

# Analyse Numérique

*Institut de Mathématiques de Bourgogne (IMB)  
Université de Bourgogne Franche-Comté*

## Organisation

- CM: 24 heures, 12 séances. TD: 18 heures, 9 séances. TP: 8 heures, 4 séances.
- Page web :  
<http://jaramillo.perso.math.cnrs.fr/Courses/AnalyseNumerique/AnalyseNumerique.html>

## Intervenants

- J.L. Jaramillo (resp. Bureau : A323), CM ([Jose-Luis.Jaramillo@u-bourgogne.fr](mailto:Jose-Luis.Jaramillo@u-bourgogne.fr)).
- X. Dupuis (Bureau: A329), TD et TP ([Xavier.Dupuis@u-bourgogne.fr](mailto:Xavier.Dupuis@u-bourgogne.fr)).



# Table des matières

<b>I</b>	<b>Problèmes linéaires</b>	<b>5</b>
<b>1</b>	<b>Introduction: problèmes linéaires</b>	<b>7</b>
1.1	Le problème et les méthodes . . . . .	7
1.2	Diagonalisation de matrices réelles symétriques . . . . .	7
1.3	Exemples de matrices symétriques . . . . .	10
<b>2</b>	<b>Méthodes directes</b>	<b>15</b>
2.1	Décomposition $LU$ (méthode de Gauss/ d'échelonnement) . . . . .	15
2.1.1	Algorithme LU, sans permutation . . . . .	18
2.1.2	Algorithme LU, avec permutation . . . . .	20
2.2	Décomposition de Choleski . . . . .	27
<b>3</b>	<b>Normes et conditionnement d'une matrice</b>	<b>35</b>
3.1	Norme matricielle, norme matricielle induite et rayon spectral . . . . .	35
3.1.1	Rappel sur les normes . . . . .	35
3.1.2	Normes matricielles . . . . .	36
3.2	Relation entre norme matricielle et rayon spectral . . . . .	39
3.2.1	Convergence et rayon spectral . . . . .	39
3.3	Conditionnement d'une matrice . . . . .	41
3.3.1	Conditionnement et ses propriétés élémentaires . . . . .	41
3.3.2	Majoration des erreurs . . . . .	42
<b>4</b>	<b>Méthodes itératives</b>	<b>45</b>
<b>5</b>	<b>Méthodes de descente</b>	<b>47</b>
<b>II</b>	<b>Interpolation</b>	<b>49</b>
<b>6</b>	<b>Interpolation polynomiale</b>	<b>51</b>
<b>III</b>	<b>Annexes</b>	<b>53</b>
<b>A</b>	<b>Espaces linéaires</b>	<b>55</b>
A.1	Rappel sur les espaces linéaires . . . . .	55
A.2	Application linéaire . . . . .	56
A.3	Produit scalaire . . . . .	56

<b>B</b>	<b>Matrices</b>	<b>59</b>
B.1	Définitions . . . . .	59
B.1.1	Opérations sur les matrices . . . . .	60
B.1.2	Rang et noyau d'une matrice . . . . .	61
B.1.3	Inverse d'une matrice . . . . .	61
B.2	Applications linéaires et matrice d'une application linéaire . . . . .	62
B.2.1	Matrices semblables . . . . .	62
B.3	Matrices et géométrie . . . . .	63
<b>C</b>	<b>Factorisation de matrices</b>	<b>65</b>
<b>D</b>	<b>Diagonalisation</b>	<b>67</b>
D.1	Notion de diagonalisation . . . . .	67
D.2	Diagonalisation de matrices normales . . . . .	68

## Part I

# Problèmes linéaires



# Chapitre 1

## Introduction: problèmes linéaires

### 1.1 Le problème et les méthodes

Le problème fondamental que nous allons aborder est la résolution de l'équation

$$A\mathbf{x} = \mathbf{b} , \tag{1.1}$$

avec  $A$  une matrice et  $\mathbf{b}$  un vecteur donnés dans des espaces linéaires de dimension fini. Un problème lié qui va aussi nous occuper est la résolution du problème spectral de la matrice  $A$

$$A\mathbf{x} = \lambda\mathbf{x} . \tag{1.2}$$

Les appendices **A-C** offrent un rappel des éléments de base sur les espaces linéaires et sujets liés.

La résolution “efficace” de l'équation (1.1) pour des dimensions élevées pose un problème difficile dont la résolution demande le concours de différentes techniques mathématiques. L'esprit ici est de souligner l'enchevêtrement de l'Algèbre Linéaire, la Géométrie et l'Analyse Mathématique dans l'Analyse Numérique. Le problème abordé dans la section suivante illustre ce point.

### 1.2 Diagonalisation de matrices réelles symétriques

Les matrices carrées réelles et symétriques sont diagonalisables, avec des valeurs et vecteurs propres réels. D'un côté, ceci est un résultat fondamental dans le contexte du problème abordé dans ce cours. D'un autre côté, la preuve de ce résultat nous permet d'illustrer concrètement la puissance de la combinaison d'Analyse Mathématique et de l'Algèbre Linéaire, un leitmotiv dans ce cours.

**Lemme 1.1.** (Une matrice symétrique réelle est diagonalisable sur  $\mathbb{R}$ ). *Soit  $V$  un espace vectoriel sur  $\mathbb{R}$  de dimension finie:  $\dim(V) = n, n \in \mathbb{N}^*$ , muni d'un produit scalaire  $\langle \cdot, \cdot \rangle; V \times V \rightarrow \mathbb{R}$ , avec norme associée  $\|\mathbf{x}\| = \langle \mathbf{x}, \mathbf{x} \rangle^{\frac{1}{2}}, \forall \mathbf{x} \in V$ . Soit  $T$  une application linéaire  $T : V \rightarrow V$ , symétrique*

$$\langle T(\mathbf{x}), \mathbf{y} \rangle = \langle T(\mathbf{y}), \mathbf{x} \rangle, \tag{1.3}$$

*Alors, il existe une base orthonormée  $\mathbf{e}_1, \dots, \mathbf{e}_n$  de  $V$  (c.à.d  $\langle \mathbf{e}_i, \mathbf{e}_j \rangle = \delta_{ij}$ ) et  $(\lambda_1, \dots, \lambda_n) \in \mathbb{R}^n$  tels que  $T\mathbf{e}_i = \lambda_i\mathbf{e}_i, \forall i \in \{1, \dots, n\}$ .*

**Preuve :** Par récurrence sur  $n = \dim(V)$ .

i) Nous supposons d'abord  $\dim(V) = 1$ .

Soit  $\mathbf{v} \in V$ ,  $\mathbf{v} \neq 0$ , alors  $V = \text{Lin}_{\mathbb{R}}(\mathbf{v}) = \text{Lin}_{\mathbb{R}}(\mathbf{e}_1)$  avec  $\mathbf{e}_1 = \frac{\mathbf{v}}{\|\mathbf{v}\|}$ . En particulier  $\langle \mathbf{e}_1, \mathbf{e}_1 \rangle = 1$ . Soit  $T : V \rightarrow V$  linéaire (symétrique), on a  $T\mathbf{e}_1 \in V = \text{Lin}_{\mathbb{R}}(\mathbf{e}_1) \implies \exists \lambda_1 \in \mathbb{R}$  tel que  $T\mathbf{e}_1 = \lambda_1 \mathbf{e}_1$ .

Notez que le raisonnement peut être étendu au cas complexe<sup>1</sup>.

ii) On suppose que le résultat est vrai pour  $\dim(V) < n$  (hypothèse de récurrence). On le montre pour  $\dim(V) = n$ .

Soit  $V$  un espace linéaire avec produit scalaire, avec  $\dim(V) = n < \infty$  et  $T : V \rightarrow V$  linéaire symétrique. On définit :

$$\begin{aligned} \varphi : V &\rightarrow \mathbb{R} \\ x &\mapsto \langle T\mathbf{x}, \mathbf{x} \rangle . \end{aligned}$$

Cette application est continue. En particulier, sa restriction à la sphère unité  $S_1 = \{\mathbf{x} \in V ; \|\mathbf{x}\| = 1\}$ , qui est compacte (fermée et bornée, parce  $\dim(V) = n < \infty$ ), est aussi continue. L'image de  $\varphi$  est alors compacte dans  $\mathbb{R}$  (intervalle fermé) alors  $\varphi$  atteint son maximum. C'est-à-dire  $\exists \mathbf{v} \in S_1$  tel que

$$\varphi(\mathbf{x}) \leq \varphi(\mathbf{v}) = \langle T\mathbf{v}, \mathbf{v} \rangle := \lambda, \forall \mathbf{x} \in S_1. \quad (1.4)$$

On va montrer que  $T\mathbf{v} = \lambda\mathbf{v}$ .

Soit  $\mathbf{y} \in V \setminus \{0\}$  et soit  $t \in ]0, \frac{1}{\|\mathbf{y}\|}[$ , alors  $\mathbf{v} + t\mathbf{y} \neq 0$ <sup>2</sup> On en déduit :

$$\frac{\mathbf{v} + t\mathbf{y}}{\|\mathbf{v} + t\mathbf{y}\|} \in S_1 \quad (1.5)$$

et donc :

$$\underbrace{\varphi(\mathbf{v})}_{\lambda} \geq \langle T \frac{\mathbf{v} + t\mathbf{y}}{\|\mathbf{v} + t\mathbf{y}\|}, \frac{\mathbf{v} + t\mathbf{y}}{\|\mathbf{v} + t\mathbf{y}\|} \rangle \Leftrightarrow \lambda \|\mathbf{v} + t\mathbf{y}\|^2 \geq \langle T(\mathbf{v} + t\mathbf{y}), \mathbf{v} + t\mathbf{y} \rangle . \quad (1.6)$$

En développant à gauche :

$$\lambda \|\mathbf{v} + t\mathbf{y}\|^2 = \lambda \langle (\mathbf{v} + t\mathbf{y}), \mathbf{v} + t\mathbf{y} \rangle = \lambda \left( \underbrace{\|\mathbf{v}\|^2}_{=1} + 2t\langle \mathbf{v}, \mathbf{y} \rangle + t^2 \|\mathbf{y}\|^2 \right) , \quad (1.7)$$

et à droite :

$$\langle T(\mathbf{v} + t\mathbf{y}), \mathbf{v} + t\mathbf{y} \rangle = \underbrace{\langle T\mathbf{v}, \mathbf{v} \rangle}_{\lambda} + t \underbrace{\langle T\mathbf{v}, \mathbf{y} \rangle}_{\langle \mathbf{y}, T\mathbf{v} \rangle = \langle T\mathbf{v}, \mathbf{y} \rangle} + t^2 \langle T\mathbf{y}, \mathbf{y} \rangle , \quad (1.8)$$

et on conclut (en utilisant  $t > 0$  pour pouvoir simplifier par  $t$ ) :

$$\lambda (2\langle \mathbf{v}, \mathbf{y} \rangle + t\|\mathbf{y}\|^2) \geq 2\langle T\mathbf{v}, \mathbf{y} \rangle + t\langle T\mathbf{y}, \mathbf{y} \rangle . \quad (1.9)$$

<sup>1</sup>Il faut utiliser une application linéaire  $T$  symétrique (hermitien) par rapport au produit scalaire hermitien (voir Appendix A.3).

<sup>2</sup>Avec  $\|\|\mathbf{x}\| - \|\mathbf{y}\|\| \leq \|\mathbf{x} - \mathbf{y}\|$ , on a  $\|\mathbf{v} + t\mathbf{y}\| \geq \|\mathbf{v}\| - \|\mathbf{t}\mathbf{y}\| = |1 - t\|\mathbf{y}\|| \neq 0$ .



Si on prend la limite lorsque  $t \rightarrow 0^+$ , on a :

$$2\lambda\langle \mathbf{v}, \mathbf{y} \rangle \geq 2\langle T\mathbf{v}, \mathbf{y} \rangle \Rightarrow 0 \geq \langle T\mathbf{v} - \lambda\mathbf{v}, \mathbf{y} \rangle, \forall \mathbf{y} \in V \setminus \{0\}. \quad (1.10)$$

Si on prend maintenant  $\mathbf{z} = -\mathbf{y}$ , on dérive aussi :

$$\begin{aligned} 0 \geq \langle T\mathbf{v} - \lambda\mathbf{v}, \mathbf{z} \rangle &= \langle T\mathbf{v} - \lambda\mathbf{v}, -\mathbf{y} \rangle = -\langle T\mathbf{v} - \lambda\mathbf{v}, \mathbf{y} \rangle \\ &\Rightarrow \langle T\mathbf{v} - \lambda\mathbf{v}, \mathbf{y} \rangle \geq 0, \forall \mathbf{y} \in V \setminus \{0\} \end{aligned} \quad (1.11)$$

On a alors :

$$\langle T\mathbf{v} - \lambda\mathbf{v}, \mathbf{y} \rangle = 0, \forall \mathbf{y} \in V \setminus \{0\} \Rightarrow T\mathbf{v} - \lambda\mathbf{v} = 0 \Leftrightarrow T\mathbf{v} = \lambda\mathbf{v} \quad (1.12)$$

On pose  $\mathbf{e}_n = \mathbf{v}$  et  $\lambda_n = \lambda$ .

Maintenant, on utilise l'hypothèse de récurrence. Soit  $W = \{\mathbf{x} \in V; \langle \mathbf{x}, \mathbf{e}_n \rangle = 0\}$ . Alors  $V \neq W$  et  $V = W \oplus \text{Lin}_{\mathbb{R}}(\mathbf{e}_n)$ . On peut décomposer  $\mathbf{x} \in V$

$$\mathbf{x} = \underbrace{\mathbf{x} - \langle \mathbf{x}, \mathbf{e}_n \rangle \mathbf{e}_n}_{\in W} + \langle \mathbf{x}, \mathbf{e}_n \rangle \mathbf{e}_n. \quad (1.13)$$

Si on note  $T|_W$  la restriction de  $T$  à  $W$ , l'application  $T|_W : W \rightarrow W$  est linéaire et symétrique (exercice). On peut alors faire usage de l'hypothèse de récurrence sur  $T|_W$  :

$$\exists \{\mathbf{e}_1, \dots, \mathbf{e}_{n-1}\} \text{ et } (\lambda_1, \dots, \lambda_{n-1}) \in \mathbb{R}^{n-1}, \quad (1.14)$$

tels que :

$$T|_W \mathbf{e}_i = T\mathbf{e}_i = \lambda_i \mathbf{e}_i, \forall i \in \{1, \dots, n-1\}, \langle \mathbf{e}_i, \mathbf{e}_j \rangle = 0. \quad (1.15)$$

On choisit finalement :

$$\{\mathbf{e}_1, \dots, \mathbf{e}_{n-1}, \mathbf{e}_n\} \text{ et } (\lambda_1, \dots, \lambda_{n-1}, \lambda_n), \quad (1.16)$$

qui montrent le théorème.

□

**Corollaire 1.1.** Une matrice  $A \in M_n(\mathbb{R})$  symétrique est diagonalisable sur  $\mathbb{R}$ .

**Preuve :**

On considère  $V = \mathbb{R}^n$  avec le produit scalaire canonique

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^t \cdot \mathbf{y} = \sum_{i=1}^n x^i y^i,$$

où  $\mathbf{x} = x^i \mathbf{e}_i$  avec  $B = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  la base canonique. Si  $A \in M_n(\mathbb{R})$  est une matrice symétrique, l'application

$$\begin{aligned} T : \mathbb{R}^n &\rightarrow \mathbb{R}^n \\ \mathbf{x} &\mapsto A \cdot \mathbf{x}, \text{ où } \mathbf{x} = x^i \mathbf{e}_i, B = \{\mathbf{e}_{i \in \mathbb{N}}\} \text{ base canonique} \end{aligned} \quad (1.17)$$

est linéaire et symétrique. Notez que  $A$  est la matrice de  $T$  dans la base canonique

$$A = M(T, B, B) \quad (1.18)$$

Pour la symétrie,

$$\begin{aligned} \langle \mathbf{x}, T\mathbf{y} \rangle &= \mathbf{x}^t \cdot A \cdot \mathbf{y} = (\mathbf{x}^t \cdot A \cdot \mathbf{y})^t = (A \cdot \mathbf{y})^t \cdot (\mathbf{x}^t)^t = \mathbf{y}^t \cdot A \cdot \mathbf{x} = (\mathbf{y}^t \cdot A \cdot \mathbf{x})^t \\ &= (A \cdot \mathbf{x})^t \cdot \mathbf{y} = \langle T\mathbf{x}, \mathbf{y} \rangle . \end{aligned} \quad (1.19)$$

Par le lemme précédent,  $\exists B' = \{\mathbf{e}'_1, \dots, \mathbf{e}'_n\}$  et  $\{\lambda_1, \dots, \lambda_n\}$  tels que

$$A\mathbf{e}'_i = T\mathbf{e}'_i = \lambda_i \mathbf{e}'_i, \text{ et } \langle \mathbf{e}'_i, \mathbf{e}'_j \rangle = \delta_{ij} \quad (1.20)$$

Par la définition [D.2](#) on conclut que  $A$  est diagonalisable □

*Remarque 1.1.* La matrice  $A$  est semblable à une matrice diagonale. En effet, on a

$$A = M(T, B) = M(T, B, B) \text{ et } D = M(T, B') = M(T, B', B'), \quad (1.21)$$

et on peut écrire

$$\begin{aligned} A &= M(T, B, B) = M(I, B', B)M(T, B', B')M(I, B, B') \\ &= M(B' \rightarrow B)M(T, B', B')M(B \rightarrow B') = PDP^{-1} \end{aligned} \quad (1.22)$$

où  $P$  est la matrice de changement de base  $P = M(B' \rightarrow B)$ . On a

$$A = PDP^{-1}, \quad D = P^{-1}AP \quad (1.23)$$

Étant donné que  $P$  est la matrice de changement de base entre des bases orthonormées, la matrice  $P$  est orthogonale, c.à.d  $P^{-1} = P^t$ .

### 1.3 Exemples de matrices symétriques

**Exemple 1.1** (Opérateur Laplacien: discrétisation). Le premier exemple de matrice symétrique est placé dans le contexte des équations aux dérivées partielles (EDP), en particulier dans une classe de problèmes qui amènent à une équation

$$A\mathbf{x} = \mathbf{b}, \quad (1.24)$$

avec  $A$  symétrique et le vecteur  $\mathbf{b}$  donnés.

