

Chapitre I : Résolution de systèmes linéaires

Introduction

Deux problèmes fondamentaux :

- (i) Étant donnée une matrice inversible A et un vecteur b , déterminer la solution x du système linéaire : $Ax = b$.
- (ii) Étant donnée une matrice carrée A , déterminer les éléments propres de A .

Un objectif : proposer des méthodes pratiques de résolution de ces deux problèmes.

I Rappels d'algèbre linéaire

Soit $N \in \mathbb{N}^*$.

1. Normes subordonnées

Soit $\|\cdot\|$, une norme sur l'espace vectoriel \mathbb{C}^N .

Définition : La norme subordonnée à la norme $\|\cdot\|$ sur l'espace $M_N(\mathbb{C})$ est l'application $\| \cdot \|$ définie par :

$$\forall A \in M_N(\mathbb{C}), \|A\| = \sup \{ \|Ax\|, x \in \mathbb{C}^N \text{ s.t. } \|x\| \leq 1 \}$$

Exemples : (i) Soit $\forall x \in \mathbb{C}^N, \|x\|_2 = \sqrt{\sum_{j=1}^N |x_j|^2}$. L'application $\|\cdot\|_2$ est une norme sur l'espace vectoriel \mathbb{C}^N , et sa norme subordonnée

donnée vaut:

$$\forall A \in \mathcal{M}_n(\mathbb{K}), \|A\|_\infty = \max_{1 \leq j \leq n} \sum_{i=1}^n |A_{ij}|.$$

(ii) Soit $\forall x \in \mathbb{K}^n$, $\|x\|_\infty = \max\{|x_j|, 1 \leq j \leq n\}$. L'application $\|\cdot\|_\infty$ est une norme sur l'espace vectoriel \mathbb{K}^n , et sa norme subordonnée vaut:

$$\forall A \in \mathcal{M}_n(\mathbb{K}), \|A\|_{\infty, \infty} = \max_{1 \leq i \leq n} \sum_{j=1}^n |A_{ij}|.$$

Propriétés: (i) La norme subordonnée $\|\cdot\|_\infty$ à la norme $\|\cdot\|_\infty$ est bien définie et est une norme sur l'espace $\mathcal{M}_n(\mathbb{K})$.

(ii) La norme subordonnée $\|\cdot\|_\infty$ est définie de manière équivalente par les caractérisations:

$$\begin{aligned} \forall A \in \mathcal{M}_n(\mathbb{K}), \|A\|_\infty &= \max\{\|A(x)\|_\infty, x \in \mathbb{K}^n \text{ t.q. } \|x\|_\infty = 1\} \\ &= \max\{\|A(x)\|_\infty, x \in \mathbb{K}^n \text{ t.q. } \|x\|_\infty = 1\} \\ &= \max\{\alpha \in \mathbb{R}_+, \text{ t.q. } \forall x \in \mathbb{K}^n, \|A(x)\|_\infty \leq \alpha \|x\|_\infty\}. \end{aligned}$$

(iii) En particulier, la norme subordonnée $\|\cdot\|_\infty$ satisfait l'inégalité:

$$\forall x \in \mathbb{K}^n, \|A(x)\|_\infty \leq \|A\|_\infty \|x\|_\infty.$$

(iv) La norme subordonnée $\|\cdot\|_\infty$ est une norme d'algèbre sur l'algèbre $\mathcal{M}_n(\mathbb{K})$:

$$\|I_n\|_\infty = 1 \text{ et } \forall (A, B) \in \mathcal{M}_n(\mathbb{K})^2, \|AB\|_\infty \leq \|A\|_\infty \|B\|_\infty.$$

Preuve:

Immédiate.

2. Relation avec le rayon spectral

Rappel: Soit $A \in \mathcal{M}_n(\mathbb{K})$.

(i) Le spectre $\sigma(A)$ de la matrice A est non vide.

