

Analyse numérique

Francis Filbet

Analyse numérique

Algorithme
et étude mathématique

2^e édition

DUNOD

Tout le catalogue sur
www.dunod.com



Illustration de couverture : © vege - Fotolia.com

Le pictogramme qui figure ci-contre mérite une explication. Son objet est d'alerter le lecteur sur la menace que représente pour l'avenir de l'écrit, particulièrement dans le domaine de l'édition technique et universitaire, le développement massif du photocopillage.

Le Code de la propriété intellectuelle du 1^{er} juillet 1992 interdit en effet expressément la photocopie à usage collectif sans autorisation des ayants droit. Or, cette pratique s'est généralisée dans les établissements

d'enseignement supérieur, provoquant une baisse brutale des achats de livres et de revues, au point que la possibilité même pour

les auteurs de créer des œuvres nouvelles et de les faire éditer correctement est aujourd'hui menacée. Nous rappelons donc que toute reproduction, partielle ou totale, de la présente publication est interdite sans autorisation de l'auteur, de son éditeur ou du Centre français d'exploitation du droit de copie (CFC, 20, rue des Grands-Augustins, 75006 Paris).



© Dunod, Paris, 2009, 2013
ISBN 978-2-10-059910-3

Le Code de la propriété intellectuelle n'autorisant, aux termes de l'article L. 122-5, 2^o et 3^o a), d'une part, que les « copies ou reproductions strictement réservées à l'usage privé du copiste et non destinées à une utilisation collective » et, d'autre part, que les analyses et les courtes citations dans un but d'exemple et d'illustration, « toute représentation ou reproduction intégrale ou partielle faite sans le consentement de l'auteur ou de ses ayants droit ou ayants cause est illicite » (art. L. 122-4).

Cette représentation ou reproduction, par quelque procédé que ce soit, constituerait donc une contrefaçon sanctionnée par les articles L. 335-2 et suivants du Code de la propriété intellectuelle.

TABLE DES MATIÈRES

Avant-propos	ix
Chapitre 1. Les systèmes linéaires	1
1.1 Exemple d'un système linéaire	1
1.2 Rappels sur les matrices	2
1.2.1 Cas des matrices carrées	3
1.2.2 Quelques matrices particulières	6
1.2.3 Conditionnement de matrices	8
1.3 Méthodes directes	15
1.3.1 Méthodologie générale	15
1.3.2 Méthode de Gauss avec et sans pivot	15
1.3.3 Factorisation de Cholesky	26
1.3.4 Factorisation QR	30
1.4 Méthodes itératives	34
1.4.1 Méthodologie générale	35
1.4.2 Méthode de Jacobi	39
1.4.3 Méthode de Gauss-Seidel	42
Exercices	44
Chapitre 2. Calcul numérique de valeurs propres	57
2.1 Quelques exemples de problèmes aux valeurs propres	57
2.1.1 Algorithme de Google	57
2.1.2 Mouvement de ressorts	59
2.2 Localisation des valeurs propres	61
2.2.1 Approximation des valeurs propres	61
2.2.2 Ce qu'il ne faut pas faire!	63
2.3 Méthode de la puissance	64
2.3.1 La méthode	64
2.3.2 Un résultat de convergence	65

Table des matières

2.4	Méthode de Jacobi	68
2.4.1	Cas de la dimension deux	68
2.4.2	Cas général	69
2.5	La méthode QR pour le calcul des valeurs propres	72
	Exercices	74
Chapitre 3. Les systèmes non linéaires		85
3.1	Introduction aux problèmes non linéaires	85
3.1.1	Modèle de coagulation et fragmentation	85
3.1.2	Résultats généraux et définitions	87
3.2	Méthode de point fixe	88
3.2.1	Méthode de Héron	88
3.2.2	Méthode générale	89
3.3	Vers la méthode de Newton-Raphson	93
3.3.1	Méthode de dichotomie	93
3.3.2	Méthode de la sécante	94
3.3.3	Méthode de Newton-Raphson	95
3.3.4	Combinaison de méthodes	98
3.4	Méthode de Newton-Raphson dans \mathbb{R}^n	99
3.4.1	Quelques rappels de calcul différentiel	99
3.4.2	Méthode de Newton-Raphson	103
	Exercices	109
Chapitre 4. Optimisation		119
4.1	Introduction	119
4.2	Optimisation sans contrainte	121
4.2.1	Rappels sur les fonctions convexes	121
4.2.2	Résultat d'existence et unicité	123
4.2.3	Méthodes d'optimisation sans contrainte	127
4.2.4	Cas d'une fonctionnelle quadratique : la méthode du gradient conjugué	130
4.3	Optimisation sous contraintes	135
4.3.1	Existence et unicité du problème sous contraintes	136
4.3.2	Conditions d'optimalité	137
4.3.3	Méthodes d'optimisation sous contraintes	143
	Exercices	146

Chapitre 5. Les polynômes	153
5.1 Introduction	153
5.1.1 Un exemple en analyse	153
5.1.2 Rappels sur les polynômes	154
5.2 Interpolation de Lagrange et de Hermite	156
5.2.1 Construction et convergence de l'interpolation de Lagrange	156
5.2.2 Interpolation composée	161
5.2.3 Applications : formules de quadrature pour l'approximation d'intégrales	163
5.2.4 Interpolation de Hermite	165
5.3 Méthode des moindres carrés	167
5.3.1 Rappel du théorème d'approximation	168
5.3.2 Résolution du problème des moindres carrés discrets	170
5.4 Polynômes orthogonaux	172
5.4.1 Construction de polynômes orthogonaux	172
5.4.2 Méthode des moindres carrés continue	176
5.4.3 Application : formules de quadrature pour l'approximation d'intégrales	178
5.4.4 Transformée de Fourier rapide	182
Exercices	190
Chapitre 6. Les équations différentielles	203
6.1 Introduction et rappels théoriques	203
6.1.1 Le pont de Tacoma aux États-Unis	203
6.1.2 Rappels théoriques	205
6.2 Schémas à un pas explicites	210
6.2.1 Les schémas de Runge-Kutta	211
6.2.2 Convergence du schéma d'Euler explicite	215
6.2.3 Consistance et stabilité des schémas à un pas	217
6.3 Les équations différentielles raides	223
6.3.1 Convergence du schéma d'Euler implicite	227
6.3.2 Stabilité des schémas implicites	227
6.4 Les systèmes hamiltoniens	230
Exercices	239