#### a) EDP et Laplacien:

L'équation de Poisson est

$$-\Delta u = f, \text{ où } \Delta = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2}. \quad (1.25)$$

L'opérateur  $\Delta$  est le Laplacien et cette équation est un exemple d'EDP de type dit elliptique. Elle joue un rôle fondamental dans plusieurs contextes mathématiques et physiques. Voici quelques exemples :

- Détermination du potentiel gravitationnel de Newton

$$-\Delta\phi_{\text{grav}} = -4\pi G\rho_{\text{masse}} . \quad (1.26)$$

- Potentiel électrostatique

$$-\Delta\phi_{\text{lec}} = \frac{\rho_{\text{charge}}}{\epsilon} . \quad (1.27)$$

- états stationnaires de l'équation de la chaleur, avec une densité de flux de chaleur  $q$

$$\rho c_p \frac{\partial T}{\partial t} = k\Delta T + q \quad (1.28)$$

Celle-ci est une EDP de type *parabolique*.

- états stationnaires de l'équation d'ondes avec forçage  $g$

$$-\frac{\partial^2 u}{\partial t^2} + c^2 \Delta u = g . \quad (1.29)$$

L'équation d'ondes est le prototype d'EDP de type hyperbolique.

**b) Le Laplacien: un opérateur symétrique, (en fait, auto-adjoint <sup>3</sup>).** Nous considérons l'espace linéaire  $L^2(D \subset \mathbb{R}^n)$  de fonctions de carré intégrable, c.à.d  $f \in L^2(D)$  ssi  $\int_D f^2 d^n x < \infty$ .  $L^2(D \subset \mathbb{R}^n)$  est un espace de Hilbert avec produit scalaire

$$\langle f, g \rangle = \int_D fg d^n x . \quad (1.30)$$

*Remarque 1.2.* espace de Hilbert complexe. Nous pouvons aussi considérer des fonctions  $\varphi : D \subset \mathbb{R}^n \rightarrow \mathbb{C}$ , avec  $\int_D |\varphi|^2 d^n x < \infty$ , qui est aussi un espace de Hilbert avec produit hermitien (sesquilinéaire)

$$\langle \varphi, \phi \rangle = \int_D \bar{\varphi} \phi d^n x \quad (1.31)$$

On va étudier une discrétisation de cet opérateur, notamment ce qu'on appelle des différences finies. Pour simplifier, nous allons étudier le cas en dimension 1, avec des conditions aux bords dites "Dirichlet homogènes". L'équation (1.25) devient

$$\begin{cases} -\frac{d^2 u}{dx^2} = f(x) \\ u(0) = u(1) = 0 \end{cases} \quad (1.32)$$

avec  $f \in C([0, 1], \mathbb{R})$ . On peut montrer que  $u \in C^2([0, \mathbb{R}])$ .

Étant donné  $N \in \mathbb{N}$ , on définit le pas de la discrétisation  $h = \frac{1}{N+1}$ , et on introduit les points :

$$x_i = ih = \frac{i}{N+1} , \quad \forall i \in \{0, 1, \dots, N+1\} \quad (1.33)$$

---

<sup>3</sup>La notion d'opérateur auto-adjoint dans un espace linéaire de dimension infinie implique des aspects subtils sur le domaine de l'opérateur. Dans ce cours, nous ne discutons pas de ces questions.

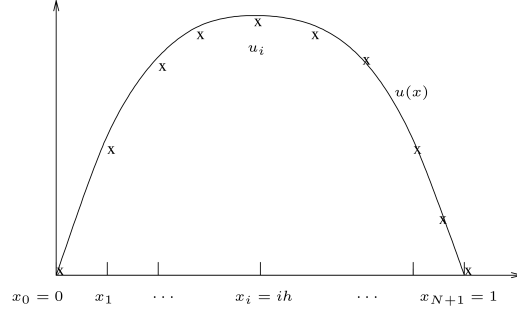


Figure 1.1: Schéma de la discrétisation par différences finies.

On évalue l'équation (1.32) dans les points  $x_i \in \{1, \dots, N\}$ . Notons que  $u(x_0) = u(x_{N+1}) = 0$ . Ainsi :

$$-u''(x_i) = f(x_i) \quad , \quad \forall i \in \{0, 1, \dots, N+1\} \quad (1.34)$$

où  $f(x_i)$  sont les *valeurs exactes* de la solution  $u(x)$  évalué sur  $x_i$ . On va supposer que  $u \in C^4([0, 1], \mathbb{R})$ . On peut alors écrire le développement de Taylor-Lagrange :

$$u(x_{i+1}) = u(x_i + h) = u(x_i) + hu'(x_i) + \frac{h^2}{2}u''(x_i) + \frac{h^3}{6}u^{(3)}(\xi_i) + \frac{h^4}{24}u^{(4)}(\xi_i) \quad , \quad \xi_i \in ]x_i, x_{i+1}[ \quad (1.35)$$

De la même manière :

$$u(x_{i-1}) = u(x_i - h) = u(x_i) - hu'(x_i) + \frac{h^2}{2}u''(x_i) - \frac{h^3}{6}u^{(3)}(\eta_i) + \frac{h^4}{24}u^{(4)}(\eta_i) \quad , \quad \eta_i \in ]x_{i-1}, x_i[ \quad (1.36)$$

Si on ajoute les deux expressions, on obtient :

$$u(x_{i+1}) + u(x_{i-1}) = 2u(x_i) + h^2u''(x_i) + \frac{h^4}{24}(u^{(4)}(\xi_i) + u^{(4)}(\eta_i)) \quad . \quad (1.37)$$

Alors, on peut écrire :

$$-u''(x_i) - \frac{-u(x_{i+1}) + 2u(x_i) - u(x_{i-1}))}{h^2} = \frac{h^2}{24}(u^{(4)}(\xi_i) + u^{(4)}(\eta_i)) \quad (1.38)$$

Si on définit maintenant *l'erreur de consistance*  $R_i$  par :

$$R_i = -u''(x_i) - \frac{-u(x_{i+1}) + 2u(x_i) - u(x_{i-1}))}{h^2} \quad , \quad (1.39)$$

on a

$$R_i = \frac{h^2}{24}(u^{(4)}(\xi_i) + u^{(4)}(\eta_i)) \quad , \quad (1.40)$$

de telle manière que :

$$\begin{aligned} |R_i| &= \frac{h^2}{24} |u^{(4)}(\xi_i) + u^{(4)}(\eta_i)| \leq \frac{h^2}{24} \left( \underbrace{|u^{(4)}(\xi_i)|}_{\leq \sup_{x \in ]0,1[} |u^{(4)}(x)|} + \underbrace{|u^{(4)}(\eta_i)|}_{\leq \sup_{x \in ]0,1[} |u^{(4)}(x)|} \right) \\ &= \frac{h^2}{12} \sup_{x \in ]0,1[} |u^{(4)}(x)| = \frac{h^2}{12} \underbrace{\|u^{(4)}\|_{\infty}}_{< \infty} \quad . \end{aligned} \quad (1.41)$$

Cela nous dit que l'erreur de consistance tend vers 0 comme  $h^2$  (on dit que le schéma est consistant à l'ordre 2).

L'idée maintenant est de se rapprocher du problème original par un problème discret, en faisant

$$\frac{d^2u}{dx^2} \sim \frac{-u(x_{i+1}) + 2u(x_i) - u(x_{i-1}))}{h^2}. \quad (1.42)$$

Néanmoins, cette expression utilise les valeurs exactes  $u(x_i)$ , mais ces valeurs sont précisément les inconnues qu'on cherche à trouver. À ce point, nous introduisons des valeurs  $u_i$  qui vont représenter des *approximations* des vraies valeurs de  $u(x_i)$ . Si on a aussi  $b_i = f(x_i)$  (valeurs exactes du membre à droite!), nous pouvons définir le problème discret

$$\begin{cases} \frac{-u_{i+1} + 2u_i - u_{i-1}}{h^2} = b_i, & i \in \{1, \dots, N\} \\ u_0 = u_{N+1} = 0 \end{cases} \quad (1.43)$$

Si on définit la matrice  $\Delta_N$ , par ses éléments matriciels :

$$(\Delta_N)_{ij} = \begin{cases} \frac{2}{h^2} & , i = j, i, j \in \{1, \dots, N\} \\ -\frac{1}{h^2} & , |i - j| = 1, i, j \in \{1, \dots, N\} \\ 0 & , |i - j| > 1, i, j \in \{1, \dots, N\} \end{cases} \quad (1.44)$$

on peut écrire le problème (1.43) comme

$$\Delta_N \mathbf{u} = \mathbf{b} \quad (1.45)$$

Pour avoir une idée, les premières matrices  $\Delta_N$  sont de la forme :

$$\Delta_2 = \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix}, \quad \Delta_3 = \begin{pmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{pmatrix}, \quad \Delta_4 = \begin{pmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{pmatrix},$$

On verra que cette matrice est inversible (en fait elle symétrique définie positive), ce qui nous permet de trouver la solution approchée comme :

$$\mathbf{u} = (\Delta_N)^{-1} \mathbf{b}. \quad (1.46)$$

Cette équation différentielle nous amène à notre premier problème matriciel de type  $A\mathbf{x} = \mathbf{b}$ , et on peut déjà voir que si  $N$  devient très grand, le problème d'inversion de  $\Delta_N$  devient non-trivial. Ce problème d'inversion des matrices de grande taille, et en particulier des matrices symétriques, sera un des sujets fondamentaux de ce cours.

Un autre aspect important, déjà présent dans ce problème? concerne la convergence de la méthode. Notamment, l'attente est que, lorsque  $h \rightarrow 0$ , la solution approchée  $u_i$  tend vers  $u(x_i)$ . C'est-à-dire que, si on définit l'erreur comme :

$$e_i = u(x_i) - u_i, \quad i \in \{1, \dots, N\}, \quad (1.47)$$

on peut voir en appliquant les résultats précédents que :

$$\Delta_n \mathbf{e} = \mathbf{R}, \quad (1.48)$$

où  $R_i$  est donné par (1.40). Alors,

$$\mathbf{e} = (\Delta_N)^{-1} \mathbf{R} . \tag{1.49}$$

Le fait que  $|R_i|$  tende vers 0 quand  $h \rightarrow 0$  est bon, mais ce n'est pas suffisant pour garantir que  $\mathbf{e} \rightarrow 0$ . Pour ça, il faut montrer que  $\Delta_N$ , qui dépend de  $h$ , est telle que l'on puisse donner un sens à la *norme*  $\|(\Delta_N)^{-1}\|$  et elle qu'elle soit, notamment bornée indépendamment de  $h$ . Ce type de question liant la convergence d'un schéma avec la norme d'une matrice approprié sera aussi l'objet de notre cours.

## Chapitre 2

# Méthodes directes

**Définition 2.1** (méthode directe). On appelle méthode directe de résolution de

$$A\mathbf{x} = \mathbf{b}, \quad A \in M_n(\mathbb{R}), \quad \mathbf{x}, \mathbf{b} \in \mathbb{R}^n \quad (2.1)$$

une méthode qui donne exactement  $\mathbf{x}$  ( $A$  et  $\mathbf{b}$  étant connus), solution de (3.34) après un nombre fini d'opérations élémentaires.

### 2.1 Décomposition $LU$ (méthode de Gauss/ d'échelonnement)

Le principe de la méthode est de transformer un système linéaire, par des combinaisons linéaires et des permutations de lignes, dans un système triangulaire équivalent. On va rappeler la procédure de cette méthode, étudiée dans les cours précédents avec un exemple.

**Exemple 2.1.** On considère le problème (3.34) avec :

$$A = \begin{pmatrix} 4 & -2 & 6 \\ 2 & -4 & 9 \\ -4 & 5 & -9 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 12 \\ -24 \\ 12 \end{pmatrix}. \quad (2.2)$$

On construit la matrice augmentée:

$$\tilde{A} = (A \quad \mathbf{b}) = \left( \begin{array}{ccc|c} 4 & -2 & 6 & 12 \\ 2 & -4 & 9 & -24 \\ -4 & 5 & -9 & 12 \end{array} \right), \quad (2.3)$$

Et on procède avec la méthode de pivot déjà étudiée pour transformer  $A$  en une matrice triangulaire supérieure. Simplement, dans ce process, on va garder systématiquement l'information des lignes qu'on retranche à chaque pas de l'itération. Concrètement, vu que  $\ell_1(1) = 4 \neq 0$ , nous pouvons l'utiliser comme *pivot*. Si on note  $\tilde{A}_0 = \tilde{A}$

$$\tilde{A}_0 \xrightarrow{\ell_2 \rightarrow \ell_2 - \frac{2}{4}\ell_1} \tilde{A}_1 = \left( \begin{array}{ccc|c} 4 & -2 & 6 & 12 \\ 0 & -3 & 6 & -30 \\ -4 & 5 & -9 & 12 \end{array} \right). \quad (2.4)$$

Notons que cette transformation de  $\tilde{A}$  et  $\mathbf{b}$ , qui change la deuxième ligne en lui retranchant la première ligne multipliée par un facteur  $\frac{1}{2}$  est le résultat de multiplier  $\tilde{A}_0$  par la gauche par la

matrice :

$$E_1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} - \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad (2.5)$$

On procède de la même manière pour échelonner la matrice :

$$\tilde{A}_1 \xrightarrow[E_2]{\ell_3 \rightarrow \ell_3 - \frac{-4}{4}\ell_1} \tilde{A}_2 = \left( \begin{array}{ccc|c} 4 & -2 & 6 & 12 \\ 0 & -3 & 6 & -30 \\ 0 & 3 & -3 & 24 \end{array} \right) \xrightarrow[E_3]{\ell_3 \rightarrow \ell_3 - \frac{-3}{-3}\ell_2} \tilde{A}_3 = \left( \begin{array}{ccc|c} 4 & -2 & 6 & 12 \\ 0 & -3 & 6 & -30 \\ 0 & 0 & 3 & -6 \end{array} \right), \quad (2.6)$$

avec

$$E_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix}, \quad E_3 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix} \quad (2.7)$$

Comme résultat de la procédure, autant que tous les pivots qu'on trouve soient différents de zéro, on construit un système linéaire équivalent  $\tilde{A}_3 \mathbf{x} = \tilde{\mathbf{b}}$  échelonné qu'on peut résoudre facilement par remontée. Le résultat est

$$\mathbf{x} = \begin{pmatrix} 9 \\ 6 \\ -2 \end{pmatrix} \quad (2.8)$$

Si on a besoin de résoudre ce système pour différents vecteurs  $\mathbf{b}$ , dans chaque cas on doit répéter le calcul du vecteur  $\tilde{\mathbf{b}}$  (dans le process d'échelonnement de  $A$ ). Si le nombre d'équations à résoudre est élevé, cette répétition est clairement à éviter en cherchant une méthode plus efficace. Si on regarde notre exemple, on note que c'est ce qu'on a fait, on peut l'écrire :

$$A\mathbf{x} = \mathbf{b} \longrightarrow \underbrace{E_3 E_2 E_1 A}_{U} \mathbf{x} = E_3 E_2 E_1 \mathbf{b}, \quad (2.9)$$

où  $U$  est une matrice triangulaire supérieure (c.à.d.  $u_{ij} = 0$ , si  $i > j$ ). Les matrices  $E_i$  sont inversibles et on peut alors multiplier l'expression à gauche

$$A\mathbf{x} = (E_3 E_2 E_1)^{-1} U \mathbf{x} = \mathbf{b} = \underbrace{E_1^{-1} E_2^{-1} E_3^{-1}}_L U \mathbf{x} = LU \mathbf{x} = \mathbf{b} \quad (2.10)$$

avec

$$(E_1)^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad E_2^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -1 & 0 & 1 \end{pmatrix}, \quad E_3^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -1 & 1 \end{pmatrix} \quad (2.11)$$

Ce qui amène à :

$$L = \begin{pmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ -1 & -1 & 1 \end{pmatrix} \quad (2.12)$$

$L$  est une matrice triangulaire inférieure ( $l_{ij} = 0, i < j$ ) avec des 1 dans la diagonale principale et les éléments non nuls sont les coefficients qu'on a utilisé pour retrancher chaque ligne dont



le pivot est non nul dans le process d'échelonnement <sup>1</sup>. En particulier nous avons décomposé la matrice  $A$  comme le produit d'une matrice triangulaire inférieure  $L$  et d'une matrice triangulaire supérieure  $U$

$$A = LU . \quad (2.15)$$

En particulier, dans notre exemple

$$A = \begin{pmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ -1 & -1 & 1 \end{pmatrix} \begin{pmatrix} 4 & -2 & 6 \\ 0 & -3 & 6 \\ 0 & 0 & 3 \end{pmatrix} \quad (2.16)$$

Notez que la décomposition de  $A$  comme produit d'une matrice triangulaire inférieure et d'une matrice triangulaire supérieure n'est pas unique. Par exemple, on peut multiplier la première par 2 et la deuxième par  $\frac{1}{2}$

$$A = \begin{pmatrix} 2 & 0 & 0 \\ 1 & 2 & 0 \\ -2 & -2 & 2 \end{pmatrix} \begin{pmatrix} 2 & -1 & 3 \\ 0 & -\frac{3}{2} & 3 \\ 0 & 0 & \frac{3}{2} \end{pmatrix} , \quad (2.17)$$

et on peut aussi trouver d'autres paires de matrices

$$A = \begin{pmatrix} 2 & 0 & 0 \\ 1 & -3 & 0 \\ -2 & 3 & 1 \end{pmatrix} \begin{pmatrix} 2 & -1 & 3 \\ 0 & 1 & -2 \\ 0 & 0 & 3 \end{pmatrix} . \quad (2.18)$$

Ce qui est particulier de la décomposition (2.16) c'est que les éléments de la diagonale principale de la matrice  $L$  sont seulement des 1. Comme va le voir plus tard, une telle décomposition est unique et on l'appelle décomposition  $LU$  de la matrice  $A$ .

Si on utilise cette décomposition, on peut résoudre le système original

$$\begin{aligned} Ax &= \mathbf{b} \\ LUx &= \mathbf{b} \end{aligned} \quad (2.19)$$

en deux pas, en introduisant un vecteur  $\mathbf{y}$  intermédiaire

$$\begin{aligned} Ly &= \mathbf{b} \\ Ux &= \mathbf{y} \end{aligned} \quad (2.20)$$

---

<sup>1</sup>En général, la structure de matrice sera toujours

$$\begin{aligned} (E_1) &= \begin{pmatrix} 1 & 0 & 0 \\ -a & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} , & E_2 &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -b & 0 & 1 \end{pmatrix} , & E_3 &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -c & 1 \end{pmatrix} \\ (E_1)^{-1} &= \begin{pmatrix} 1 & 0 & 0 \\ a & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} , & E_2^{-1} &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ b & 0 & 1 \end{pmatrix} , & E_3^{-1} &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & c & 1 \end{pmatrix} \end{aligned} \quad (2.13)$$

Ce qui amène à :

$$L = E_1^{-1}E_2^{-1}E_3^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ a & 1 & 0 \\ b & c & 1 \end{pmatrix} \quad (2.14)$$

Étant donné  $\mathbf{b}$ , la résolution pour  $\mathbf{y}$  est directe car  $L$  est échelonnée. En particulier, vu que  $L$  est triangulaire inférieure, on va résoudre par descente. Une fois  $\mathbf{y}$  calculé, on résout pour  $\mathbf{x}$  par remontée. La procédure peut être répétée en commençant toujours avec  $\mathbf{b}$ .

