible, on peut choisir  $B = M^{-1}N$  et  $c = M^{-1}b$ . Dans la suite on supposera toujours que  $B \in \mathbb{M}_{n \times n}(\mathbb{K})$  et  $c \in \mathbb{K}^n$  sont choisis tels que  $\mathbb{I} - B$  inversible et  $c = (\mathbb{I} - B)A^{-1}b$  (e.g., méthode itérative *consistante*). On se donne alors un vecteur  $x^{(0)} \in \mathbb{K}^n$  et on construit une suite de vecteurs  $x^{(k)} \in \mathbb{K}^n$  à l'aide du schéma itératif

$$x^{(k+1)} = Bx^{(k)} + c, \quad k = 1, 2, \dots \quad (4.2)$$

Si la suite  $(x^{(k)})_{k \in \mathbb{N}}$  est convergente, alors elle converge vers la solution  $A^{-1}b$  de (S). En effet, si elle existe, la limite  $x^*$  est un point fixe de la fonction  $x \mapsto Bx + c$ , i.e., vérifie  $x^* = Bx^* + c$  qui est équivalent à  $Ax^* = b$  d'après (4.1).

La mise en oeuvre pratique d'une méthode itérative de la forme (4.2) nécessite la donnée d'un point de départ  $x^{(0)}$  (en général, sauf si l'on possède des informations *a priori* sur la solution, on choisit le vecteur nul) et d'une tolérance sur la solution que l'on cherche à calculer. On calcule ensuite les itérés  $x^{(k)}$ ,  $k = 1, 2, \dots$  en utilisant la formule (4.2) jusqu'à ce que le résidu  $b - Ax^{(k)}$  soit plus petit que la tolérance.

## 4.2.2 Convergence

**Définition 4.1.** La méthode itérative (4.2) pour résoudre  $Ax = b$  est dite convergente si pour toute valeur initiale  $x^{(0)} \in \mathbb{K}^n$ , on a  $\lim_{k \rightarrow +\infty} x^{(k)} = A^{-1}b$ .

**Lemme 4.2.** Si la méthode itérative (4.2) est convergente et si on note  $x = A^{-1}b$  la solution, alors

$$x^{(k)} - x = B^k(x^{(0)} - x).$$

*Démonstration.* On a  $c = (\mathbb{I}_n - B)A^{-1}b = (\mathbb{I}_n - B)x$  d'où  $x^{(k+1)} = Bx^{(k)} + (\mathbb{I}_n - B)x$  ou encore  $x^{(k+1)} - x = B(x^{(k)} - x)$  d'où le résultat.  $\square$

Remarquons que  $x^{(k)} - x$  représente l'erreur à la k-ième itération de sorte que la formule ci-dessus permet d'estimer cette erreur en fonction de l'erreur initiale.

Le résultat suivant nous donne des critères pour tester la convergence de la méthode itérative (4.2).

**Théorème 4.3.** Les assertions suivantes sont équivalentes :

- (i) La méthode itérative (4.2) est convergente ;

(ii) Pour tout  $y \in \mathbb{K}^n$ ,  $\lim_{k \rightarrow +\infty} B^k y = 0$  ;

(iii) Pour toute norme matricielle  $\|\cdot\|$  sur  $\mathbb{M}_{n \times n}(\mathbb{K})$ , on a  $\lim_{k \rightarrow +\infty} \|B^k\| = 0$ .

*Démonstration.* Admis pour ce cours. □

En pratique, les caractérisations précédentes de la convergence d'une méthode itérative ne sont pas faciles à vérifier. On utilise plutôt le résultat suivant :

**Théorème 4.4.** *Les assertions suivantes sont équivalentes :*

(i) La méthode itérative (4.2) est convergente ;

(ii)  $\rho(B) < 1$ , où  $\rho(B)$  désigne le rayon spectral de la matrice  $B$ , i.e., le maximum des modules des valeurs propres de  $B$  ;

(iii) Il existe une norme matricielle  $\|\cdot\|$  sur  $\mathbb{M}_{n \times n}(\mathbb{K})$  subordonnée à une norme vectorielle sur  $\mathbb{K}^n$  telle que  $\|B\| < 1$ .

*Démonstration.* Admis pour ce cours. □

### 4.2.3 Vitesse de convergence

L'égalité  $x^{(k)} - x = B^k(x^{(0)} - x)$  donnée précédemment implique que c'est la norme des puissances de la matrice  $B$  qui va nous renseigner sur la vitesse de convergence de la méthode itérative. Nous définissons ici les outils permettant de comparer les vitesses de convergence de différentes méthodes itératives.

**Définition 4.5.** *Considérons le schéma itératif (4.2) convergent. Soit  $\|\cdot\|$  une norme matricielle sur  $\mathbb{M}_{n \times n}(\mathbb{K})$  et soit  $k$  un entier tel que  $\|B^k\| < 1$  (l'existence d'un tel  $k$  découle du théorème 4.3). On appelle taux moyen de convergence associé à la norme  $\|\cdot\|$  pour  $k$  itérations de (4.2) le nombre positif*

$$R_k(B) = -\ln \left( \left[ \|B^k\| \right]^{\frac{1}{k}} \right).$$

**Définition 4.6.** *Considérons deux méthodes itératives consistantes et convergentes :*

$$(1) \quad x^{(k+1)} = B_1 x^{(k)} + c_1, \quad k = 1, 2, \dots,$$

$$(2) \quad x^{(k+1)} = B_2 x^{(k)} + c_2, \quad k = 1, 2, \dots$$

*Soit  $k$  un entier tel que  $\|B_1^k\| < 1$  et  $\|B_2^k\| < 1$ . On dit que (1) est plus rapide que (2) relativement à la norme  $\|\cdot\|$  si  $R_k(B_1) \geq R_k(B_2)$ .*

En pratique le calcul des  $R_k(B)$  est trop coûteux car il nécessite l'évaluation des  $B^k$ . On préfère donc estimer le taux asymptotique de convergence.

**Définition 4.7.** *On appelle taux asymptotique de convergence le nombre*

$$R_\infty(B) = \lim_{k \rightarrow +\infty} R_k(B) = -\ln(\rho(B)).$$

**Théorème 4.8.** *Avec les notations précédentes, une méthode itérative est d'autant plus rapide que son taux asymptotique de convergence est grand c'est-à-dire que  $\rho(B)$  est petit.*

## 4.3 Les méthodes itératives classiques

### 4.3.1 Principe

On considère un système linéaire  $(S) : Ax = b$  avec  $A$  inversible. L'idée est de déduire un schéma itératif de la décomposition de  $A$  sous la forme  $A = M - N$  où  $M$  est une matrice inversible. En pratique on suppose que les systèmes de matrice  $M$  sont « faciles » à résoudre (par exemple  $M$  diagonale, triangulaire, ...). Le système  $(S)$  s'écrit alors  $Mx = Nx + b$  c'est-à-dire  $x = Bx + c$  avec  $B = M^{-1}N$  et  $c = M^{-1}b$  et on considère le schéma itératif associé :

$$x^{(0)} \in \mathbb{K}^n, \quad Mx^{(k+1)} = Nx^{(k)} + b.$$

Nous allons maintenant considérer trois exemples classiques : les méthodes de Jacobi, Gauss-Seidel et de relaxation. Le point de départ de chacune de ces méthodes est l'unique décomposition de la matrice  $A = (a_{i,j})_{1 \leq i,j \leq n}$  sous la forme  $A = D - E - F$  avec :

- $D = (d_{i,j})_{1 \leq i,j \leq n}$  diagonale, telle que  $d_{i,i} = a_{i,i}$  et  $d_{i,j} = 0$  pour  $i \neq j$  ;
- $E = (e_{i,j})_{1 \leq i,j \leq n}$  triangulaire inférieure **stricte** telle que  $e_{i,j} = -a_{i,j}$  si  $i > j$  et  $e_{i,j} = 0$  si  $i \leq j$  ;
- $F = (f_{i,j})_{1 \leq i,j \leq n}$  triangulaire supérieure **stricte** telle que  $f_{i,j} = -a_{i,j}$  si  $i < j$  et  $f_{i,j} = 0$  si  $i \geq j$  ;

**Exemple :** Considérons la matrice

$$A = \begin{pmatrix} 2 & -1 & 1 \\ 2 & 2 & 2 \\ -1 & -1 & 2 \end{pmatrix}. \quad (4.3)$$

La décomposition de  $A$  sous la forme  $A = D - E - F$  décrite ci-dessus s'écrit alors

$$\underbrace{\begin{pmatrix} 2 & -1 & 1 \\ 2 & 2 & 2 \\ -1 & -1 & 2 \end{pmatrix}}_A = \underbrace{\begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{pmatrix}}_D - \underbrace{\begin{pmatrix} 0 & 0 & 0 \\ -2 & 0 & 0 \\ 1 & 1 & 0 \end{pmatrix}}_E - \underbrace{\begin{pmatrix} 0 & 1 & -1 \\ 0 & 0 & -2 \\ 0 & 0 & 0 \end{pmatrix}}_F.$$

On supposera de plus que  $D$  est inversible et on distingue les trois méthodes suivantes :

- Méthode de Jacobi :  $M = D, N = E + F$  ;
- Méthode de Gauss-Seidel :  $M = D - E, N = F$  ;
- Méthode de relaxation :  $M = \frac{1}{\omega}(D - \omega E), N = \left(\frac{1-\omega}{\omega}\right) D + F$  avec  $\omega$  paramètre réel non nul.

On remarque que la méthode de Gauss-Seidel est un cas particulier de la méthode relaxation pour  $\omega = 1$ .

### 4.3.2 Méthode de Jacobi

**Description :** On considère un système linéaire  $(S) : Ax = b$  avec  $A$  inversible. On pose  $A = M - N$  avec  $M = D$  inversible et  $N = E + F$ . Le schéma itératif s'écrit alors

$$Dx^{(k+1)} = (E + F)x^{(k)} + b \iff x^{(k+1)} = D^{-1}(E + F)x^{(k)} + D^{-1}b.$$

**Définition 4.9.** La matrice  $B_J = D^{-1}(E + F)$  s'appelle la matrice de Jacobi associée à  $A$ .

**Mise en œuvre et complexité arithmétique :** On se propose d'estimer le nombre d'opérations à virgule flottante nécessaires pour calculer  $x^{(k+1)}$  à partir de  $x^{(k)}$ . On a  $Dx^{(k+1)} = (E + F)x^{(k)} + b$  donc pour tout  $i = 1, \dots, n$ ,  $(Dx^{(k+1)})_i = ((E + F)x^{(k)})_i + b_i$  c'est-à-dire

$$a_{i,i}x_i^{(k+1)} = - \sum_{\substack{j=1 \\ j \neq i}}^n a_{i,j}x_j^{(k)} + b_i \iff x_i^{(k+1)} = \frac{1}{a_{i,i}} \left[ - \sum_{\substack{j=1 \\ j \neq i}}^n a_{i,j}x_j^{(k)} + b_i \right].$$

Pour calculer  $x_i^{(k+1)}$  à partir de  $x^{(k)}$ , on a donc besoin de  $n - 1$  multiplications,  $n - 1$  additions et 1 division soit  $2n - 1$  opérations à virgule flottante. Par conséquent il nous faudra  $n(2n - 1)$  opérations à virgule flottante pour calculer  $x^{(k+1)}$  à partir de  $x^{(k)}$  et pour  $K$  itérations, on aura besoin de  $Kn(2n - 1)$  opérations à virgule flottante. Pour comparaison, pour  $n = 1000$ , l'élimination de Gauss coûte environ  $\frac{2}{3}n^3 = 6,6 \cdot 10^8$  opérations à virgule flottante alors que par exemple  $K = 100$  itérations de la méthode de Jacobi coûtent approximativement  $2Kn^2 = 2 \cdot 10^8$  opérations à virgule flottante.

**Convergence :** D'après le théorème 4.4, on a le résultat suivant :

**Théorème 4.10.** La méthode de Jacobi converge si et seulement si  $\rho(B_J) < 1$ .

**Exemple :** Pour la matrice  $A$  donnée par (4.3), on obtient :

$$B_J = D^{-1}(E + F) = \begin{pmatrix} \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{2} \end{pmatrix} \begin{pmatrix} 0 & 1 & -1 \\ -2 & 0 & -2 \\ 1 & 1 & 0 \end{pmatrix} = \begin{pmatrix} 0 & \frac{1}{2} & -\frac{1}{2} \\ -1 & 0 & -1 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix}.$$

Les valeurs propres de la matrices  $B_J$  sont 0 et  $\pm \frac{i\sqrt{5}}{2}$ . On a donc  $\rho(B_J) = \frac{\sqrt{5}}{2} > 1$  et la méthode de Jacobi diverge.

### 4.3.3 Méthode de Gauss-Seidel

**Description** : On considère un système linéaire  $(S) : Ax = b$  avec  $A$  inversible. On pose  $A = M - N$  avec  $M = D - E$  inversible et  $N = F$ . Le schéma itératif s'écrit alors

$$(D - E)x^{(k+1)} = Fx^{(k)} + b \iff x^{(k+1)} = (D - E)^{-1}Fx^{(k)} + (D - E)^{-1}b.$$

**Définition 4.11.** La matrice  $B_{GS} = (D - E)^{-1}F$  s'appelle la matrice de Gauss-Seidel associée à  $A$ .

**Mise en œuvre et complexité arithmétique** : On se propose d'estimer le nombre d'opérations à virgule flottante nécessaires pour calculer  $x^{(k+1)}$  à partir de  $x^{(k)}$ . On a  $(D - E)x^{(k+1)} = Fx^{(k)} + b$  donc pour tout  $i = 1, \dots, n$ ,  $((D - E)x^{(k+1)})_i = (Fx^{(k)})_i + b_i$  c'est-à-dire

$$a_{i,i}x_i^{(k+1)} + \sum_{j=1}^{i-1} a_{i,j}x_j^{(k+1)} = - \sum_{j=i+1}^n a_{i,j}x_j^{(k)} + b_i,$$

ce qui entraîne

$$x_1^{(k+1)} = \frac{1}{a_{1,1}} \left[ - \sum_{j=2}^n a_{1,j}x_j^{(k)} + b_1 \right],$$

et pour  $i = 2, \dots, n$ ,

$$x_i^{(k+1)} = \frac{1}{a_{i,i}} \left[ - \sum_{j=1}^{i-1} a_{i,j}x_j^{(k)} + b_i - \sum_{j=i+1}^n a_{i,j}x_j^{(k)} \right].$$

La complexité arithmétique de la méthode de Gauss-Seidel est la même que celle de la méthode de Jacobi. Cependant, on peut remarquer que la méthode de Gauss-Seidel est plus intéressante en ce qui concerne la gestion de la mémoire. En effet, on peut écraser au fur et à mesure la valeur de  $x_i^{(k)}$  et ne stocker au cours des calculs qu'un seul vecteur de taille  $n$ , *e.g.*, le vecteur  $(x_1^{(k+1)} \dots x_i^{(k+1)} x_{i+1}^{(k)} \dots x_n^{(k)})^T$ , au lieu de deux vecteurs pour la méthode de Jacobi.

**Convergence** : D'après le théorème 4.4, on a le résultat suivant :

**Théorème 4.12.** La méthode de Gauss-Seidel converge si et seulement si  $\rho(B_{GS}) < 1$ .

**Exemple** : Pour la matrice  $A$  donnée par (4.3), on obtient :

$$B_{GS} = (D - E)^{-1}F = \begin{pmatrix} 2 & 0 & 0 \\ 2 & 2 & 0 \\ -1 & -1 & 2 \end{pmatrix}^{-1} \begin{pmatrix} 0 & 1 & -1 \\ 0 & 0 & -2 \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} & 0 & 0 \\ -\frac{1}{2} & \frac{1}{2} & 0 \\ 0 & \frac{1}{4} & \frac{1}{2} \end{pmatrix} \begin{pmatrix} 0 & 1 & -1 \\ 0 & 0 & -2 \\ 0 & 0 & 0 \end{pmatrix},$$

d'où

$$B_{GS} = \begin{pmatrix} 0 & \frac{1}{2} & -\frac{1}{2} \\ 0 & -\frac{1}{2} & -\frac{1}{2} \\ 0 & 0 & -\frac{1}{2} \end{pmatrix}.$$

Les valeurs propres de la matrices  $B_{GS}$  sont 0 et  $-\frac{1}{2}$  (de multiplicité 2). On a donc  $\rho(B_{GS}) = \frac{1}{2} < 1$  donc la méthode de Gauss-Seidel converge.