Table des matières

Chapitre 7. Approximation numérique des équations dérivées partielles	249
7.1 Modélisation de la répartition de chaleur	249
7.2 La méthode aux différences finies	252
7.2.1 Approximation de la dérivée	252
7.2.2 L'équation de Poisson	254
7.2.3 Étude de l'erreur	256
7.2.4 Conditions aux limites mixtes	260
7.2.5 Problèmes plus généraux	262
7.2.6 Cas multidimensionnel	264
7.3 La méthode des éléments finis	265
7.3.1 Formulation variationnelle	265
7.3.2 Méthodologie générale	267
7.3.3 Cas uni-dimensionnel	269
7.4 Les équations d'évolution	272
7.4.1 L'équation de la chaleur	275
7.4.2 L'équation des ondes	283
7.5 La méthode des volumes finis pour les équations de transport	288
7.5.1 L'équation de transport	288
7.5.2 Les volumes finis	297
Bibliographie	307
Index	309

AVANT-PROPOS

L'objet de cet ouvrage est de fournir au lecteur une présentation et une analyse des méthodes numériques les plus couramment utilisées pour la résolution de problèmes de mathématiques appliquées. Chaque chapitre fait l'objet d'une introduction et d'exemples puisés dans des applications à divers domaines des mathématiques. Il s'agit ici d'illustrer l'utilisation d'outils mathématiques pour des problèmes actuels (modèles économiques, physiques, algorithme de Google, etc).

Cet ouvrage s'adresse tout particulièrement aux étudiants de Licence 3, Master ou préparation à l'agrégation ainsi qu'aux élèves des écoles d'Ingénieurs intéressés par des problèmes issus de la physique, de la biologie ou encore de l'économie, et qui souhaitent avoir une idée des méthodes de bases utilisées aujourd'hui et être en mesure de les mettre en œuvre sur ordinateur.

Nous avons volontiers glissé plusieurs rappels théoriques qui devraient aider le lecteur à mieux appréhender les difficultés mathématiques de certaines démonstrations. Notons toutefois que certains résultats sont présentés dans un cadre simplifié de manière à ne pas noyer le lecteur dans des difficultés techniques.

L'ensemble des sujets abordés ici à déjà largement été traité dans la littérature et ce cours doit beaucoup à d'autres ouvrages et cours dispensés en France, comme les cours d'analyse numérique de Jean-Marie Thomas, Michel Pierre et Abderrahmane Bendali.

Lors de la rédaction de cet ouvrage, de nombreux collègues et étudiants ont bien voulu me faire part de leurs avis, commentaires, suggestions. À ce propos, je tiens tout particulièrement à exprimer ma reconnaissance à Jean-Pierre Bourgade, Nicolas Crouseilles, Thomas Lepoutre, Céline Parzani, Miguel Rodrigues, qui m'ont fait part de leurs conseils éclairés et pour certains ont bien voulu effectuer une relecture de ce document.

Je dédie ce livre à Raphaël et Flavien.

La résolution d'un système linéaire algébrique est au cœur de la plupart des calculs en analyse numérique. Il paraît donc naturel de débiter un cours de calcul scientifique par là. Ici, nous décrivons les algorithmes de résolution les plus populaires qui sont appliqués à des systèmes généraux. Nous considérons le problème suivant : trouver le vecteur x solution de

$$Ax = b,$$

où A est une matrice carrée et b un vecteur donné à coefficients réels ou complexes. La discrétisation d'équations différentielles ordinaires ou d'équations aux dérivées partielles, la modélisation de problèmes en physique, chimie ou économie conduit souvent à la résolution de systèmes linéaires de grande taille avec plusieurs milliers d'inconnues et il devient pratiquement impossible de résoudre ces systèmes d'équations sans l'aide d'un ordinateur. Il s'agit alors de trouver des algorithmes de résolution efficaces où le nombre d'opérations et donc le temps de calcul, n'est pas prohibitif. C'est un problème classique mais difficile en analyse numérique.

L'objectif de ce chapitre est de proposer différentes méthodes numériques de résolution des systèmes linéaires et de sensibiliser le lecteur à l'importance du choix de la méthode en fonction des propriétés du système. Nous distinguerons deux types de méthodes : *les méthodes directes* où nous calculons exactement la solution et *les méthodes itératives* où nous calculons une solution approchée.

1.1 EXEMPLE D'UN SYSTÈME LINÉAIRE

Commençons par présenter un problème provenant de l'économie, où nous sommes amenés à la résolution d'un système linéaire.

Exemple d'une analyse de l'offre et de la demande

Imaginons que plusieurs artisans décident de coopérer pour fabriquer différents produits utiles à chacun d'eux. Afin d'éviter le coût du stockage des produits fabriqués, nous recherchons une situation d'équilibre entre l'offre et la demande. Pour cela, considérons n artisans, sachant que chacun fabrique un produit spécifique. Ces artisans ont choisi de coopérer, c'est-à-dire qu'ils souhaitent adapter leur production, d'une part, à leurs propres besoins déterminés par la quantité de produit nécessaire aux autres artisans pour qu'ils puissent fabriquer leur propre produit et, d'autre part, aux besoins du marché.