### 2.1.1 Algorithme LU, sans permutation

Nous présentons maintenant les pas pour résoudre le système par factorisation, dans le cas que tous les pivots qu'on trouve soient différents de zéro.

#### 1. Factorisation.

On passe en  $n - 1$  pas de la matrice  $A^{(1)}$  à une matrice  $A^{(n)} = U$  triangulaire supérieure. Dans le process, on construit la matrice  $L$ .

- i) Pas  $k = 1$ : On initialise la matrice  $A^{(1)} = A$  et on définit la première ligne de  $U$ . Si  $a_{11} \neq 0$  (pivot non-nul), on fait

$$\begin{aligned} a_{ij}^{(1)} &= a_{ij}, \quad i, j \in \{1, \dots, n\} \\ u_{1j} &= a_{1j}^{(1)}, \quad j \in \{1, \dots, n\} \end{aligned} \quad (2.21)$$

- ii) Pas  $k + 1 < n$ : On construit la colonne  $k$  de  $L$ , la matrice  $A^{(k+1)}$  et la ligne  $k + 1$  de  $U$ . Si  $a_{k,k} \neq 0$  (pivot non-nul)

- a) Pour  $L$ :

$$l_{i,k} = 0, \quad i < k, \quad l_{kk} = 1, \quad l_{ik} = \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}}, \quad i > k. \quad (2.22)$$

- b) Pour  $A^{(k+1)}$ :

$$\begin{aligned} a_{ij}^{(k+1)} &= a_{ij}^{(k)}, \quad i \in \{1, \dots, k\}, j \in \{1, \dots, n\} \\ a_{ij}^{(k+1)} &= a_{ij}^{(k)} - l_{ik} a_{kj}^{(k)}, \quad i \in \{k+1, \dots, n\}, j \in \{1, \dots, n\}. \end{aligned} \quad (2.23)$$

- c) Pour  $U$ :

$$u_{k+1,j} = a_{k+1,j}^{(k+1)}, \quad j \in \{1, \dots, n\} \quad (2.24)$$

- iii) Pas  $k = n$ : On agit comme dans le cas  $k + 1 < n$  et on finit la dernière colonne de  $L$

$$l_{i,n} = 0, \quad i < n, \quad l_{nn} = 1. \quad (2.25)$$

On constate que  $U = A^{(n)}$ .

2. Descente: On calcule  $\mathbf{y}$  (où  $L\mathbf{y} = \mathbf{b}$ )

$$y_i = b_i - \sum_{k=1}^{i-1} l_{ik} y_k, \quad i \in \{1, \dots, n\} \quad (2.26)$$

3. Remontée: On calcule  $\mathbf{x}$  (où  $U\mathbf{x} = \mathbf{y}$ )

$$x_i = \frac{1}{u_{ii}} \left( y_i - \sum_{k=i+1}^n u_{ik} x_k \right), \quad i \in \{n, \dots, 1\} \quad (2.27)$$

*Remarque 2.1* (Factorisation, version matricielle). Le pas 1 de la procédure précédente peut s'écrire d'une façon plus compacte :

i) Pas  $k = 1$ :

$$A^{(1)} = A. \quad (2.28)$$

ii) Pas  $k + 1 \leq n$ : On introduit la matrice  $E_k$  avec des zéros partout sauf dans la diagonale principale (diag[1...1]) et en dessous de la diagonale principale dans la colonne  $k$ :

$$E_k = \begin{pmatrix} 1 & \dots & 0 & \dots & \dots & 0 \\ 0 & \ddots & 0 & & & \vdots \\ \vdots & & 1 & 0 & & \vdots \\ \vdots & & \frac{-a_{k+1,k}^{(k)}}{a_{kk}^{(k)}} & 1 & & \vdots \\ \vdots & & \vdots & & \ddots & 0 \\ 0 & \dots & \frac{-a_{n,k}^{(k)}}{a_{kk}^{(k)}} & \dots & 0 & 1 \end{pmatrix}, \quad E_k^{-1} = \begin{pmatrix} 1 & \dots & 0 & \dots & \dots & 0 \\ 0 & \ddots & 0 & & & \vdots \\ \vdots & & 1 & 0 & & \vdots \\ \vdots & & \frac{a_{k+1,k}^{(k)}}{a_{kk}^{(k)}} & 1 & & \vdots \\ \vdots & & \vdots & & \ddots & 0 \\ 0 & \dots & \frac{a_{n,k}^{(k)}}{a_{kk}^{(k)}} & \dots & 0 & 1 \end{pmatrix} \quad (2.29)$$

On a alors

$$A^{(k+1)} = E_k A^{(k)} \quad (2.30)$$

iii) Une fois le process fini, on construit :

$$U = A^{(n)}$$

$$L = E_1^{-1} \cdot \dots \cdot E_n^{-1} = \begin{pmatrix} 1 & \dots & 0 & \dots & \dots & 0 \\ \frac{a_{2,1}^{(1)}}{a_{11}^{(1)}} & \ddots & 0 & & & \vdots \\ \vdots & & 1 & 0 & & \vdots \\ \vdots & & \frac{a_{k+1,k}^{(k)}}{a_{kk}^{(k)}} & 1 & & \vdots \\ \vdots & & \vdots & & \ddots & 0 \\ \frac{a_{n,1}^{(1)}}{a_{11}^{(1)}} & \dots & \frac{a_{n,k}^{(k)}}{a_{kk}^{(k)}} & \dots & \frac{a_{n,n-1}^{(n-1)}}{a_{n-1,n-1}^{(n-1)}} & 1 \end{pmatrix} \quad (2.31)$$

Notez:

- a) Comme on peut le voir dans la notation adoptée dans le pas ii), on a regroupé dans la matrice  $E_k$  tous les coefficients utilisés pour obtenir des zéros sous la diagonale principale à la colonne  $k$ . En général, une matrice  $E_k$  est une matrice triangulaire inférieure, avec des 1 dans la diagonale principale et des éléments non nuls seulement dans la colonne  $k$  en dessous de la diagonale.
- b) La matrice  $L$  a une forme très simple: elle est triangulaire inférieure, avec pour diagonale diag(1, ..., 1) et l'élément  $l_{ij}$  est le coefficient par lequel on a multiplié la ligne de pivot  $a_{jj}^{(j)}$  avant de la retrancher de la ligne  $i$ .

### 2.1.2 Algorithme LU, avec permutation

Dans la discussion précédente, on a supposé qu'à chaque pas de l'itération, on trouve un pivot non nul. Ceci n'est pas le cas en général. Pour illustrer la manière de procéder quand on trouve un pivot nul, on commence avec un exemple.

**Exemple 2.2** (Une première décomposition LU avec permutation). On considère l'exemple de la matrice :

$$A = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & -1 \\ 2 & 1 & 1 \end{pmatrix}. \quad (2.32)$$

Le premier pas de la méthode LU :

$$A^{(1)} = A \begin{array}{l} \xrightarrow{\ell_2 \rightarrow \ell_2 - \ell_1} \\ \xrightarrow{\ell_3 \rightarrow \ell_3 - 2\ell_1} \end{array} A^{(2)} = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 0 & -2 \\ 0 & -1 & -1 \end{pmatrix} \quad (2.33)$$

$$E_1 = \begin{pmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ -2 & 0 & 1 \end{pmatrix}$$

produit un zéro pour le deuxième pivot. Pour pouvoir continuer, on interchange simplement la deuxième et la troisième ligne. En effet, cette opération ne change le système linéaire associé. Concrètement :

$$A^{(2)} \begin{array}{l} \xrightarrow{\ell_2 \leftrightarrow \ell_3} \\ \xrightarrow{\ell_3 \leftrightarrow \ell_2} \end{array} \tilde{A}^{(3)} = \begin{pmatrix} 1 & 1 & 1 \\ 0 & -1 & -1 \\ 0 & 0 & -2 \end{pmatrix} \quad (2.34)$$

$$P_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$$

Dans ce cas, l'itération a déjà abouti à une matrice triangulaire supérieure. On peut écrire les pas suivants :

$$P_2 E_1 A = U \quad (2.35)$$

Notons que  $P_2 E_1$  n'est pas triangulaire inférieure. En effet,

$$P_2 E_1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ -2 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ -2 & 0 & 1 \\ -1 & 1 & 0 \end{pmatrix} \neq L \quad (2.36)$$

Notez, néanmoins, qu'on peut écrire :

$$P_2 E_1 = \begin{pmatrix} 1 & 0 & 0 \\ -2 & 0 & 1 \\ -1 & 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ -1 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} = E'_1 P_2 \quad (2.37)$$

$\underbrace{\hspace{10em}}_{E'_1}$

où  $E'_1$  est obtenue à partir de  $E_1$  en interchangeant, dans la colonne 1, deux éléments (dans ce cas, les deuxième et troisième). Ainsi, on peut réécrire (2.35)

$$\begin{aligned} E'_1 P_2 A &= U \\ P_2 A &= \underbrace{(E'_1)^{-1}}_L U = LU. \end{aligned} \quad (2.38)$$

Autrement dit,

$$\underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}}_{P_2} \underbrace{\begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & -1 \\ 2 & 1 & 1 \end{pmatrix}}_A = \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix}}_L \underbrace{\begin{pmatrix} 1 & 1 & 1 \\ 0 & -1 & -1 \\ 0 & 0 & -2 \end{pmatrix}}_U \quad (2.39)$$

On a pas réussi<sup>2</sup> à écrire  $A = LU$ . Par contre, on a trouvé un réarrangement des lignes de  $A$ , donné par  $P_2A$ , qui admet en fait la décomposition  $P_2A = LU$ . Dans le théorème 2.1, on va montrer qu'on peut toujours trouver une telle décomposition.

La matrice  $P_2$  dans l'exemple précédent est un cas de matrice de permutation  $P$ , qui nous a servi à réordonner les lignes dans la procédure de Gauss pour trouver un pivot non nul.

### Erreur d'arrondi est choix de pivot.

Même si nous ne sommes pas obligés de réarranger les lignes pour obtenir un pivot non nul, nous pouvons nous intéresser à un tel changement d'ordre. En général, étant donnée une matrice  $A$ , il n'y a pas une seule matrice de permutation  $P$  telle qu'on puisse écrire  $PA = LU$ . Différents  $P$  sont possibles et chaque  $P$  amène génériquement a un couple  $(L, U)$  différent.

Un choix privilégié du point de vue numérique consiste à arranger les lignes de telle manière qu'on choisisse comme pivot, à chaque pas de l'itération de la méthode de Gauss, le coefficient avec la plus grande valeur absolue. La raison est la limitation des erreurs d'arrondi: si on peut garder seulement un nombre fini de chiffres significatifs, comme c'est le cas avec un ordinateur, c'est mieux de travailler avec des nombres avec des petites valeurs absolues plutôt qu'avec des nombres des grandes valeurs absolues. Vu que, dans la méthode de Gauss, il faut diviser à chaque pas le pivot, ça nous amène au choix signalé.

On reprend l'exemple 2.2 pour illustrer ce point.

**Exemple 2.3** (Non-unicité de  $P$  dans la décomposition  $PA = LU$ ). Si on considère à nouveau la matrice (2.32), en regardant la colonne 1, on voit que, dans le premier pas de l'itération, on peut interchanger les lignes 1 et 3 de manière à ce que le pivot choisi soit celui avec la valeur absolue la plus grande. Explicitement, cela donne :

$$\begin{aligned} A^{(1)} &= A = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & -1 \\ 2 & 1 & 1 \end{pmatrix} \xrightarrow{P_1 = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}} \tilde{A}^{(1)} = \begin{pmatrix} 2 & 1 & 1 \\ 1 & 1 & -1 \\ 1 & 1 & 1 \end{pmatrix} \\ E_1 &= \begin{pmatrix} 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \\ -\frac{1}{2} & 0 & 1 \end{pmatrix} \xrightarrow{\begin{matrix} \ell_2 \rightarrow \ell_2 - \frac{1}{2}\ell_1 \\ \ell_3 \rightarrow \ell_3 - \frac{1}{2}\ell_1 \end{matrix}} A^{(2)} = \begin{pmatrix} 2 & 1 & 1 \\ 0 & \frac{1}{2} & -\frac{3}{2} \\ 0 & \frac{1}{2} & \frac{1}{2} \end{pmatrix} \\ E_2 &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -1 & 1 \end{pmatrix} \xrightarrow{\ell_3 \rightarrow \ell_3 - \ell_2} A^{(3)} = \begin{pmatrix} 2 & 1 & 1 \\ 0 & \frac{1}{2} & -\frac{3}{2} \\ 0 & 0 & 2 \end{pmatrix} \end{aligned} \quad (2.40)$$

<sup>2</sup>En fait, on peut montrer que pour cette matrice, ce n'est possible à faire (voir caractérisation dans le théorème 2.1).



Nous pouvons maintenant présenter le procédé de factorisation dans la décomposition  $LU$  dans le cas général d'un possible arrangement des lignes à chaque pas de l'itération.

**Factorisation dans l'algorithme  $LU$  avec permutations** (version matricielle).

Nous suivons les trois pas dans 2.1:

i) Pas  $k = 1$ :

$$A^{(1)} = A . \quad (2.46)$$

ii) Pas  $k + 1 \leq n$ : D'abord, on arrange les lignes de la matrice en multipliant à gauche par une permutation  $P_k$  (notamment une transposition) qui change seulement l'ordre des lignes  $\ell(i)$  avec  $i \geq k$ . Après, on multiplie (à gauche) par la matrice  $E_k$  de (2.29). On a alors

$$A^{(k+1)} = E_k P_k A^{(k)} \quad (2.47)$$

iii) Une fois le procédé fini, on construit une matrice triangulaire supérieure

$$(E_{n-1} P_{n-1}) \dots (E_2 P_2) (E_1 P_1) A = A^{(n)} = U . \quad (2.48)$$

La différence fondamentale avec le cas sans permutation est que les pas précédents n'amènent pas directement, en général, à une matrice triangulaire inférieure comme  $L = (E_1)^{-1} \dots (E_{n-1})^{-1}$ . En effet, en général, les matrices  $E_k$  et  $P_{k+l}$  (avec  $l \geq 1$ ) ne commutent pas, ce qui empêche de factoriser la matrice  $E_{n-1} \dots E_1$  à gauche dans (2.48). On peut illustrer ceci avec un exemple. Si on considère le cas  $n = 4$ , on peut écrire (2.48) comme :

$$(E_3 P_3)(E_2 P_2)(E_1 P_1)A = U \quad (2.49)$$

Comme on l'a dit, les matrices  $E_i$  et  $P_j$  ne commutent pas en général, alors on ne peut pas écrire  $(E_3 E_2 E_1)(P_3 P_2 P_1)A = U$ . En effet, si on considère par exemple le troisième pas de l'algorithme, on peut écrire :

$$E_2 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & a & 1 & 0 \\ 0 & b & 0 & 1 \end{pmatrix} , \quad P_3 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix} , \quad (2.50)$$

pour certains  $a$  et  $b$ . Si on calcule  $P_3 E_2$  et  $E_2 P_3$ , on obtient :

$$P_3 E_2 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & b & 0 & 1 \\ 0 & a & 1 & 0 \end{pmatrix} , \quad E_2 P_3 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & a & 0 & 1 \\ 0 & b & 1 & 0 \end{pmatrix} \quad (2.51)$$

et on constate que  $P_3 E_2 - E_2 P_3 \neq 0$ . Néanmoins on peut observer que si dans  $P_3 E_2$ , on interchange la troisième et la quatrième colonne (ce qu'on fait en multipliant à droite par  $P_3$ ), on récupère une matrice "type  $E_k$ ". Plus spécifiquement, on a :

$$P_3 E_2 = P_3 E_2 \underbrace{(P_3 P_3^{-1})}_{I_4} = E_2' P_3 , \quad (2.52)$$

où (on a aussi utilisé  $P_3^{-1} = P_3$ )

$$E'_2 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & b & 1 & 0 \\ 0 & a & 0 & 1 \end{pmatrix} \quad (2.53)$$

est, en effet, une matrice type  $E_k$  dans (2.29), obtenue de  $E_2$  précisément en interchangeant les coefficients 3 et 4 (ceux interchangeés par  $P_3$ ) dans la colonne 2.

Cette propriété de commutation n'est pas valable pour n'importe quelle permutation et matrice "type  $E_k$ ", mais il n'est pas difficile à se convaincre qu'elle est en fait valide pour les matrices spécifiques qu'on est en train de considérer: notamment, les permutations  $P_{k+l}$  ( $l \geq 1$ ) correspondant à l'arrangement de lignes après la ligne dont le pivot correspond à  $E_k$ . Dans ce cas, on a

$$P_{k+l}E_k = E'_k P_{k+l} \quad , \quad l \geq 1 \quad , \quad (2.54)$$

où  $E'_k$  est obtenue à partir de  $E_k$  en interchangeant les éléments dans la colonne  $k$  sous la diagonale principale, d'accord avec la permutation  $P_{k+l}$  (cf.  $E_2$  et  $E'_2$  en (2.50) et (2.53), respectivement). Pour l'expression générale, pour le cas des transpositions, on a :

$$P_{k+l}^{(i_o, j_o)} E_k = E'_k P_{k+l}^{(i_o, j_o)} \quad , \quad (2.55)$$

avec

$$\begin{aligned} (E'_k)_{i_o, k} &= (E_k)_{j_o, k} \\ (E'_k)_{j_o, k} &= (E_k)_{i_o, k} \\ (E'_k)_{i, j} &= (E_k)_{i, j} \quad \text{pour le reste, c.à.d. } j \neq k \text{ ou } (i \neq i_o, j_o) . \end{aligned} \quad (2.56)$$

Nous sommes maintenant en condition pour énoncer et démontrer le théorème de décomposition  $LU$  avec permutations.