#### 4.3.4 Méthode de relaxation

On considère un système linéaire  $(S) : Ax = b$  avec  $A$  inversible. Soit  $\omega$  un paramètre réel non nul. On pose  $A = M - N$  avec  $M = \frac{1}{\omega}(D - \omega E)$  inversible et  $N = \left(\frac{1-\omega}{\omega}\right) D + F$ . Le schéma itératif s'écrit alors

$$\frac{1}{\omega}(D - \omega E)x^{(k+1)} = \left( \left( \frac{1-\omega}{\omega} \right) D + F \right) x^{(k)} + b,$$

qui est équivalent à :

$$x^{(k+1)} = (D - \omega E)^{-1} [(1 - \omega) D + \omega F] x^{(k)} + \omega (D - \omega E)^{-1} b.$$

**Définition 4.13.** La matrice  $B_R(\omega) = (D - \omega E)^{-1} [(1 - \omega) D + \omega F]$  s'appelle la matrice de relaxation associée à  $A$  et  $\omega$  est le facteur de relaxation. Si  $\omega < 1$ , on parle de sous-relaxation, si  $\omega = 1$ , on retrouve la méthode de Gauss-Seidel et si  $\omega > 1$ , on parle de sur-relaxation.

D'après le théorème 4.4, on a le résultat suivant :

**Théorème 4.14.** La méthode de relaxation converge si et seulement si  $\rho(B_R(\omega)) < 1$ .

**Exemple :** Pour la matrice  $A$  donnée par (4.3), on obtient :

$$B_R(\omega) = \begin{pmatrix} 1 - \omega & \frac{1}{2}\omega & -\frac{1}{2}\omega \\ \omega(\omega - 1) & -\frac{1}{2}\omega^2 + 1 - \omega & \frac{1}{2}\omega^2 - \omega \\ \frac{1}{2}\omega(\omega - 1)^2 & -\frac{1}{4}\omega^3 - \frac{1}{4}\omega^2 + \frac{1}{2}\omega & \frac{1}{4}\omega^3 - \frac{3}{4}\omega^2 + 1 - \omega \end{pmatrix}.$$

Les valeurs propres de la matrice  $B_R(\omega)$  dépendent en général de  $\omega$  donc la convergence de la méthode de relaxation dépendra aussi de la valeur de  $\omega$ .

#### 4.3.5 Résultats de convergence dans des cas particuliers

On s'intéresse tout d'abord au cas des matrices symétriques définies positives.

**Théorème 4.15.** Soit  $A$  une matrice symétrique définie positive et écrivons  $A = M - N$  avec  $M$  inversible et  $M^T + N$  définie positive. Alors la méthode itérative

$$x^{(0)} \in \mathbb{K}^n, \quad x^{(k+1)} = M^{-1} N x^{(k)} + M^{-1} b,$$

converge.

*Démonstration.* Admis pour ce cours. □

**Corollaire 4.16.** *Soit  $A$  une matrice symétrique définie positive. Alors la méthode de Gauss-Seidel converge.*

*Démonstration.* Pour la méthode de Gauss-Seidel, on a  $M = D - E$  et  $N = F$ . La matrice  $M$  est inversible car  $A$  est supposée définie positive et est donc inversible (voir la section 2.3). On a de plus  $M^T + N = D - E^T + F$ . Or  $A$  étant supposée symétrique, on a  $E^T = F$  d'où  $M^T + N = D$ . La matrice  $M^T + N$  est donc définie positive car, pour tout  $i = 1, \dots, n$ ,  $\langle D e_i, e_i \rangle = a_{i,i}$  et  $a_{i,i} > 0$  puisque  $A$  est définie positive (voir la section 2.3). Le théorème 4.15 précédent permet alors de conclure. □

Considérons maintenant le cas des matrices à diagonale strictement dominante.

**Définition 4.17.** *Une matrice  $A = (a_{i,j})_{1 \leq i,j \leq n}$  est dite à diagonale strictement dominante si :*

$$\forall i = 1, \dots, n, \quad |a_{i,i}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{i,j}|.$$

Par exemple, la matrice du système linéaire obtenu à la fin de la section 4.1 pour la résolution d'un système d'équations aux dérivées partielles est à diagonale strictement dominante.

**Théorème 4.18.** *Soit  $A$  une matrice à diagonale strictement dominante. Alors  $A$  est inversible et les méthodes de Jacobi et de Gauss-Seidel convergent toutes les deux.*

*Démonstration.* Admis pour ce cours. □

## 4.4 Méthode du gradient conjugué

Il s'agit d'une méthode itérative qui permet de résoudre un système linéaire  $(S) : Ax = b$  lorsque  $A$  est une matrice symétrique et définie positive. Dans cette méthode, la matrice  $A$  du système intervient une seule fois à chaque itération, lorsqu'on calcule son produit par un vecteur. Par conséquent, la méthode du gradient conjugué est particulièrement bien adaptée aux systèmes creux et de grande taille. Pour ce type de systèmes, la méthode du gradient conjugué est souvent plus efficace que les méthodes déjà présentées, à la fois en termes de complexité arithmétique et en termes d'espace mémoire nécessaire : elle est donc très utilisée en pratique.

Soit  $(S) : Ax = b$  avec  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  symétrique et définie positive, et  $b \in \mathbb{R}^n$ . La méthode du gradient conjugué construit une suite de vecteurs  $(x^{(k)})_{k=0,1,\dots}$  telle que  $x^{(m)} = A^{-1}b$  pour un indice  $m \leq n$ . En principe, il s'agit donc d'une méthode exacte. En pratique, à cause des erreurs numériques, le gradient conjugué est considéré comme une méthode itérative. Dans les applications, le nombre d'itérations nécessaires pour atteindre la précision voulue



est significativement plus petit que la taille  $n$  du système, en particulier dans le cas où on utilise des techniques de préconditionnement.

Dans la suite, le gradient conjugué sera présenté comme un cas particulier d'une famille de méthodes itératives connue sous le nom de *méthodes du gradient*.

**Définition 4.19.** Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  symétrique et définie positive. On définit la fonction  $\|\cdot\|_A : \mathbb{R}^n \rightarrow \mathbb{R}_+$  par  $\|x\|_A = \sqrt{x^T A x}$ .

**Proposition 4.20.** La fonction  $\|\cdot\|_A$  est une norme vectorielle.

*Démonstration.* En utilisant le fait que  $A$  est symétrique et définie positive, on montre que la fonction  $\|\cdot\|_A$  vérifie les trois propriétés de la définition 3.1, avec  $E = \mathbb{R}^n$  et  $\mathbb{K} = \mathbb{R}$  :

- $\forall x \in \mathbb{R}^n, \|x\|_A = 0 \Rightarrow x^T A x = 0 \Rightarrow x = 0$ , car  $A$  est définie positive ;
- $\forall \lambda \in \mathbb{R}, \forall x \in \mathbb{R}^n, \|\lambda x\|_A = \sqrt{(\lambda x)^T A (\lambda x)} = \sqrt{\lambda^2 x^T A x} = |\lambda| \sqrt{x^T A x} = |\lambda| \|x\|_A$  ;
- $\forall x, y \in \mathbb{R}^n, \|x+y\|_A = \sqrt{(x+y)^T A (x+y)} = \sqrt{x^T A x + y^T A y} \leq \sqrt{x^T A x} + \sqrt{y^T A y} = \|x\|_A + \|y\|_A$ .

□

**Définition 4.21.** Soit  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  symétrique et définie positive. On dit que les vecteurs  $u$  et  $v$  de  $\mathbb{R}^n$  sont  $A$ -conjugués si  $u^T A v = 0$ .

On remarque que l'application  $(u, v) \in (\mathbb{R}^n)^2 \mapsto u^T A v \in \mathbb{R}$  est un produit scalaire sur  $\mathbb{R}^n$ . Deux vecteurs  $u$  et  $v$  de  $\mathbb{R}^n$  sont donc  $A$ -conjugués s'ils sont orthogonaux (pour ce produit scalaire).

#### 4.4.1 Méthodes du gradient

On considère le problème suivant : minimiser sur  $\mathbb{R}^n$  la fonction  $\phi$  définie par

$$\phi(x) = \frac{1}{2} x^T A x - b^T x,$$

où la matrice  $A$  est symétrique et définie positive. L'intérêt de la question vient du fait que le minimum de  $\phi$  est atteint pour  $x^* = A^{-1} b$ , et cette solution est unique. En effet, le gradient de  $\phi$  en  $x = (x_1 \dots x_n)^T \in \mathbb{R}^n$  est le vecteur

$$\nabla \phi(x) = \left( \frac{\partial \phi}{\partial x_1} \quad \frac{\partial \phi}{\partial x_2} \quad \dots \quad \frac{\partial \phi}{\partial x_n} \right)^T = \frac{1}{2} A x + \frac{1}{2} A^T x - b = A x - b,$$

qui s'annule seulement pour  $x^* = A^{-1} b$ . Notons que dans ce calcul on a utilisé le fait que  $A$  est une matrice symétrique. Le vecteur  $x^* = A^{-1} b$  est donc l'unique point critique de  $\phi$  et, puisque  $A$  est définie positive, il s'agit d'un minimum global (admis).

On peut donc conclure de ce qui précède que minimiser  $\phi$  sur  $\mathbb{R}^n$  et résoudre le système linéaire  $Ax = b$  sont deux problèmes équivalents.

**Définition 4.22.** Soit  $(S) : Ax = b$  avec  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  symétrique et définie positive, et  $b \in \mathbb{R}^n$ . La quantité  $r(x)$  définie par

$$r(x) = b - Ax = -\nabla\phi(x),$$

est appelée résidu du système  $(S)$  en  $x$ . En particulier, on notera

$$r^{(k)} = b - Ax^{(k)} = -\nabla\phi(x^{(k)}),$$

le résidu à l'itération  $k$ .

**Remarque :** on rappelle que la valeur  $\nabla\phi(x)$  du gradient de  $\phi$  en  $x$  donne la direction de plus forte pente pour la fonction  $\phi$  au point  $x$ .

Les méthodes du gradient procèdent généralement en choisissant à l'étape  $k$  une *direction de descente* pour  $\phi$  c'est-à-dire un vecteur  $p^{(k)} \in \mathbb{R}^n$  tel que  $p^{(k)T} \nabla\phi(x^{(k)}) < 0$ . Pour calculer  $x^{(k+1)}$  à partir de  $x^{(k)}$ , on minimise alors la fonction  $\phi$  sur la droite de vecteur directeur  $p^{(k)}$  et passant par  $x^{(k)}$  : on choisit donc

$$x^{(k+1)} = x^{(k)} + \alpha_k p^{(k)}, \quad (4.4)$$

où  $\alpha_k \in \mathbb{R}$  est tel que

$$\phi(x^{(k+1)}) = \min_{\alpha \in \mathbb{R}} \phi(x^{(k)} + \alpha p^{(k)}).$$

En dérivant  $\phi(x^{(k)} + \alpha p^{(k)})$  par rapport à  $\alpha$ , on obtient

$$\begin{aligned} \frac{\partial}{\partial \alpha} (\phi(x^{(k)} + \alpha p^{(k)})) &= \frac{\partial}{\partial \alpha} \left( \frac{1}{2} (x^{(k)} + \alpha p^{(k)})^T A (x^{(k)} + \alpha p^{(k)}) - (x^{(k)} + \alpha p^{(k)})^T b \right), \\ &= \frac{\partial}{\partial \alpha} \left( \frac{1}{2} (x^{(k)T} A x^{(k)} + \alpha p^{(k)T} A x^{(k)} + \alpha x^{(k)T} A p^{(k)} + \alpha^2 p^{(k)T} A p^{(k)}) - x^{(k)T} b - \alpha p^{(k)T} b \right), \\ &= \frac{1}{2} (p^{(k)T} A x^{(k)} + x^{(k)T} A p^{(k)} + 2\alpha p^{(k)T} A p^{(k)}) - p^{(k)T} b, \\ &= (x^{(k)} + \alpha p^{(k)})^T A p^{(k)} - b^T p^{(k)}, \end{aligned}$$

car  $A$  étant symétrique,  $p^{(k)T} A x^{(k)} = x^{(k)T} A p^{(k)}$ . En imposant  $\frac{\partial}{\partial \alpha} (\phi(x^{(k)} + \alpha p^{(k)})) = 0$ , on trouve alors

$$\alpha_k = \frac{(b - Ax^{(k)})^T p^{(k)}}{p^{(k)T} A p^{(k)}} = \frac{r^{(k)T} p^{(k)}}{p^{(k)T} A p^{(k)}}. \quad (4.5)$$

En particulier, on observe que  $\alpha_k > 0$  puisque  $r^{(k)T} p^{(k)} = -p^{(k)T} \nabla\phi(x^{(k)}) > 0$ .

**Proposition 4.23.** À chaque itération, le résidu  $r^{(k+1)}$  est orthogonal à la direction de descente  $p^{(k)}$  utilisée à l'étape précédente, i.e.,  $r^{(k+1)T} p^{(k)} = 0$ .

*Démonstration.* D'après (4.4), pour  $k \in \mathbb{N}$ , on a  $b - Ax^{(k+1)} = b - Ax^{(k)} - \alpha_k A p^{(k)}$  et donc  $r^{(k+1)} = r^{(k)} - \alpha_k A p^{(k)}$ . D'où, à partir de (4.5),  $r^{(k+1)T} p^{(k)} = (r^{(k)} - \alpha_k A p^{(k)})^T p^{(k)} = r^{(k)T} p^{(k)} - \alpha_k p^{(k)T} A p^{(k)} = 0$ .  $\square$

## 4.4.2 Méthode de la plus forte pente

Dans cette méthode, on choisit à chaque itération la direction de descente

$$p^{(k)} = r^{(k)} = -\nabla\phi(x^{(k)}),$$

c'est-à-dire la direction de plus forte pente pour  $\phi$  au point  $x^{(k)}$ . La proposition 4.23 implique alors qu'à chaque itération, on choisit une direction de descente orthogonale à la précédente, *i.e.*,  $p^{(k+1)T} p^{(k)} = 0$ .

**Théorème 4.24.** *Pour la méthode de la plus forte pente on a, à l'itération  $k$  :*

$$\|x^* - x^{(k)}\|_A \leq \left( \frac{\text{Cond}_2(A) - 1}{\text{Cond}_2(A) + 1} \right)^k \|x^* - x^{(0)}\|_A.$$

*Démonstration.* Admis pour ce cours. □

Le choix  $p^{(k)} = -\nabla\phi(x^{(k)})$  peut paraître intuitivement assez efficace pour minimiser  $\phi$  vu qu'on se déplace le long de la direction de plus forte décroissance de la fonction. Pourtant, le théorème 4.24 suggère que dans certains cas la convergence de cette méthode pourrait être lente, notamment si la matrice  $A$  est mal conditionnée, *i.e.*, lorsque  $\text{Cond}_2(A) \gg 1$  (Voir Chapitre 3). Soient  $\lambda_1 \geq \dots \geq \lambda_n > 0$  les valeurs propres de  $A$  : géométriquement, le fait que le nombre de conditionnement  $\text{Cond}_2(A) = \lambda_1/\lambda_n$  (voir Proposition 3.6) soit grand est équivalent au fait que les courbes de niveau de  $\phi$  soient des hyperellipsoïdes très allongés.

**Exemple :** Soit  $A$  une matrice symétrique de taille  $2 \times 2$  ayant pour valeurs propres  $\lambda_1 = e^4 \geq \lambda_2 = e^{-2} > 0$ , et  $b$  le vecteur  $(1 \ -1)^T$ . Le conditionnement de  $A$  est alors  $\text{Cond}_2(A) \approx 403,43$  de sorte que  $A$  est mal conditionnée. La figure 4.1 montre la surface définie dans  $\mathbb{R}^3$  par la fonction  $\phi$  : on voit que cette surface ressemble à une vallée au fond assez plat.