Puisque chaque artisan fabrique un produit différent, nous notons par $i \in \{1, \dots, n\}$ le produit fabriqué par un artisan, x_i désigne le nombre total de produits i fabriqués par un artisan et b_i la demande du marché en ce produit. D'autre part $(c_{i,j})_{1 \leq i, j \leq n}$ exprime la quantité de produit i nécessaire à la confection d'une unité de produit j . En supposant que la relation qui lie les différents produits est linéaire, nous recherchons l'équilibre entre les besoins et la production, c'est-à-dire que nous imposons que la quantité de produit i fabriquée soit égale à la somme des besoins des autres artisans en produit i pour fabriquer leur propre produit et des besoins du marché en i

$$x_i = \sum_{j=1}^n C_{i,j} x_j + b_i, \quad i = 1, \dots, n$$

ou encore

$$x = C x + b,$$

où la matrice C est formée par les coefficients $(c_{i,j})_{1 \leq i, j \leq n}$ tandis que le vecteur b correspond à la demande du marché $b := (b_1, \dots, b_n)^T$. Par conséquent, la production totale $x = (x_1, \dots, x_n)^T$ est la solution du système linéaire $A x = b$, où la matrice A est donnée par $A = I_n - C$ et I_n est la matrice identité composée de 1 sur sa diagonale et de 0 ailleurs.

Cette modélisation pose plusieurs difficultés mathématiques. Tout d'abord nous nous interrogeons sur les propriétés de la matrice A permettant d'assurer que ce problème a bien une solution. Il vient ensuite la question du calcul de cette solution. En pratique, lorsque le nombre d'artisans considérés devient grand, il n'est plus possible de calculer l'inverse de la matrice A , il faudra donc fournir des algorithmes qui ne nécessitent pas l'inversion de cette matrice. C'est l'objectif de ce premier chapitre.

1.2 RAPPELS SUR LES MATRICES

Avant de s'intéresser à la résolution numérique de systèmes linéaires, présentons quelques rappels d'algèbre linéaire.

Par la suite, nous noterons $\mathcal{M}_{m,n}(\mathbb{K})$ l'ensemble des matrices à m lignes, n colonnes et à coefficient dans le corps des réels $\mathbb{K} = \mathbb{R}$ ou des complexes $\mathbb{K} = \mathbb{C}$. Pour un vecteur colonne $u \in \mathbb{K}^n$, la valeur u_i avec $1 \leq i \leq n$, désigne la i -ème composante dans la base canonique de \mathbb{K}^n du vecteur u . Nous appelons *adjoint* du vecteur colonne u , le vecteur ligne u^* de \mathbb{K}^n tel que $u^* = (\bar{u}_1, \dots, \bar{u}_n)$ où \bar{u}_i désigne le conjugué de u_i et transposé de u est le vecteur ligne $u^T = (u_1, \dots, u_n)$.

Rappelons également le produit matrice-vecteur. Soient A une matrice à m lignes et n colonnes et u un vecteur de \mathbb{K}^n , nous définissons $v = A u \in \mathbb{K}^m$ le vecteur dont les composantes sont données par

$$v_i = \sum_{j=1}^n a_{i,j} u_j, \quad i = 1, \dots, m$$

1.2. Rappels sur les matrices

et pour B une matrice à n lignes et p colonnes, le produit AB fournit une matrice $C \in \mathcal{M}_{m,p}$ donnée par

$$c_{i,j} = \sum_{k=1}^n a_{i,k} b_{k,j}, \quad i = 1, \dots, m, \quad j = 1, \dots, p.$$

Introduisons également la notion de matrice adjointe et transposée.

Définition 1.1

Soit $A \in \mathcal{M}_{m,n}(\mathbb{K})$, la matrice adjointe de A , notée A^* , à n lignes et m colonnes, est donnée par $A^* = (a_{i,j}^*)_{\substack{1 \leq i \leq n \\ 1 \leq j \leq m}}$ et $a_{i,j}^* = \bar{a}_{j,i}$, pour $1 \leq i \leq n$ et $1 \leq j \leq m$.

La matrice transposée de A , notée A^T , à n lignes et m colonnes, est donnée par $A^T = (a_{i,j}^T)_{\substack{1 \leq i \leq n \\ 1 \leq j \leq m}}$ et $a_{i,j}^T = a_{j,i}$, pour $1 \leq i \leq n$ et $1 \leq j \leq m$.

Dans la base canonique de \mathbb{K}^n , définissons le produit scalaire ou produit hermitien entre deux vecteurs u et $v \in \mathbb{K}^n$, le scalaire de \mathbb{K} donné par

$$\langle u, v \rangle = u^* v = \sum_{i=1}^n u_i^* v_i.$$

1.2.1 Cas des matrices carrées

Considérons maintenant le cas particulier des matrices carrées, pour lesquelles le nombre de lignes est égal au nombre de colonnes. Nous rappelons les définitions suivantes :

Définition 1.2

Une matrice carrée $A \in \mathcal{M}_{n,n}(\mathbb{K})$ est dite inversible, ou régulière, ou encore non singulière, s'il existe une matrice $B \in \mathcal{M}_{n,n}(\mathbb{K})$ telle que $AB = BA = I_n$. Dans ce cas, la matrice B est unique et s'appelle la matrice inverse de A , elle est notée A^{-1} .

Définition 1.3

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$, alors

- A est une matrice normale si $A^* A = A A^*$.
- A est une matrice unitaire si $A^* A = A A^* = I_n$, où I_n désigne la matrice identité. Dans le cas où $\mathbb{K} = \mathbb{R}$, nous parlons de matrice orthogonale et $A^T A = A A^T = I_n$, c'est-à-dire $A^T = A^{-1}$.
- A est une matrice hermitienne si $A^* = A$. Dans le cas où $\mathbb{K} = \mathbb{R}$, nous parlons de matrice symétrique et $A^T = A$.

Il est alors facile de faire le lien entre les matrices hermitiennes et normales.

Proposition 1.4

Toute matrice hermitienne est une matrice normale.

Comme nous l'avons vu en début de chapitre, l'objectif de cette partie est de mettre au point des algorithmes de résolution numérique pour un système de la forme

$$Ax = b, \tag{1.1}$$

où $A \in \mathcal{M}_{n,n}(\mathbb{K})$ et $b \in \mathbb{K}^n$ sont donnés et $x \in \mathbb{K}^n$ est l'inconnue. Nous donnons d'abord une condition nécessaire et suffisante sur la matrice A pour que ce système admette une solution unique. Le système linéaire (1.1) a une solution unique dès que la matrice A est inversible et cette solution est donnée par $x = A^{-1}b$. Il est alors important de connaître quelques propriétés des matrices inversibles [16].