**Théorème 2.1** (Décomposition  $LU$ ). *Soit  $A \in M_n(\mathbb{R})$  une matrice inversible. Alors il existe une permutation  $P$  (non unique) telle que, pour ce  $P$ , il existe un et un seul couple de matrices  $(L, U)$  où  $L$  est triangulaire inférieure de termes diagonaux égaux à 1 et  $U$  triangulaire supérieure, avec*

$$PA = LU \quad . \quad (2.57)$$

**Preuve :**

- *Existence.*

L'existence découle directement de l'expression (2.48) comme résultat de la factorisation dans l'algorithme  $LU$  avec permutations. On peut alors écrire :

$$E_{n-1}P_{n-1} \dots E_2P_2 E_1P_1 A = U. \quad (2.58)$$

Si on utilise maintenant (2.54) pour passer les permutations  $P_i$  vers la droite, on a :

$$\begin{aligned} E_{n-1}E'_{n-2}P_{n-1} \dots E'_2P_3E'_1P_2P_1 A &= U \\ \dots & \\ (E_{n-1}E'_{n-2} \dots E_2^{n-3}E_1^{n-2}) (P_{n-1}P_{n-2} \dots P_3P_2P_1)A &= U . \end{aligned} \quad (2.59)$$



Si on définit maintenant  $L$  et  $P$  comme :

$$\begin{aligned} L &= (E_{n-1}E'_{n-2} \cdots E_2^{n-3}E_1^{n-2})^{-1} = (E_1^{n-2})^{-1}(E_2^{n-3})^{-1}(E'_{n-2})^{-1}(E_{n-1})^{-1} \\ P &= P_{n-1}P_{n-2} \cdots P_3P_2P_1, \end{aligned} \quad (2.60)$$

on peut écrire

$$PA = LU. \quad (2.61)$$

Notez que la preuve d'existence n'utilise pas le fait que  $A$  soit inversible.

- *Unicité.*

On suppose que pour  $P$  donné, il existe deux couples  $(L_1, U_1)$  et  $(L_2, U_2)$ , tels que:

$$PA = L_1U_1, \quad PA = L_2U_2, \quad (2.62)$$

alors

$$L_1U_1 = L_2U_2 \quad (2.63)$$

Si  $A$  est inversible, vu que  $L_1$  et  $L_2$  sont inversibles, on conclut que  $U_1$  et  $U_2$  sont inversibles. Ainsi, on peut multiplier (2.63) à gauche par  $L_2^{-1}$  et à droite par  $U_1^{-1}$ . On obtient

$$L_2^{-1}L_1 = U_2U_1^{-1} \quad (2.64)$$

Maintenant, on utilise<sup>3</sup> le fait que  $L_2^{-1}L_1$  est une matrice triangulaire inférieure avec  $\text{diag}(1, \dots, 1)$  pour diagonale principale et  $U_2U_1^{-1}$  est une matrice triangulaire supérieure. Alors, nécessairement, on a :

$$L_2^{-1}L_1 = I_n = U_2U_1^{-1}, \quad (2.65)$$

et on conclut

$$L_1 = L_2, \quad U_1 = U_2 \quad (2.66)$$

□

*Remarque 2.3.* Faisons quelques remarques générales sur le résultat de la décomposition  $PA = LU$ :

- i) Dans le résultat de la décomposition  $PA = LU$ , on affirme l'existence de  $P$ , mais pas son unicité.
- ii) étant donnée  $P$ , si  $A$  est inversible,  $L$  et  $U$  sont uniques. Mais si  $A$  n'est pas inversible, la décomposition  $PA = LU$  peut toujours se faire, mais elle n'est pas unique.

---

<sup>3</sup>Exercice! (Matrices triangulaires). Montrer que pour  $L_1, L_2$  matrices triangulaires inférieures (respectivement  $U_1, U_2$  matrices triangulaires supérieures), les matrices  $L_1^{-1}, L_2^{-1}, L_1 \cdot L_2$  sont triangulaires inférieures (respectivement,  $U_1^{-1}, U_2^{-1}, U_1 \cdot U_2$  sont triangulaires supérieures).

*Exemple 2.4* (Non-unicité de  $PA = LU$ , pour  $A$  non-inversible). Si on considère

$$A = \begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix} \quad (2.67)$$

on peut vérifier  $\forall \lambda \in \mathbb{R}$

$$\begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ \lambda & 1 \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 0 & 1 - \lambda \end{pmatrix} \quad (2.68)$$

Cette non-unicité est consistant avec la non-inversibilité de  $A$

- i) *Pivot partiel versus pivot total.* Dans l'algorithme qui amène à la construction des matrices  $L$  et  $U$  (partie d'existence), on interchange seulement les lignes, et pas les colonnes, pour trouver le pivot avec la valeur absolue la plus grande. On parle de *pivot partiel*. Pour le pivot total, on chercherait le plus grand pivot en changeant aussi les colonnes.
- ii) L'algorithme du pivot partiel dans le problème  $A\mathbf{x} = \mathbf{b}$ , ne réarrange pas les composantes de  $\mathbf{x}$  vu qu'on ne change pas les colonnes de  $A$ . Le vecteur  $\mathbf{b}$ , par contre, est en général réarrangé.

### Caractérisation $A = LU$ sans permutation

Avant de considérer plus spécifiquement le cas de matrices symétriques définies positives, nous allons donner une caractérisation de la possibilité d'écrire  $A = LU$  sans matrice de permutation  $P$ . D'abord, nous devons introduire la définition suivante.

**Définition 2.3** (matrice principale, mineur principal). Soit  $A \in M_n(\mathbb{R})$  et  $k \in \{1, \dots, k\}$ . On appelle:

- i) Matrice principale d'ordre  $k$  de  $A$ , la matrice  $A_k \in M_k(\mathbb{R})$  définie par  $(A_k)_{ij} = A_{ij}, \forall i, j \in \{1, \dots, k\}$ .
- ii) Mineur principal d'ordre  $k$ , le déterminant de la matrice principale d'ordre  $k$ .

Avec ces éléments, nous pouvons énoncer la caractérisation suivante.

**Proposition 2.1** (Caractérisation de  $A = LU$ , sans permutation). Soit  $A \in M_n(\mathbb{R})$ . Les deux propriétés suivantes sont équivalentes:

- i) Il existe une unique couple  $(L, U)$ , avec  $L$  matrice triangulaire inférieure dont les coefficients diagonaux sont égaux à 1 et  $U$  une matrice inversible triangulaire supérieure, tel que  $A = LU$ .
- ii) Les mineurs principaux de  $A$  sont tous non nuls.

Avant de démontrer ce résultat, on va introduire le lemme technique suivant.

**Lemme 2.1.** Soit  $A \in M_n(\mathbb{R})$  et  $k \in \{1, \dots, n\}$ . On suppose que pour chaque matrice principale  $A_k$  de  $A$ , il existe une matrice  $L_k \in M_k(\mathbb{R})$  triangulaire inférieure avec des coefficients diagonaux tous égaux à 1 et une matrice triangulaire supérieure  $U_k \in M_k(\mathbb{R})$  inversible, telles que  $A_k = L_k U_k$ . On a alors:

a) On peut écrire  $A$  sous la forme :

$$A = \begin{pmatrix} L_k & 0_{k \times (n-k)} \\ C_k & I_{n-k} \end{pmatrix} \begin{pmatrix} U_k & B_k \\ 0_{(n-k) \times k} & D_k \end{pmatrix}, \quad (2.69)$$

où  $B_k \in M_{k, n-k}(\mathbb{R})$ ,  $C_k \in M_{n-k, k}(\mathbb{R})$  et  $D_k \in M_{n-k, n-k}(\mathbb{R})$ .

b) En particulier, la matrice principale d'ordre  $k+1$  s'écrit sous la forme

$$A_{k+1} = \begin{pmatrix} L_k & 0_{1 \times k} \\ c_k & 1 \end{pmatrix} \begin{pmatrix} U_k & b_k \\ 0_{k \times 1} & d_k \end{pmatrix}, \quad (2.70)$$

où  $(b_k) \in M_{k,1}(\mathbb{R})$  est la première colonne de la matrice  $B_k$ ,  $(c_k) \in M_{k,1}(\mathbb{R})$  est la première ligne de la matrice  $C_k$ , et  $d_k$  est le coefficient de la ligne 1 et colonne 1 de  $D_k$ .

**Preuve:** À compléter. □

**Preuve** (de la proposition 2.1): À compléter. □

## 2.2 Décomposition de Choleski

Dans cette section, nous allons nous concentrer sur le cas des matrices symétriques définies positives.

**Définition 2.4** (Matrice définie positive).  $A \in M_n(\mathbb{R})$  est symétrique définie positive si :

- i)  $A = A^t$  (symétrique).
- ii)  $Ax \cdot x > 0, \forall x \in \mathbb{R}^n \setminus \{0\}$  (définie positive).

Un exemple important de ce type de matrice est la discrétisation du Laplacien qu'on a étudié dans l'exemple 1.1.

Une caractérisation très utile des matrices symétriques définies positives est donnée en termes de mineurs principaux.

**Proposition 2.2** (Caractérisation d'une matrice symétrique définie positive par ses mineurs principaux).  $A \in M_n(\mathbb{R})$  est symétrique définie positive si et seulement si tous ses mineurs principaux sont strictement positifs.

En particulier, comme conséquence de la proposition 2.1, une matrice  $A$  symétrique définie positive admet toujours une (unique) décomposition  $A = LU$  sans permutation. Dans cette section nous allons discuter une autre factorisation spécialement adaptée au caractère symétrique et défini positif de ces matrices. Nous commençons avec un exemple qui part de la décomposition  $LU$ .

**Exemple 2.5** (Premier contact avec la factorisation de Choleski). Nous considérons la matrice

$$A = \begin{pmatrix} 1 & 1 & 2 \\ 1 & 5 & 6 \\ 2 & 6 & 17 \end{pmatrix}. \quad (2.71)$$

La matrice  $A$  est symétrique et, comme on peut le vérifier à partir e.g. du calcul de ses mineurs principaux, elle est aussi définie positive. Alors, elle admet une décomposition  $A = LU$  unique. On va construire explicitement une telle décomposition.

$$\begin{pmatrix} 1 & 1 & 2 \\ 1 & 5 & 6 \\ 2 & 6 & 17 \end{pmatrix} \xrightarrow[\substack{\ell_2 \rightarrow \ell_2 - \ell_1 \\ \ell_3 \rightarrow \ell_3 - 2\ell_1}]{(E_1)^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 2 & 0 & 1 \end{pmatrix}} \begin{pmatrix} 1 & 1 & 2 \\ 0 & 4 & 4 \\ 0 & 4 & 13 \end{pmatrix} \xrightarrow[\ell_3 \rightarrow \ell_3 - \ell_2]{(E_2)^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix}} \begin{pmatrix} 1 & 1 & 2 \\ 0 & 4 & 4 \\ 0 & 0 & 9 \end{pmatrix}. \quad (2.72)$$

Alors, on peut écrire

$$\begin{pmatrix} 1 & 1 & 2 \\ 1 & 5 & 6 \\ 2 & 6 & 17 \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 2 & 1 & 1 \end{pmatrix}}_L \underbrace{\begin{pmatrix} 1 & 1 & 2 \\ 0 & 4 & 4 \\ 0 & 0 & 9 \end{pmatrix}}_U \quad (2.73)$$

Notez qu'on peut écrire :

$$U = \begin{pmatrix} 1 & 1 & 2 \\ 0 & 4 & 4 \\ 0 & 0 & 9 \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 9 \end{pmatrix}}_D \underbrace{\begin{pmatrix} 1 & 1 & 2 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix}}_{L^t} \quad (2.74)$$

où  $D$  est la diagonale principale de  $U$ . Alors,

$$A = LDL^t \quad (2.75)$$

Nous notons que tous les coefficients de  $D$  sont strictement positifs. Alors, nous pouvons introduire la matrice racine carrée

$$D^{\frac{1}{2}} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix} \quad (2.76)$$

pour écrire :

$$\begin{aligned} A = LD^{\frac{1}{2}}D^{\frac{1}{2}}L^t &= \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 2 & 1 & 1 \end{pmatrix}}_{\tilde{L}} \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix}}_D \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix}}_{\tilde{L}^t} \underbrace{\begin{pmatrix} 1 & 1 & 2 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix}}_{L^t} \\ &= \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 1 & 2 & 0 \\ 2 & 2 & 3 \end{pmatrix}}_{\tilde{L}} \underbrace{\begin{pmatrix} 1 & 1 & 2 \\ 0 & 2 & 2 \\ 0 & 0 & 3 \end{pmatrix}}_{\tilde{L}^t} \end{aligned} \quad (2.77)$$

C'est-à-dire, on a écrit

$$A = \tilde{L}\tilde{L}^t, \quad \text{avec } \tilde{L} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 2 & 0 \\ 2 & 2 & 3 \end{pmatrix}. \quad (2.78)$$

La matrice  $\tilde{L}$  est triangulaire inférieure avec coefficients positifs dans la diagonale principale. L'expression (2.78) s'appelle la décomposition de Choleski de  $A$ .

Le pas clé, dans l'exemple précédent, est la possibilité de prendre (dans les réels) la racine carrée de la matrice  $D$  définie pour la diagonale principale de  $U$ . Dans l'exemple suivant, le cas de la décomposition de Choleski dans  $M_2(\mathbb{R})$ , on voit que ceci est une conséquence du caractère symétrique défini positif des matrices considérées. Après l'exemple, ceci est énoncé en toute généralité.

**Exemple 2.6** (Décomposition de Choleski: cas  $n = 2$ ). On considère une matrice  $A \in M_2(\mathbb{R})$  symétrique définie positive, on peut toujours écrire :

$$A = \begin{pmatrix} a & b \\ b & c \end{pmatrix}. \quad (2.79)$$

Si on prends  $\mathbf{x} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ , le caractère défini positif de  $A$  implique

$$A\mathbf{x} \cdot \mathbf{x} = \mathbf{x}^t A \mathbf{x} = \begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} a & b \\ b & c \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = a > 0. \quad (2.80)$$

Autrement dit, le mineur principal  $A_1 = a > 0$ , comme il en résulte de la caractérisation 2.2. Ainsi, étant donné que  $a > 0$ , on peut procéder à la décomposition  $LU$

$$A = \begin{pmatrix} a & b \\ b & c \end{pmatrix} \xrightarrow[L = \begin{pmatrix} 1 & 0 \\ \frac{b}{a} & 1 \end{pmatrix}]{\ell_2 \rightarrow \ell_2 - \frac{b}{a} \ell_1} \begin{pmatrix} a & b \\ 0 & c - \frac{b^2}{a} \end{pmatrix}. \quad (2.81)$$

Alors, on peut écrire :

$$A = \underbrace{\begin{pmatrix} 1 & 0 \\ \frac{b}{a} & 1 \end{pmatrix}}_L \underbrace{\begin{pmatrix} a & b \\ 0 & c - \frac{b^2}{a} \end{pmatrix}}_U = \underbrace{\begin{pmatrix} 1 & 0 \\ \frac{b}{a} & 1 \end{pmatrix}}_L \underbrace{\begin{pmatrix} a & 0 \\ 0 & c - \frac{b^2}{a} \end{pmatrix}}_D \underbrace{\begin{pmatrix} 1 & \frac{b}{a} \\ 0 & 1 \end{pmatrix}}_{L^t} \quad (2.82)$$

Si on utilise maintenant le reste de l'information sur  $A$ , notamment la positivité du mineur  $A_2 = ac - b^2 > 0$ , on conclut (avec  $a > 0$ )  $c - \frac{b^2}{a} > 0$ . On peut alors écrire :

$$A = \underbrace{\begin{pmatrix} 1 & 0 \\ \frac{b}{a} & 1 \end{pmatrix}}_L \underbrace{\begin{pmatrix} \sqrt{a} & 0 \\ 0 & \sqrt{c - \frac{b^2}{a}} \end{pmatrix}}_{D^{\frac{1}{2}}} \underbrace{\begin{pmatrix} \sqrt{a} & 0 \\ 0 & \sqrt{c - \frac{b^2}{a}} \end{pmatrix}}_{D^{\frac{1}{2}}} \underbrace{\begin{pmatrix} 1 & \frac{b}{a} \\ 0 & 1 \end{pmatrix}}_{L^t} \quad (2.83)$$

$$= \underbrace{\begin{pmatrix} \sqrt{a} & 0 \\ \frac{b}{\sqrt{a}} & \sqrt{c - \frac{b^2}{a}} \end{pmatrix}}_{\tilde{L}} \underbrace{\begin{pmatrix} \sqrt{a} & \frac{b}{\sqrt{a}} \\ 0 & \sqrt{c - \frac{b^2}{a}} \end{pmatrix}}_{\tilde{L}^t} \quad (2.84)$$

Comme on l'a annoncé ci-dessus, la proposition suivante nous garantit que  $D^{\frac{1}{2}}$  peut toujours être construite comme matrice réelle, pour des matrices symétriques définies positives.

**Proposition 2.3** (Caractérisation d'une matrice symétrique définie positive à partir de sa décomposition  $LU$ ). Soit  $A \in M_n(\mathbb{R})$  une matrice symétrique admettant une décomposition  $LU$  sans permutation. Alors  $A$  est symétrique définie positive si et seulement si tous les pivots (c.à.d. les coefficients diagonaux de la matrice  $U$ ) sont strictement positifs.

**Preuve:** À compléter. □

Pour continuer, nous allons énoncer et prouver le théorème de décomposition de Choleski, qui présente une partie d'existence et une autre d'unicité. En particulier, la partie d'existence suit l'esprit du schéma de l'exemple 2.6 pour le cas  $n = 2$ , dont le contenu est essentiellement celui de la proposition 2.3. La partie de l'unicité, qui utilise d'une manière basique le résultat d'existence, proportionne l'algorithme de Choleski proprement dit.

**Théorème 2.2** (Décomposition de Choleski). *Soit  $A \in M_n(\mathbb{R})$ , une matrice symétrique définie positive. Alors il existe une unique matrice  $\tilde{L} \in M_n(\mathbb{R})$  telle que:*

- i)  $\tilde{L}$  est triangulaire inférieure.
- ii)  $\tilde{l}_{ii} > 0, \forall i \in \{1, \dots, n\}$ .
- iii)  $A = \tilde{L}\tilde{L}^t$ .

**Preuve:**

- *Existence.* Par récurrence sur la dimension  $n$ .

1) Cas  $n = 1$ .