La solution du système est approximativement  $x^* \approx (8,7361 \ -2,3478)^T$ . On choisit un vecteur initial  $x^{(0)} = (8 \ -2)^T$  assez proche de  $x^*$  et on effectue 10 itérations de la méthode de la plus forte pente. La figure 4.2 montre, en échelle logarithmique, les résidus relatifs  $\frac{\|r^{(k)}\|_2}{\|b\|_2}$  obtenus à chaque itération : on constate que 10 itérations ne suffisent pas à atteindre un résidu de norme inférieure à  $10^{-6}$ . Géométriquement, ceci peut s'expliquer par le fait qu'à chaque itération de la méthode de la plus forte pente on choisit une direction orthogonale à la précédente et donc on va rebondir sur les parois de la vallée, ce qui fait qu'on s'approche très lentement du minimum situé sur le fond.

Par contre, pour ce même exemple, la méthode du gradient conjugué converge en 2 itérations avec un résidu relatif comparable à la précision machine (voir la section suivante).

## 4.4.3 Gradient conjugué

Pour la méthode du gradient conjugué, le choix de  $p^{(k)}$  dans le schéma itératif (4.4) se fait en tenant compte des directions  $p^{(j)}$ ,  $j = 0, 1, \dots, k-1$ , calculées aux itérations précédentes.

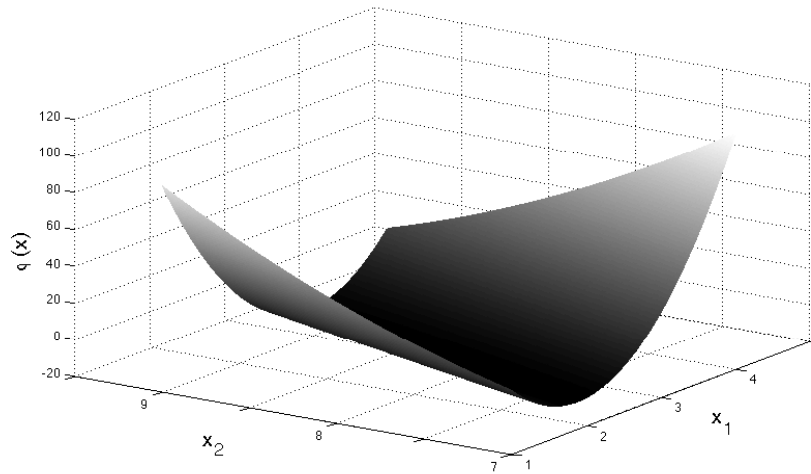


FIGURE 4.1 : Surface définie par la fonction  $\phi$ .

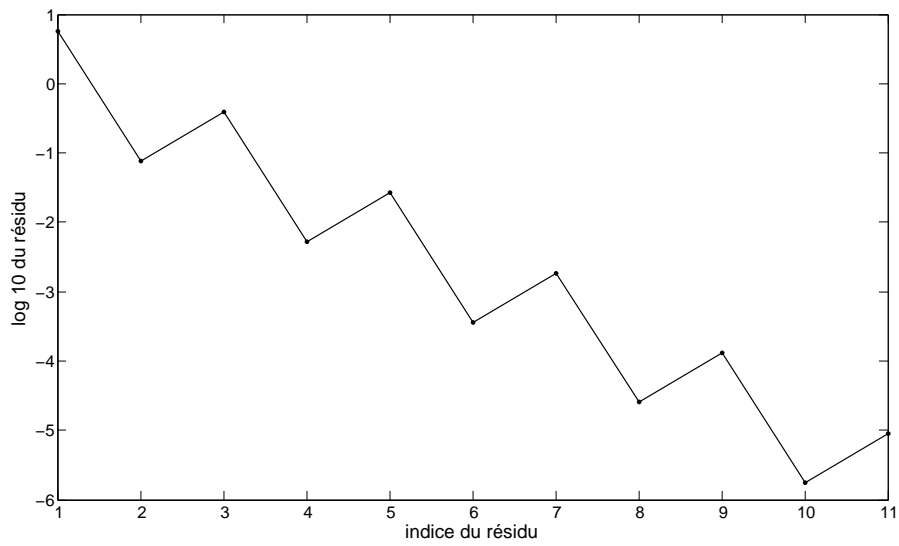


FIGURE 4.2 : Résidus relatifs calculés à chaque itération de la méthode de la plus forte pente.

On définit

$$p^{(k)} = \begin{cases} r^{(0)} & \text{si } k = 0, \\ r^{(k)} + \beta_k p^{(k-1)} & \text{si } k \geq 1, \end{cases} \quad (4.6)$$

où  $\beta_k \in \mathbb{R}$  est tel que

$$p^{(k)T} A p^{(k-1)} = 0. \quad (4.7)$$

La condition (4.7) revient à dire que les directions de descente  $p^{(k-1)}$  et  $p^{(k)}$  calculées à deux itérations consécutives de la méthode du gradient conjugué sont  $A$ -conjuguées.

En utilisant (4.6) et (4.7), on peut écrire  $\beta_k$  sous la forme

$$\beta_k = -\frac{r^{(k)T} A p^{(k-1)}}{p^{(k-1)T} A p^{(k-1)}}. \quad (4.8)$$

On vérifie alors que  $p^{(k)}$  est effectivement une direction de descente pour  $\phi$  : si  $r^{(k)} \neq 0$ , c'est-à-dire  $x^{(k)} \neq x^*$ , on a

$$p^{(k)T} \nabla \phi(x^{(k)}) = -p^{(k)T} r^{(k)} = -r^{(k)T} r^{(k)} - \beta_k p^{(k-1)T} r^{(k)} = -r^{(k)T} r^{(k)} < 0.$$

En particulier, on constate que  $p^{(k)T} r^{(k)} = r^{(k)T} r^{(k)}$  donc l'expression (4.5) pour  $\alpha_k$  devient

$$\alpha_k = \frac{r^{(k)T} r^{(k)}}{p^{(k)T} A p^{(k)}}.$$

**Lemme 4.25.** *À chaque itération, le résidu  $r^{(k)}$  est orthogonal au résidu  $r^{(k-1)}$  calculé à l'itération précédente, i.e.,  $r^{(k)T} r^{(k-1)} = 0$ .*

*Démonstration.* On a à la fois  $r^{(k)T} r^{(k-1)} = r^{(k)T} p^{(k-1)} - \beta_{k-1} r^{(k)T} p^{(k-2)} = -\beta_{k-1} r^{(k)T} p^{(k-2)}$  et  $r^{(k)T} p^{(k-2)} = r^{(k-1)T} p^{(k-2)} - \alpha_k p^{(k-1)T} A p^{(k-2)} = 0$ , d'où  $r^{(k)T} r^{(k-1)} = 0$ .  $\square$

**Lemme 4.26.** *La quantité  $\beta_k$  peut s'écrire sous la forme*

$$\beta_k = \frac{r^{(k)T} r^{(k)}}{r^{(k-1)T} r^{(k-1)}}.$$

*Démonstration.* Pour démontrer cette formule, on écrit la quantité  $p^{(k)T} r^{(k-1)}$  de deux manières différentes. On a

$$p^{(k)T} r^{(k-1)} = r^{(k)T} r^{(k-1)} + \beta_k p^{(k-1)T} r^{(k-1)} = \beta_k p^{(k-1)T} r^{(k-1)} = \beta_k r^{(k-1)T} r^{(k-1)},$$

et

$$p^{(k)T} r^{(k-1)} = p^{(k)T} r^{(k)} + \alpha_{k-1} p^{(k)T} A p^{(k-1)} = p^{(k)T} r^{(k)} = r^{(k)T} r^{(k)},$$

d'où le résultat.  $\square$

**Théorème 4.27.** *Soit  $\mathcal{S}_k$  le sous-espace vectoriel de  $\mathbb{R}^n$  engendré par les vecteurs  $p^{(0)}, \dots, p^{(k-1)}$ . Alors le vecteur  $x^{(k)}$  défini par la méthode du gradient conjugué à l'itération  $k$  minimise la fonction  $\phi$  sur  $\mathcal{S}_k$  :*

$$\phi(x^{(k)}) = \min_{x \in \mathcal{S}_k} \phi(x), \quad k \geq 1.$$

*Démonstration.* Admis pour ce cours. □

**Théorème 4.28.** Soit  $r^{(0)} \neq 0$  et  $h \geq 1$  tels que  $r^{(k)} \neq 0$  pour tout  $k \leq h$ . Alors pour  $k, j \in \{0, \dots, h\}$  avec  $k \neq j$ , on a :

$$r^{(k)T} r^{(j)} = 0 \quad \text{et} \quad p^{(k)T} A p^{(j)} = 0.$$

Autrement dit, dans la méthode du gradient conjugué, les résidus forment un ensemble de vecteurs orthogonaux et les directions de descente  $p^{(k)}$  forment un ensemble de vecteurs  $A$ -conjugués.

*Démonstration.* Par récurrence sur  $h$  (exercice). □

**Corollaire 4.29.** Il existe  $m \leq n$  tel que  $r^{(m)} = 0$ . Autrement dit, le gradient conjugué calcule la solution  $x^* = A^{-1}b$  en au plus  $n$  itérations.

Pour récapituler, l'algorithme du gradient conjugué se déroule de la manière suivante :

**Entrée :**  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  symétrique et définie positive,  $b \in \mathbb{R}^n$ , et  $x^{(0)} \in \mathbb{R}^n$ .

**Sortie :**  $x^* \in \mathbb{R}^n$  tel que  $Ax^* = b$ .

1.  $k = 0$  ;
2.  $r^{(0)} = b - Ax^{(0)}$  ;
3. Tant que  $r^{(k)} \neq 0$ , faire :
  - Si  $k = 0$ , alors faire :
    - $p^{(0)} = r^{(0)}$  ;
    - Sinon faire :
      - $\beta_k = r^{(k)T} r^{(k)} / r^{(k-1)T} r^{(k-1)}$  ;
      - $p^{(k)} = r^{(k)} + \beta_k p^{(k-1)}$  ;
  - $\alpha_k = r^{(k)T} r^{(k)} / p^{(k)T} A p^{(k)}$  ;
  - $x^{(k+1)} = x^{(k)} + \alpha_k p^{(k)}$  ;
  - $r^{(k+1)} = r^{(k)} - \alpha_k A p^{(k)}$  ;
  - $k = k + 1$  ;
4. Retourner  $x^* = x^{(k)}$ .

Les propriétés de convergence de la méthode du gradient conjugué sont données par le résultat suivant (à comparer au théorème 4.24) :

**Théorème 4.30.** Pour la méthode du gradient conjugué, on a, à l'itération  $k$  :

$$\|x^* - x^{(k)}\|_A \leq 2 \left( \frac{\sqrt{\text{Cond}_2(A)} - 1}{\sqrt{\text{Cond}_2(A)} + 1} \right)^k \|x^* - x^{(0)}\|_A$$

*Démonstration.* Admis pour ce cours. □

Quelques remarques :

- Dans la pratique, le critère d'arrêt  $r^{(k)} = 0$  est remplacé par  $\|r^{(k)}\|_2 < \epsilon_M \|b\|_2$ , et  $k$  est borné par  $k_{\max} \ll n$ . En effet, la méthode du gradient conjugué est souvent appliquée à des systèmes de grande taille, où on espère atteindre une bonne approximation de la solution  $x^*$  après un nombre d'itérations significativement plus petit que  $n$ .
- En principe, le vecteur initial  $x^{(0)}$  peut être choisi de manière arbitraire, par exemple comme le vecteur nul. Évidemment, le choix de  $x^{(0)}$  a un effet sur le nombre d'itérations nécessaires pour atteindre la solution.
- À chaque itération, l'opération la plus coûteuse du point de vue de la complexité arithmétique est la multiplication matrice-vecteur  $Ap^{(k)}$ , qui nécessite en général  $n^2$  opérations. On peut donc estimer qu'asymptotiquement le coût de la méthode du gradient conjugué est de l'ordre de  $k_{\max} n^2$  opérations. De plus, si la matrice  $A$  est creuse, comme c'est souvent le cas dans les applications (*e.g.*, résolution d'équations aux dérivées partielles par discrétisation- voir Section 4.1), le produit matrice-vecteur  $Ap^{(k)}$  peut se faire en seulement  $n$  opérations, ce qui rend la méthode du gradient conjugué plus avantageuse qu'une méthode directe comme celle de Cholesky (voir la proposition 2.18).
- Si on veut résoudre un système linéaire  $Ax = b$  où la matrice  $A$  est inversible mais n'est pas symétrique et définie positive, on peut toujours appliquer la méthode du gradient conjugué au système équivalent  $A^T Ax = A^T b$  (*méthode des équations normales*). Cependant, cette approche n'est pas recommandée lorsque  $A$  est mal conditionnée car le passage aux équations normales élève le conditionnement de la matrice au carré. Dans ce cas, il existe des versions de la méthode du gradient conjugué spécialement adaptées aux matrices non symétriques : l'une des plus utilisées est la méthode GMRES.

#### 4.4.4 Gradient conjugué avec préconditionnement

Le préconditionnement est une technique qui vise à accélérer la convergence d'une méthode itérative. On rappelle que, dans le cas de la méthode du gradient conjugué, la convergence est très rapide si la matrice  $A$  est proche de la matrice identité, ou si ses valeurs propres sont bien regroupées (voir Théorème 4.30).

**Description :** On considère un système linéaire  $(S) : Ax = b$  avec  $A$  symétrique et définie positive. Étant donnée une matrice  $C \in \mathbb{M}_{n \times n}(\mathbb{R})$  inversible, on définit le système transformé

$$(\tilde{S}) : \tilde{A} \tilde{x} = \tilde{b}, \quad \text{avec} \quad \tilde{A} = C^{-1} A (C^{-1})^T, \quad \tilde{x} = C^T x, \quad \text{et} \quad \tilde{b} = C^{-1} b.$$

On remarque que  $\tilde{A}$  est aussi symétrique et définie positive donc on peut appliquer la méthode du gradient conjugué à  $(\tilde{S})$ .

**Définition 4.31.** La matrice  $M = C C^T$  est appelée préconditionneur du système  $(\tilde{S})$ .

Le choix du préconditionneur  $M$  est généralement fait de sorte que :

- la matrice  $\tilde{A}$  soit mieux conditionnée que  $A$ , ou idéalement proche de la matrice identité, pour que la méthode du gradient conjugué appliquée à  $(\tilde{S})$  converge rapidement,
- $M$  soit « facilement » inversible car dans l'algorithme il faut résoudre un système linéaire de matrice  $M$  à chaque itération (voir ci-dessous), donc cette opération doit pouvoir être effectuée de manière stable et rapide, idéalement en un nombre d'opérations de l'ordre de  $n$ .

Pour écrire l'algorithme du gradient conjugué préconditionné, on observe que, si on note  $s^{(k)}$  le résidu de la méthode préconditionnée à l'itération  $k$ , on a, avec les notations précédentes,  $s^{(k)} = C^{-1} r^{(k)}$ , et donc  $s^{(k)T} s^{(k)} = r^{(k)T} M^{-1} r^{(k)}$ . Si  $z^{(k)}$  est tel que  $M z^{(k)} = r^{(k)}$ , alors  $s^{(k)T} s^{(k)} = z^{(k)T} r^{(k)}$ . On obtient donc l'algorithme suivant :

**Entrée :**  $A \in \mathbb{M}_{n \times n}(\mathbb{R})$  symétrique et définie positive,  $b \in \mathbb{R}^n$ ,  $x^{(0)} \in \mathbb{R}^n$ , et un préconditionneur  $M \in \mathbb{M}_{n \times n}(\mathbb{R})$  symétrique et défini positif.

**Sortie :**  $x^* \in \mathbb{R}^n$  tel que  $A x^* = b$ .

1.  $k = 0$  ;
2.  $r^{(0)} = b - A x^{(0)}$  ;
3. Tant que  $r^{(k)} \neq 0$ , faire :
  - résoudre le système linéaire  $M z^{(k)} = r^{(k)}$  ;
  - Si  $k = 0$ , alors faire :
 
$$p^{(0)} = z^{(0)} ;$$
 Sinon faire :
 
$$\beta_k = z^{(k)T} r^{(k)} / z^{(k-1)T} r^{(k-1)} ;$$

$$p^{(k)} = z^{(k)} + \beta_k p^{(k-1)} ;$$
  - $\alpha_k = z^{(k)T} r^{(k)} / p^{(k)T} A p^{(k)}$  ;
  - $x^{(k+1)} = x^{(k)} + \alpha_k p^{(k)}$  ;
  - $r^{(k+1)} = r^{(k)} - \alpha_k A p^{(k)}$  ;
  - $k = k + 1$  ;
4. Retourner  $x^* = x^{(k)}$ .