(ii) La matrice A est trigonalisable: il existe une matrice inversible $P \in \mathcal{GL}_n(\mathbb{K})$ telle que:

$$A = P^{-1} \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_N \end{pmatrix} P, \text{ où } \sigma(A) = \{\lambda_1, \dots, \lambda_N\}.$$

Definition: Soit $A \in M_N(\mathbb{C})$. Le rayon spectral $\rho(A)$ de la matrice A est défini par:

$$\rho(A) = \max \{ |\lambda|, \lambda \in \sigma(A) \}.$$

Exemples: (i) $\rho(0) = 0$ et $\rho(I_N) = 1$

(ii) Si une matrice $U \in M_N(\mathbb{C})$ est unitaire, alors: $\rho(U) = 1$.

Propriété: Soit $\|\cdot\|$ la norme subordonnée à la norme $\|\cdot\|$ de \mathbb{C}^N .

(i) Si $A \in M_N(\mathbb{C})$, alors:

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

(ii) La réciproque est fautive.

Preuve:

(i) Soit $\lambda \in \sigma(A)$ telle que $|\lambda| = \rho(A)$, et $v \in \mathbb{C}^N \setminus \{0\}$ tel que $A(v) = \lambda v$

$$\text{d'alors: } \rho(A) \|v\| = \|A(v)\| \leq \|A\| \|v\| \Rightarrow \|A\| \geq \rho(A).$$

(ii) $\rho \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} = 0 < \left\| \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \right\|$.

Lemme: Soit $\varepsilon \in]0, +\infty[$ et $A \in M_N(\mathbb{C})$. Il existe une norme subordonnée $\|\cdot\|_A$ à une norme $\|\cdot\|_A$ de \mathbb{C}^N telle que:

$$\|A\|_A \leq \rho(A) + \varepsilon.$$

Preuve:

Soit $P \in GL_N(\mathbb{C})$ et $(\lambda_1, \dots, \lambda_N) \in \mathbb{C}^N$ tels que:

$$A = P^{-1} \begin{pmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \ddots & \\ 0 & & & \lambda_N \end{pmatrix} P$$

Pour $\delta > 0$, introduisons la matrice $D_\delta = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & \delta^{-1} \end{pmatrix}$. Il vient:

$$(D_\delta^{-1} P)^{-1} A (D_\delta^{-1} P)^{-1} = \begin{pmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \ddots & \\ 0 & & & \lambda_N \end{pmatrix},$$

De sorte que :

$$\| (D\delta^{-2\epsilon}) A (D\delta^{-2\epsilon})^t \|_\infty = \max_{z \leq i \leq N} |\lambda_i| + \sum_{j=i+1}^N |b_{ij}| \delta^{j-i}$$

Prions alors le choix de δ tel que :

$$\forall z \leq i \leq N, \sum_{j=i+1}^N |b_{ij}| \delta^{j-i} \leq \epsilon$$

Dans ce cas, nous obtenons :

$$\| (D\delta^{-2\epsilon}) A (D\delta^{-2\epsilon})^t \|_\infty \leq \rho(A) + \epsilon.$$

Enfin, nous savons que :

$$\begin{aligned} \| A \|_A &:= \| (D\delta^{-2\epsilon}) A (D\delta^{-2\epsilon})^{-2} \|_\infty \\ &= \max \{ \alpha \in \mathbb{R}_+ \text{ t.q. } \forall u \in \mathbb{S}^N, \| (D\delta^{-2\epsilon}) A (D\delta^{-2\epsilon})^{-2}(u) \|_\infty \leq \alpha \| u \|_\infty \} \\ &= \max \{ \alpha \in \mathbb{R}_+ \text{ t.q. } \forall y = (D\delta^{-2\epsilon})^{-2}(u) \in \mathbb{S}^N, \| (D\delta^{-2\epsilon}) A(y) \|_\infty \leq \alpha \| (D\delta^{-2\epsilon})^{-2}(u) \|_\infty \} \end{aligned}$$

De sorte que l'application $\| \cdot \|_A$ est la norme subordonnée à la norme de \mathbb{S}^N définie par :

$$\forall x \in \mathbb{S}^N, \| x \|_A = \| A \delta^{-2\epsilon}(x) \|_\infty.$$

Ce qui conclut la preuve du théorème.

Alébrème : Soit $\forall x \in \mathbb{S}^N, \| x \|_2 = \left(\sum_{j=1}^N |x_j|^2 \right)^{\frac{1}{2}}$.