On a  $A = (a_{11})$ .  $A$  est automatiquement symétrique et elle est définie positive si et seulement si  $a_{11} > 0$ . On définit  $\tilde{L} = (\tilde{l}_{11})$ , avec  $\tilde{l}_{11} = \sqrt{a_{11}}$  et on a

$$A = \tilde{L}\tilde{L}^t . \quad (2.85)$$

2) On suppose  $A = \tilde{L}\tilde{L}^t$ , avec  $\tilde{L}$  triangulaire inférieure et  $\tilde{l}_{ii} > 0, \forall A \in M_p(\mathbb{R}), 1 \leq p \leq n$ . On va démontrer que, pour  $A \in M_{n+1}(\mathbb{R})$  symétrique et définie positive, il existe une décomposition de Choleski. On considère  $A \in M_{n+1}(\mathbb{R})$ , que l'on peut écrire :

$$A = \left( \begin{array}{c|c} B & \mathbf{a} \\ \mathbf{a}^t & \alpha \end{array} \right) \quad (2.86)$$

où  $\mathbf{a} \in \mathbb{R}^n$  est un vecteur colonne,  $B$  est symétrique et  $\alpha \in \mathbb{R}$ . D'abord, on montre que  $B$  est définie positive à partir du fait que  $A$  est définie positive, i.e.  $\mathbf{x}^t A \mathbf{x} > 0, \forall \mathbf{x} \in \mathbb{R}^{n+1}, \mathbf{x} \neq 0$ . Si on prend en particulier  $\mathbf{x} = \begin{pmatrix} \mathbf{y} \\ 0 \end{pmatrix}$ , avec  $\mathbf{y} \in \mathbb{R}^n$ , on a :

$$0 < \mathbf{x}^t A \mathbf{x} = (\mathbf{y} \ 0) \begin{pmatrix} B & \mathbf{a} \\ \mathbf{a}^t & \alpha \end{pmatrix} \begin{pmatrix} \mathbf{y} \\ 0 \end{pmatrix} = \mathbf{y}^t B \mathbf{y} , \quad \forall \mathbf{y} \in \mathbb{R}^n . \quad (2.87)$$

Alors  $B$  est symétrique définie positive. Par hypothèse de récurrence  $\exists M \in M_n(\mathbb{R})$  triangulaire inférieure telle que

$$B = M M^t , \quad m_{ii} > 0 , \quad \forall i \in \{1, \dots, n\} . \quad (2.88)$$

On cherche une matrice  $\tilde{L}$  avec l'hypothèse (Ansatz) de la forme :

$$\tilde{L} = \begin{pmatrix} M & 0 \\ \mathbf{b}^t & \lambda \end{pmatrix} \quad (2.89)$$

où  $\mathbf{b} \in \mathbb{R}^n$  et  $\lambda \in \mathbb{R}$  (on veut que  $\lambda > 0$ ). On veut imposer :

$$A = \tilde{L}\tilde{L}^t . \quad (2.90)$$

D'abord, nous évaluons :

$$\tilde{L}\tilde{L}^t = \left( \begin{array}{c|c} M & 0 \\ \mathbf{b}^t & \lambda \end{array} \right) \left( \begin{array}{c|c} M^t & \mathbf{b} \\ 0 & \lambda \end{array} \right) = \left( \begin{array}{c|c} MM^t & M\mathbf{b} \\ \mathbf{b}^t M^t & \mathbf{b}^t \mathbf{b} + \lambda^2 \end{array} \right) = \left( \begin{array}{c|c} B & M\mathbf{b} \\ (M\mathbf{b})^t & \mathbf{b}^t \mathbf{b} + \lambda^2 \end{array} \right) \quad (2.91)$$

où on a utilisé (2.88). Maintenant, si on impose (2.90) en utilisant (2.86), on obtient :

$$\begin{aligned} M\mathbf{b} &= \mathbf{a} \\ \mathbf{b}^t \mathbf{b} + \lambda^2 &= \alpha . \end{aligned} \quad (2.92)$$

Par hypothèse de récurrence,  $M$  est inversible ( $m_{ii} > 0$ ), alors :

$$\mathbf{b} = M^{-1}\mathbf{a} , \quad (2.93)$$

et

$$\alpha = (M^{-1}\mathbf{a})^t (M^{-1}\mathbf{a}) + \lambda^2 = \dots = \mathbf{a}^t \underbrace{(MM^t)^{-1}}_B \mathbf{a} + \lambda^2 = \mathbf{a}^t B^{-1} \mathbf{a} + \lambda^2 , \quad (2.94)$$

et on peut écrire :

$$\lambda^2 = \alpha - \mathbf{a}^t B^{-1} \mathbf{a} . \quad (2.95)$$

On va montrer que  $\alpha - \mathbf{a}^t B^{-1} \mathbf{a} > 0$ , en utilisant la pièce d'information qu'on n'a pas encore utilisé, notamment que  $A$  est aussi définie positive pour un vecteur  $\mathbf{x}$  avec la composante  $x^{n+1} \neq 0$ . Pour ça, on choisit:

$$\mathbf{x} = \begin{pmatrix} B^{-1}\mathbf{a} \\ -1 \end{pmatrix} \quad (2.96)$$

On a  $\mathbf{x} \neq 0$ , alors :

$$\begin{aligned} 0 < \mathbf{x}^t A \mathbf{x} &= ((B^{-1}\mathbf{a})^t \quad -1) \left( \begin{array}{c|c} B & \mathbf{a} \\ \mathbf{a}^t & \alpha \end{array} \right) \begin{pmatrix} B^{-1}\mathbf{a} \\ -1 \end{pmatrix} = ((B^{-1}\mathbf{a})^t \quad -1) \left( \begin{array}{c} \overbrace{BB^{-1}}^{I_n} \mathbf{a} - \mathbf{a} \\ \mathbf{a}^t B^{-1} \mathbf{a} - \alpha \end{array} \right) \\ &= ((B^{-1}\mathbf{a})^t \quad -1) \begin{pmatrix} 0 \\ \mathbf{a}^t B^{-1} \mathbf{a} - \alpha \end{pmatrix} = \alpha - \mathbf{a}^t B^{-1} \mathbf{a} \end{aligned} \quad (2.97)$$

Alors, de (2.95) on peut écrire

$$\lambda = \sqrt{\alpha - \mathbf{a}^t B^{-1} \mathbf{a}} , \quad (2.98)$$

et, finalement

$$\tilde{L} = \left( \begin{array}{c|c} M & 0 \\ (M^{-1}\mathbf{a})^t & \sqrt{\alpha - \mathbf{a}^t B^{-1} \mathbf{a}} \end{array} \right) , \quad (2.99)$$

est une matrice triangulaire inférieure, avec  $\tilde{l}_{ii}, \forall i \in \{1, \dots, n+1\}$  et telle que (2.90) est satisfaite.

- *Unicité.*

Dans la deuxième de partie de la preuve, on utilise de manière fondamentale l'existence montrée dans la première partie. Cela nous garantit, en particulier, de prendre la racine carrée des expressions qu'on va trouver dans l'algorithme.

Étant donnée  $A$  symétrique définie positive, dans la première partie on a montré l'existence d'une matrice  $\tilde{L}$  telle que :

$$A = \tilde{L}\tilde{L}^t ; \tilde{l}_{ij} = 0 , j > i ; \tilde{l}_{ii} > 0 \forall i . \quad (2.100)$$

Alors, si on écrit le produit de matrices en termes des coefficients :

$$a_{ij} = \sum_{k=1}^n \tilde{l}_{ik}\tilde{l}_{kj}^t = \sum_{k=1}^n \tilde{l}_{ik}\tilde{l}_{jk} \quad (2.101)$$

On va employer cette expression pour déterminer d'une manière constructive l'expression des coefficients  $\tilde{l}_{ij}$  en termes des données  $a_{ij}$ , en particulier *de forme unique*.

On va procéder colonne par colonne.

- i) On commence avec la première colonne. Si on fait  $j = 1$  dans (2.101) on a, pour la première ligne

$$a_{11} = \sum_{k=1}^n \tilde{l}_{1k}\tilde{l}_{1k} = \tilde{l}_{11}\tilde{l}_{11} + \sum_{k=2}^n \tilde{l}_{1k} \underbrace{\tilde{l}_{1k}}_{=0} = (\tilde{l}_{11})^2 . \quad (2.102)$$

( $k > 1$ )

Alors

$$l_{11} = \sqrt{a_{11}} > 0 . \quad (2.103)$$

On a la garantie qu'on peut prendre la racine carrée de  $a_{11}$  parce qu'on a montré dans la première partie de la preuve l'existence d'un  $l_{11} > 0$ : dans le cas de trouver  $a_{11} \leq 0$ , il y aurait donc une contradiction avec le résultat déjà prouvé.

On continue avec le reste des lignes dans la première colonne :

$$\begin{aligned} a_{21} &= \sum_{k=1}^n \tilde{l}_{2k}\tilde{l}_{1k} = \tilde{l}_{21}\tilde{l}_{11} + \sum_{k=2}^n \tilde{l}_{2k} \underbrace{\tilde{l}_{1k}}_{=0} = \tilde{l}_{21}\tilde{l}_{11} \\ \tilde{l}_{21} &= \frac{a_{21}}{\tilde{l}_{11}} = \frac{a_{21}}{\sqrt{a_{11}}} \\ &\vdots \\ a_{i1} &= \sum_{k=1}^n \tilde{l}_{ik}\tilde{l}_{1k} = \dots = \tilde{l}_{i1}\tilde{l}_{11} \\ \tilde{l}_{i1} &= \frac{a_{i1}}{\sqrt{a_{11}}} \end{aligned} \quad (2.104)$$

Ceci fixe la première colonne de manière unique.



- ii) Nous continuons avec le reste des colonnes, de façon récursive. Si on a déjà calculé les  $q$  premières colonnes, on calcule la colonne  $q + 1$ . Pour ça, on commence dans la ligne  $q + 1$ , dans la diagonale principale, car dans le résultat d'existence, on a montré que la matrice est diagonale inférieure. Ainsi

$$\begin{aligned} a_{q+1,q+1} &= \sum_{k=1}^n \tilde{l}_{q+1,k} \tilde{l}_{q+1,k} = \sum_{k=1}^{q+1} (\tilde{l}_{q+1,k})^2 = (\tilde{l}_{q+1,q+1})^2 + \sum_{k=1}^q (\tilde{l}_{q+1,k})^2 \\ \tilde{l}_{q+1,q+1} &= \left( a_{q+1,q+1} - \sum_{k=1}^q (\tilde{l}_{q+1,k})^2 \right)^{\frac{1}{2}} > 0 \end{aligned} \quad (2.105)$$

On note que l'expression sous la racine carrée est positive: sinon on serait en contradiction avec le résultat d'existence de la matrice  $\tilde{L}$ . On continue avec le reste de la colonne. Pour  $i \in \{q + 2, \dots, n\}$

$$\begin{aligned} a_{i,q+1} &= \sum_{k=1}^n \tilde{l}_{i,k} \tilde{l}_{q+1,k} = \sum_{k=1}^{q+1} \tilde{l}_{i,k} \tilde{l}_{q+1,k} = \tilde{l}_{i,q+1} \underbrace{\tilde{l}_{q+1,q+1}}_{>0} + \sum_{k=1}^q \tilde{l}_{i,k} \tilde{l}_{q+1,k} \\ \tilde{l}_{i,q+1} &= \frac{1}{\tilde{l}_{q+1,q+1}} \left( a_{i,q+1} - \sum_{k=1}^q \tilde{l}_{i,k} \tilde{l}_{q+1,k} \right) \end{aligned} \quad (2.106)$$

De cette manière, nous pouvons calculer toutes les colonnes de  $\tilde{L}$  et celle-ci est complètement fixée de manière unique. □

*Remarque 2.4* (Décomposition de Choleski: quelques points). Nous soulignons les points suivants:

- i) En général, toute matrice symétrique admet une décomposition unique  $PA = LU$ . Si  $A$  est symétrique définie positive, on peut choisir  $P = I_n$ , alors on a  $A = LU = \tilde{L}\tilde{L}^t$ , avec  $\tilde{L} = L\sqrt{D}$  où  $D = \text{diag}(U)$ .
- ii) Coût de la méthode Choleski: la méthode de Choleski (avec l'algorithme donné dans la partie d'unicité du théorème) est deux fois plus rapide que celui de la méthode de Gauss (décomposition  $LU$ ).
- iii) *Conservation du profil de  $A$* . Étant donnée la matrice symétrique  $A$  on peut définir pour chaque ligne  $i$  l'entier  $j_i$

$$j_i = \min \{j \in \{1, \dots, n\}\} \text{ tels que } a_{ij} \neq 0. \quad (2.107)$$

L'ensemble  $\{j_i\}_{i \in \{1, \dots, n\}}$  détermine le profil de la matrice. En particulier (pour  $A$  symétrique), seuls les éléments  $a_{ij}$  avec  $j \in \{j_i, \dots, i\}$  doivent être stockés. Une propriété importante de la décomposition de Choleski est qu'elle préserve le profil  $\{j_i\}_{i \in \{1, \dots, n\}}$  de la matrice  $A$ . Ceci peut être démontré à partir de la preuve d'unicité de la décomposition de Choleski en appliquant à chaque ligne  $i$  un raisonnement par l'absurde ou par récurrence sur les colonnes.

- iv) Si on étudie le problème  $A\mathbf{x} = \mathbf{b}$  avec  $A$  inversible mais non-symétrique, on peut utiliser la méthode de Choleski en notant que la matrice  $AA^t$  est symétrique définie positive:

$$A\mathbf{x} = \mathbf{b} \Leftrightarrow A^t(A\mathbf{x} = \mathbf{b}) \Leftrightarrow A^tA\mathbf{x} = A^t\mathbf{b} \quad (2.108)$$

La matrice  $A^tA$ , admet une décomposition de Choleski  $A^tA = \tilde{L}\tilde{L}$ .



## Chapitre 3

# Normes et conditionnement d'une matrice

Dans ce chapitre, nous introduisons des éléments de base autour les notions de norme matricielle et de conditionnement d'une matrice. Le but fondamental est la formulation de conditions suffisantes pour la convergence des suites définies par l'itération de l'action d'une matrice. Un objectif secondaire est la discussion de la majoration des erreurs d'arrondis dans la résolution du problème

$$A\mathbf{x} = \mathbf{b}, \quad (3.1)$$

en conséquence d'une erreur dans  $A$  ou dans  $\mathbf{b}$ .

### 3.1 Norme matricielle, norme matricielle induite et rayon spectral

#### 3.1.1 Rappel sur les normes

Nous commençons par rappeler la notion de norme dans un espace linéaire (sur  $\mathbb{C}$ ).

**Définition 3.1** (Norme). Étant donné un espace linéaire  $V$  sur  $\mathbb{C}$ , une norme est une application  $\|\cdot\| : V \rightarrow \mathbb{R}$  telle que

- i)  $\|\alpha\mathbf{x}\| = |\alpha|\|\mathbf{x}\|$ ,  $\forall \alpha \in \mathbb{C}$ ,  $\mathbf{x} \in V$ .
- ii)  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ ,  $\forall \mathbf{x}, \mathbf{y} \in V$
- iii)  $\|\mathbf{x}\| = 0 \Leftrightarrow \mathbf{x} = 0$ .

Nous allons travailler avec des espaces linéaires réels de dimension fini.

**Exemple 3.1** (Exemples de normes sur  $\mathbb{R}^n$ ). Nous donnons trois exemples des normes en  $\mathbb{R}^n$ . Étant donnée  $\mathbf{x} = (x_1, \dots, x_n)$  dans la base canonique, on définit :

- i) La norme  $\|\mathbf{x}\|_\infty$ :  $\|\mathbf{x}\|_\infty = \max\{|x_1|, \dots, |x_n|\}$ .
- ii) La norme  $\|\mathbf{x}\|_1$ :  $\|\mathbf{x}\|_1 = |x_1| + \dots + |x_n| = \sum_{i=1}^n |x_i|$ .

iii) La norme  $\|\mathbf{x}\|_2$ :  $\|\mathbf{x}\|_2 = (|x_1|^2 + \dots + |x_n|^2)^{\frac{1}{2}} = \left( \sum_{i=1}^n |x_i|^2 \right)^{\frac{1}{2}}$ .

iv) Les deux dernières sont cas particuliers de la norme  $\|\mathbf{x}\|_p$ , avec  $p \in \mathbb{R}, p \geq 1$ :

$$\|\mathbf{x}\|_p = (|x_1|^p + \dots + |x_n|^p)^{\frac{1}{p}} = \left( \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}.$$

**Définition 3.2** (Normes équivalentes). Deux normes  $\|\cdot\|_1$  et  $\|\cdot\|_2$  sur  $V$  sont équivalentes s'il existe deux constantes réelles  $C_1 > 0$  et  $C_2 > 0$  telles que

$$C_1 \|\mathbf{x}\|_1 \leq \|\mathbf{x}\|_2 \leq C_2 \|\mathbf{x}\|_2, \quad \forall \mathbf{x} \in V. \quad (3.2)$$

*Remarque 3.1.* Dans la suite, on va utiliser les deux propriétés suivantes :

- Deux normes équivalentes induisent la même topologie sur  $V$ .
- Si  $V$  est de dimension finie, toutes les normes sont équivalentes.

### 3.1.2 Normes matricielles

**Définition 3.3** (Norme matricielle). Étant donnée l'espace vectoriel  $M_n(\mathbb{R})$ :

i) On appelle *em norme matricielle* sur  $M_n(\mathbb{R})$  une norme  $\|\cdot\|$  sur  $\mathbb{R}$  telle que

$$\|AB\| \leq \|A\| \|B\|, \quad \forall A, B \in M_n(\mathbb{R}). \quad (3.3)$$

ii) Si  $\mathbb{R}^n$  est muni d'une norme  $\|\cdot\|$ , on appelle *norme matricielle induite (ou subordonnée)* sur  $M_n(\mathbb{R})$  par  $\|\cdot\|$  à la norme sur  $M_n(\mathbb{R})$  définie par

$$\|A\| = \sup_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|=1} \{\|A\mathbf{x}\|\}, \quad \forall A \in M_n(\mathbb{R}). \quad (3.4)$$

**Proposition 3.1** (propriété de la norme induite). . Soit  $M_n(\mathbb{R}^n)$  muni d'une norme induite  $\|\cdot\|$ . Alors, pour tout  $A \in M_n(\mathbb{R})$  on a:

i)  $\|A\mathbf{x}\| \leq \|A\| \|\mathbf{x}\|, \quad \forall \mathbf{x} \in \mathbb{R}^n$ .

ii)  $\|A\| = \max_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|=1} \{\|A\mathbf{x}\|\}$ .

iii)  $\|A\| = \max \left\{ \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|}; \mathbf{x} \in \mathbb{R}^n \setminus \{0\} \right\}$

iv)  $\|\cdot\|$  est une norme matricielle.