Proposition 1.5

Soient $A, B \in \mathcal{M}_{n,n}(\mathbb{K})$ inversibles, nous avons alors

- pour tout $\alpha \in \mathbb{K}^*$, $(\alpha A)^{-1} = \frac{1}{\alpha} A^{-1}$.
- $(AB)^{-1} = B^{-1}A^{-1}$.
- $(A^*)^{-1} = (A^{-1})^*$.

Énonçons quelques propriétés des matrices inversibles, qui sont autant d'outils permettant de s'assurer que le problème (1.1) admet bien une solution [16].

Théorème 1.6

Théorème des matrices inversibles

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$, les propositions suivantes sont équivalentes :

- A est inversible,
- le déterminant de A est non nul,
- le rang de A est égal à n , c'est-à-dire que les n vecteurs colonnes forment une famille libre,
- le système homogène $Ax = 0$ a pour unique solution $x = 0$,
- pour tout b dans \mathbb{K}^n , le système linéaire $Ax = b$ a exactement une solution.

Pour conclure cette partie, rappelons qu'une valeur propre de A est donnée par $\lambda \in \mathbb{K}$ telle que $\det(A - \lambda I_n) = 0$. Ainsi, il existe au moins un vecteur v non nul, dit *vecteur propre*, vérifiant $Av = \lambda v$. Définissons alors le spectre d'une matrice.

Définition 1.7

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$. Nous appelons spectre de A l'ensemble des valeurs propres de A :

$$\text{Sp}(A) = \{\lambda \in \mathbb{K}; \exists v \in \mathbb{K}^n \quad v \neq 0 \quad Av = \lambda v\}$$

Nous appelons rayon spectral de A le nombre réel positif $\rho(A)$ tel que

$$\rho(A) = \max \{|\lambda|, \lambda \in \mathbb{C}, v \in \mathbb{C} \setminus \{0\} Av = \lambda v\}$$

À partir de la définition d'une valeur propre, nous vérifions facilement qu'une matrice est inversible si et seulement si 0 n'est pas valeur propre.

La conséquence de ces propriétés est que l'ensemble des matrices carrées inversibles forme un groupe, appelé le groupe linéaire et noté habituellement $GL_n(\mathbb{K})$. En général, « presque toutes » les matrices sont inversibles. Sur le corps \mathbb{K} , cela peut être formulé de façon plus précise : l'ensemble des matrices non inversibles, considéré comme sous-ensemble de $\mathcal{M}_{n,n}(\mathbb{K})$, est un ensemble négligeable, c'est-à-dire de mesure de Lebesgue nulle. Intuitivement, cela signifie que si vous choisissez au hasard une matrice carrée à coefficients réels, la probabilité pour qu'elle soit non inversible est égale à zéro. La raison est que des matrices non inversibles peuvent être considérées comme racines d'une fonction polynôme donnée par le déterminant [16].

Par la suite, nous nous intéresserons au calcul proprement dit de la solution de (1.1). Proposons d'abord un premier algorithme naïf qui permet de calculer la solution d'un système linéaire ; cet algorithme est dû à Cramer¹ [23].

Théorème 1.8

Résolution de Cramer

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$ une matrice inversible. Alors la solution du système $Ax = b$ est donnée par $x_i = \det(A_i)/\det(A)$, pour tout $i = 1, \dots, n$, où A_i est la matrice A pour laquelle la i -ème colonne est remplacée par le vecteur b .

Cette méthode bien que très élégante est très coûteuse puisqu'elle nécessite plus de $n!$ opérations où $n!$ est la factorielle de n , donnée par $n! = n \times (n-1) \times \dots \times 2 \times 1$. Elle n'est donc jamais utilisée en pratique sauf en dimension $n = 2$. Cet algorithme revient à calculer explicitement la matrice A^{-1} . Hélas, ce calcul est souvent long et fastidieux, même pour un ordinateur, c'est pourquoi nous avons recours à des algorithmes de résolution exacte d'une complexité moindre (nous parlons alors d'une méthode directe), ou des méthodes itératives qui consistent à construire une suite de solutions approchées qui converge vers la solution exacte. Avant de présenter de tels algorithmes, nous introduirons des matrices dont la structure particulière est bien adaptée à la résolution du système linéaire (1.1). Puis nous verrons comment se ramener à ces cas particuliers.

1. En référence au mathématicien suisse Gabriel Cramer (1704-1752).

1.2.2 Quelques matrices particulières

Cette partie constitue la base théorique permettant de mettre au point des méthodes pour la résolution exacte de (1.1).

• Matrices triangulaires

Nous introduisons la notion de matrice triangulaire et montrons que pour les matrices de ce type, la résolution du système linéaire (1.1) devient très facile.

Définition 1.9

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$,

- la matrice A est une matrice triangulaire inférieure si $A = (a_{i,j})_{1 \leq i, j \leq n}$ avec $a_{i,j} = 0, \quad 1 \leq i < j \leq n$.
- A est une matrice triangulaire supérieure lorsque $a_{i,j} = 0, \quad 1 \leq j < i \leq n$.
- A est une matrice diagonale dès lors que $a_{i,j} = 0$ pour $i \neq j$.

Nous vérifions d'abord que l'ensemble des matrices triangulaires supérieures (resp. inférieures) est stable par la somme et le produit. En outre, nous avons le résultat suivant.

Proposition 1.10

Soit $L \in \mathcal{M}_{n,n}(\mathbb{K})$ une matrice triangulaire inférieure inversible, ce qui signifie que tous les éléments diagonaux sont non nuls. Alors L^{-1} est aussi une matrice triangulaire inférieure.

Soit U une matrice triangulaire supérieure inversible. Alors U^{-1} est aussi une matrice triangulaire supérieure.