**Remarque :** la matrice  $C$  est utilisée dans la description théorique de la méthode mais n'apparaît pas explicitement dans l'algorithme.

Le choix d'un préconditionneur est un problème délicat et il existe une vaste littérature à ce sujet. Ici nous présentons deux exemples simples qui peuvent s'appliquer de manière assez générale.

**Exemple 1 :** Préconditionnement diagonal. On note  $A = (a_{i,j})_{1 \leq i,j \leq n}$ , et on choisit le conditionneur  $M = (m_{i,j})_{1 \leq i,j \leq n}$  défini par

$$m_{i,j} = \begin{cases} a_{i,i} & \text{si } i = j, \\ 0 & \text{si } i \neq j. \end{cases}$$

**Exemple 2 :** Préconditionnement de *Cholesky incomplet*. On choisit  $M = LL^T$ , où  $L = (\ell_{i,j})_{1 \leq i,j \leq n}$  est une matrice triangulaire inférieure définie de la manière suivante :

$$\ell_{i,i} = \sqrt{a_{i,i} - \sum_{r=1}^{i-1} \ell_{i,r}^2}, \quad i = 1, \dots, n,$$

$$\ell_{i,j} = \begin{cases} 0 & \text{si } a_{i,j} = 0, \\ \frac{1}{\ell_{i,j}} \left( a_{i,j} - \sum_{r=1}^{j-1} \ell_{i,r} \ell_{j,r} \right) & \text{si } a_{i,j} \neq 0, \end{cases} \quad j = 1, \dots, i-1, \quad i = 2, \dots, n.$$

On remarque que la matrice  $L$  est définie de manière à préserver l'éventuelle structure creuse de  $A$ . Les éléments non nuls de  $L$  sont calculés comme pour le facteur de Cholesky de  $A$  (voir l'algorithme de Cholesky dans la section 2.3).



# Chapitre 5

## Interpolation polynomiale

### 5.1 Le problème considéré

Dans ce chapitre, on notera  $\mathcal{P}_n = \mathbb{R}_n[x]$  l'ensemble des polynômes de degré inférieur ou égal à  $n$  et à coefficients dans  $\mathbb{R}$ . On rappelle que  $\mathcal{P}_n$  est un espace vectoriel de dimension  $n + 1$  sur  $\mathbb{R}$ . Soit  $(a, b) \in \mathbb{R}^2$  avec  $a < b$  et  $f : [a, b] \rightarrow \mathbb{R}$  une fonction que l'on suppose continue sur  $[a, b]$ . On considère  $n + 1$  points  $x_0, \dots, x_n$  de l'intervalle  $[a, b]$  tels que  $a \leq x_0 \leq x_1 \leq \dots \leq x_n \leq b$ . Le problème d'interpolation polynomiale  $(I)_{m,n}^f$  consiste à chercher s'il existe un polynôme  $P_m \in \mathcal{P}_m$  qui coïncide avec  $f$  aux nœuds  $(x_i)_{0 \leq i \leq n}$ , i.e., tel que, pour tout  $i = 0, \dots, n$ ,  $P_m(x_i) = f(x_i)$ . Si on pose  $P_m(x) = \lambda_0 + \lambda_1 x + \dots + \lambda_m x^m$  avec les  $\lambda_i$  dans  $\mathbb{R}$ , alors le problème se ramène à trouver  $\lambda_0, \dots, \lambda_m$  tels que :

$$(S) : \underbrace{\begin{pmatrix} 1 & x_0 & x_0^2 & \dots & x_0^m \\ 1 & x_1 & x_1^2 & \dots & x_1^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^m \end{pmatrix}}_V \begin{pmatrix} \lambda_0 \\ \lambda_1 \\ \vdots \\ \lambda_m \end{pmatrix} = \begin{pmatrix} f(x_0) \\ f(x_1) \\ \vdots \\ f(x_m) \end{pmatrix}. \quad (5.1)$$

Nous avons donc un système de  $n + 1$  équations linéaires en  $m + 1$  inconnues.

**Proposition 5.1.** *Le problème d'interpolation  $(I)_{m,n}^f$  admet une unique solution si et seulement si  $m = n$  et les nœuds  $(x_i)_{0 \leq i \leq n}$  sont deux à deux distincts.*

*Démonstration.* Cela découle des propriétés des matrices de Vandermonde comme  $V$ , e.g., lorsque  $n = m$ ,  $\det(V) = \prod_{0 \leq i < j \leq n} (x_j - x_i)$ .  $\square$

Dans la suite, on s'intéresse au cas où le problème d'interpolation admet une unique solution et on le notera  $(I)_n^f$ .

**Définition 5.2.** *Soit  $f : [a, b] \rightarrow \mathbb{R}$  et  $n + 1$  nœuds  $(x_i)_{0 \leq i \leq n}$  deux à deux distincts. La solution (unique) du problème  $(I)_n^f$  est appelée polynôme d'interpolation de  $f$  aux nœuds  $(x_i)_{0 \leq i \leq n}$  : ce polynôme est noté  $P_n(x; f)$ .*

Notons que même si cela n'apparaît pas dans la notation,  $P_n(x; f)$  dépend des nœuds choisis.

Ce problème d'interpolation apparaît dans un contexte expérimental. Supposons que l'on souhaite calculer les valeurs d'une fonction  $f$  mais qu'on ne connaisse pas la fonction  $f$  et qu'il soit difficile d'en calculer des valeurs. Par exemple lorsque les valeurs de  $f$  s'obtiennent par intégration ou par résolution d'une équation non linéaire ou en sommant une série compliquée. On va donc faire une table de valeurs de  $f$  pour des valeurs du paramètre  $x$  convenablement choisies et construire une règle de calcul qui permette d'approcher les valeurs de  $f$  en dehors de ces valeurs discrètes. Cette règle de calcul peut par exemple être une formule d'interpolation. Le même problème apparaît lorsque la fonction  $f$  est obtenue comme résultat de mesures. Dans ce cas, on fait un nombre fini de mesures pour des valeurs distinctes de  $x$  et on souhaite avoir une formule pour calculer les valeurs de  $f$  en dehors des valeurs pour lesquelles on a fait des mesures.

Notons que pour espérer la réussite de ces techniques, il est naturel de supposer que l'on connaisse un minimum d'information sur la fonction  $f$  à interpoler.

## 5.2 La méthode d'interpolation de Lagrange

**Définition 5.3.** Pour  $j \in \{0, \dots, n\}$ , le polynôme  $L_j^{(n)}$  défini par

$$L_j^{(n)}(x) = \prod_{\substack{i=0 \\ i \neq j}}^n \frac{x - x_i}{x_j - x_i} = \frac{(x - x_0) \cdots (x - x_{j-1})(x - x_{j+1}) \cdots (x - x_n)}{(x_j - x_0) \cdots (x_j - x_{j-1})(x_j - x_{j+1}) \cdots (x_j - x_n)},$$

est appelé interpolant de base de Lagrange ou polynôme de base de Lagrange associé à la suite  $(x_i)_{0 \leq i \leq n}$  et relatif au point  $x_j$ .

**Proposition 5.4.** Pour  $n \in \mathbb{N}$  fixé, les  $(L_j^{(n)}(x))_{0 \leq j \leq n}$  forment une base de l'espace vectoriel  $\mathcal{P}_n$  que l'on appelle base de Lagrange.

**Proposition 5.5.** Les interpolants de base de Lagrange vérifient les propriétés suivantes :

1. Pour tout  $j = 0, \dots, n$ , si on note  $g_j$  la fonction de  $[a, b]$  dans  $\mathbb{R}$  définie par  $\forall i = 0, \dots, n$ ,  $g_j(x_i) = \delta_{ij}$ , alors  $P_n(x; g_j) = L_j^{(n)}(x)$  ;

2. Si on pose

$$\pi_{n+1}(x) = \prod_{j=0}^n (x - x_j) \in \mathcal{P}_{n+1},$$

alors, pour tout  $j = 0, \dots, n$ ,

$$L_j^{(n)}(x) = \frac{\pi_{n+1}(x)}{(x - x_j) \pi'_{n+1}(x_j)}.$$

3. Pour tout  $k = 0, \dots, n$ ,  $x^k = \sum_{j=0}^n x_j^k L_j^{(n)}(x)$ .

*Démonstration.* Exercice. □

La méthode d'interpolation de Lagrange consiste à écrire le polynôme d'interpolation sur la base de Lagrange.

**Théorème 5.6.** Soit  $f : [a, b] \rightarrow \mathbb{R}$  et  $n + 1$  nœuds  $(x_i)_{0 \leq i \leq n}$  deux à deux distincts. Le polynôme d'interpolation de  $f$  aux nœuds  $(x_i)_{0 \leq i \leq n}$  s'écrit alors :

$$P_n(x; f) = \sum_{j=0}^n f(x_j) L_j^{(n)}(x).$$

*Démonstration.* On doit montrer que le polynôme  $Q(x) = \sum_{j=0}^n f(x_j) L_j^{(n)}(x)$  ainsi défini est solution du problème d'interpolation  $(I)_n^f$ . Il est clair que  $Q \in \mathcal{P}_n$ , i.e.,  $Q$  est de degré inférieur ou égal à  $n$ . De plus, pour tout  $i = 0, \dots, n$ ,  $L_j^{(n)}(x_i) = \delta_{ij}$  où  $\delta_{ij} = 1$  si  $i = j$  et  $\delta_{ij} = 0$  si  $i \neq j$  d'où  $Q(x_i) = f(x_i)$ . □

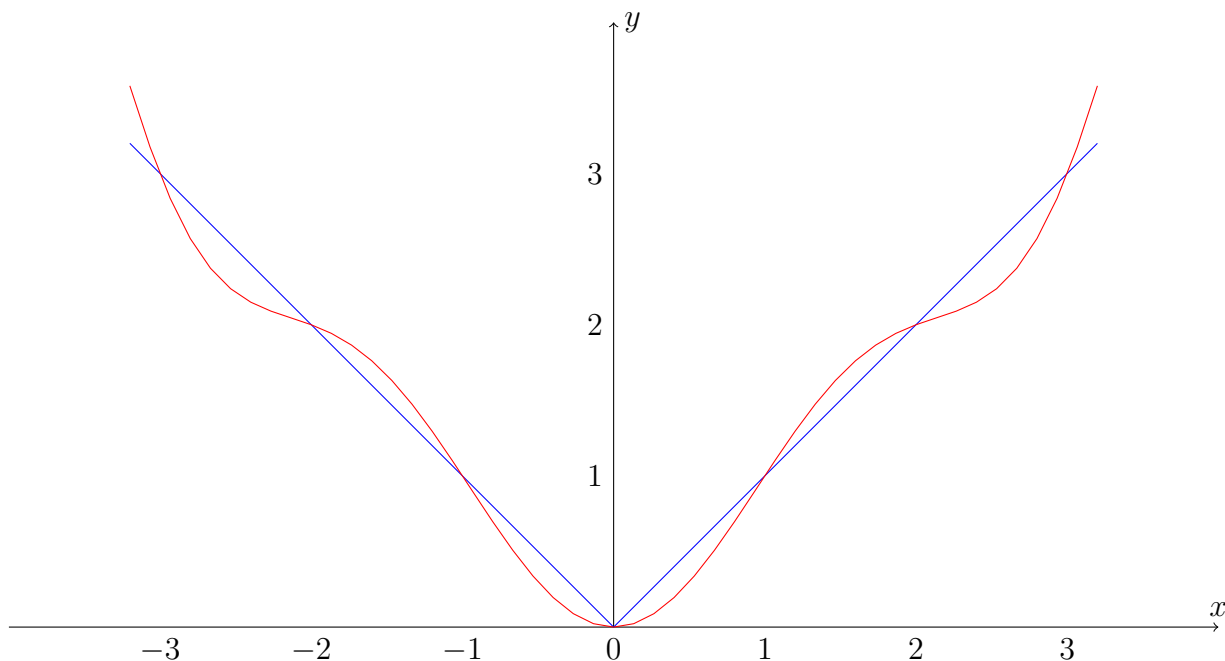
L'intérêt de la base de Lagrange par rapport à la base constituée des monômes est que nous n'avons pas besoin de résoudre un système linéaire de la forme (5.1) pour écrire le polynôme d'interpolation. L'expression de celui-ci dans la base de Lagrange s'écrit facilement. Par exemple si on choisit les nœuds  $-1, 0, 1$ , on obtient :

$$\begin{aligned} P_2(x; f) &= f(-1) \frac{(x-0)(x-1)}{(-1-0)(-1-1)} + f(0) \frac{(x+1)(x-1)}{(0+1)(0-1)} + f(1) \frac{(x+1)(x-0)}{(1+1)(1-0)}, \\ &= f(-1) \frac{x(x-1)}{2} - f(0) (x^2 - 1) + f(1) \frac{x(x+1)}{2}, \\ &= \frac{f(-1) - 2f(0) + f(1)}{2} x^2 + \frac{f(1) - f(-1)}{2} x + f(0). \end{aligned}$$

**Exemple :** Si on considère la fonction  $f : [-4, 4] \rightarrow \mathbb{R}$ ,  $x \mapsto |x|$ . Alors le polynôme d'interpolation de  $f$  relatif aux nœuds  $(x_j)_{0 \leq j \leq 8} = (-4, -3, -2, -1, 0, 1, 2, 3, 4)$  se décompose sur la base de Lagrange sous la forme  $P_8(x; f) = \sum_{j=0}^8 |x_j| L_j^{(8)}(x)$  et en développant on obtient :

$$P_8(x; f) = \frac{533}{420} x^2 - \frac{43}{144} x^4 + \frac{11}{360} x^6 - \frac{1}{1008} x^8.$$

Le graphe ci-dessous montre en bleu la courbe de  $f$  et en rouge celle de l'interpolant de Lagrange  $P_8(x; f)$ .



D'un point de vue efficacité pratique, on ne développe les  $L_j^{(n)}(x)$  pour écrire le polynôme d'interpolation dans la base monomiale. Si l'on veut évaluer le polynôme d'interpolation en un point, on utilise la formule

$$P_n(x; f) = \pi_{n+1}(x) \sum_{j=0}^n \frac{f(x_j)}{\pi'_{n+1}(x_j) (x - x_j)},$$

qui nécessite moins d'opérations à virgule flottante. Ce calcul demeure tout de même coûteux.

Le principal inconvénient de la méthode d'interpolation de Lagrange est que le fait de rajouter un nœud change complètement les interpolants de base de Lagrange et on doit donc recalculer entièrement le polynôme  $P_n(x; f)$ . On va donc considérer une autre approche qui se comporte mieux lorsqu'on rajoute des nœuds.

Notons enfin que cette méthode permet aussi d'interpoler un nuage de points. Au lieu de se donner une fonction  $f$ , on se donne une suite de valeurs discrètes  $(b_i)_{0 \leq i \leq n}$  aux nœuds  $(x_i)_{0 \leq i \leq n}$  et on cherche un polynôme  $P_n$  tel que  $P_n(x_i) = b_i$  pour  $i = 0, \dots, n$ . On obtient exactement les mêmes résultats en remplaçant  $f(x_i)$  par  $b_i$  dans tout ce qui précède.

## 5.3 Effectivité de l'interpolation : interpolant de Newton

### 5.3.1 Base d'interpolation de Newton

**Définition 5.7.** Les polynômes  $N_j^{(n)}$  définis pour  $j = 0, \dots, n$  par :

$$\left\{ \begin{array}{l} N_0^{(n)}(x) = 1, \\ N_1^{(n)}(x) = (x - x_0), \\ N_2^{(n)}(x) = (x - x_0)(x - x_1), \\ \vdots \\ N_j^{(n)}(x) = (x - x_0)(x - x_1) \cdots (x - x_{j-1}), \\ \vdots \\ N_n^{(n)}(x) = (x - x_0)(x - x_1) \cdots (x - x_{n-1}), \end{array} \right.$$

sont appelés interpolants de base de Newton ou polynômes de base de Newton relatifs à la suite de points  $(x_i)_{i=0, \dots, n-1}$ .