**Preuve:**

i) Soit  $\mathbf{x} \in \mathbb{R}^n, \mathbf{x} \neq 0$ . On prend  $\mathbf{y} = \frac{\mathbf{x}}{\|\mathbf{x}\|}$ , alors  $\|\mathbf{y}\| = 1$ . De la définition de norme induite on a  $A\mathbf{y} \leq \|A\| (= \sup\{\|A\mathbf{y}\|; \|\mathbf{y}\| = 1\})$ . Alors

$$\|A\| \geq \left\| A \frac{\mathbf{x}}{\|\mathbf{x}\|} \right\| = \frac{1}{\|\mathbf{x}\|} \|A\mathbf{x}\| \Leftrightarrow \|A\| \|\mathbf{x}\| \geq \|A\mathbf{x}\|, \quad \forall \mathbf{x} \in \mathbb{R}^n \setminus \{0\}. \quad (3.5)$$

Si  $\mathbf{x} = 0$ , il se suit  $A\mathbf{x} = 0$  et  $\|\mathbf{x}\| = 0$ , alors  $\|A\mathbf{x}\| \leq \|A\| \|\mathbf{x}\|$  est vérifié.

ii) On définit  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$ , avec  $\varphi(\mathbf{x}) = \|\mathbf{Ax}\|$ . L'application  $\varphi$  est continue sur  $S^1 = \{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\| = 1\}$ , que c'est un compact dans  $\mathbb{R}^n$ . Alors, l'image de  $S^1$  par  $\varphi$  est un compacte, donc  $\varphi$  est bornée et atteint ses bornes. On conclut qu'il existe  $\mathbf{x}_{\max}$ , tel que  $\|A\| = \|\mathbf{Ax}_{\max}\|$ .

iii) On prend  $\|\mathbf{x} \neq 0\|$ . Alors

$$\frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} = A \frac{\|\mathbf{x}\|}{\|\mathbf{x}\|}, \mathbf{x} \in S^1 \Rightarrow \max_{\mathbf{x} \in \mathbb{R}^n \setminus \{0\}} \left\{ \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} \right\} = \max_{\mathbf{x} \in S^1} \{\|\mathbf{Ax}\|\} = \|A\|. \quad (3.6)$$

iv) Étant donné  $A$  et  $B \in M_n(\mathbb{R})$ , on a  $\|AB\| = \max\{\|AB\mathbf{x}\| : \|\mathbf{x}\| = 1, \mathbf{x} \in \mathbb{R}^n\}$ . Or, en utilisant deux fois la propriété dans le point i):

$$\|AB\mathbf{x}\| \leq \|A\| \|\mathbf{Bx}\| \leq \|A\| \|B\| \|\mathbf{x}\| = \|A\| \|B\| \quad (3.7)$$

Alors

$$\|AB\| = \max\{\|AB\mathbf{x}\| : \|\mathbf{x}\| = 1, \mathbf{x} \in \mathbb{R}^n\} \leq \|A\| \|B\|. \quad (3.8)$$

On conclut que  $\|\cdot\|$  est une norme matricielle

□

**Définition 3.4** (Rayon spectral). Soit  $A \in M_n(\mathbb{R})$  inversible. On appelle *rayon spectral*  $\rho(A)$  de  $A$  à:

$$\rho(A) = \max\{|\lambda|; \lambda \in \mathbb{C}, \lambda \text{ valeur propre de } A\} \quad (3.9)$$

**Proposition 3.2** (Caractérisation de de normes induites; calculs effectifs)). Soit  $A = (a_{ij}) \in M_n(\mathbb{R})$ . Montrer:

i) On munit  $\mathbb{R}^n$  de la norme  $\|\cdot\|_\infty$  et  $M_n(\mathbb{R})$  de la norme induite correspondante; notée aussi  $\|\cdot\|_\infty$ . Alors

$$\|A\|_\infty = \max_{i \in \{1, \dots, n\}} \sum_{j=1}^n |a_{ij}|$$

ii) On munit  $\mathbb{R}^n$  de la norme  $\|\cdot\|_1$  et  $M_n(\mathbb{R})$  de la norme induite correspondante; notée aussi  $\|\cdot\|_1$ . Alors

$$\|A\|_1 = \max_{j \in \{1, \dots, n\}} \sum_{i=1}^n |a_{ij}|$$

iii) On munit  $\mathbb{R}^n$  de la norme  $\|\cdot\|_2$  et  $M_n(\mathbb{R})$  de la norme induite correspondante; notée aussi  $\|\cdot\|_2$ . Alors

$$\|A\|_2 = (\rho(A^t A))^{\frac{1}{2}}$$

En particulier; si  $A$  est symétrique,  $\|A\|_2 = \rho(A)$ .

**Preuve:** i) et ii) Exercice.

iii) Par définition

$$\begin{aligned} \|A\|_2^2 &= \left( \max_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|=1} \{\|A\mathbf{x}\|_2\} \right)^2 = \left( \max_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|=1} \{(A\mathbf{x} \cdot A\mathbf{x})^{\frac{1}{2}}\} \right)^2 \\ &= \max_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|=1} \{A\mathbf{x} \cdot A\mathbf{x}\} = \max_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|=1} \{A^t A \mathbf{x} \cdot \mathbf{x}\} \end{aligned} \quad (3.10)$$

La matrice  $A^t A$  est symétrique de “non-negative”. En effet

$$\mathbf{x}^t A^t A \mathbf{x} = \mathbf{x} \cdot A^t A \mathbf{x} = A \mathbf{x} \cdot A \mathbf{x} \geq 0 . \quad (3.11)$$

(Noter que  $A\mathbf{x}=0$  n'implique pas  $\mathbf{x}$  en générale, au moins que  $A$  soit inversible). Alors, on peut diagonaliser  $A^t A$  sur une base orthonormée  $B = \{\mathbf{e}_i, i \in \{1, \dots, n\}\}$  avec valeurs propres non-négatifs :

$$\begin{aligned} A^t A \mathbf{e}_i &= \lambda_i \mathbf{e}_i \quad , \quad \mathbf{e}_i \cdot \mathbf{e}_j \delta_{ij} \\ 0 &\leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n . \end{aligned} \quad (3.12)$$

On considère  $\mathbf{x} = \sum_{i=1}^n x^i \mathbf{e}_i$ . Alors

$$\begin{aligned} A^t A \mathbf{x} \cdot \mathbf{x} &= \left( A^t A \sum_i x^i \mathbf{e}_i \right) \cdot \left( \sum_j x^j \mathbf{e}_j \right) = \left( \sum_i x^i \lambda_i \mathbf{e}_i \right) \cdot \left( \sum_j x^j \mathbf{e}_j \right) \\ &= \sum_{i,j} \lambda_i x^i x^j \underbrace{\mathbf{e}_i \cdot \mathbf{e}_j}_{\delta_{ij}} = \sum_i \lambda_i (x^i)^2 \leq \sum_i \lambda_n (x^i)^2 = \lambda_n \underbrace{\sum_i (x^i)^2}_{\|\mathbf{x}\|_2^2} \end{aligned} \quad (3.13)$$

Avec ceci, on conclut

$$\left\langle A^t A \frac{\mathbf{x}}{\|\mathbf{x}\|}, \frac{\mathbf{x}}{\|\mathbf{x}\|} \right\rangle = \frac{A^t A \mathbf{x} \cdot \mathbf{x}}{\|\mathbf{x}\|^2} \leq \lambda_n = \max\{|\lambda|, \lambda \text{ valeur propre de } A^t A\} = \rho(A^t A) , \quad (3.14)$$

Avec, (3.10) on a

$$\|A\|_2^2 \leq \rho(A^t A) . \quad (3.15)$$

Pour l'égalité, on considère  $\mathbf{x} = \mathbf{e}_n$ ,  $\|\mathbf{e}_n\| = 1$ , et on a

$$\|A\mathbf{e}_n\|_2^2 = \langle A\mathbf{e}_n, A\mathbf{e}_n \rangle = \langle A^t A \mathbf{e}_n, \mathbf{e}_n \rangle = \langle \lambda_n \mathbf{e}_n, \mathbf{e}_n \rangle = \lambda_n = \rho(A^t A) . \quad (3.16)$$

Alors

$$\|A\|_2^2 = \rho(A^t A) , \quad (3.17)$$

qui montre le résultat. En particulier, si  $A$  est symétrique elle diagonalisable avec des valeurs propres  $\lambda_i^A$ , on a  $A^t A = A^2$  et  $\lambda_i = (\lambda_i^A)^2, \forall i \in \{1, \dots, n\}$ . En particulier, pour  $i = n$  on a  $\lambda_n = (|\lambda_n^A|)^2$ , où  $|\lambda_n^A| = \rho(A)$  et on conclut

$$\|A\|_2^2 = \rho(A^t A) = \lambda_n = (\rho(A))^2 \Rightarrow \|A\|_2 = \rho(A) . \quad (3.18)$$

□

## 3.2 Relation entre norme matricielle et rayon spectral

**Théorème 3.1** (Approximation spectral du rayon spectral par une norme induite). *On considère  $\mathbb{R}^n$  et  $M_n(\mathbb{R})$ :*

i) *Soit  $\|\cdot\|$  une norme induite en  $M_n(\mathbb{R})$ . Alors*

$$\rho(A) \leq \|A\|, \quad \forall A \in M_n(\mathbb{R}). \quad (3.19)$$

ii) *Étant donné  $A \in M_n(\mathbb{R}^n)$  et  $\varepsilon > 0$ , il existe une norme sur  $\mathbb{R}^n$  (qui dépend de  $A$  et de  $\varepsilon$ ) telle que la correspondante norme induite sur  $M_n(\mathbb{R})$ , notée  $\|\cdot\|_{A,\varepsilon}$ , vérifie*

$$\|A\|_{A,\varepsilon} \leq \rho(A) + \varepsilon, \quad \forall A \in M_n(\mathbb{R}). \quad (3.20)$$

**Preuve:** Pour i), étant donnée  $A$  on considère  $\lambda \in \mathbb{C}$  valeur propre de  $A$ , c'est-à-dire,  $\exists \mathbf{x} \neq 0$  tel que  $A\mathbf{x} = \lambda\mathbf{x}$ . Alors

$$\|A\mathbf{x}\| = \|\lambda\mathbf{x}\| = |\lambda|\|\mathbf{x}\|, \quad \|A\mathbf{x}\| \leq \|A\|\|\mathbf{x}\|, \quad (3.21)$$

donc  $|\lambda| \leq \|A\|$  et ceci  $\forall \lambda$  valeur propre de  $A$ . Alors  $\rho(A) \leq \|A\|$ .

On laisse ii) comme exercice. □

*Remarque 3.2.* Après le théorème précédent, pour une matrice  $A \in M_n(\mathbb{R})$  donnée, on peut considérer  $\rho(A)$  comme le infimum sur l'ensemble  $\{\|A\|;$  avec  $\|\cdot\|$  norme induite sur  $M_n(\mathbb{R})\}$ .

### 3.2.1 Convergence et rayon spectral

**Corollaire 3.1** (Convergence et rayon spectral). *Soit  $A \in M_n(\mathbb{R})$ . A condition nécessaire et suffisante pour la convergence à zéro de la suite de matrices  $\{A_k\}_{k \in \mathbb{N}}$  est:*

$$\rho(A) < 1 \Leftrightarrow A^k \rightarrow 0. \quad (3.22)$$

**Preuve:** On montre d'abord la direction  $\Rightarrow$  (condition suffisante). Si  $\rho(A) < 1$ , on peut toujours choisir  $\varepsilon < 0$  tel que  $\rho(A) < 1 - 2\varepsilon$ . Par le théorème 3.1 d'approximation du rayon spectral, il existe une norme induite  $\|\cdot\|_{A,\varepsilon}$  telle que:

$$\|A\|_{A,\varepsilon} \equiv \mu \leq \rho(A) + \varepsilon < 1 - 2\varepsilon + \varepsilon = 1 - \varepsilon < 1. \quad (3.23)$$

Étant donné que la norme induite est une norme matricielle, on a

$$\|A^k\|_{A,\varepsilon} \leq \|A\|_{A,\varepsilon} \|A^{k-1}\|_{A,\varepsilon} \leq \dots \leq \|A\|_{A,\varepsilon}^k = \mu^k \rightarrow 0 \text{ (quand } k \rightarrow \infty \text{)}. \quad (3.24)$$

Notons que dans l'espace  $M_n(\mathbb{R})$ , de dimension finie, toutes les normes sont équivalentes (gènèrent la même topologie, alors la même convergence) alors, pour n'importe quelle norme  $\|A^k\| \rightarrow 0$  ( $k \rightarrow \infty$ )  $\Rightarrow A^k \rightarrow 0$ .

Pour la condition nécessaire ( $\Leftarrow$ ), supposons  $A^k \rightarrow 0$  ( $k \rightarrow \infty$ ). On considère  $\lambda$  une valeur propre de  $A$  et  $\mathbf{x} \neq 0$  un vecteur propre associé. Alors

$$A^k \mathbf{x} = \lambda^k \mathbf{x}. \quad (3.25)$$

Si  $A^k \rightarrow 0$ , on a  $A^k \mathbf{x} \rightarrow 0$  et alors  $\lambda^k \mathbf{x}$ . Vu que  $\mathbf{x} \neq 0$ ,  $\lambda^k \rightarrow 0$ , ce qui est vrai si et seulement si  $|\lambda| < 1$ . Ceci est vrai pour n'importe quel valeur propre  $\lambda$ . Alors

$$\rho(A) = \max\{|\lambda|, \lambda \text{ valeur propre de } A\} < 1. \quad (3.26)$$

□

*Remarque 3.3* (Convergence de suites de vecteurs). Du corollaire précédent on déduit que la suite  $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$  définie par  $\mathbf{x}^{(0)} = \mathbf{x}_o$  et  $\mathbf{x}^{(k+1)} = A\mathbf{x}^{(k)}$ , on a

$$\mathbf{x}^k \rightarrow 0, \forall \mathbf{x}^{(0)} = \mathbf{x}_o \Leftrightarrow \rho(A) < 1. \quad (3.27)$$

Cette caractérisation sera utilisé dans l'étude des méthodes itératifs pour la résolution d'un problème linéaire.

Nous venons de voir que pour montrer la convergence d'une suite de vecteurs  $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$  donnée par  $\mathbf{x}^{(k+1)} = A\mathbf{x}^{(k)}$ , il faut contrôler la valeur du rayon spectrale. Le calcul de  $\rho(A)$  peut être difficile mais, après le théorème 3.1, il suffit de trouver une norme induite  $\|\cdot\|$  pour laquelle  $\|A\| < 1$ . Nous allons donner maintenant une condition suffisante plus générale, en montrant qu'on peut relaxer la condition sur le caractère *induit* de la norme matricielle employée. Pour montrer ça nous donnons d'abord la caractérisation suivante du rayon spectrale (valable aussi en dimension infinie).

**Proposition 3.3** (Caractérisation du rayon spectrale). *On munit  $M_n(\mathbb{R})$  d'une norme générale, notée  $\|\cdot\|$ . Soit  $A \in M_n(\mathbb{R})$ . Alors*

$$\rho(A) = \lim_{k \rightarrow \infty} \|A^k\|^{\frac{1}{k}}. \quad (3.28)$$

**Preuve:** Exercice.

*Remarque 3.4.* Noter que la norme dans la proposition est une norme générale, pas nécessairement matricielle.

Nous pouvons maintenant prouver le résultat annoncé.

**Corollaire 3.2** (Comparaison rayon spectrale et norme matricielle). *On munit  $M_n(\mathbb{R})$  d'une norme matricielle générale, notée  $\|\cdot\|$ . Soit  $A \in M_n(\mathbb{R})$ . Alors*

$$\rho(A) \leq \|A\| \quad (3.29)$$

**Preuve:** Si  $\|\cdot\|$  est matricielle, on a  $\|A^k\| \leq \|A\|^k$  et donc, en utilisant la caractérisation (3.28) du rayon spectral, on a

$$\rho(A) = \lim_{k \rightarrow \infty} \|A^k\|^{\frac{1}{k}} \leq \lim_{k \rightarrow \infty} \left(\|A\|^k\right)^{\frac{1}{k}} = \|A\| \quad (3.30)$$

.

□

*Remarque 3.5.* C'est fondamental que la norme  $\|\cdot\|$  soit matricielle. Ça ne marche pas pour une norme générique non-matricielle (exercice !).

Nous finissons cette section en énonçant un résultat sur les matrices "perturbation" de l'identité.

**Théorème 3.2** (Matrices de la forme  $I_n + A$ ). *Soit  $\|\cdot\|$  une norme matricielle induite sur  $M_n(\mathbb{R})$ .*

*i) Étant donné  $I_n \in M_n(\mathbb{R})$  la matrice identité et  $A \in M_n(\mathbb{R})$  telle que  $\|A\| < 1$ , on a que la matrice  $I_n + A$  est inversible et*

$$\|(I_n + A)^{-1}\| \leq \frac{1}{1 - \|A\|}. \quad (3.31)$$



ii) Si une matrice de la forme  $I_n + A \in M_n(\mathbb{R})$  est singulière, alors  $\|A\| \geq 1$  pour toute norme matricielle.

**Preuve:**

i) Exercice.

ii) Si  $I_n + A \in M_n(\mathbb{R})$  est singulière, alors  $\exists \mathbf{x} \neq 0$  tel que

$$(I_n + A)\mathbf{x} = 0 \Leftrightarrow A\mathbf{x} = -\mathbf{x} \Rightarrow \lambda = -1 \text{ est une valeur propre.} \quad (3.32)$$

Alors  $\rho(A) \geq 1$ . Si on utilise (3.29), étant donnée que  $\|\cdot\|$  est matricielle, on a

$$\|A\| \geq \rho(A) \geq 1. \quad (3.33)$$

□

### 3.3 Conditionnement d'une matrice

Quand on étudie le problème

$$A\mathbf{x} = \mathbf{b}, \quad (3.34)$$

en particulier sa résolution numérique par un ordinateur, un problème crucial à considérer est posé par la question suivante : quel est le changement de  $\mathbf{x}$  si on fait une petite erreur sur  $A$  ou  $\mathbf{b}$  par rapport à ses vraies valeurs ?

Si un petit changement  $\delta A$  dans  $A$  ou une petite erreur  $\delta \mathbf{b}$  entraînent un grand changement  $\delta \mathbf{x}$  dans la valeur de  $\mathbf{x}$ , alors le problème est très mal posé pour être résolu par un ordinateur, étant donné que la résolution dans  $A$  et  $\mathbf{b}$  par l'ordinateur sera toujours finie.