Les matrices triangulaires jouent un rôle important en analyse numérique car elles sont facilement inversibles ou du moins, nous pouvons facilement trouver la solution $x \in \mathbb{K}^n$ du système linéaire $Ax = b$. En effet, considérons le cas d'une matrice triangulaire supérieure, alors la solution x se calcule par un algorithme dit de remontée. Nous observons d'abord qu'une matrice triangulaire inversible a tous ses éléments diagonaux non nuls, c'est pourquoi nous pouvons écrire $x_n = b_n/a_{n,n}$ puis pour tout $i = n - 1, n - 2, \dots, 1$, nous pouvons calculer x_i de la manière suivante :

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

Exemple

Soit $A \in \mathcal{M}_{3,3}(\mathbb{R})$ donnée par

$$A = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 2 & 2 \\ 0 & 0 & 3 \end{pmatrix},$$

1.2. Rappels sur les matrices

nous cherchons une solution de $Ax = b$ avec $b = (3, 4, 3)^T$. En commençant par la dernière ligne, nous trouvons $x_3 = 1$. Puis en injectant cette valeur dans l'avant-dernière équation, cela donne $x_2 = 1$. Finalement, connaissant x_2 et x_3 , la première équation donne directement $x_1 = 1$.

Nous allons voir par la suite qu'une méthode directe calcule la solution exacte de (1.1) et consiste le plus souvent à trouver un système triangulaire équivalent.

C'est pourquoi nous énonçons le Théorème de Shur [7] qui sera justement utile pour rendre certaines matrices triangulaires par simple changement de base.

Théorème 1.11 *Théorème de Shur*

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$ une matrice quelconque. Alors il existe une matrice unitaire $U \in \mathcal{M}_{n,n}(\mathbb{K})$, c'est-à-dire $U^* U = I_n$ et une matrice $T \in \mathcal{M}_{n,n}(\mathbb{K})$ triangulaire dont la diagonale est composée par l'ensemble des valeurs propres de A telles que $T = U^* A U$.

• *Matrices hermitiennes*

Dans le cas des matrices hermitiennes, ce dernier résultat peut être sensiblement amélioré.

Corollaire 1.12

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$ une matrice hermitienne. Alors il existe une matrice unitaire $U \in \mathcal{M}_{n,n}(\mathbb{K})$ et une matrice diagonale $D \in \mathcal{M}_{n,n}(\mathbb{K})$ dont la diagonale est composée par l'ensemble des valeurs propres de A , telles que $D = U^* A U$.

Nous rappelons enfin quelques propriétés de base des matrices hermitiennes, définies positives. Nous introduisons d'abord la notion de sous-matrice principale.

Définition 1.13

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$ une matrice quelconque. Nous appelons sous-matrice principale d'ordre i (pour $1 \leq i \leq n$) de A , la matrice $A_i \in \mathcal{M}_{i,i}(\mathbb{K})$ obtenue en ne gardant que les i premières lignes et les i premières colonnes de A .

Définition 1.14

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$, nous disons que A est positive si pour tout $x \in \mathbb{K}^n$, elle vérifie $x^T A x \geq 0$ et A est définie positive si pour tout $x \in \mathbb{K}^n \setminus \{0\}$, elle vérifie $x^T A x > 0$ et $x^T A x = 0$ implique $x = 0$.

Nous avons alors les résultats suivants dont la preuve est laissée en exercice.

Proposition 1.15

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$ une matrice hermitienne définie positive. Alors elle vérifie

- toute sous-matrice principale A_i , $1 \leq i \leq n$ est hermitienne définie positive ;
- tous les coefficients diagonaux de A sont des réels strictement positifs ;
- A est diagonalisable et ses valeurs propres sont strictement positives ;
- le déterminant de A est strictement positif, c'est-à-dire que A est inversible ;
- il existe une constante $\alpha > 0$, telle que $x^* A x \geq \alpha \|x\|^2$, pour n'importe qu'elle norme vectorielle $\|\cdot\|$.

• *Matrices de permutations*

Par la suite, nous utiliserons souvent les matrices de permutations pour réordonner les lignes ou les colonnes d'une matrice quelconque.

Définition 1.16

Une matrice de permutations est une matrice carrée qui ne possède que des 0 et des 1 comme coefficients, telle qu'il y ait un seul 1 par ligne et par colonne.

Une matrice de permutations vérifie alors

Proposition 1.17

Soient σ et τ deux permutations des indices $\{1, \dots, n\}$, nous avons alors

- la matrice de permutations P_σ correspondant à la permutation σ s'écrit comme

$$(P_\sigma)_{i,j} = \delta_{i,\sigma(j)} = \begin{cases} 1, & \text{si } i = \sigma(j), \\ 0, & \text{sinon.} \end{cases}$$

- $P_\sigma P_\tau = P_{\sigma \circ \tau}$.
- $(P_\sigma)^T P_\sigma = I_n$, c'est une matrice orthogonale.

1.2.3 Conditionnement de matrices

Avant de décrire des algorithmes de résolution de systèmes linéaires (1.1), nous mettons en évidence une difficulté supplémentaire : la sensibilité de la solution par rapport aux perturbations des données $b \in \mathbb{K}^n$ ou des coefficients de la matrice A . En effet, considérons par exemple la matrice de Hilbert donnée par

$$a_{i,j} = \frac{1}{i+j-1}, \quad i, j = 1, \dots, n.$$

Nous avons pour $n = 4$,

$$\begin{pmatrix} 1 & 1/2 & 1/3 & 1/4 \\ 1/2 & 1/3 & 1/4 & 1/5 \\ 1/3 & 1/4 & 1/5 & 1/6 \\ 1/4 & 1/5 & 1/6 & 1/7 \end{pmatrix}.$$

Ainsi, pour un vecteur donné $b \in \mathbb{R}^4$ tel que

$$b^T = \left(\frac{25}{12}, \frac{77}{60}, \frac{57}{60}, \frac{319}{420} \right) \simeq (2,08, 1,28, 0,95, 0,76),$$

nous calculons $x = A^{-1}b$ et vérifions facilement que $x^T = (1, 1, 1, 1)$. Ensuite, si nous modifions légèrement la source $\tilde{b}^T = (2,1, 1,3, 1, 0,8)$, la solution \tilde{x} du système $\tilde{x} = A^{-1}\tilde{b}$ est donnée par $\tilde{x}^T = (5,6, -48, 114, -70)$, ce qui est très loin de la solution de départ x .