(i) L'application $\| \cdot \|_2$ est une norme sur l'espace vectoriel \mathbb{S}^N .

(ii) La norme subordonnée $\| \cdot \|_2$ vaut :

$$\forall A \in \mathcal{M}_N(\mathbb{C}), \| A \|_2 = \rho(A^*A)^{\frac{1}{2}} = \rho(AA^*)^{\frac{1}{2}}.$$

En particulier :

(iii) $\forall A \in \mathcal{M}_N(\mathbb{C}), \| A^* \|_2 = \| A \|_2$

(iv) $\forall V \in \mathcal{U}_N(\mathbb{C}), \| V \|_2 = 1.$

(v) $\forall A \in \mathcal{M}_N(\mathbb{C}), \forall V \in \mathcal{U}_N(\mathbb{C}), \| V^* A V \|_2 = \| A \|_2 = \| A V \|_2 = \| A \|_2.$

(vi) En particulier, si la matrice $N \in \mathcal{M}_N(\mathbb{C})$ est normale, alors :

$$\| N \|_2 = \rho(N).$$

Preuve :

(i) Transmissité.

(ii) Par définition, nous savons que :

$$\| A \|_2^2 = \max \{ \| A(x) \|_2^2, x \in \mathbb{S}^N \text{ t.q. } \| x \|_2 \leq 1 \}.$$

Comme

$$\forall x \in \mathbb{C}^N, \|A(x)\|_2^2 = \langle A^*A(x), x \rangle_2,$$

il vient:

$$\|A\|_2^2 = \max \{ \langle A^*A(x), x \rangle_2, x \in \mathbb{C}^N \text{ s.t. } \|x\|_2 \leq 1 \}.$$

La matrice A^*A est hermitienne positive de sorte qu'il existe une matrice unitaire $U \in \mathcal{U}_N(\mathbb{C})$ et des nombres positifs $\lambda_1 \leq \dots \leq \lambda_N$ tels que:

$$A^*A = U^* \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_N \end{pmatrix} U.$$

Comme l'application $y \mapsto U^*y$ est une isométrie de l'espace vectoriel hermitien $(\mathbb{C}^N, \langle \cdot, \cdot \rangle_2)$, nous aurons au fait que:

$$\begin{aligned} \|A\|_2^2 &= \max \{ \langle \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_N \end{pmatrix} (y), y \rangle_2, y \in \mathbb{C}^N \text{ s.t. } \|y\|_2 \leq 1 \} \\ &= \lambda_N \\ &= \rho(A^*A) \end{aligned}$$

Deux cas se présentent alors:

- si $\rho(A^*A) = 0$, alors,

$$\|A\|_2 = 0 \Rightarrow A = 0 \Rightarrow \rho(AA^*) = 0 = \rho(A^*A) = \|A\|_2^2$$

- sinon, il existe un vecteur $x \in \mathbb{C}^N \setminus \{0\}$ tel que:

$$A^*A(x) = \rho(A^*A)x \Rightarrow (AA^*)(A(x)) = \rho(A^*A)A(x).$$

Comme

$$\|A(x)\|_2^2 = \langle A^*A(x), x \rangle_2 = \rho(A^*A) \|x\|_2^2 > 0,$$

la matrice hermitienne positive AA^* a pour valeur propre $\rho(A^*A)$, de sorte que:

$$\rho(A^*A) \leq \rho(AA^*).$$

De même, il vient:

$$\rho(AA^*) = \rho((A^*)^*A^*) \leq \rho(A^*A),$$

De sorte que l'égalité a lieu.

(iii) Immédiat

(iv) Immédiat.

(v) Cette propriété découle en particulier du fait que l'application $y \mapsto U^*y$ est une isométrie de l'espace hermitien $(\mathbb{C}^N, \langle \cdot, \cdot \rangle_2)$.