D'une manière plus précise, on veut estimer l'erreur  $\delta \mathbf{x}$  en  $\mathbf{x}$ , si on fait des erreurs  $\delta A$  en  $A$  et  $\delta \mathbf{b}$  en  $\mathbf{b}$ , en résolvant l'équation

$$(A + \delta A)(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} + \delta \mathbf{b}. \quad (3.35)$$

Si  $\delta A$  est suffisamment petite, la matrice  $A + \delta A$  est inversible et le problème posé par (3.51) a du sens. D'une manière plus précise, si  $\|A^{-1}\delta A\| < 1$ , le théorème 3.2 garantit que la matrice  $I_n + A^{-1}\delta A$  est inversible, donc  $A + \delta A = A(I_n + A^{-1}\delta A)$  est aussi inversible. Dans ces conditions on peut estimer  $\delta \mathbf{x}$  dans (3.51) en fonction de  $\delta A$  et  $\delta \mathbf{b}$ .

#### 3.3.1 Conditionnement et ses propriétés élémentaires

**Définition 3.5.** Soit  $\mathbb{R}^n$  muni d'une norme  $\|\cdot\|$  et  $M_n(\mathbb{R})$  muni de la norme matricielle induite. Étant donnée  $A \in M_n(\mathbb{R})$  inversible, on appelle conditionnement de  $A$  par rapport à la norme  $\|\cdot\|$  le nombre réel positif

$$\text{cond}(A) = \|A\| \|A^{-1}\|. \quad (3.36)$$

**Proposition 3.4** (propriétés générales du conditionnement). . Soit  $\mathbb{R}^n$  muni d'une norme  $\|\cdot\|$  et  $M_n(\mathbb{R})$  muni de la norme matricielle induite, et  $A, B \in M_n(\mathbb{R})$  inversibles et  $\alpha \in \mathbb{R}^*$ . On a :

i)  $\text{cond}(A) \geq 1$ .

$$ii) \text{ cond}(\alpha A) = \text{cond}(A).$$

$$iii) \text{ cond}(AB) \leq \text{cond}(A) \text{ cond}(B).$$

**Preuve:**

i) On utilise que  $\|\cdot\|$  est matricielle. Alors

$$\|I_n\| = \|AA^{-1}\| \leq \|A\| \|A^{-1}\| = \text{cond}(A). \quad (3.37)$$

D'autre part

$$\|I_n\| = \sup_{\|\mathbf{x}\|=1} \{\|I_n \mathbf{x}\|\} = 1 \quad (3.38)$$

Alors  $1 \leq \text{cond}(A)$  (noter qu'on a utilisé que la norme est induite, et pas seulement matricielle).

ii) En utilisant la définition (3.36)

$$\text{cond}(\alpha A) = \|\alpha A\| \|(\alpha A)^{-1}\| = |\alpha| \|A\| \cdot \frac{1}{|\alpha|} \|A^{-1}\| = \|A\| \|A^{-1}\| = \text{cond}(A) \quad (3.39)$$

iii) Étant donné que  $A$  et  $B$  sont inversibles, alors  $AB$  est aussi inversible et, utilisant le fait que  $\|\cdot\|$  est matricielle, on a

$$\begin{aligned} \text{cond}(AB) &= \|AB\| \|(AB)^{-1}\| = \|AB\| \|B^{-1}A^{-1}\| \\ &\leq \|A\| \|B\| \|B^{-1}\| \|A^{-1}\| = \text{cond}(A) \text{ cond}(B). \end{aligned} \quad (3.40)$$

□

### 3.3.2 Majoration des erreurs

Nous allons étudier le cas d'une perturbation seulement par  $\delta A$  et le cas d'une perturbation seulement par  $\delta \mathbf{b}$ . Le cas plus général d'une perturbation simultanée par  $\delta A$  et  $\delta \mathbf{b}$  suit une stratégie pareille, mais avec une expression plus compliquée.

**Proposition 3.5** (majoration de l'erreur relative par rapport à une erreur dans le deuxième membre). *Soit  $A \in M_n(\mathbb{R})$  inversible et  $\mathbf{b} \in \mathbb{R}^n$ , avec  $\mathbf{b} \neq 0$ . On munit  $\mathbb{R}^n$  d'une norme  $\|\cdot\|$  et  $M_n(\mathbb{R})$  de la norme matricielle induite associée. Étant donné  $\delta \mathbf{b} \in \mathbb{R}^n$ , si  $\mathbf{x}$  est solution de*

$$A\mathbf{x} = \mathbf{b}, \quad \mathbf{b} \neq 0, \quad (3.41)$$

et  $\mathbf{x} + \delta \mathbf{x}$  est solution de

$$A(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} + \delta \mathbf{b}, \quad (3.42)$$

on a la majoration suivante de l'erreur relative en  $\mathbf{x}$

$$\frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} \leq \text{cond}(A) \frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|} \quad (3.43)$$

**Preuve:** On considère

$$\begin{aligned} A\mathbf{x} &= \mathbf{b} \\ A(\mathbf{x} + \delta\mathbf{x}) &= \mathbf{b} + \delta\mathbf{b}. \end{aligned} \quad (3.44)$$

Si on retranche la première à la deuxième, on a :

$$A\delta\mathbf{x} = \delta\mathbf{b} \Leftrightarrow \delta\mathbf{x} = A^{-1}\delta\mathbf{b}, \quad (3.45)$$

donc

$$\|\delta\mathbf{x}\| \leq \|A^{-1}\| \|\delta\mathbf{b}\|. \quad (3.46)$$

D'autre part, de  $\mathbf{b} = A\mathbf{x}$ , on a

$$\|\mathbf{b}\| = \|A\mathbf{x}\| \leq \|A\| \|\mathbf{x}\|. \quad (3.47)$$

Si on utilise que  $\mathbf{b} \neq 0$ , on a  $\mathbf{x} \neq 0$ , donc on peut écrire

$$\frac{1}{\|\mathbf{x}\|} \leq \frac{\|A\|}{\|\mathbf{b}\|}. \quad (3.48)$$

Si on multiplie, membre à membre, cette expression par celle dans (3.46) on obtient

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \underbrace{\|A\| \|A^{-1}\|}_{\text{cond}(A)} \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} = \text{cond}(A) \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|}. \quad (3.49)$$

□

*Remarque 3.6.* L'inégalité (3.52) est optimale, c'est à dire, il y a des choix que la rendre une égalité.

**Proposition 3.6** (majoration de l'erreur relative par rapport à une erreur sur la matrice). *Soit  $A \in M_n(\mathbb{R})$  inversible et  $\mathbf{b} \in \mathbb{R}^n$ , avec  $\mathbf{b} \neq 0$ . On munit  $\mathbb{R}^n$  d'une norme  $\|\cdot\|$  et  $M_n(\mathbb{R})$  de la norme matricielle induite associée. Étant donné  $\delta A \in M_n(\mathbb{R})$ , on suppose que  $A + \delta A$  est inversible. Si  $\mathbf{x}$  est solution de*

$$A\mathbf{x} = \mathbf{b}, \quad \mathbf{b} \neq 0, \quad (3.50)$$

et  $\mathbf{x} + \delta\mathbf{x}$  est solution de

$$(A + \delta A)(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b}, \quad (3.51)$$

alors

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x} + \delta\mathbf{x}\|} \leq \text{cond}(A) \frac{\|\delta A\|}{\|A\|}. \quad (3.52)$$

**Preuve:** On part de

$$(A + \delta A)(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} . \quad (3.53)$$

En le développant

$$A\mathbf{x} + A\delta\mathbf{x} + \delta A\mathbf{x} + \delta A\delta\mathbf{x} = \mathbf{b} . \quad (3.54)$$

Si on retranche  $A\mathbf{x} = \mathbf{b}$  on trouve

$$A\delta\mathbf{x} + \delta A\mathbf{x} + \delta A\delta\mathbf{x} = 0 , \quad (3.55)$$

c'est à dire

$$\begin{aligned} A\delta\mathbf{x} &= -A\delta A(\mathbf{x} + \delta\mathbf{x}) \\ \delta\mathbf{x} &= A^{-1}\delta A(\mathbf{x} + \delta\mathbf{x}) . \end{aligned} \quad (3.56)$$

Alors,

$$\|\delta\mathbf{x}\| \leq \|A^{-1}\| \|\delta A\| \|\mathbf{x} + \delta\mathbf{x}\| . \quad (3.57)$$

Si on utilise maintenant que  $\mathbf{b} \neq 0$  et  $A + \delta A$  est inversible, on conclut  $\mathbf{x} + \delta\mathbf{x} \neq 0$  et on divise par  $\|\mathbf{x} + \delta\mathbf{x}\|$  pour trouver

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x} + \delta\mathbf{x}\|} \leq \|A^{-1}\| \|\delta A\| . \quad (3.58)$$

Si on multiplie maintenant le deuxième par  $\frac{\|A\|}{\|A\|} = 1$

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x} + \delta\mathbf{x}\|} \leq \underbrace{\|A^{-1}\| \|A\|}_{\text{cond}(A)} \frac{\|\delta A\|}{\|A\|} \leq \text{cond}(A) \frac{\|\delta A\|}{\|A\|} . \quad (3.59)$$

□

## Chapitre 4

# Méthodes itératives



## Chapitre 5

# Méthodes de descente





Part II

Interpolation



## Chapitre 6

# Interpolation polynomiale



Part III  
**Annexes**



# Appendix A

## Espaces linéaires

### A.1 Rappel sur les espaces linéaires

**Définition A.1** (*espace linéaire ou vectoriel*). Un espace linéaire sur un corps  $\mathbb{K}$  (ici  $\mathbb{K} = \mathbb{R}$  ou  $\mathbb{K} = \mathbb{C}$ ) est un ensemble non vide  $V$ , dont les éléments sont appelés vecteurs et où on a deux opérations définies l'*addition* et la *multiplication par un scalaire*, avec les propriétés:

1. L'addition est commutative et associative.
2. Il existe un élément  $\mathbf{0}$  (*vecteur zéro ou nul*) tel que  $\mathbf{v} + \mathbf{0} = \mathbf{v}$ ,  $\forall \mathbf{v} \in V$ .
3.  $0 \cdot \mathbf{v} = \mathbf{0}$ ,  $1 \cdot \mathbf{v} = \mathbf{v}$ ,  $\forall \mathbf{v} \in V$ , où 0 et 1 sont le zéro et l'unité de  $\mathbb{K}$ .
4. Pour chaque  $\mathbf{v} \in V$  il existe un unique opposé,  $-\mathbf{v} \in V$ , tel que  $\mathbf{v} - \mathbf{v} = \mathbf{0}$ .
5. Propriété distributive:

$$\alpha(\mathbf{v} + \mathbf{w}) = \alpha\mathbf{v} + \alpha\mathbf{w}, \quad \forall \alpha \in \mathbb{K}, \forall \mathbf{v}, \mathbf{w} \in V \quad (\text{A.1})$$

$$(\alpha + \beta)\mathbf{v} = \alpha\mathbf{v} + \beta\mathbf{v}, \quad \forall \alpha, \beta \in \mathbb{K}, \forall \mathbf{v} \in V \quad (\text{A.2})$$

6. Propriété associative:

$$(\alpha\beta)\mathbf{v} = \alpha(\beta\mathbf{v}) \quad \forall \alpha, \beta \in \mathbb{K}, \forall \mathbf{v} \in V \quad (\text{A.3})$$

**Exemple A.1.** On donne plusieurs exemples d'espace linéaire.

- i)  $V = \mathbb{R}^n$ ,  $V = \mathbb{R}^n$ ,  $M_{m \times n}(\mathbb{R}) \equiv \mathbb{R}^{m \times n}$ .
- ii)  $V = \mathbb{P}_n$ , ensemble de polynômes  $p_n(x) = \sum_{i=0}^n a_i x^i$
- iii)  $V = C^p([a, b])$ , ensemble des fonctions réelles (complexes) continues jusqu'à la dérivée p-ième.

**Définition A.2** (*espace linéaire ou vectoriel*). Un sous-espace non vide  $W$  de  $V$  est un sous-espace linéaire de  $V$  si  $W$  est un espace linéaire sur  $\mathbb{K}$ .

**Exemple A.2.**  $V = \mathbb{P}_n$  est un sous espace linéaire de  $V = C^\infty(\mathbb{R})$ .

**Définition A.3** (système de générateurs). Étant donné un ensemble  $\{\mathbf{v}_1, \dots, \mathbf{v}_p\}$ , le sous-espace

$$W = \text{Lin}\{\mathbf{v}_1, \dots, \mathbf{v}_p\} = \langle \mathbf{v}_1, \dots, \mathbf{v}_p \rangle = \{\mathbf{w} = \alpha_1 \mathbf{v}_1 + \dots + \alpha_p \mathbf{v}_p, \text{ avec } \alpha_i \in \mathbb{K}\} \quad (\text{A.4})$$

est un espace linéaire et  $\{\mathbf{v}_1, \dots, \mathbf{v}_p\}$  est appelé *système de générateurs* de  $W$ . Si  $W_1, \dots, W_m$  sont des sous-espaces linéaires de  $V$ , et  $W_i \cap W_j = \{\mathbf{0}\} \forall i \neq j$ , l'ensemble

$$S = \{\mathbf{w} \in V, \text{ tel que } \mathbf{w} = \mathbf{w}_1 + \dots + \mathbf{w}_m, \text{ avec } \mathbf{w}_i \in W_i\} \quad (\text{A.5})$$

est un espace linéaire appelé *somme directe* et on note  $S = W_1 \oplus \dots \oplus W_m$ .

**Définition A.4** (Indépendance linéaire). Un système de vecteurs  $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$  es dit *linéairement indépendant* si

$$\alpha_1 \mathbf{v}_1 + \dots + \alpha_m \mathbf{v}_m = \mathbf{0}, \quad \alpha_i \in \mathbb{K} \quad (\text{A.6})$$

implique  $\alpha_1 = \dots = \alpha_m = 0$ . Dans le cas contraire, le système est dit *linéairement dépendant*.

**Définition A.5** (base). Une base de  $V$  est un système linéairement indépendant de générateurs de  $V$ .

**Propriétés** (base). Les bases d'un espace linéaire satisfont:

- i) Si  $B = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  (avec  $n < \infty$ ) est une base de  $V$ , l'expression de  $\mathbf{v} \in V$   $\mathbf{v} = v_1 \mathbf{e}_1 + \dots + v_n \mathbf{e}_n$  est appelée la décomposition de  $\mathbf{v}$  en  $B$ . Les *composants*  $v_i$  sont uniques.
- ii) Dans le cas où  $n < \infty$ , chaque système de vecteurs linéairement indépendants a un maximum de  $n$  vecteurs. En particulier, toutes les bases de  $V$  ont le même nombre des vecteurs  $n < \infty$ , appelé dimension de  $V$ :  $\dim(V) = n$ .
- iii) Si pour chaque  $n$  il y a un système de vecteurs linéairement indépendant, l'espace linéaire  $V$  est dit de dimension infinie. Ici, on va se focaliser sur des  $V$  de dimension finie.

## A.2 Application linéaire

**Définition A.6** (application linéaire). Étant donnés deux espaces linéaires  $V$  et  $W$ , une application  $\mathcal{A} : V \rightarrow W$  est dite linéaire si elle respecte la structure de espace linéaire, ceci signifie que

$$\mathcal{A}(\alpha \mathbf{v} + \beta \mathbf{w}) = \alpha \mathcal{A}(\mathbf{v}) + \beta \mathcal{A}(\mathbf{w}), \quad \mathbf{v}, \mathbf{w} \in V, \alpha, \beta \in \mathbb{K}. \quad (\text{A.7})$$

## A.3 Produit scalaire

**Définition** (produit scalaire). Un produit scalaire est une application bilinéaire  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$ , symétrique et positivement définie. Cela veut dire

i) Application bilinéaire:

$$\begin{aligned} \langle \alpha \mathbf{v} + \beta \mathbf{w}, \mathbf{x} \rangle &= \alpha \langle \mathbf{v}, \mathbf{x} \rangle + \beta \langle \mathbf{w}, \mathbf{x} \rangle \\ \langle \mathbf{x}, \alpha \mathbf{v} + \beta \mathbf{w} \rangle &= \alpha \langle \mathbf{x}, \mathbf{v} \rangle + \beta \langle \mathbf{x}, \mathbf{w} \rangle, \quad \mathbf{v}, \mathbf{w}, \mathbf{x} \in V, \alpha, \beta \in \mathbb{K} \end{aligned} \quad (\text{A.8})$$



ii) Symétrique:

$$\langle \mathbf{v}, \mathbf{w} \rangle = \langle \mathbf{w}, \mathbf{v} \rangle, \quad \mathbf{v}, \mathbf{w} \in V \quad (\text{A.9})$$

iii) Définie positive:

$$\langle \mathbf{v}, \mathbf{v} \rangle \geq 0 \quad \text{et} \quad \langle \mathbf{v}, \mathbf{v} \rangle = 0 \quad \text{sii} \quad \mathbf{v} = \mathbf{0}. \quad (\text{A.10})$$

**Définition** (produit scalaire hermitien). Un produit scalaire est une application  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{C}$  qui satisfait:

i) Linéaire dans la deuxième variable et anti-linéaire dans la première variable:

$$\begin{aligned} \langle \alpha \mathbf{v} + \beta \mathbf{w}, \mathbf{x} \rangle &= \bar{\alpha} \langle \mathbf{v}, \mathbf{x} \rangle + \bar{\beta} \langle \mathbf{w}, \mathbf{x} \rangle \\ \langle \mathbf{x}, \alpha \mathbf{v} + \beta \mathbf{w} \rangle &= \alpha \langle \mathbf{x}, \mathbf{v} \rangle + \beta \langle \mathbf{x}, \mathbf{w} \rangle, \quad \mathbf{v}, \mathbf{w}, \mathbf{x} \in V, \alpha, \beta \in \mathbb{K} \end{aligned} \quad (\text{A.11})$$

ii) Symétrique hermitienne:

$$\langle \mathbf{v}, \mathbf{w} \rangle = \overline{\langle \mathbf{w}, \mathbf{v} \rangle}, \quad \mathbf{v}, \mathbf{w} \in V \quad (\text{A.12})$$

iii) Définie positive:

$$\langle \mathbf{v}, \mathbf{v} \rangle \geq 0 \quad \text{et} \quad \langle \mathbf{v}, \mathbf{v} \rangle = 0 \quad \text{sii} \quad \mathbf{v} = \mathbf{0}. \quad (\text{A.13})$$



# Appendix B

## Matrices

### B.1 Définitions

**Définition B.1** (matrice). Étant donnés  $m$  et  $n$  deux entiers positifs, on appelle matrice  $A$  à  $m$  lignes et  $n$  colonnes l'ensemble ordonné de  $m \cdot n$  éléments de  $\mathbb{K}$

$$A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix} \quad (\text{B.1})$$

On note par  $M_{m,n}(\mathbb{K})$  l'ensemble de matrices de  $m$  lignes et  $n$  colonnes, avec des coefficients dans  $\mathbb{K}$ .