En définitive, nous constatons qu'une petite perturbation de la donnée b conduit à une grande modification de la solution x , ce qui engendre des instabilités lors de la résolution du système.

• Norme matricielle

Pour mesurer l'ampleur de cette instabilité, nous introduisons la notion de *norme matricielle*. En effet, l'ensemble $\mathcal{M}_{n,n}(\mathbb{K})$ peut être considéré comme étant un \mathbb{K} -espace vectoriel muni d'une norme $\|\cdot\|$.

Définition 1.18

Nous appelons *norme matricielle* toute application $\|\cdot\|$ de $\mathcal{M}_{n,n}(\mathbb{K})$ à valeur dans $\mathbb{R}^+ := [0, +\infty[$ qui vérifie les propriétés suivantes :

- pour toute matrice $A \in \mathcal{M}_{n,n}(\mathbb{K})$, $\|A\| = 0 \Rightarrow A = 0_{\mathbb{K}^{n \times n}}$,
- pour toute matrice $A \in \mathcal{M}_{n,n}(\mathbb{K})$, et pour tout $\alpha \in \mathbb{K}$, $\|\alpha A\| = |\alpha| \|A\|$,
- pour toutes matrices A et $B \in \mathcal{M}_{n,n}(\mathbb{K})$, $\|A + B\| \leq \|A\| + \|B\|$, c'est l'inégalité triangulaire,
- pour toutes matrices A et $B \in \mathcal{M}_{n,n}(\mathbb{K})$, $\|A B\| \leq \|A\| \cdot \|B\|$.

Notons bien que définir une norme matricielle requiert une condition supplémentaire par rapport à la définition d'une norme vectorielle (la dernière propriété de la définition). Il n'est donc pas évident *a priori* de pouvoir construire une telle application seulement à partir d'une norme vectorielle. À titre d'exemple la norme de Froebenius $\|\cdot\|_F$ donnée par

$$\|A\|_F = \left(\sum_{i,j=1}^n |a_{i,j}|^2 \right)^{1/2} \quad (1.2)$$

est bien une norme matricielle. Les trois premières propriétés sont connues puisque $\|\cdot\|_F$ est une norme vectorielle de \mathbb{K}^{n^2} . Il reste à démontrer que pour toutes matrices A et $B \in \mathcal{M}_{n,n}(\mathbb{K})$, $\|A B\|_F \leq \|A\|_F \|B\|_F$. En effet,

$$\|A B\|_F^2 = \|C\|_F^2 = \sum_{i,j=1}^n c_{i,j}^2 = \sum_{i,j=1}^n \left(\sum_{k=1}^n a_{i,k} b_{k,j} \right)^2.$$

En utilisant l'inégalité de Cauchy-Schwarz, nous avons

$$\begin{aligned} \|A B\|_F^2 &\leq \sum_{i,j=1}^n \left(\sum_{k=1}^n a_{i,k}^2 \right) \left(\sum_{k=1}^n b_{k,j}^2 \right) = \left(\sum_{i,k=1}^n a_{i,k}^2 \right) \left(\sum_{k,j=1}^n b_{k,j}^2 \right), \\ &= \|A\|_F^2 \|B\|_F^2. \end{aligned}$$

Pour construire une norme matricielle à partir d'une norme vectorielle quelconque, nous introduisons la définition suivante :

Définition 1.19

Soit $\|\cdot\|$ une norme vectorielle sur \mathbb{K}^n . Nous définissons la norme matricielle $\|\cdot\|$ subordonnée à la norme vectorielle $\|\cdot\|$ comme étant l'application donnée par

$$A \in \mathcal{M}_{n,n}(\mathbb{K}) \mapsto \|A\| := \sup_{\substack{v \in \mathbb{K}^n \\ v \neq 0}} \frac{\|A v\|}{\|v\|}.$$

Nous vérifions facilement que cette application définit bien une norme matricielle.

Par exemple pour $1 \leq p \leq \infty$, nous savons que l'application qui à $v \in \mathbb{K}^n$ fait correspondre le réel positif

$$\|v\|_p = \left(\sum_{i=1}^n |v_i|^p \right)^{1/p}$$

ou $\|v\|_\infty = \max_{1 \leq i \leq n} |v_i|$ pour $p = \infty$, est une norme vectorielle sur \mathbb{K}^n . Posons alors pour $p \in [1, +\infty]$ et $A \in \mathcal{M}_{n,n}(\mathbb{K})$

$$\|A\|_p = \sup_{\substack{v \in \mathbb{K}^n \\ v \neq 0}} \frac{\|A v\|_p}{\|v\|_p},$$

qui est une norme matricielle subordonnée à la norme vectorielle $\|\cdot\|_p$. En général nous ne pouvons pas calculer $\|A\|_p$ directement en fonction de $(a_{i,j})_{1 \leq i,j \leq n}$ sauf pour $p = 1$ et $p = \infty$.

Proposition 1.20 *Caractérisation de $\|\cdot\|_1$ et $\|\cdot\|_\infty$*

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$ une matrice quelconque. Alors nous avons

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

et

$$\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{i,j}|.$$

DÉMONSTRATION. Soit $v \in \mathbb{K}^n$. D'une part, nous avons

$$\begin{aligned} \|A v\|_1 &= \sum_{i=1}^n \left| \sum_{j=1}^n a_{i,j} v_j \right| \leq \sum_{i=1}^n \sum_{j=1}^n |a_{i,j}| |v_j|, \\ &\leq \sum_{j=1}^n \left(\sum_{i=1}^n |a_{i,j}| \right) |v_j| \leq \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{i,j}| \|v\|_1. \end{aligned}$$