(vi) Si la matrice N est normale, alors, il existe une matrice unitaire $U \in \mathcal{U}_N(\mathbb{C})$ et

Des nombres complexes $\lambda_1, \dots, \lambda_n$ tels que:

$$N = U^* \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix} U,$$

On note que:

$$\|N\|_2 = \left\| \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix} \right\|_2 = \max_{1 \leq j \leq n} |\lambda_j| = \rho(N).$$

3. Convergence de la suite des puissances itérées

Rappel: Soit $(A_n)_{n \in \mathbb{N}}$, une suite de $M_n(\mathbb{C})$ et $A_0 \in M_n(\mathbb{C})$.

La suite $(A_n)_{n \in \mathbb{N}}$ converge vers la matrice A_0 dans $M_n(\mathbb{C})$ si il existe une norme $\|\cdot\|$ sur l'espace vectoriel $M_n(\mathbb{C})$ telle que:

$$\|A_n - A_0\| \xrightarrow{n \rightarrow +\infty} 0.$$

En particulier, cette propriété ne dépend pas du choix de la norme

$\|\cdot\|$ sur $M_n(\mathbb{C})$.

Lemme: Soit $A \in M_n(\mathbb{C})$.

(i) La suite des puissances itérées $(A^n)_{n \in \mathbb{N}}$ converge vers la matrice nulle dans $M_n(\mathbb{C})$ si le rayon spectral de la matrice A satisfait l'inégalité: $\rho(A) < 1$.

(ii) Dans ce cas, la matrice $I_n - A$ est inversible, et son inverse est donné par la série convergente:

$$(I_n - A)^{-1} = \sum_{n=0}^{+\infty} A^n.$$

Preuve:

(i) Soit $A^n \xrightarrow{n \rightarrow +\infty} 0$ dans $M_n(\mathbb{C})$ et si λ désigne une valeur propre de la matrice A telle que:

$$|\lambda| = \rho(A),$$

alors, il existe un vecteur propre $v \in \mathbb{C}^n \setminus \{0\}$ tel que:

$$A(v) = \lambda v \Rightarrow \forall n \in \mathbb{N}, A^n(v) = \lambda^n v.$$

Il s'ensuit que:

$$\rho(A)^m \|r\|_2 = |\lambda|^m \|r\|_2 = \|A^m(r)\|_2 \leq \|A^m\|_2 \|r\|_2 \xrightarrow{m \rightarrow +\infty} 0,$$

ce qui entraîne que: $\rho(A) < 1$.

Réciproquement, si $\rho(A) < 1$, alors, il existe une norme subordonnée $\|\cdot\|$ sur $\mathcal{M}_N(\mathbb{C})$ telle que:

$$\|A\| < 1.$$

Comme

$$\forall n \in \mathbb{N}, \|A^n\| \leq \|A\|^n,$$

il est clair que:

$$A^n \xrightarrow{m \rightarrow +\infty} 0 \text{ dans } \mathcal{M}_N(\mathbb{C}).$$

(ii) Soit $\|\cdot\|$, une norme subordonnée sur $\mathcal{M}_N(\mathbb{C})$ telle que:

$$\|A\| < 1.$$

Comme

$$\forall n \in \mathbb{N}, \|A^n\| \leq \|A\|^n,$$

la série $\sum_{n=0}^{+\infty} A^n$ est absolument convergente, donc convergente, et:

$$(I_N - A) \left(\sum_{n=0}^{+\infty} A^n \right) = \left(\sum_{n=0}^{+\infty} A^n \right) (I_N - A) = I_N.$$

La matrice $I_N - A$ est donc inversible et l'inverse égal à la somme $\sum_{n=0}^{+\infty} A^n$.

Corollaire (Formule du rayon spectral): Soit $A \in \mathcal{M}_N(\mathbb{C})$ et $\|\cdot\|$, une norme d'algèbre sur l'algèbre $\mathcal{M}_N(\mathbb{C})$. Alors:

$$\rho(A) = \lim_{n \rightarrow +\infty} \|A^n\|^{1/n}.$$

Preuve:

Soit $\lambda \in \sigma(A)$ telle que $|\lambda| = \rho(A)$, et $v \in \mathbb{C}^N \setminus \{0\}$ tel que $A(v) = \lambda v$. Alors:

$$\rho(A)^m \|v\| = \|\lambda^m v\| = \|A^m(v)\| \leq \|A^m\| \|v\|.$$

Comme $\|v\| \neq 0$, il vient:

$$\rho(A) \leq \|A^m\|^{1/m} \Rightarrow \rho(A) \leq \lim_{m \rightarrow +\infty} \|A^m\|^{1/m}.$$