On remarque que là où on avait besoin de  $n+1$  points pour définir les  $L_j^{(n)}(x)$ ,  $j = 0, \dots, n$ , la définition des  $N_j^{(n)}(x)$ ,  $j = 0, \dots, n$ , ne nécessite que  $n$  points.

**Proposition 5.8.** Pour  $n \in \mathbb{N}$  fixé, les  $(N_j^{(n)}(x))_{0 \leq j \leq n}$  forment une base de l'espace vectoriel  $\mathcal{P}_n$  que l'on appelle base de Newton.

### 5.3.2 Expression de l'interpolant de Newton

Soit  $f : [a, b] \rightarrow \mathbb{R}$  et  $n$  nœuds  $(x_i)_{0 \leq i \leq n-1}$ . Essayons d'écrire le polynôme d'interpolation sur la base de Newton. Autrement dit, cherchons des nombres  $\alpha_i$ ,  $i = 0, \dots, n$  tels que  $P_n(x; f) = \sum_{i=0}^n \alpha_i N_i^{(n)}(x)$ . On a :

$$P_n(x_0; f) = \alpha_0 = f(x_0) \implies \alpha_0 = f(x_0)$$

$$P_n(x_1; f) = f(x_0) + \alpha_1(x_1 - x_0) = f(x_1) \implies \alpha_1 = \frac{f(x_0) - f(x_1)}{x_0 - x_1}$$

$$P_n(x_2; f) = f(x_0) + \frac{f(x_0) - f(x_1)}{x_0 - x_1}(x_2 - x_0) + \alpha_2(x_2 - x_0)(x_2 - x_1) = f(x_2)$$

$$\implies \alpha_2 = \frac{\frac{f(x_0) - f(x_2)}{x_0 - x_2} - \frac{f(x_0) - f(x_1)}{x_0 - x_1}}{x_2 - x_1}$$

En posant

$$f[u, v] = \frac{f(u) - f(v)}{u - v},$$

on a alors

$$\alpha_1 = f[x_0, x_1], \quad \alpha_2 = \frac{f[x_0, x_2] - f[x_0, x_1]}{x_2 - x_1} = \frac{f[x_0, x_1] - f[x_1, x_2]}{x_0 - x_2}.$$

**Définition 5.9.** Pour tout  $k \in \mathbb{N}$ , on appelle différence divisée d'ordre  $k$  de  $f$  associée à la suite de points deux à deux distincts  $(x_j)_{j \in \mathbb{N}}$  la quantité  $f[x_0, x_1, \dots, x_k]$  définie par :

$$f[x_0] = f(x_0), \quad \forall k \in \mathbb{N}^*, f[x_0, x_1, \dots, x_k] = \frac{f[x_0, x_1, \dots, x_{k-1}] - f[x_1, x_2, \dots, x_k]}{x_0 - x_k}.$$

**Théorème 5.10.** Avec les notations précédentes, on a

$$P_n(x; f) = \sum_{k=0}^n f[x_0, x_1, \dots, x_k] N_k^{(n)}(x).$$

*Démonstration.* Admis pour ce cours. □

Même si la définition de la base d'interpolation de Newton de  $\mathcal{P}_n$  ne nécessite que la donnée de  $n$  nœuds, le coefficient  $f[x_0, x_1, \dots, x_n]$  de  $N_n^{(n)}(x)$  fait intervenir le nœud  $x_n$ .

**Corollaire 5.11.** Avec les notations précédentes, on a

$$P_n(x; f) = P_{n-1}(x; f) + f[x_0, x_1, \dots, x_n] N_n^{(n)}(x).$$

La formule de ce corollaire montre que si l'on écrit le polynôme d'interpolation sur la base de Newton, alors il est facile de rajouter un point.

L'évaluation du polynôme d'interpolation dans la base de Newton est relativement simple comparé à celui sur la base de Lagrange (schéma d'évaluation type Horner).

**Proposition 5.12.** Pour tout  $k \in \mathbb{N}$ , on a :

$$f[x_0, x_1, \dots, x_k] = \sum_{j=0}^k \frac{f(x_j)}{\prod_{\substack{l=0 \\ l \neq j}}^k (x_j - x_l)} = \sum_{j=0}^k \frac{f(x_j)}{\pi'_{k+1}(x_j)}$$

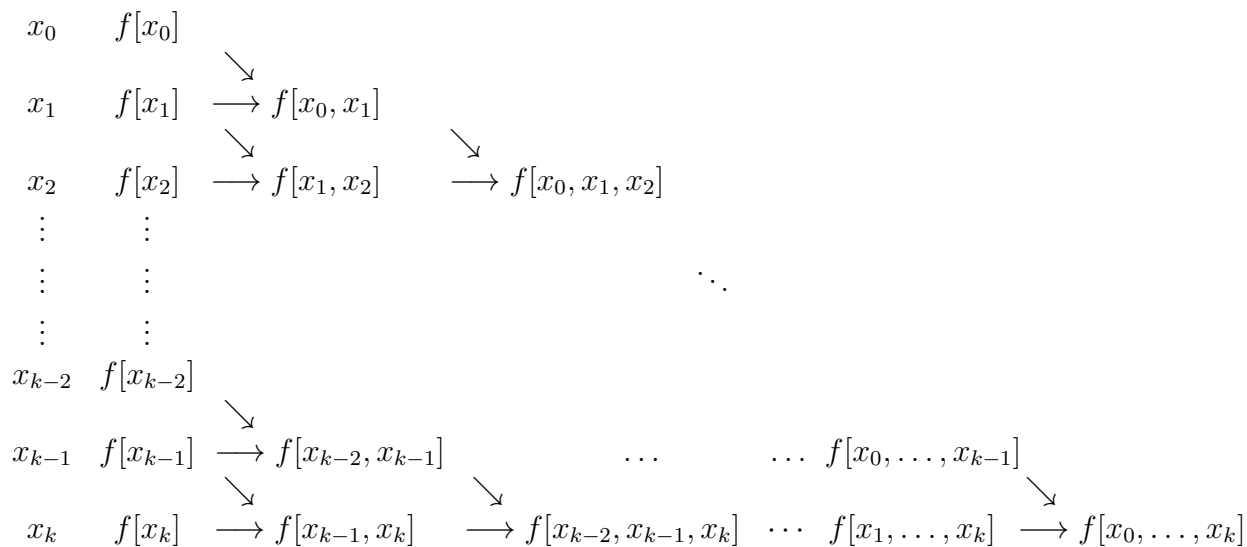
*Démonstration.* Admis pour ce cours. □

**Corollaire 5.13.** Soit  $\mathcal{S}_{k+1}$  l'ensemble des permutations sur  $\{0, 1, \dots, k+1\}$ . Pour tout  $\sigma \in \mathcal{S}_{k+1}$ , on a  $f[x_{\sigma(0)}, x_{\sigma(1)}, \dots, x_{\sigma(k)}] = f[x_0, x_1, \dots, x_k]$ .



### 5.3.3 Algorithme de calcul des différences divisées

Le calcul pratique des différences divisées est basé sur la formule de récurrence de la définition 5.9. Nous avons le tableau suivant :



On en déduit immédiatement un algorithme de calcul des différences divisées.

On note que contrairement à ce qu'il se passait pour l'interpolation de Lagrange, l'ajout d'un nouveau nœud n'oblige pas à recalculer toutes les différences divisées. Plus précisément, passer de  $n$  à  $n+1$  nœuds demande simplement le calcul de  $n$  différences divisées. Par exemple, pour passer de 3 à 4 points, on rajoute une ligne  $x_3$  en dessous de la ligne  $x_2$  dans notre tableau et on doit calculer  $(f[x_3])$ ,  $f[x_2, x_3]$ ,  $f[x_1, x_2, x_3]$  et  $f[x_0, x_1, x_2, x_3]$  soit 3 différences divisées.

## 5.4 Erreur d'interpolation

On va maintenant voir comment estimer l'erreur ponctuelle  $|f(x) - P_n(x; f)|$ .

**Lemme 5.14.** Soit  $(x_i)_{0 \leq i \leq n}$  tels que, pour tout  $i = 0, \dots, n$ ,  $x_i \in [a, b]$  et soit  $P_n(x; f)$  le polynôme d'interpolation de  $f$  aux nœuds  $(x_i)_{0 \leq i \leq n}$ . Alors, avec les notations précédentes, pour tout  $x \in [a, b]$  tel que, pour tout  $i = 0, \dots, n$ ,  $x \neq x_i$ , on a :

$$f(x) - P_n(x; f) = f[x_0, x_1, \dots, x_n, x] N_{n+1}^{(n+1)}(x).$$

*Démonstration.* Admis pour ce cours. □

**Lemme 5.15.** Si  $f \in C^n([a, b])$ , alors :

$$\exists \xi \in ]a, b[, \quad f[x_0, x_1, \dots, x_n] = \frac{1}{n!} f^{(n)}(\xi).$$

*Démonstration.* Admis pour ce cours. □

Les deux lemmes précédents mènent alors au résultat suivant et à son corollaire :

**Théorème 5.16.** *Soit  $(x_i)_{0 \leq i \leq n}$  tels que, pour tout  $i = 0, \dots, n$ ,  $x_i \in [a, b]$  et soit  $P_n(x; f)$  le polynôme d'interpolation de  $f$  aux nœuds  $(x_i)_{0 \leq i \leq n}$ . Si  $f \in \mathcal{C}^{n+1}([a, b])$ , alors :*

$$\forall x \in [a, b], \exists \xi_x \in ]a, b[, \quad f(x) - P_n(x; f) = \frac{1}{(n+1)!} f^{(n+1)}(\xi_x) N_{n+1}^{(n+1)}(x).$$

**Corollaire 5.17.** *Avec les mêmes hypothèses, on a :*

$$\forall x \in [a, b], \quad |f(x) - P_n(x; f)| \leq \frac{|N_{n+1}^{(n+1)}(x)|}{(n+1)!} \sup_{y \in [a, b]} |f^{(n+1)}(y)|.$$

# Chapitre 6

## Intégration numérique

Dans ce chapitre, nous allons voir comment approcher de façon numérique la valeur d'intégrales de la forme  $I(f) = \int_a^b f(x) dx$ . Ce problème est d'autant plus intéressant qu'en pratique on ne connaît pas forcément l'expression symbolique de  $f$  et que même si c'est le cas, la plupart des fonctions n'admettent pas de primitives pouvant s'exprimer à l'aide de fonctions élémentaires.

### 6.1 Introduction et méthodes classiques

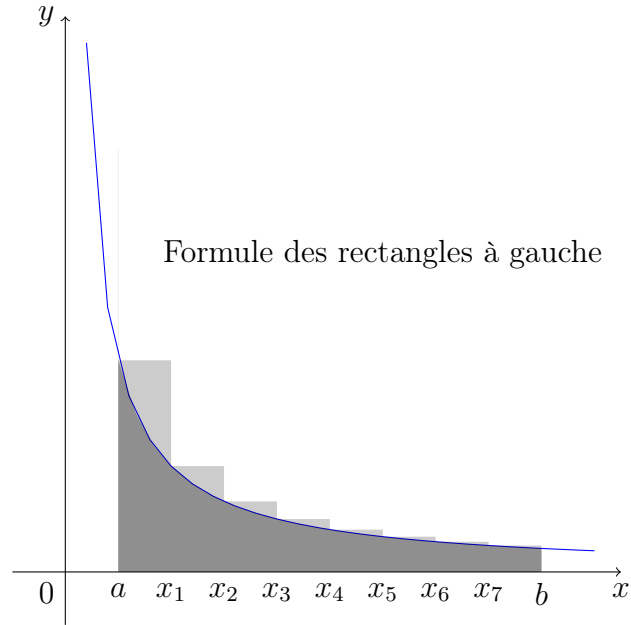
On supposera dans ce qui suit que les fonctions que l'on cherche à intégrer numériquement sont continues sur l'intervalle  $[a, b]$ . Soit  $x_0 = a < x_1 < x_2 < \dots < x_{n-1} < x_n = b$  une subdivision de l'intervalle  $[a, b]$ . La théorie élémentaire de l'intégration implique :

$$I(f) = \int_a^b f(x) dx = \lim_{n \rightarrow +\infty} \underbrace{\sum_{j=0}^{n-1} f(\xi_j)(x_{j+1} - x_j)}_{\text{Somme de Riemann}}, \quad \text{avec } \forall j, \xi_j \in [x_j, x_{j+1}].$$

Différents choix des  $\xi_j$  mènent aux méthodes classiques :

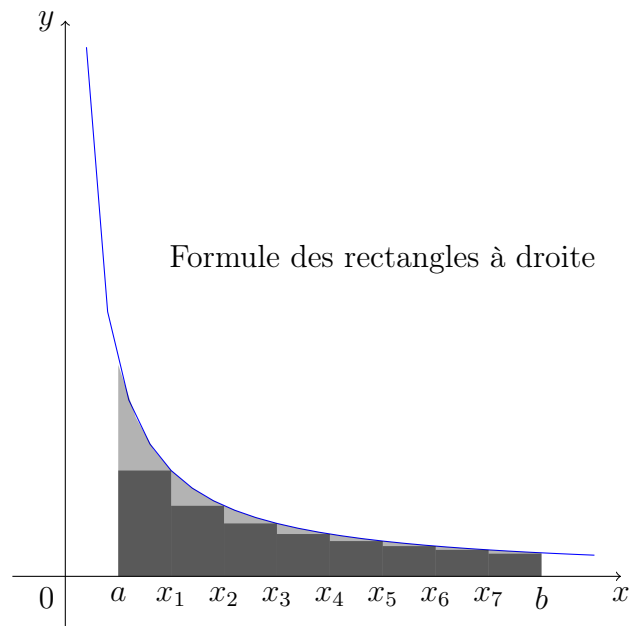
1.  $\xi_j = x_j$  donne la *formule des rectangles à gauche*  $I_{rg}(f)$  :

$$I_{rg}(f) = \sum_{j=0}^{n-1} f(x_j)(x_{j+1} - x_j).$$



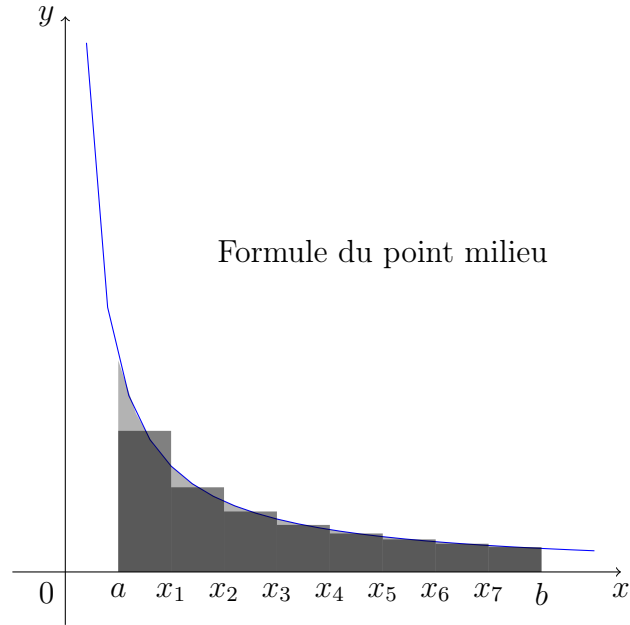
2.  $\xi_j = x_{j+1}$  donne la *formule des rectangles à droite*  $I_{rd}(f)$  :

$$I_{rd}(f) = \sum_{j=0}^{n-1} f(x_{j+1}) (x_{j+1} - x_j).$$



3.  $\xi_j = \frac{x_j + x_{j+1}}{2}$  donne la *formule du point milieu*  $I_{pm}(f)$  :

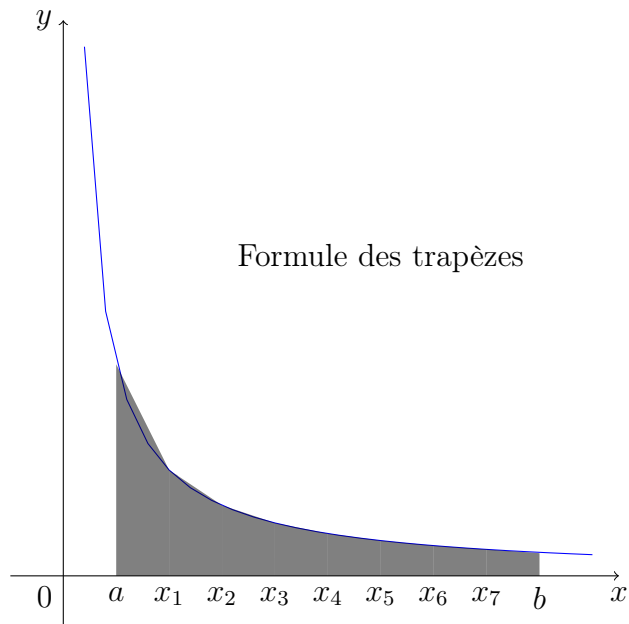
$$I_{pm}(f) = \sum_{j=0}^{n-1} f\left(\frac{x_j + x_{j+1}}{2}\right) (x_{j+1} - x_j).$$



Remarquons que les méthodes précédentes reviennent à interpoler  $f$  sur chaque intervalle  $[x_j, x_{j+1}]$  par le polynôme d'interpolation de degré 0 relatif à l'unique nœud  $\xi_j$ . Ces formules seront donc exactes pour les fonctions constantes sur  $[a, b]$  et en particulier pour  $f \in \mathcal{P}_0$ .