D'où

$$\|A\|_1 \leq \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{i,j}|.$$

D'autre part, nous montrons qu'il existe $w \in \mathbb{K}^n$ tel que $\|w\|_1 = 1$ et

$$\|A w\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{i,j}|.$$

En effet, choisissons $j_0 \in \{1, \dots, n\}$ tel que

$$\sum_{i=1}^n |a_{i,j_0}| = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{i,j}|,$$

puis prenons $w \in \mathbb{K}^n$ tel que $w_i = 0$ pour $i \neq j_0$ et $w_{j_0} = 1$. Alors

$$\|A w\|_1 = \sum_{i=1}^n |(A w)_i| = \sum_{i=1}^n |a_{i,j_0}| = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{i,j}|.$$

Par un raisonnement analogue, nous montrons le résultat pour la norme matricielle $\|\cdot\|_\infty$. \square

Dans le cas particulier de la norme subordonnée à la norme euclidienne de \mathbb{K}^n définie pour $A \in \mathcal{M}_{n,n}(\mathbb{K})$ par

$$\|A\|_2 = \sup_{\substack{v \in \mathbb{K}^n \\ v \neq 0}} \frac{\|A v\|_2}{\|v\|_2},$$

nous démontrons le résultat suivant.

Proposition 1.21

Pour $A \in \mathcal{M}_{n,n}(\mathbb{K})$, nous avons $\|A\|_2 = \sqrt{\rho(A^* A)} = \sqrt{\rho(A A^*)}$.

DÉMONSTRATION. Prenons $\mathbb{K} = \mathbb{C}$ ou \mathbb{R} , par définition de la norme subordonnée à la norme euclidienne, il vient

$$\begin{aligned} \|A\|_2^2 &= \sup_{\substack{v \in \mathbb{K}^n \\ v \neq 0}} \frac{\|A v\|_2^2}{\|v\|_2^2} = \sup_{\substack{v \in \mathbb{K}^n \\ v \neq 0}} \frac{(A v)^* (A v)}{v^* v}, \\ &= \sup_{\substack{v \in \mathbb{K}^n \\ v \neq 0}} \frac{v^* A^* A v}{v^* v}. \end{aligned}$$

Or, nous vérifions que

$$(A^* A)^* = A^* (A^*)^* = A^* A.$$

Ainsi, d'après le Corollaire 1.12 puisque la matrice $A^* A$ est hermitienne, elle est diagonalisable et il existe $U \in \mathcal{M}_{n,n}(\mathbb{K})$ unitaire telle que $U A^* A U = \text{diag}(\mu_k)$, où $\text{diag}(\mu_k)$ représente une matrice diagonale formée à partir des valeurs μ_k sur la diagonale. Les scalaires $(\mu_k)_{1 \leq k \leq n}$ sont les valeurs propres de la matrice hermitienne $A^* A$. Notons que $\mu_k \geq 0$ puisque de la relation $A^* A p_k = \mu_k p_k$, nous déduisons que $(A p_k)^* A p_k = \mu_k p_k^* p_k$, c'est-à-dire $\mu_k = \|A p_k\|_2^2 / \|p_k\|_2^2 \geq 0$. Il vient ensuite

$$\sup_{\substack{v \in \mathbb{K}^n \\ v \neq 0}} \frac{v^* A^* A v}{v^* v} = \sup_{\substack{v \in \mathbb{K}^n \\ v \neq 0}} \frac{v^* U^* U A^* A U^* U v}{v^* U^* U v}.$$

Puis, comme U est inversible, nous avons par changement de variable $w = U v$

$$\|A\|_2^2 = \sup_{\substack{w \in \mathbb{K}^n \\ w \neq 0}} \frac{w^* U A^* A U^* w}{w^* w} = \sup_{\substack{w \in \mathbb{K}^n \\ w \neq 0}} \frac{\sum_{k=1}^n \mu_k |w_k|^2}{\sum_{k=1}^n |w_k|^2},$$

où les μ_k sont les valeurs propres de $A^* A$. Enfin, en prenant le vecteur de la base canonique dont toutes les composantes sont nulles exceptée la k_0 -ème qui correspond à la plus grande valeur propre μ_{k_0} en module, nous obtenons $\|A\|_2^2 = \rho(A A^*)$. \square

• Conditionnement d'une matrice

Essayons maintenant de comprendre de manière plus générale le phénomène d'instabilité par rapport à la donnée $b \in \mathbb{K}^n$. Soient $A \in \mathcal{M}_{n,n}(\mathbb{K})$ inversible et b un vecteur de \mathbb{K}^n , non nul. Nous désignons par x la solution du système $Ax = b$ et pour une perturbation δb du vecteur b , le vecteur $x + \delta x$ est la solution de $A(x + \delta x) = b + \delta b$.

Pour une norme vectorielle $\|\cdot\|$ de \mathbb{K}^n , nous cherchons à contrôler l'erreur relative $\|\delta x\| / \|x\|$ en fonction de l'erreur relative $\|\delta b\| / \|b\|$ et de la norme matricielle subordonnée $\|A\|$.

Par linéarité et puisque A est inversible, nous avons d'une part $A\delta x = \delta b \Rightarrow \|\delta x\| \leq \|A^{-1}\| \|\delta b\|$ et d'autre part $Ax = b$ implique que $\|b\| \leq \|A\| \|x\|$, ou encore de manière équivalente puisque b est non nul

$$\frac{1}{\|x\|} \leq \|A\| \frac{1}{\|b\|}.$$

Ainsi, nous obtenons l'estimation suivante :

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

et proposons la définition suivante :

Définition 1.22

Nous appelons conditionnement de la matrice A relativement à la norme matricielle $\|\cdot\|$ subordonnée à la norme vectorielle $\|\cdot\|$, le nombre

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

Observons bien que le conditionnement sert à mesurer la sensibilité du système aux perturbations de b et de A . Il permet d'assurer le contrôle de $\|\delta x\|$ en fonction de la perturbation $\|\delta b\|$.

Définition 1.23

Soit $A \in \mathcal{M}_{n,n}(\mathbb{K})$ une matrice inversible. Nous disons que

- un système linéaire est bien conditionné, si $\text{cond}(A)$ n'est pas trop grand par rapport à un, ce qui correspond au conditionnement de l'identité ;
- un système linéaire est mal conditionné, si $\text{cond}(A)$ est grand par rapport à un.