Soit alors $\varepsilon > 0$. Il vient:

$$\rho \left(\frac{A}{\rho(A) + \varepsilon} \right) < 1,$$

D'où suit que

$$\frac{A^m}{(\rho(A) + \varepsilon)^m} \xrightarrow{m \rightarrow +\infty} 0.$$

En particulier, il existe un entier $m \in \mathbb{N}$ tel que:

$$\forall m > m_0, \|A^m\|_{\infty}^{\frac{1}{m}} \leq \rho(A) + \varepsilon \Rightarrow \lim_{m \rightarrow \infty} \|A^m\|_{\infty}^{\frac{1}{m}} = \rho(A).$$

En conclusion, nous avons :

$$\rho(A) = \lim_{m \rightarrow \infty} \|A^m\|_{\infty}^{\frac{1}{m}}$$

II Conditionnement de problèmes linéaires

1. Position du problème

Considérons la matrice : $A := \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix}$, et le vecteur : $b := \begin{pmatrix} 32 \\ 23 \\ 33 \\ 32 \end{pmatrix}$.

La matrice A est symétrique et inversible, et son inverse A^{-1} est égale à :

$$A^{-1} = \begin{pmatrix} 25 & -42 & 10 & -6 \\ -42 & 68 & -23 & 20 \\ 10 & -23 & 5 & -3 \\ -6 & 20 & -3 & 2 \end{pmatrix}$$

L'unique solution x du système linéaire $Ax = b$ est égale à : $x = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$.

Les calculs numériques sont approchés et doivent prendre en compte les problèmes d'arrondis : de légères modifications de la matrice A ou des second membre b ne doivent pas entraîner des erreurs trop importantes sur le calcul des vecteurs solution x . Qu'en est-il ici ?

Lorsque le vecteur b est remplacé par le vecteur $b' := \begin{pmatrix} 32,1 \\ 23,5 \\ 33,1 \\ 32,9 \end{pmatrix}$, l'unique solution x' est désormais égale à : $x' = \begin{pmatrix} 9,6 \\ -24,6 \\ 4,5 \\ -4,2 \end{pmatrix}$.

Une erreur relative de l'ordre de $1/200$ entraîne une erreur relative sur la solution de l'ordre de $10/1$: l'amplification des erreurs est de l'ordre de 2000, ce qui est colossal.

Lorsque la matrice A est remplacée par la matrice $A' := \begin{pmatrix} 10 & 7 & 3,2 & 7,2 \\ 7,02 & 5,04 & 6 & 5 \\ 8 & 5,98 & 3,92 & 9 \\ 6,92 & 4,94 & 9 & 3,92 \end{pmatrix}$

l'unique solution x' devient : $x' := \begin{pmatrix} -82 \\ 237 \\ -34 \\ 82 \end{pmatrix}$. L'amplification des erreurs est à nouveau colossale, et ce bien que les matrices A , A^{-1} , et le second membre b est l'un des plus raisonnables.

La notion de conditionnement permet de quantifier cette difficulté, et parfois d'y remédier à travers l'introduction de méthodes de préconditionnement.

Notons pour terminer qu'un problème identique apparaît dans la résolution de problèmes de valeurs propres. Le spectre de la matrice $B = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$ est réduit au singleton $\{0\}$, tandis que le spectre de la matrice $B' = \begin{pmatrix} 0 & 0 & 0 & 10^{-4} \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$ est donné par : $\sigma(B') = \{\pm 10^{-4}, \pm 10^{-2} i\}$. L'amplification de l'erreur relative est dans ce cas de l'ordre de 10000 !

2. Conditionnement d'un système linéaire

Soit $N \in \mathbb{N}^*$, et $\|\cdot\|$, une norme sur l'espace vectoriel \mathbb{C}^N . Désignons par la notation $\|\cdot\|_{\|\cdot\|}$, la norme subordonnée à la norme $\|\cdot\|$ sur l'espace $\mathcal{L}_N(\mathbb{C})$.