L'autre méthode classique est la *méthode des trapèzes* basée sur la formule :

$$I_t(f) = \sum_{j=0}^{n-1} \frac{f(x_j) + f(x_{j+1})}{2} (x_{j+1} - x_j).$$



La méthode des trapèzes revient à interpoler  $f$  sur chaque intervalle  $[x_j, x_{j+1}]$  par le polynôme d'interpolation de degré 1. Cette formule sera donc exacte pour  $f \in \mathcal{P}_1$ .

## 6.2 Formalisation de l'intégration approchée

Nous allons maintenant définir un cadre d'étude général au problème de l'intégration approchée et donner un certain nombre de résultats généraux qui seront admis pour ce cours.

Soit  $\mathcal{C}([a, b])$  l'espace vectoriel des fonctions continues sur l'intervalle  $[a, b]$  de  $\mathbb{R}$  et  $f$  une fonction de  $\mathcal{C}([a, b])$ . On suppose que l'on connaît au moins les valeurs de  $f$  en certains points  $x_0, x_1, \dots, x_n$  de l'intervalle  $[a, b]$ . On cherche alors une *formule d'intégration approchée* de la forme

$$I(f) = \int_a^b f(x) dx \approx \sum_{k=0}^n \lambda_k f(x_k) = \tilde{I}^{(n)}(f),$$

où les  $\lambda_k$  sont à déterminer. On parle aussi de *méthode d'intégration numérique* ou *formule de quadrature*.

**Définition 6.1.** Une méthode d'intégration numérique est dite d'ordre  $N$  ( $N \in \mathbb{N}$ ) si elle est exacte sur  $\mathcal{P}_N$ .

Par exemple, la méthode des rectangles à gauche ou à droite et la méthode du point milieu sont d'ordre 0 et celle des trapèzes est d'ordre 1.

En pratique, connaissant les valeurs de  $f$  aux points  $x_0, \dots, x_n$ , on remplace  $f$  par le polynôme d'interpolation  $\sum_{k=0}^n f(x_k) L_k^{(n)}(x)$  écrit dans la base de Lagrange. On a alors la formule d'intégration approchée :

$$\tilde{I}^{(n)}(f) = \sum_{k=0}^n A_k^{(n)} f(x_k), \quad A_k^{(n)} = \int_a^b L_k^{(n)}(x) dx, \quad (6.1)$$

qui est exacte sur  $\mathcal{P}_n([a, b])$  (admis).

**Théorème 6.2.** Soit  $f \in \mathcal{C}^{n+1}([a, b])$  et  $\tilde{I}^{(n)}(f)$  donnée par (6.1). Alors on a la majoration suivante de l'erreur d'intégration :

$$|I(f) - \tilde{I}^{(n)}(f)| \leq \frac{M_{n+1}}{(n+1)!} \int_a^b |\pi_{n+1}(x)| dx, \quad M_{n+1} = \sup_{x \in [a, b]} |f^{(n+1)}(x)|, \quad \pi_{n+1}(x) = \prod_{j=0}^n (x - x_j).$$

### 6.3 Formules de Newton-Côtes

D'après ce qui précède, pour obtenir notre formule d'intégration approchée, on doit donc calculer les

$$A_k^{(n)} = \int_a^b L_k^{(n)}(x) dx.$$

Pour ceci, nous allons supposer que les points d'interpolation sont équidistants, *i.e.*,  $x_{j+1} - x_j$  ne dépend pas de  $j$ , que  $x_0 = a$ ,  $x_n = b$  et  $n \geq 1$ .

**Proposition 6.3.** *Pour  $k = 0, 1, \dots, n$ , on a :*

(i)

$$A_k^{(n)} = \frac{(b-a)}{n} \frac{(-1)^{n-k}}{k!(n-k)!} \int_0^n \prod_{\substack{j=0 \\ j \neq k}}^n (y-j) dy.$$

(ii)  $A_{n-k}^{(n)} = A_k^{(n)}$ .

*Démonstration.* Le (i) s'obtient par un changement de variable  $x = a + y \frac{b-a}{n}$ . Pour le (ii), on fait le changement d'indice  $k = n - k$  en remarquant que  $n - (n - k) = k$ .  $\square$

Pour  $n = 1$ , on obtient  $A_0^{(1)} = A_1^{(1)} = \frac{b-a}{2}$  d'où  $\tilde{I}^{(1)}(f) = \frac{b-a}{2}(f(a) + f(b))$  et on retrouve la formule des trapèzes.

Pour  $n = 2$ , on trouve  $A_0^{(2)} = A_2^{(2)} = \frac{b-a}{6}$  et  $A_1^{(2)} = \frac{4(b-a)}{6}$  d'où

$$\tilde{I}^{(2)}(f) = \frac{b-a}{6} \left( f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right) \text{ (formule de Simpson).}$$

**Calcul pratique des coefficients :** Pour calculer les coefficients  $A_k^{(n)}$ , on peut utiliser le fait que  $\tilde{I}^{(n)}$  est exacte sur  $\mathcal{P}_n([a, b])$ . Par exemple pour  $n = 1$ ,  $a = -1$  et  $b = 1$ , on a  $\tilde{I}^{(1)}(f) = A_0^{(1)} f(-1) + A_1^{(1)} f(1)$ . Or  $\tilde{I}^{(1)}$  est exacte sur  $\mathcal{P}_1([-1, 1])$  donc  $\tilde{I}^{(1)}(1) = I(1) = \int_{-1}^1 1 dx = 2$  et  $\tilde{I}^{(1)}(x) = I(x) = \int_{-1}^1 x dx = 0$ . On obtient donc le système linéaire

$$\begin{cases} A_0^{(1)} + A_1^{(1)} = 2, \\ -A_0^{(1)} + A_1^{(1)} = 0, \end{cases}$$

d'où  $A_0^{(1)} = A_1^{(1)} = 1$ . Remarquons que cette formule n'est pas exacte sur  $\mathcal{P}_2([-1, 1])$  puisque  $I(x^2) = \int_{-1}^1 x^2 dx = \frac{2}{3}$  alors que  $\tilde{I}^{(1)}(x^2) = 2$ .

**Théorème 6.4.** *Considérons l'erreur  $\mathcal{E}_n(f) = I(f) - \sum_{i=0}^n A_i^{(n)} f(x_i)$ . Alors :*

1. Si  $n$  est impair et si  $f \in \mathcal{C}^{n+1}([a, b])$ , alors il existe  $\xi \in [a, b]$  tel que :

$$\mathcal{E}_n(f) = \left(\frac{b-a}{n}\right)^{n+2} \frac{f^{(n+1)}(\xi)}{(n+1)!} \int_0^n t(t-1)\cdots(t-n) dt.$$

2. Si  $n$  est pair et si  $f \in \mathcal{C}^{n+2}([a, b])$ , alors il existe  $\xi \in [a, b]$  tel que :

$$\mathcal{E}_n(f) = \left(\frac{b-a}{n}\right)^{n+3} \frac{f^{(n+2)}(\xi)}{(n+2)!} \int_0^n t^2(t-1)\cdots(t-n) dt.$$

Pour  $n = 2$ , si  $f \in \mathcal{C}^4([a, b])$ , alors l'erreur d'approximation commise par la formule de Simpson vaut  $-h^5 \frac{f^{(4)}(\xi)}{90}$  où  $h = \frac{b-a}{2}$  et  $\xi \in [a, b]$ .

## 6.4 Stabilité des méthodes d'intégration

La stabilité d'une méthode numérique mesure la « sensibilité de la méthode aux erreurs de calculs ». Considérons une formule d'intégration approchée  $\tilde{I}^{(n)}(f) = \sum_{k=0}^n A_k^{(n)} f(x_k)$ . Supposons maintenant que les valeurs calculées des  $f(x_k)$  ne soient pas exactes. On a :

$\sum_{k=0}^n A_k^{(n)} (f(x_k) + \epsilon_k) - \sum_{k=0}^n A_k^{(n)} f(x_k) = \sum_{k=0}^n A_k^{(n)} \epsilon_k$ . Par conséquent,

$$\left| \sum_{k=0}^n A_k^{(n)} \epsilon_k \right| \leq \left( \max_{0 \leq k \leq n} |\epsilon_k| \right) \sum_{k=0}^n |A_k^{(n)}|,$$

et le terme  $\sum_{k=0}^n |A_k^{(n)}|$  dépend de la méthode.

**Définition 6.5.** La formule d'intégration numérique  $\tilde{I}^{(n)}(f) = \sum_{k=0}^n A_k^{(n)} f(x_k)$  est dite stable s'il existe  $M \in \mathbb{R}$  tel que :  $\forall n \in \mathbb{N}, \forall (\epsilon_0, \dots, \epsilon_n) \in \mathbb{R}^{n+1}, \left| \sum_{k=0}^n A_k^{(n)} \epsilon_k \right| \leq M \max_{0 \leq k \leq n} |\epsilon_k|$ .

**Théorème 6.6.** Avec les notations précédentes, une condition nécessaire et suffisante de stabilité est qu'il existe  $M \in \mathbb{R}$  (indépendant de  $n$ ) tel que  $\sum_{k=0}^n |A_k^{(n)}| \leq M$ .

Concernant les formules de Newton-Côtes, on peut montrer que pour certaines valeurs de  $k$ ,  $\lim_{n \rightarrow \infty} |A_k^{(n)}| = +\infty$  de sorte que pour de grandes valeurs de  $n$  ces formules ne sont pas stables.

## 6.5 Formules d'intégration composées

Les formules d'intégration composées sont les plus utilisées en pratique. Le principe de ces formules consiste à décomposer l'intervalle  $[a, b]$  en  $k$  intervalles  $[a_i, a_{i+1}]$ ,  $i = 0, \dots, k-1$ . Grâce à la relation de Chasles, on écrit alors :

$$I(f) = \int_a^b f(x) dx = \sum_{i=0}^{k-1} \underbrace{\int_{a_i}^{a_{i+1}} f(x) dx}_{I_i(f)},$$



et on approche chaque  $I_i(f)$  par une formule d'intégration numérique vu précédemment. Notons que pour la stabilité, il est judicieux de choisir une formule avec un  $n$  petit comme par exemple celle de Simpson ( $n = 2$ ).

Ces méthodes composées sont d'autant plus intéressantes que l'erreur d'approximation diminue lorsque la taille de l'intervalle diminue. Par exemple, avec la formule de Simpson obtenue précédemment, si l'on subdivise l'intervalle  $[a, b]$  en  $k$  sous intervalles avec  $k$  pair, on obtient la formule d'intégration suivante :

$$\frac{h}{3} \left( f(a_0) + 2 \sum_{i=1}^{k/2-1} f(a_{2i}) + 4 \sum_{i=1}^{k/2} f(a_{2i-1}) + f(a_k) \right) \text{ (formule de Simpson composée),}$$

avec  $h = \frac{b-a}{k}$ ,  $a_0 = a$ ,  $a_k = b$  et  $a_i = a_{i-1} + h$ . Lorsque  $f \in \mathcal{C}^4([a, b])$ , l'erreur d'approximation de cette formule composée est alors de  $-k h^5 \frac{f^{(4)}(\xi)}{180}$  où  $h = \frac{b-a}{k}$  et  $\xi \in [a, b]$ .



# Chapitre 7

## Résolution d'équations et de systèmes d'équations non linéaires

Le dernier chapitre de ce cours est consacré à la résolution d'équations et de systèmes d'équations non linéaires.

On considère tout d'abord une fonction  $f : \mathbb{R} \rightarrow \mathbb{R}$  d'une seule variable réelle et on cherche à résoudre l'équation  $f(x) = 0$  c'est-à-dire trouver une valeur approchée  $\bar{x}$  d'un réel  $\tilde{x}$  vérifiant  $f(\tilde{x}) = 0$ .

La mise en oeuvre pratique des méthodes que nous allons voir nécessite la donnée d'une tolérance sur la solution que l'on cherche à calculer. L'algorithme numérique utilisé doit alors avoir un *critère d'arrêt* dépendant de cette tolérance et nous assurant que la solution calculée a bien la précision recherchée. En fonction de la méthode utilisée, on peut parfois savoir à l'avance combien d'étapes de l'algorithme sont nécessaires pour obtenir la précision recherchée (méthode de dichotomie) ou alors il nous faut à chaque étape vérifier une condition nous permettant d'arrêter le processus lorsque l'on est certain d'avoir obtenu la précision requise sur la solution (méthodes de points fixes).

Pour comparer les différentes méthodes de résolution que l'on va considérer, on utilise les notions suivantes de *vitesse de convergence* d'une suite :

**Définition 7.1.** Soit  $(x_n)_{n \in \mathbb{N}}$  une suite convergente et soit  $\tilde{x}$  sa limite.

1. On dit que la convergence de  $(x_n)_{n \in \mathbb{N}}$  est linéaire de facteur  $K \in ]0, 1[$  s'il existe  $n_0 \in \mathbb{N}$  tel que, pour tout  $n \geq n_0$ ,  $|x_{n+1} - \tilde{x}| \leq K |x_n - \tilde{x}|$ .
2. On dit que la convergence de  $(x_n)_{n \in \mathbb{N}}$  est superlinéaire d'ordre  $p \in \mathbb{N}$ ,  $p > 1$  s'il existe  $n_0 \in \mathbb{N}$  et  $K > 0$  tels que, pour tout  $n \geq n_0$ ,  $|x_{n+1} - \tilde{x}| \leq K |x_n - \tilde{x}|^p$ . Si  $p = 2$ , on parle de convergence quadratique et si  $p = 3$  on parle de convergence cubique.

Remarquons que  $K$  n'est pas unique et qu'en pratique il peut être difficile de prouver la

convergence d'une méthode d'autant plus qu'il faut tenir compte des erreurs d'arrondis. On utilise en général les notions de convergence plus faible suivantes :

**Définition 7.2.** Soit  $(x_n)_{n \in \mathbb{N}}$  une suite convergent vers une limite  $\tilde{x}$ . On dit que la convergence de  $(x_n)_{n \in \mathbb{N}}$  est linéaire de facteur  $K$  (resp. superlinéaire d'ordre  $p$ ) s'il existe une suite  $(y_n)_{n \in \mathbb{N}}$  convergent vers 0, linéaire de facteur  $K$  (resp. superlinéaire d'ordre  $p$ ) au sens de la définition 7.1 telle que  $|x_n - \tilde{x}| \leq y_n$ .

Soit  $d_n = -\log_{10}(|x_n - \tilde{x}|)$  une « mesure » du nombre de décimales exactes de  $x_n$ . Si la convergence est d'ordre  $p$ , alors asymptotiquement, on a  $|x_{n+1} - \tilde{x}| \sim K |x_n - \tilde{x}|^p$  d'où  $-d_{n+1} \sim \log_{10}(K) - p d_n$  et donc asymptotiquement  $x_{n+1}$  a  $p$  fois plus de décimales exactes que  $x_n$ . Ainsi, l'ordre  $p$  de la convergence représente asymptotiquement le facteur multiplicatif du nombre de décimales exactes que l'on gagne à chaque itération. Nous avons donc intérêt à ce qu'il soit le plus grand possible.