Vérifions sur l'exemple précédent la cohérence de cette définition. Nous avons pour la norme $\|\cdot\|_1$, $\text{cond}_1(A) = \|A\|_1 \|A^{-1}\|_1 \simeq 28\,375 \gg 1$ ou pour la norme $\|\cdot\|_2$, $\text{cond}_2(A) = \|A\|_2 \|A^{-1}\|_2 \simeq 15\,514 \gg 1$. Les valeurs du conditionnement de la matrice A semblent être assez élevées indépendamment de la norme choisie.

Pour conclure cette partie, nous donnons quelques propriétés sur le conditionnement.

Proposition 1.24

Soient $A \in \mathcal{M}_{n,n}(\mathbb{K})$ une matrice carrée inversible et $\|\cdot\|$ une norme matricielle subordonnée à la norme vectorielle $\|\cdot\|$. Alors nous avons

- $\text{cond}(A^{-1}) = \text{cond}(A)$.
- Soit $\alpha \in \mathbb{K} \setminus \{0\}$, $\text{cond}(\alpha A) = \text{cond}(A)$.
- $\text{cond}(I_n) = 1$.
- $\text{cond}(A) \geq 1$.

L'inconvénient de la définition du conditionnement est qu'il fait apparaître $\|A^{-1}\|$, lequel n'est pas facile à calculer d'autant plus que nous ne connaissons pas la forme explicite de la matrice A^{-1} . Dans le cas particulier d'une matrice hermitienne et pour la norme matricielle $\|\cdot\|_2$, nous avons néanmoins le résultat suivant qui ne nécessite pas le calcul de A^{-1} mais seulement la connaissance des valeurs propres de A .

Proposition 1.25

Soit A une matrice hermitienne ($A^* = A$). Alors $\|A\|_2 = \rho(A)$. De plus, si A est une matrice hermitienne inversible et $(\lambda_i)_{1 \leq i \leq n}$ ses valeurs propres. Alors,

$$\text{cond}_2(A) = \frac{\max\{|\lambda_i|, i = 1, \dots, n\}}{\min\{|\lambda_i|, i = 1, \dots, n\}}.$$

DÉMONSTRATION. Appliquons la Proposition 1.21 et puisque A est hermitienne, nous montrons que $\|A\|_2^2 = \rho(A)^2$.

Supposons ensuite que A est une matrice hermitienne inversible. Pour calculer $\text{cond}_2(A)$, il suffit de remarquer que $1/\lambda_i$, où $\lambda_i \neq 0$, est valeur propre de A^{-1} et donc en appliquant le résultat précédent à A^{-1} (qui est également une matrice hermitienne), nous avons

$$\|A^{-1}\|_2 = \rho(A^{-1}) = \frac{1}{\min\{|\lambda_i|, i = 1, \dots, n\}}.$$

Nous en déduisons le résultat

$$\text{cond}_2(A) = \|A\|_2 \|A^{-1}\|_2 = \frac{\max\{|\lambda_i|, i = 1, \dots, n\}}{\min\{|\lambda_i|, i = 1, \dots, n\}}. \quad \square$$

- *Préconditionnement d'un système linéaire*

Pour remédier au problème du mauvais conditionnement d'une matrice, nous pouvons appliquer une méthode de preconditionnement. En effet, en vue de résoudre $Ax = b$ nous multiplions ce système d'équation à gauche par une matrice inversible P , il vient alors $PAx = Pb$, avec P choisie de manière à ce que la matrice PA soit mieux conditionnée que A (dans le cas le plus favorable, nous aurions $P = A^{-1}$). Cependant, il n'y a pas de méthode standard pour trouver la matrice P , le plus souvent nous chercherons une matrice à la fois facile à inverser et « assez proche » de A^{-1} .

Présentons maintenant les deux types de méthodes pour la résolution d'un système linéaire : les méthodes directes et les méthodes itératives.

1.3 MÉTHODES DIRECTES

1.3.1 Méthodologie générale

Soient $A \in \mathcal{M}_{n,n}(\mathbb{K})$ une matrice inversible et un vecteur $b \in \mathbb{K}^n$, nous recherchons $x \in \mathbb{K}^n$ solution de $Ax = b$. Pour cela, construisons des matrices M et $N \in \mathcal{M}_{n,n}(\mathbb{K})$ telles que $A = MN$, où M est facile à inverser (triangulaire ou unitaire) et N triangulaire. Le système s'écrit alors $Nx = M^{-1}b$, que nous résolvons en deux étapes de la manière suivante :

$$\begin{cases} \text{trouver } y \in \mathbb{K}^n \text{ tel que } My = b, \\ \text{trouver } x \in \mathbb{K}^n \text{ tel que } Nx = y. \end{cases}$$

1.3.2 Méthode de Gauss avec et sans pivot

Les méthodes directes permettent de calculer la solution exacte du problème (1.1) en un nombre fini d'étapes (en l'absence d'erreurs d'arrondi). La méthode directe la plus classique est la méthode d'élimination de Gauss ou Gauss-Jordan, qui consiste à décomposer la matrice A comme le produit LU où L est une matrice triangulaire inférieure et U une matrice triangulaire supérieure¹.

L'élimination de Gauss ou l'élimination de Gauss-Jordan est un algorithme d'algèbre linéaire pour déterminer les solutions d'un système d'équations linéaires, pour déterminer le rang d'une matrice ou pour calculer l'inverse d'une matrice carrée inversible.

1. Cette méthode fut nommée d'après Carl Friedrich Gauss, mathématicien allemand (1777-1855) surnommé le prince des mathématiciens. Cependant, si l'on se réfère au livre chinois « Les neuf chapitres sur l'art du calcul » sa naissance remonte à la dynastie Han au premier siècle de notre ère. Elle constitue le huitième chapitre de ce livre sous le titre de la « disposition rectangulaire » et est présentée au moyen de dix-huit exercices [6].