Définition: Soit $A \in \mathcal{L}_N(\mathbb{C})$. Le conditionnement de la matrice A relativement à la norme subordonnée $\|\cdot\|_{\|\cdot\|}$ est le nombre strictement positif :

$$\text{cond}(A) := \|A\|_{\|\cdot\|} \times \|A^{-1}\|_{\|\cdot\|}.$$

Exemples: (i) $\text{cond}(I_N) = 1$ (quel que soit le choix de la norme $\|\cdot\|$ sur \mathbb{C}^N)
 (ii) Si $\|\cdot\|_{\|\cdot\|_2}$ est la norme subordonnée à la norme hermitienne

canonique $\|\cdot\|_2$ sur \mathbb{C}^N , alors:

$$\forall V \in \mathcal{U}_N(\mathbb{C}), \text{cond}_2(V) := \|V\|_2 \|V^{-1}\|_2 = 1.$$

Propriétés: Soit $A \in \mathcal{G}_N(\mathbb{C})$. Le conditionnement de la matrice A relativement à la norme subordonnée $\|\cdot\|$ vérifie les propriétés suivantes:

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

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

(iii) $\forall \alpha \in \mathbb{C}^*$, $\text{cond}(\alpha A) = \text{cond}(A)$.

Preuve:

(i) Comme $I_N = A A^{-1}$, il vient:

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

(ii) (iii) Immédiat.

Dans le cas particulier de la norme subordonnée $\|\cdot\|_2$ à la norme hermitienne canonique $\|\cdot\|_2$ de \mathbb{C}^N , le conditionnement cond_2 vérifie les propriétés suivantes.

Propriétés: Soit $A \in \mathcal{G}_N(\mathbb{C})$.

(i) Si $U \in \mathcal{U}_N(\mathbb{C})$, alors: $\text{cond}_2(A) = \text{cond}_2(AU) = \text{cond}_2(UA) = \text{cond}_2(U^*AU)$.

(ii) Si $\mu_1(A)$ et $\mu_N(A)$ désignent les plus petite, resp. la plus grande valeur propre de la matrice hermitienne définie positive A^*A , alors: $\text{cond}_2(A) = \left(\frac{\mu_N(A)}{\mu_1(A)} \right)^{\frac{1}{2}}$.

(iii) Si la matrice A est normale, et si son spectre est donné par l'ensemble $\sigma(A) = \{\lambda_1(A), \dots, \lambda_n(A)\}$, alors:

$$\text{cond}_2(A) = \frac{\max_{1 \leq i \leq n} |\lambda_i(A)|}{\min_{1 \leq i \leq n} |\lambda_i(A)|}$$

Preuve:

(i) Immédiat.

(ii) Rappelons que:

$$\forall n \in \mathbb{N}, \forall A \in \mathcal{M}_n(\mathbb{C}), \|A\|_2 = \rho(n \cdot n)^{\frac{1}{2}} = \rho(n \cdot n)^{\frac{1}{2}}$$

De sorte que:

$$\text{cond}_2(A) = \rho(A^*A)^{\frac{1}{2}} \rho((A^*A)^{-1})^{\frac{1}{2}} = \left(\frac{\rho(A^*A)}{\rho(A^*A)} \right)^{\frac{1}{2}}$$

(ii) Lorsque la matrice $A \in \mathcal{M}_n(\mathbb{C})$ est normale, son inverse l'est également.

En particulier, il vient:

$$\text{cond}_2(A) = \rho(A) \rho(A^{-1}),$$

De sorte que:

$$\text{cond}_2(A) = \frac{\max_{\lambda \in \sigma_n} |\lambda(A)|}{\min_{\lambda \in \sigma_n} |\lambda(A)|}$$

Exemple: La formule précédente permet de calculer une valeur approchée du conditionnement relativement à la norme $\|\cdot\|_2$ de la matrice

$$A = \begin{pmatrix} 2.0 & 7 & 8.7 \\ 7 & 5 & 6.5 \\ 8 & 6 & 10.9 \\ 7 & 5 & 5 & 10 \end{pmatrix}, \text{ laquelle est égale à: } \text{cond}_2(A) \approx 2.924.$$

Comme on soulignait les deux fractions qui suivent, cette valeur élevée du conditionnement explique les problèmes d'arrondis mentionnés plus tôt.