## 7.1 Méthode de dichotomie

La méthode classique de dichotomie est une méthode de localisation des racines d'une équation  $f(x) = 0$  basée sur le théorème des valeurs intermédiaires : si  $f$  est continue sur  $[a, b]$  et  $f(a)f(b) < 0$ , alors il existe  $\tilde{x} \in ]a, b[$  tel que  $f(\tilde{x}) = 0$ . L'idée est donc de partir d'un intervalle  $[a, b]$  vérifiant la propriété  $f(a)f(b) < 0$ , de le scinder en deux intervalles  $[a, c]$  et  $[c, b]$  avec  $c = \frac{a+b}{2}$ , et de tester les bornes des nouveaux intervalles (on calcule  $f(a)f(c)$  et  $f(c)f(b)$ ) pour en trouver un (au moins) qui vérifie encore la propriété, *i.e.*,  $f(a)f(c) < 0$  ou/et  $f(c)f(b) < 0$ . On itère ensuite ce procédé un certain nombre de fois dépendant de la précision que l'on recherche sur la solution (voir Théorème 7.3 ci-dessous). On obtient l'algorithme suivant :

**Entrées** : la fonction<sup>1</sup>  $f$ ,  $(a, b) \in \mathbb{R}^2$  tels que  $f$  est continue sur  $[a, b]$  et  $f(a)f(b) < 0$  et la précision  $\epsilon$ .

**Sortie** :  $x_{k+1}$  valeur approchée de  $\tilde{x}$  solution de  $f(\tilde{x}) = 0$  à  $\epsilon$  près.

1.  $x_0 \leftarrow a, y_0 \leftarrow b$  ;
2. Pour  $k$  de 0 à  $E\left(\frac{\ln(b-a) - \ln(\epsilon)}{\ln(2)}\right)$  par pas de 1, faire :
  - Si  $f(x_k)f\left(\frac{x_k + y_k}{2}\right) > 0$ , alors  $x_{k+1} \leftarrow \frac{x_k + y_k}{2}, y_{k+1} \leftarrow y_k$  ;
  - Si  $f(x_k)f\left(\frac{x_k + y_k}{2}\right) < 0$ , alors  $x_{k+1} \leftarrow x_k, y_{k+1} \leftarrow \frac{x_k + y_k}{2}$  ;
  - Sinon retourner  $\frac{x_k + y_k}{2}$  ;

---

<sup>1</sup>Il suffit en fait de connaître un moyen d'évaluer les valeurs de la fonction

3. Retourner  $x_{k+1}$ .

Cet algorithme construit une suite de segments emboîtés contenant tous la solution  $\tilde{x}$ . À chaque passage dans la boucle nous devons calculer une évaluation de  $f$ . Remarquons qu'en pratique, avec les arrondis, «  $> 0$  » et «  $< 0$  » ne veulent rien dire ! On démontre maintenant que cet algorithme est correct dans le sens où il calcule bien une valeur approchée de la solution à  $\epsilon$  près.

**Théorème 7.3.** *Le nombre minimum d'itérations de la méthode de dichotomie nécessaire pour approcher  $\tilde{x}$  à  $\epsilon$  près est  $E\left(\frac{\ln(b-a)-\ln(\epsilon)}{\ln(2)}\right) + 1$ , où  $E(x)$  désigne la partie entière d'un réel  $x$ .*

*Démonstration.* À la première itération, la longueur de l'intervalle est  $\frac{b-a}{2}$ , ..., à la  $n$ ème itération, la longueur de l'intervalle est  $\frac{b-a}{2^n}$ . L'erreur commise à l'étape  $n$  est donc majorée par  $\frac{b-a}{2^n}$ . Le nombre  $n$  d'itérations à effectuer doit alors vérifier  $\frac{b-a}{2^n} \leq \epsilon$  qui est équivalent à  $n \geq \frac{\ln(b-a)-\ln(\epsilon)}{\ln(2)}$  d'où le résultat.  $\square$

**Proposition 7.4.** *La convergence de la méthode de dichotomie est linéaire de facteur  $\frac{1}{2}$ .*

*Démonstration.* Prendre  $y_n = \frac{b-a}{2^n}$  dans la définition 7.2.  $\square$

## 7.2 Méthode du point fixe

La méthode itérative du point fixe que nous allons décrire est aussi appelée méthode des approximations successives.

**Définition 7.5.** *Soit  $g : \mathbb{R} \rightarrow \mathbb{R}$ . On dit que  $x \in \mathbb{R}$  est un point fixe de  $g$  si  $g(x) = x$ .*

Le principe de la méthode du point fixe est d'associer à l'équation  $f(x) = 0$  une équation de point fixe  $g(x) = x$  de sorte que trouver une solution de  $f(x) = 0$  équivaut à trouver un point fixe de  $g$ . La technique pour approximer le point fixe de  $g$  est alors basée sur le résultat suivant :

**Lemme 7.6.** *Soit  $(x_n)_{n \in \mathbb{N}}$  la suite définie par  $x_0 \in \mathbb{R}$  donné et  $x_{n+1} = g(x_n)$ . Si  $(x_n)_{n \in \mathbb{N}}$  est convergente et  $g$  est continue, alors la limite de  $(x_n)_{n \in \mathbb{N}}$  est un point fixe de  $g$ .*

Nous devons donc trouver des conditions sur  $g$  pour que la suite  $(x_n)_{n \in \mathbb{N}}$  définie ci-dessus converge.

**Définition 7.7.** *Soit  $g : \Omega \subseteq \mathbb{R} \rightarrow \mathbb{R}$ . On dit que  $g$  est lipschitzienne sur  $\Omega$  de constante de lipschitz  $\gamma$  (ou  $\gamma$ -lipschitzienne) si pour tout  $(x, y) \in \Omega^2$ , on a  $|g(x) - g(y)| \leq \gamma|x - y|$ . On dit que  $g$  est strictement contractante sur  $\Omega$  si  $g$  est  $\gamma$ -lipschitzienne sur  $\Omega$  avec  $\gamma < 1$ .*

On remarque que  $g$  lipschitzienne sur  $\Omega$  implique en particulier  $g$  continue sur  $\Omega$ .

**Théorème 7.8** (Théorème du point fixe). *Soit  $g$  une application strictement contractante sur un intervalle  $[a, b] \subset \mathbb{R}$  de constante de Lipschitz  $\gamma < 1$ . Supposons que l'intervalle  $[a, b]$  soit stable sous  $g$ , i.e.,  $g([a, b]) \subseteq [a, b]$  ou encore pour tout  $x \in [a, b]$ ,  $g(x) \in [a, b]$ . Alors  $g$  admet un unique point fixe  $x^* \in [a, b]$  et la suite définie par  $x_{n+1} = g(x_n)$  converge linéairement de facteur  $\gamma$  vers  $x^*$  pour tout point initial  $x_0 \in [a, b]$ . De plus,*

$$\forall n \in \mathbb{N}, |x_n - x^*| \leq \frac{\gamma^n}{1 - \gamma} |x_1 - x_0|.$$

*Démonstration.* Admis pour ce cours. □

On remarque que l'erreur est d'autant plus petite que  $\gamma$  est proche de 0. De plus, on peut montrer que l'on a aussi

$$\forall n \in \mathbb{N}, |x_n - x^*| \leq \frac{\gamma}{1 - \gamma} |x_n - x_{n-1}|,$$

et si  $\gamma \leq \frac{1}{2}$ , alors  $|x_n - x^*| \leq |x_n - x_{n-1}|$ . Dans ce cas, on pourra utiliser le test d'arrêt  $|x_n - x_{n-1}| < \epsilon$  qui certifiera une précision  $\epsilon$  sur le résultat.

La proposition suivante donne un critère pour tester le fait qu'une fonction soit contractante sur un intervalle donné.

**Proposition 7.9.** *Soit  $g$  une fonction dérivable sur l'intervalle  $[a, b]$ . Si sa dérivée  $g'$  vérifie  $\max_{x \in [a, b]} |g'(x)| = L < 1$ , alors  $g$  est strictement contractante sur  $[a, b]$  de constante de Lipschitz  $L$ .*

*Démonstration.* Utiliser le théorème des accroissements finis. □

On en déduit alors la proposition suivante :

**Proposition 7.10.** *Soit  $x^* \in [a, b]$  un point fixe d'une fonction  $g \in \mathcal{C}^1([a, b])$ .*

- *Si  $|g'(x^*)| < 1$ , alors il existe un intervalle  $[\alpha, \beta] \subseteq [a, b]$  contenant  $x^*$  pour lequel la suite définie par  $x_0 \in [\alpha, \beta]$  et  $x_{n+1} = g(x_n)$  converge vers  $x^*$  ;*
- *Si  $|g'(x^*)| > 1$ , alors pour tout  $x_0 \neq x^*$ , la suite définie par  $x_0$  et  $x_{n+1} = g(x_n)$  ne converge pas vers  $x^*$  ;*
- *Si  $|g'(x^*)| = 1$ , on ne peut pas conclure.*

En pratique, on « estime »  $g'(x^*)$ , i.e., on en a une valeur approchée  $\overline{g'(x^*)}$ . Si  $|\overline{g'(x^*)}| > 1$ , alors on élimine la méthode et si  $|\overline{g'(x^*)}| < 1$ , on cherche un intervalle  $[\alpha, \beta] \subseteq [a, b]$  dans lequel  $\max_{x \in [\alpha, \beta]} |g'(x)| < 1$  et  $g([\alpha, \beta]) \subseteq [\alpha, \beta]$ .

Revenons maintenant à notre problème initial où on cherche à résoudre une équation  $f(x) = 0$ . Posons  $g(x) = x - f(x)$  de sorte que trouver les solutions de  $f(x) = 0$  soit équivalent

à trouver les points fixes de  $g$ . D'après le théorème du point fixe (Théorème 7.8), une condition suffisante pour que  $g$  admette un point fixe dans l'intervalle  $[a, b]$  est que  $[a, b]$  soit stable sous  $g$  et que  $g$  soit strictement contractante sur  $[a, b]$  de constante de Lipschitz  $\gamma < 1$ . On a alors (conséquence directe de la définition d'une fonction contractante)

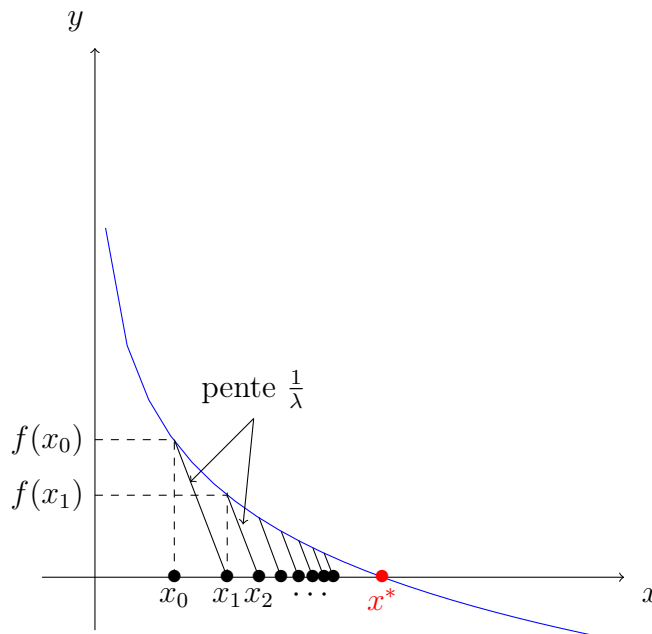
$$\forall x \in [a, b], \quad |g'(x)| < \gamma \iff |1 - f'(x)| < \gamma,$$

qui implique en particulier que  $f'$  ne s'annule pas sur  $[a, b]$  et que  $f$  est donc monotone sur  $[a, b]$ .

Étant donnée  $f$ , le choix précédent de  $g$  est restrictif et on peut choisir  $g(x) = x - \lambda f(x)$  où  $\lambda$  est une constante arbitraire non nulle. Comme précédemment, une condition suffisante pour que  $g$  ait un point fixe dans  $[a, b]$  est que  $[a, b]$  soit stable sous  $g$  et que  $g$  soit strictement contractante sur  $[a, b]$  de constante de Lipschitz  $\gamma < 1$ . On obtient alors

$$\forall x \in [a, b], \quad |1 - \lambda f'(x)| < \gamma < 1, \tag{7.1}$$

ce qui implique en particulier que  $f'$  ne change pas de signe sur  $[a, b]$  et que  $\lambda$  est du même signe que  $f'$ . Géométriquement, on construit la suite des itérés  $x_{n+1} = x_n - \lambda f(x_n)$ . En remarquant que la droite de pente  $\mu$  et passant par le point  $(x_n, f(x_n))$  a pour équation  $y = f(x_n) + \mu(x - x_n)$  et coupe l'axe des abscisses en  $x = x_n - \frac{f(x_n)}{\mu}$ , on voit que  $x_{n+1}$  s'obtient comme point d'intersection de la droite de pente  $\frac{1}{\lambda}$  passant par le point  $(x_n, f(x_n))$  avec l'axe des abscisses. En itérant ce procédé, on obtient la méthode illustrée sur le graphique suivant :



La proposition suivante donne des précisions sur l'ordre de convergence d'une méthode de point fixe.

**Proposition 7.11.** On considère l'équation  $g(x) = x$  où  $g$  est une fonction au moins  $p + 1$  fois dérivable avec  $p \geq 1$ . Supposons que les hypothèses du théorème 7.8 soient vérifiées de sorte que  $g$  admette un unique point fixe  $x^* \in [a, b]$ . Si  $g'(x^*) = g''(x^*) = \dots = g^{(p)}(x^*) = 0$  et  $g^{(p+1)}(x^*) \neq 0$ , alors la convergence de la méthode  $x_{n+1} = g(x_n)$  est superlinéaire d'ordre  $p + 1$ .

*Démonstration.* Utiliser la formule de Taylor. □

### 7.3 Méthode de Newton

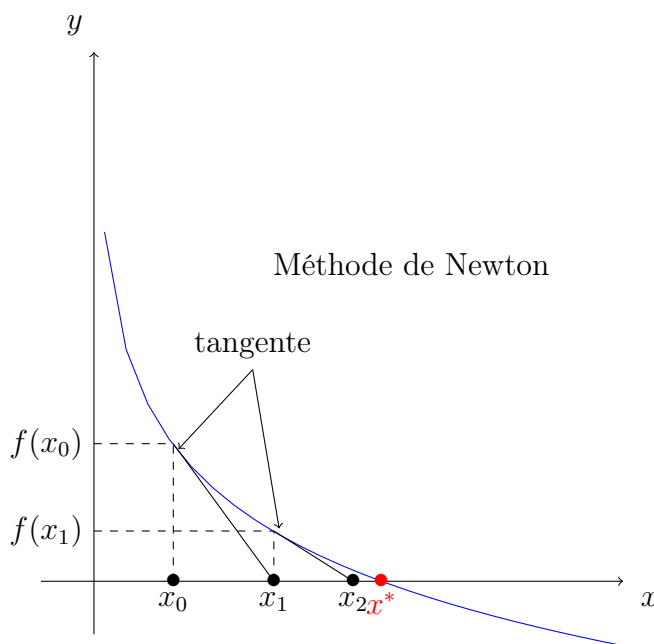
Si on regarde l'équation (7.1), on voit que la méthode convergera d'autant plus vite que la constante  $\gamma$  est petite. Ceci motive donc le fait de remplacer la constante  $\lambda$  par  $\frac{1}{f'(x)}$  et de poser  $g(x) = x - \frac{f(x)}{f'(x)}$ .

**Définition 7.12.** La fonction d'itération de Newton associée à l'équation  $f(x) = 0$  sur  $[a, b]$  est

$$\mathcal{N} : \begin{cases} [a, b] & \rightarrow \mathbb{R}, \\ x & \mapsto \mathcal{N}(x) = x - \frac{f(x)}{f'(x)}. \end{cases}$$

Cette fonction est définie pour  $f$  dérivable sur  $[a, b]$  et telle que  $f'$  ne s'annule pas sur  $[a, b]$ .