**Remarques :**

- i) Notation:  $A = (A_{ij})$  avec  $i = 1, \dots, m, j = 1, \dots, n$ .
- ii) Matrice carrée:  $m = n$ . Diagonale principale:  $\text{diag}(a_{11}, \dots, a_{nn})$ .
- iii) Une matrice avec une ligne est un *vecteur ligne* et une matrice avec une seule colonne est un *vecteur colonne*

$$\mathbf{v} = \begin{pmatrix} v^1 \\ \vdots \\ v^n \end{pmatrix}, \quad \boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n), \quad \boldsymbol{\alpha}^t = (\alpha_1 \quad \dots \quad \alpha_n) \quad (\text{B.2})$$

iv)

$$\text{Ligne-}i : \boldsymbol{\ell}_i = (A_{i1}, \dots, A_{in}), \quad \text{Colonne-}j : \mathbf{c}_j = \begin{pmatrix} A_{1j} \\ \vdots \\ A_{mj} \end{pmatrix} \quad (\text{B.3})$$

iv)  $M_{m,n}(\mathbb{K}) = \{(A_{ij}), a_{ij} \in \mathbb{K}; i = 1, \dots, m, j = 1, \dots, n\}$

### B.1.1 Opérations sur les matrices

i) Somme. Pour  $A \in M_{m \times n}(\mathbb{K})$ ,  $B \in M_{m \times n}(\mathbb{K})$ :

$$C = A + B \quad , \quad (c_{ij}) = (a_{ij}) + (b_{ij}) \quad (\text{B.4})$$

ii) Produit par un scalaire. Pour  $A \in M_{m \times n}(\mathbb{K})$  et  $\lambda \in \mathbb{K}$ :

$$C = \lambda A \quad , \quad (c_{ij}) = (\lambda a_{ij}) \quad (\text{B.5})$$

iii) Produit de matrices. Pour  $A \in M_{m \times n}(\mathbb{K})$ ,  $B \in M_{n \times p}(\mathbb{K})$ , on définit la matrice  $A \in M_{m \times p}(\mathbb{K})$

$$C = A \cdot B \quad , \quad (c_{ij}) = \left( \sum_{k=1}^n a_{ik} b_{kj} \right) \quad (\text{B.6})$$

**Exercice B.1.** Montrer que le produit de matrices est associatif est distributif par rapport à l'addition.

*Remarque B.1.* :

i) Pour des matrices carrées d'ordre  $n$ , la matrice  $I_n = (\delta_{ij})$  est la *matrice unité*, ce qui veut dire  $A \cdot I_n = I_n \cdot A = A$ .

ii) Matrice diagonale:  $D = (d_{ii} \delta_{ij})$ .

iii)  $A^p = \underbrace{A \cdot \dots \cdot A}_{p \text{ fois}}$

iv) **Définition** (matrice transposée). Étant donnée une matrice  $A \in M_{m \times n}(\mathbb{R})$ , on appelle *matrice transposée* de  $A$  la matrice  $A^t \in M_{n \times m}(\mathbb{R})$  telle que

$$(a^t_{ij}) = (a_{ji}) \quad (\text{B.7})$$

Propriétés:

$$(A^t)^t = A \quad (\text{B.8})$$

$$(A \cdot B)^t = B^t \cdot A^t \quad (\text{B.9})$$

**Définition** (matrice conjuguée transposée). Étant donnée une matrice  $A \in M_{m \times n}(\mathbb{C})$ , on appelle *matrice conjuguée transposée* de  $A$  (ou *adjointe*, voir ci-après) la matrice  $A^\dagger \in M_{n \times m}(\mathbb{C})$  telle que

$$(a^\dagger_{ij}) = (\bar{a}_{ji}) \quad (\text{B.10})$$

où  $\bar{z}$  le complexe conjugué de  $z \in \mathbb{C}$  (l'adjointe de  $A$  est souvent notée  $A^*$ .) Propriétés:

$$(A^\dagger)^\dagger = A \quad (\text{B.11})$$

$$(A \cdot B)^\dagger = B^\dagger \cdot A^\dagger \quad (\text{B.12})$$

- v) Un vecteur  $\mathbf{v} \in \mathbb{K}^n$  sera considéré comme un vecteur colonne. On peut définir un produit scalaire naturel en  $\mathbb{R}^n$ <sup>1</sup>

$$\langle \mathbf{v}, \mathbf{w} \rangle = \mathbf{v} \cdot \mathbf{w} := \mathbf{v}^t \mathbf{w} = \sum_i^n v^i w^i . \quad (\text{B.13})$$

Un vecteur ligne sera considéré comme un élément de l'espace dual  $\alpha \in (\mathbb{K}^n)^*$ , avec la relation de dualité

$$\alpha(\mathbf{v}) = \alpha \mathbf{v} = \alpha^t \cdot \mathbf{v} \quad (\text{B.14})$$

- vi) L'expression (B.6) peut être écrite  $c_{ij} = \ell_i \mathbf{c}_j$ .
- vii) Les colonnes de  $AB$  sont des combinaison linéaires des colonnes de  $A$ .  
Les lignes de  $AB$  sont de combinaisons linéaires de lignes de  $B$ .

$$\begin{aligned} \ell_i(AB) &= \sum_{k=1}^n a_{ik} \ell_k(B) \\ \mathbf{c}_j(AB) &= \sum_{k=1}^n b_{kj} \mathbf{c}_k(A) \end{aligned}$$

Rélevant pour la méthode d'élimination de Gauss.

*Exemple B.1.*

### B.1.2 Rang et noyau d'une matrice

**Définition B.2** (rang). Étant donnée  $A \in M_{m,n}$ , l'image de  $A$  (notée  $\text{Im}(A)$ ) est donnée par  $\text{Im}(A) = \{\mathbf{y} \in \mathbb{K}^m; \mathbf{y} = A\mathbf{x}, \mathbf{x} \in \mathbb{K}^n\}$ . Le rang de  $A$  est défini par  $\text{rang}(A) = \dim(\text{Im}(A))$ .

**Définition B.3** (Noyau). Étant donnée  $A \in M_{m,n}$ , le noyau de  $A$  (notée par  $\text{Ker}(A)$ ) est donné par  $\text{Ker}(A) = \{\mathbf{x} \in \mathbb{K}^n; A\mathbf{x} = 0\}$ .

**Proposition B.1.** Étant donnée  $A \in M_{m,n}$ , les relations suivantes sont satisfaites

- i)  $\text{rang}(A) = \text{rang}(A^t)$ .
- ii)  $\text{rang}(A) + \dim(\text{Ker}(A)) = n$ .

### B.1.3 Inverse d'une matrice

**Définition B.4** (matrice inverse). Une matrice carrée d'ordre  $n$  est *inversible* (ou régulière ou non-singulière) s'il existe une matrice  $B$  telle que:  $A \cdot B = B \cdot A = I_n$ . On note  $B = A^{-1}$ .

Une matrice qui n'est pas inversible est dite *singulière*. Propriétés:

$$(A^{-1})^{-1} = A \quad (\text{B.15})$$

$$(A \cdot B)^{-1} = B^{-1} \cdot A^{-1} \quad (\text{B.16})$$

---

<sup>1</sup>Cas complexe: produit scalaire hermitien

## B.2 Applications linéaires et matrice d'une application linéaire

**Définition B.5** (matrice d'une application unitaire dans bases données). Soit une application linéaire  $\mathcal{A} : V \rightarrow W$ , avec  $\dim(V) = n$  et  $\dim(W) = m$ . Soient  $B = (\mathbf{e}_1, \dots, \mathbf{e}_n)$  et  $B' = (\mathbf{f}_1, \dots, \mathbf{f}_m)$  des bases de  $V$  et  $W$ . On introduit la matrice de  $\mathcal{A}$  dans les bases  $B$  et  $B'$ ,  $A = (a_{ij})$

$$\mathcal{A}\mathbf{e}_j = \sum_i A_{ij}\mathbf{f}_i \quad (\text{B.17})$$

Alors, si les composantes de  $\mathbf{v} \in V$  dans  $B$  sont  $\mathbf{v} = v^i\mathbf{e}_i$ , les composantes de  $\mathbf{w} = \mathcal{A}\mathbf{v}$  dans  $B'$  sont ( $\mathbf{w} = \sum_i w^i\mathbf{f}_i$ )

$$\mathcal{A}\mathbf{v} = \sum_i \left( \sum_j a_{ij}v^j \right) \mathbf{f}_i \quad (\text{B.18})$$

$$w^i = \sum_j a_{ij}v^j \quad (\text{B.19})$$

$$\begin{pmatrix} w_1 \\ \vdots \\ w_m \end{pmatrix} = A \cdot \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix} \quad (\text{B.20})$$

La matrice  $A$  ainsi construite est nommée matrice de  $\mathcal{A}$  dans les bases  $B$  et  $B'$  et on peut la noter  $A = M(\mathcal{A}, B, B')$ . En particulier, l'action de la matrice  $M(\mathcal{A}, B, B')$  sur le vecteur dans  $\mathbf{R}^n$  donnée par les composantes de  $\mathbf{v}$  en  $B$ , renvoient le vecteur  $w$  dans  $\mathbf{R}^m$  avec les composantes de  $\mathbf{w} = \mathcal{A}\mathbf{v}$  dans  $B'$

$$\begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix} \mapsto \begin{pmatrix} w_1 \\ \vdots \\ w_m \end{pmatrix} = A \cdot \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix} = M(\mathcal{A}, B, B') \cdot \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix} \quad (\text{B.21})$$

En particulier, si  $V = W$  et  $\mathcal{A} = \mathbf{I}_n$ , la matrice  $M(\mathbf{I}_n, B, B')$  est la matrice de changement de base de  $B$  à  $B'$  dans le sens (passif) qu'elle prend les composantes en  $B'$  et donne les composantes en  $B$ . Ainsi, si  $\mathbf{v} = \sum_i v^i\mathbf{e}_i = \sum_i (v')^i\mathbf{e}'_i$  et si on note  $M(\mathbf{I}_n, B, B') = M(B \rightarrow B')$ , on a

$$\begin{pmatrix} v'_1 \\ \vdots \\ v'_n \end{pmatrix} = M(\mathbf{I}_n, B, B') \cdot \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix} = M(B \rightarrow B') \cdot \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix} \quad (\text{B.22})$$

où les colonnes de  $M(\mathbf{I}_n, B, B') = M(B \rightarrow B')$  sont données par les vecteurs colonnes de  $\mathbf{e}_i$  exprimées dans la base  $B' = (\mathbf{f}_1, \dots, \mathbf{f}_m)$

$$M(B \rightarrow B') = (\mathbf{e}_1 \quad \dots \quad \mathbf{e}_n) \quad (\text{B.23})$$

### B.2.1 Matrices semblables

**Définition B.6** (matrices semblables). Deux matrices carrées  $A, B \in M_n(\mathbb{K})$  sont semblables ssi  $\exists P$  inversible, telle que

$$A = P \cdot B \cdot P^{-1} \quad (\text{B.24})$$

**Proposition B.2.** *Deux matrices semblables correspondent à l'expression matricielle d'un même endomorphisme dans deux bases différentes.*

*Remarque B.2.* La notion de matrice semblable est un cas particulier dans le cas de matrices rectangulaires *équivalentes*, qui correspondent à l'expression matricielle d'une même application linéaire dans des bases différentes.

### B.3 Matrices et géométrie

En utilisant le produit de matrices, l'application suivante définit un produit scalaire en  $\mathbb{C}^n$ .

$$\langle \cdot, \cdot \rangle : \mathbb{C}^n \times \mathbb{R}^n \rightarrow \mathbb{R} \quad (\text{B.25})$$

$$(\mathbf{v}, \mathbf{v}) \mapsto \langle \mathbf{v}, \mathbf{v} \rangle = \mathbf{v}^t \cdot \mathbf{v} \quad (\text{B.26})$$

**Définition** (matrice symétrique). Une matrice  $A$  est dite symétrique si

$$A = A^t \Leftrightarrow a_{ij} = a_{ji} \quad (\text{B.27})$$

**Définition** (matrice hermitienne). Une matrice  $A$  est dite hermitienne si

$$A = A^* \Leftrightarrow a_{ij} = \bar{a}_{ji} \quad (\text{B.28})$$

**Définition** (matrice orthogonale). Une matrice  $U$  est dite orthogonale si

$$UU^t = I \quad (\text{B.29})$$

**Définition** (matrice unitaire). Une matrice  $U$  est dite unitaire si

$$UU^* = I \quad (\text{B.30})$$

**Définition** (matrice normale). Une matrice  $A$  est dite normale si

$$AA^* = A^*A \quad (\text{B.31})$$





## Appendix C

# Factorisation de matrices

**Théorème C.1** (Factorisation unitaire d'une matrice: Décomposition de Schur). *Toute matrice carrée (réelle ou complexes) peut s'écrire*

$$A = UTU^* , \tag{C.1}$$

*avec  $U$  une matrice unitaire  $U^{-1} = U^*$  et  $T$  une matrice triangulaire supérieure.*



# Appendix D

## Diagonalisation

### D.1 Notion de diagonalisation

On commence par introduire la notion de valeur propre.

**Définition D.1.** (*valeur propres*). Soit  $A \in \mathcal{M}_n(\mathbb{K})$ . On appelle valeur propre de  $A$  tout  $\lambda \in \mathbb{C}$  tel qu'il existe  $\mathbf{x} \in \mathbb{C}^n$ ,  $\mathbf{x} \neq 0$ , qui satisfait  $A\mathbf{x} = \lambda\mathbf{x}$ . L'élément  $\mathbf{x}$  est appelé vecteur propre ("à droite") de  $A$  associé à  $\lambda$ .

Ceci nous amène à la notion de diagonalisation.

**Définition D.2.** (*diagonalisation/réduction des endomorphismes*). Soit  $A \in \mathcal{M}_n(\mathbb{K})$ . On dit que  $A$  est diagonalisable dans  $\mathbb{K}$  s'il existe une base  $(\mathbf{v}_1, \dots, \mathbf{v}_n)$  avec  $v_i \in \mathbb{K}^n$  et des valeurs  $\lambda_1, \dots, \lambda_n \in \mathbb{K}$  telles que  $A\lambda_i = \lambda_i v_i$ .

**Lemme D.1.** *Une matrice diagonalisable est semblable à une matrice diagonale  $D$  donnée par*

$$D = \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix}, \quad (\text{D.1})$$

*c.à.d., s'il existe une matrice inversible  $P$ , telle que  $A = PDP^{-1}$ .*

**Lemme D.2.** *Les valeurs propres sont les racines du polynôme caractéristique  $p_A(\lambda) = \det(A - \lambda I)$ .*

**Preuve** (idée) : Le déterminant de  $A - \lambda I$  s'annule si son noyau n'est pas trivial. On montre que  $x \in \text{Ker}(A - \lambda I)$  ssi  $\lambda$  est une racine de  $p_A(\lambda)$ .  $\square$

Dans le cas  $\mathbb{K} = \mathbb{C}$ , le théorème fondamental de l'algèbre nous garantit l'existence de  $n$  solutions complexes pour le polynôme caractéristique de  $A \in \mathcal{M}_n(\mathbb{C})$ . La proposition suivante caractérise la possibilité de diagonaliser une matrice.

**Lemme D.3.** *Étant donné  $A \in \mathcal{M}_n(\mathbb{C})$ , une condition nécessaire et suffisante pour pouvoir diagonaliser  $A$  est que la dimension géométrique de chaque valeur propre  $\lambda$ , c.à.d.  $\dim E_\lambda$ , avec  $E_\lambda = \text{Ker}(A - \lambda I)$  le sous-espace propre de  $\lambda$ , coïncide avec sa dimension algébrique (multiplicité comme racine de  $p_A(\lambda)$ ).*

**Exemple D.1.** Nous considérons les exemples suivants:

i) La matrice  $\sigma_x$

$$\sigma_x = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \in \mathcal{M}_2(\mathbb{R}), \quad (\text{D.2})$$

est réelle et symétrique. Ses valeurs propres sont  $\lambda_1 = 1$  et  $\lambda_2 = -1$ . On peut vérifier que  $A$  est diagonalisable dans  $\mathbb{R}$ .

ii) La matrice  $A$

$$A = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} \in \mathcal{M}_2(\mathbb{R}), \quad (\text{D.3})$$

n'est pas diagonalisable en  $\mathbb{R}$ . En effet, ses valeurs propres sont  $\lambda_1 = i$  et  $\lambda_2 = -i$ . Elle est diagonalisable dans  $\mathbb{C}$ .

iii) La matrice  $\sigma_y$

$$\sigma_y = \begin{pmatrix} 0 & i \\ -i & 0 \end{pmatrix} \in \mathcal{M}_2(\mathbb{C}), \quad (\text{D.4})$$

est une matrice hermitienne. Ses valeurs propres sont  $\lambda_1 = 1$  et  $\lambda_2 = -1$ , alors la matrice est diagonalisable dans  $\mathbb{C}$ . Noter que  $\lambda_1$  et  $\lambda_2$  sont réels. Néanmoins, les vecteurs propres sont complexes.

iv) La matrice  $N$

$$N = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \in \mathcal{M}_2(\mathbb{R}), \quad (\text{D.5})$$

n'est pas diagonalisable dans  $\mathbb{C}$ . En effet, sa seule valeur propre est  $\lambda = 0$ . Si elle était diagonalisable, elle serait la matrice nulle, et ce n'est pas le cas. De la même manière, la matrice  $J$

$$N = \begin{pmatrix} 1 & a \\ 0 & 1 \end{pmatrix} \in \mathcal{M}_2(\mathbb{C}), \quad (\text{D.6})$$

n'est pas diagonalisable dans  $\mathbb{C}$ . En effet, la seule valeur propre est  $\lambda = 1$ . Si  $J$  était diagonalisable, elle serait semblable à la matrice unité, ce qui amène à une contradiction.

v) On peut construire des exemples de manière systématique en utilisant la *forme normale de Jordan*, valable pour n'importe quelle matrice complexe carrée (preuve par récurrence).

## D.2 Diagonalisation de matrices normales

**Théorème D.1.** Une matrice  $A$  est normale si et seulement si il existe  $U$  une matrice unitaire telle que

$$A = UDU^* \quad (\text{D.7})$$

avec  $D$  la matrice diagonale formée des valeurs propres. C'est-à-dire, une matrice normale est diagonalisable et ses vecteurs propres sont orthonormés.

Dém: (voir Nantes)