Le conditionnement d'une matrice permet d'évaluer les erreurs relatives dans la résolution d'un système linéaire de la façon suivante.

Théorème: Soit $A \in \mathcal{M}_n(\mathbb{C})$.

(i) Soit $b \in \mathbb{C}^n \setminus \{0\}$ et $\delta b \in \mathbb{C}^n$. Si x et $x + \delta x$ désignent les solutions des équations $Ax = b$, resp. $A(x + \delta x) = b + \delta b$, alors:

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

(ii) De plus, il existe des vecteurs $b \in \mathbb{C}^n \setminus \{0\}$ et $\delta b \in \mathbb{C}^n$ tels que cette inégalité soit une égalité.

Preuve:

(i) Par linéarité, il vient:

$$A(\delta x) = \delta b,$$

De sorte que:

$$\begin{cases} \| \delta x \| \leq \| A^{-1} \| \| \delta b \|, \\ \| b \| \leq \| A \| \| x \|, \end{cases}$$

Soit:

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

(ii) Cette inégalité devient une égalité lorsque les vecteurs b et δb sont choisis de façon à satisfaire les identités:

$$\begin{cases} \| A^{-1}(\delta b) \| = \| A^{-1} \| \| \delta b \|, \\ \| b \| \leq \| A \| \| A^{-1}(b) \|, \end{cases}$$

ce qui est possible.

Exercice: Soit $A \in \mathcal{GL}_n(\mathbb{C})$ et $\delta A \in \mathcal{M}_n(\mathbb{C})$ telle que $A + \delta A$ soit inversible.

(i) Soit $b \in \mathbb{C}^n \setminus \{0\}$. Si x et $x + \delta x$ désignent les solutions des équations $Ax = b$, resp. $(A + \delta A)(x + \delta x) = b$, alors:

$$\frac{\| \delta x \|}{\| x + \delta x \|} \leq \text{cond}(A) \frac{\| \delta A \|}{\| A \|}$$

(ii) De plus, il existe un vecteur $b \in \mathbb{C}^n \setminus \{0\}$ et une matrice $\delta A \in \mathcal{M}_n(\mathbb{C}) \setminus \{0\}$ tels que cette inégalité soit une égalité.

(iii) Supposons que: $\| \delta A \| \leq \frac{1}{2 \| A^{-1} \|}$. Si $b \in \mathbb{C}^n \setminus \{0\}$ et si x et $x + \delta x$ désignent les solutions des équations $Ax = b$, resp. $(A + \delta A)(x + \delta x) = b$, alors:

$$\frac{\| \delta x \|}{\| x \|} \leq 2 \text{cond}(A) \frac{\| \delta A \|}{\| A \|}$$

Preuve:

(i) Par linéarité, il vient:

$$A(\delta x) + \delta A(x + \delta x) = 0,$$

De sorte que:

$$\| \delta x \| \leq \| A^{-1} \| \| \delta A \| \| x + \delta x \| \Leftrightarrow \frac{\| \delta x \|}{\| x + \delta x \|} \leq \text{cond}(A) \frac{\| \delta A \|}{\| A \|}$$

(ii) Soit $\lambda \in \mathbb{C}^*$ tel que $A + \lambda I_n$ soit inversible et $\delta A := \lambda I_n$. Il existe un vecteur $\delta b \in \mathbb{C}^n \setminus \{0\}$ tel que:

$$\| A^{-1}(\delta b) \| = \| A^{-1} \| \| \delta b \|.$$

Si $b := \lambda^{-1}(A(\delta b) + \lambda \delta b)$, alors, l'unique solution de l'équation $Ax = b$ est égale à:

$$x = \lambda^{-2}(\delta b) + A^{-2}(\delta b),$$

De sorte que l'unique solution de l'équation $A(\delta x) + \delta A(x + \delta x) = 0$ vaut :

$$\delta x = -\lambda(A + \lambda I_N)^{-2}(x) = -A^{-2}(\delta b).$$

Il vient donc :

$$x + \delta x = \lambda^{-2}(\delta b),$$

De sorte que :

$$\|\delta x\| = \|A^{-2}\| \|\delta b\| = \lambda \|A^{-2}\| \|x + \delta x\| = \text{cond}(A) \frac{\|\delta A\|}{\|A\|} \|x + \delta x\|.$$