Géométriquement, la suite des itérés de Newton se construit comme suit : étant donné  $x_n \in [a, b]$ , on cherche à construire  $x_{n+1}$  tel que  $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$ . En reprenant, le raisonnement développé dans la sous-section précédente, on voit que  $x_{n+1}$  s'obtient comme point d'intersection de la droite tangente au graphe de  $f$  en  $(x_n, f(x_n))$  avec l'axe des abscisses. On illustre ceci sur le graphique suivant :





**Théorème 7.13** (Théorème de convergence locale). *Soit  $f$  une fonction de classe  $\mathcal{C}^2$  sur un intervalle  $[a, b]$  de  $\mathbb{R}$ . On suppose qu'il existe  $\tilde{x} \in [a, b]$  tel que  $f(\tilde{x}) = 0$  et  $f'(\tilde{x}) \neq 0$  ( $\tilde{x}$  est un zéro simple de  $f$ ). Alors il existe  $\epsilon > 0$ , tel que pour tout  $x_0 \in [\tilde{x} - \epsilon, \tilde{x} + \epsilon]$ , la suite des itérés de Newton donnée par  $x_{n+1} = \mathcal{N}(x_n)$  pour  $n \geq 1$  est bien définie, reste dans l'intervalle  $[\tilde{x} - \epsilon, \tilde{x} + \epsilon]$  et converge vers  $\tilde{x}$  quand  $n$  tend vers l'infini. De plus, cette convergence est (au moins) quadratique.*

*Démonstration.* Admis pour ce cours. □

**Exemple :** Prenons l'exemple du calcul de la racine carrée d'un nombre réel  $a > 0$ . On cherche alors à résoudre l'équation  $f(x) = 0$  avec  $f(x) = x^2 - a$ . On a alors

$$\mathcal{N}(x) = x - \frac{f(x)}{f'(x)} = x - \frac{x^2 - a}{2x} = \frac{1}{2} \left( x + \frac{a}{x} \right).$$

Dans le cas  $a = 2$ , si l'on prend comme point de départ  $x_0 = 1$  on obtient les valeurs suivantes des itérés de Newton :

$$\begin{aligned} x_0 &= 1,0000000000000000 \\ x_1 &= 1,5000000000000000 \\ x_2 &= 1,4166666666666667 \\ x_3 &= 1,414215686274510 \\ x_4 &= 1.414213562374690 \\ x_5 &= 1,414213562373095 \end{aligned}$$

La valeur de  $\sqrt{2}$  donnée par MATLAB en format `long` est 1,414213562373095. Le nombre de décimales justes double approximativement à chaque itération ce qui est bien cohérent avec le fait que la convergence de la méthode de Newton est quadratique. Notons que si on utilise la méthode de dichotomie sur l'intervalle  $[1, 2]$ , alors nous avons besoin de 51 itérations pour avoir une valeur approchée de  $\sqrt{2}$  à  $10^{-15}$  près (voir Théorème 7.3).

Le théorème 7.13 suppose que  $\tilde{x}$  est un zéro simple de  $f$ . Nous avons la généralisation suivante au cas d'un zéro de multiplicité quelconque.

**Théorème 7.14.** *Avec les notations, précédentes, si  $\tilde{x}$  est un zéro de multiplicité  $m$  de  $f$ , i.e.,  $f(x^*) = f'(x^*) = \dots = f^{(m-1)}(x^*) = 0$  et  $f^{(m)}(x^*) \neq 0$ , alors la méthode itérative définie par  $x_{n+1} = \mathcal{N}_m(x_n)$  avec  $\mathcal{N}_m(x_n) = x - m \frac{f(x)}{f'(x)}$  est d'ordre supérieure ou égal à 2.*

Finalement, nous avons le résultat « global » suivant :

**Théorème 7.15** (Théorème de convergence globale). *Soit  $f$  une fonction de classe  $\mathcal{C}^2$  sur un intervalle  $[a, b]$  de  $\mathbb{R}$  vérifiant :*

- $f(a)f(b) < 0$ ,
- $\forall x \in [a, b], f'(x) \neq 0$  ( $f$  strictement monotone),

- $\forall x \in [a, b], f''(x) \neq 0$  (concavité dans le même sens sur  $[a, b]$ ).

Alors, en choisissant  $x_0 \in [a, b]$  tel que  $f(x_0) f''(x_0) > 0$ , la suite  $(x_n)_{n \in \mathbb{N}}$  définie par  $x_0$  et  $x_{n+1} = \mathcal{N}(x_n)$  converge vers l'unique solution de  $f(x) = 0$  dans  $[a, b]$ .

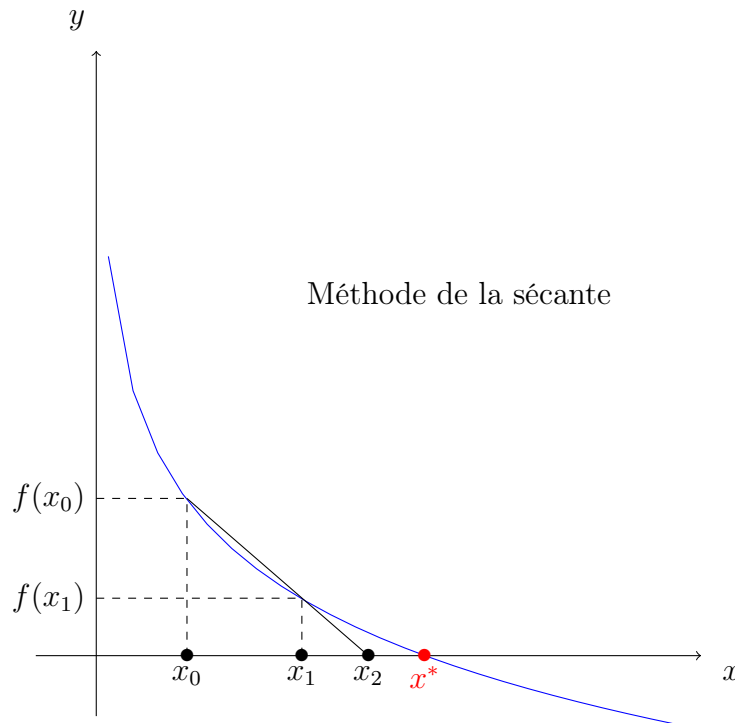
*Démonstration.* Admis pour ce cours. □

## 7.4 Méthode de la sécante

La méthode de Newton étudiée dans la sous-section précédente présente le désavantage de nécessiter le calcul de la dérivée de la fonction  $f$  qui peut s'avérer difficile. Gardons en tête qu'on ne connaît pas nécessairement une expression symbolique de  $f$ . Partant de ce constat, l'idée de la *méthode de la sécante* est de remplacer la dérivée  $f'$  de  $f$  qui apparaît dans la méthode de Newton par une différence divisée (voir la définition 5.9). La méthode de la sécante est alors donnée par l'itération suivante :

$$x_{n+1} = x_n - \frac{f(x_n)}{f[x_n, x_{n-1}]}, \quad \text{où} \quad f[x_n, x_{n-1}] = \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}. \quad (7.2)$$

Cette méthode doit alors être initialisée par deux points  $x_0$  et  $x_1$  à partir desquels on peut utiliser la formule précédente pour calculer  $x_2$  et les itérés suivants.



**Théorème 7.16.** Soit  $f$  une fonction de classe  $\mathcal{C}^2$  sur un intervalle  $[a, b]$  de  $\mathbb{R}$ . On suppose qu'il existe  $\tilde{x} \in [a, b]$  tel que  $f(\tilde{x}) = 0$  et  $f'(\tilde{x}) \neq 0$  ( $\tilde{x}$  est un zéro simple de  $f$ ). Alors il existe

$\epsilon > 0$ , tel que pour tout  $x_0, x_1 \in [\tilde{x} - \epsilon, \tilde{x} + \epsilon]$ , la suite des itérés de la méthode de la sécante donnée par (7.2) pour  $n \geq 1$  est bien définie, reste dans l'intervalle  $[\tilde{x} - \epsilon, \tilde{x} + \epsilon]$  et converge vers  $\tilde{x}$  quand  $n$  tend vers l'infini. De plus, cette convergence est d'ordre  $p = \frac{1+\sqrt{5}}{2} \approx 1,618$  (nombre d'or).

*Démonstration.* Admis pour ce cours. □

**Exemple :** Reprenons l'exemple du calcul de la racine carrée d'un réel  $a > 0$  en résolvant  $f(x) = 0$  pour  $f(x) = x^2 - a$ . La suite des itérés de la méthode de la sécante est donnée par :

$$x_{n+1} = x_n - \frac{x_n^2 - a}{\frac{(x_n^2 - a) - (x_{n-1}^2 - a)}{x_n - x_{n-1}}} = x_n - \frac{x_n^2 - a}{x_n + x_{n-1}}.$$

Dans le cas  $a = 2$ , si l'on prend comme points de départ  $x_0 = x_1 = 1$  on obtient les valeurs suivantes des itérés de la méthode de la sécante :

$$\begin{aligned} x_0 &= 1,0000000000000000 \\ x_1 &= 1,0000000000000000 \\ x_2 &= 1,5000000000000000 \\ x_3 &= 1,4000000000000000 \\ x_4 &= 1,413793103448276 \\ x_5 &= 1,414215686274510 \\ x_6 &= 1,414213562057320 \\ x_7 &= 1,414213562373095 \end{aligned}$$

On voit donc que l'on a besoin de plus d'itérations que pour la méthode de Newton pour avoir la même précision. Ceci vient du fait que l'ordre de cette méthode est plus petit que celui de la méthode de Newton. Par contre, on n'a pas besoin de calculer la dérivée.

## 7.5 Systèmes d'équations non linéaires

On considère maintenant un système d'équations non linéaires donné par une fonction  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $x = (x_1 \dots x_n)^T \mapsto f(x) = (f_1(x_1, \dots, x_n), \dots, f_n(x_1, \dots, x_n))^T$ . On cherche donc un vecteur  $x = (x_1 \dots x_n)^T \in \mathbb{R}^n$  tel que

$$f(x) = 0_{\mathbb{R}^n} \iff \begin{cases} f_1(x_1, \dots, x_n) = 0, \\ \vdots \\ f_n(x_1, \dots, x_n) = 0. \end{cases}$$

La méthode décrite à la fin de la section 7.2 se généralise immédiatement à ce cadre en définissant l'itération

$$x^{(n+1)} = x^{(n)} + M^{-1} f(x^{(n)}), \tag{7.3}$$

où  $M$  est une certaine matrice, et nous avons les mêmes résultats de convergence que dans le cas d'une seule équation.

**Définition 7.17.** La matrice jacobienne d'une fonction  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  notée  $J_f$  est définie (lorsqu'elle existe) par :

$$\forall x = (x_1 \dots x_n)^T \in \mathbb{R}^n, \quad J_f(x) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1}(x) & \frac{\partial f_1}{\partial x_2}(x) & \dots & \frac{\partial f_1}{\partial x_n}(x) \\ \frac{\partial f_2}{\partial x_1}(x) & \frac{\partial f_2}{\partial x_2}(x) & \dots & \frac{\partial f_2}{\partial x_n}(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1}(x) & \frac{\partial f_n}{\partial x_2}(x) & \dots & \frac{\partial f_n}{\partial x_n}(x) \end{pmatrix}.$$

La méthode de Newton se généralise naturellement au cas des systèmes d'équations non linéaires de la manière suivante : on choisit  $x^{(0)} \in \mathbb{R}^n$  et on utilise la formule d'itération

$$x^{(n+1)} = x^{(n)} - J_f(x^{(n)})^{-1} f(x^{(n)}), \quad (7.4)$$

où  $J_f(x^{(n)})^{-1}$  désigne l'inverse de la matrice jacobienne de  $f$  évaluée en  $x^{(n)}$ .

**Théorème 7.18.** Soit  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  une fonction de classe  $\mathcal{C}^2$  sur une boule fermée  $B$  de  $\mathbb{R}^n$ . On suppose qu'il existe un zéro  $\tilde{x}$  de  $f$  dans  $B$  et que  $J_f(\tilde{x})$  est inversible. Alors il existe  $\epsilon > 0$  tel que pour tout  $x^{(0)} \in B$  tel que  $\|x^{(0)} - \tilde{x}\| \leq \epsilon$ , la suite des itérés de la méthode de Newton définie par (7.4) est bien définie et converge vers  $\tilde{x}$  quand  $n$  tend vers l'infini.

Calculer l'itéré  $n+1$  à partir de l'itéré  $n$  en utilisant la formule (7.4) nécessite d'inverser la matrice  $J_f(x^{(n)})$ . Or, calculer l'inverse d'une matrice peut s'avérer coûteux. Par conséquent, nous ré-écrivons la formule d'itération (7.4) sous la forme

$$J_f(x^{(n)}) (x^{(n+1)} - x^{(n)}) = -f(x^{(n)}), \quad (7.5)$$

de sorte qu'à chaque itération, le calcul de l'inverse d'une matrice est remplacé par la résolution d'un système d'équations linéaires ce qui est asymptotiquement moins coûteux en nombre d'opérations à virgule flottante (voir la sous-section 2.2.4 du chapitre 2).

**Exemple :** Considérons le système d'équations non linéaires :

$$(S) : \begin{cases} x_1^2 + 2x_1 - x_2^2 - 2 = 0, \\ x_1^3 + 3x_1x_2^2 - x_2^3 - 3 = 0. \end{cases}$$

Avec les notations précédentes, cela correspond à  $n = 2$ ,  $f_1(x_1, x_2) = x_1^2 + 2x_1 - x_2^2 - 2$ , et  $f_2(x_1, x_2) = x_1^3 + 3x_1x_2^2 - x_2^3 - 3$ . La matrice jacobienne de  $f$  est alors :

$$J_f(x_1, x_2) = \begin{pmatrix} 2x_1 + 2 & -2x_2 \\ 3(x_1^2 + x_2^2) & 6x_1x_2 - 3x_2^2 \end{pmatrix}.$$

Partant du point  $x^{(0)} = (1 \quad -1)^T$ , calculons le premier itéré de la méthode de Newton pour résoudre le système (S). Pour  $n = 1$ , la formule d'itération (7.5) s'écrit :

$$J_f(x^{(0)}) (x^{(1)} - x^{(0)}) = -f(x^{(0)}),$$

c'est-à-dire

$$\begin{pmatrix} 4 & 2 \\ 6 & -9 \end{pmatrix} \begin{pmatrix} x_1^{(1)} - 1 \\ x_2^{(1)} + 1 \end{pmatrix} = - \begin{pmatrix} 0 \\ 2 \end{pmatrix}.$$

En résolvant ce système linéaire, on trouve  $x_1^{(1)} - 1 = -\frac{1}{12}$  et  $x_2^{(1)} + 1 = \frac{1}{6}$  de sorte que  $x^{(1)} = \left(\frac{11}{12} \quad -\frac{5}{6}\right)^T$ .

La méthode de la sécante ne se généralise pas facilement au cas de plusieurs équations. En pratique, pour résoudre un système de plusieurs équations non linéaires, soit on utilise la méthode de Newton, soit on utilise une méthode type (7.3) mais en ajustant la matrice  $M$  au bout d'un certain nombre d'itérations. En général, on prend pour  $M$  une matrice assez proche de la jacobienne  $J_f$  de  $f$  pour que la convergence soit d'un ordre supérieur à 1. On obtient ainsi des *méthodes de Newton généralisées* qui sont surtout utilisées dans le cadre de l'optimisation.

## Remerciements

Ce document a été rédigé avec l'aide de Samir Adly, Paola Boito et Marc Rybowicz, enseignants-chercheurs à la faculté des sciences et techniques de l'Université de Limoges. Je remercie en particulier Paola Boito pour la rédaction de la section [4.4](#).

# Bibliographie

- [1] Jean-Claude Bajard et Jean-Michel Muller. Calcul et arithmétique des ordinateurs. Informatique et Systèmes d'Information, Hermes Science publications, Lavoisier, 2004.
- [2] Jean-Pierre Demailly. Analyse numérique et équations différentielles. Collections Grenoble Sciences, 1991.
- [3] Nicolas J. Higham. Accuracy and Stability of Numerical Algorithms. Second Edition, SIAM, 2002.
- [4] Jean-Michel Muller. Arithmétique des ordinateurs. Masson Paris, 1989.
- [5] Michelle Schatzman. Analyse numérique, cours et exercices pour la licence. InterEditions, Paris, 1991.
- [6] Alain Yger et Jacques-Arthur Weil. Mathématiques L3 - Mathématiques appliquées (Cours complet avec 500 tests et exercices corrigés, 890p et Dvd). Pearson, 2009.