(iii) Lorsque $\|\delta A\| \leq \frac{1}{2\|A^{-2}\|}$, la matrice $(\delta A)A^{-2}$ satisfait la condition :

$$\|(\delta A)A^{-2}\| \leq \frac{1}{2}.$$

La matrice $I_N + (\delta A)A^{-2}$ est donc inversible et son inverse vaut :

$$(I_N + (\delta A)A^{-2})^{-2} = \sum_{n=0}^{+\infty} (-(\delta A)A^{-2})^n.$$

En particulier, il vient :

$$\|(I_N + (\delta A)A^{-2})^{-2}\| \leq \sum_{n=0}^{+\infty} \|(\delta A)A^{-2}\|^n \leq 2.$$

Par linéarité, nous savons par ailleurs que :

$$(A + \delta A)\delta x = -(\delta A)x,$$

De sorte que :

$$\delta x = -A^{-2}(I_N + (\delta A)A^{-2})^{-2}(\delta A)x,$$

ce qui implique :

$$\|\delta x\| \leq \|A^{-2}\| \|(I_N + (\delta A)A^{-2})^{-2}\| \|\delta A\| \|x\|,$$

Soit :

$$\frac{\|\delta x\|}{\|x\|} \leq 2 \text{cond}(A) \frac{\|\delta A\|}{\|A\|}.$$

Corollaire : Soit $A \in \mathcal{G}\mathcal{L}_N(\mathbb{R})$ et $\delta A \in \mathcal{M}_N(\mathbb{R})$ telle que : $\|\delta A\| \leq \frac{1}{2\|A^{-2}\|}$. Pour $b \in$

$\mathbb{R}^N \setminus \{0\}$ et $\delta b \in \mathbb{R}^N$, désignons par x et $x + \delta x$ les solutions des équations

$Ax = b$, resp. $(A + \delta A)(x + \delta x) = b + \delta b$. Alors :

$$\frac{\|\delta x\|}{\|x\|} \leq 2 \text{cond}(A) \left(\frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right).$$

Preuve :

Sous les hypothèses du corollaire, la matrice $I_N + (\delta A)A^{-2}$ est inversible et son inverse satisfait l'inégalité :

$$\|(I_N + (\delta A)A^{-2})^{-2}\| \leq 2.$$

En particulier, la matrice $A + \delta A$ est également inversible, de sorte que le vecteur δx est bien défini. Il vient alors:

$$(A + \delta A) (\delta x) = \delta b - \delta A (x),$$

De sorte que:

$$\begin{aligned} \|\delta x\| &\leq \|A^{-1}\| \|(\mathbb{I}_n + (\delta A)A^{-1})^{-1}\| (\|\delta b\| + \|\delta A\| \|x\|) \\ &\leq 2 \|A^{-1}\| (\|\delta b\| + \|\delta A\| \|x\|). \end{aligned}$$

Pour $b = A(x)$, nous savons que:

$$\|b\| \leq \|A\| \|x\|,$$

ce qui suffit à montrer que:

$$\|\delta x\| \leq 2 \operatorname{cond}(A) \left[\frac{\|\delta b\|}{\|b\|} + \frac{\|\delta A\|}{\|A\|} \right] \|x\|.$$

En conclusion, les erreurs relatives dans la résolution du système linéaire $A(x) = b$ sont contrôlées par le conditionnement de la matrice A . Afin de minimiser les erreurs d'arrondi, il est pertinent de réduire la valeur du conditionnement de la matrice considérée. Ce préconditionnement des systèmes linéaires à résoudre repose (en général) sur la multiplication par une matrice $D \in \mathbb{R}^{n,n}(\mathbb{C})$ ad hoc (par exemple, diagonale) de façon à se ramener à un système linéaire équivalent: $DA(x) = D(b)$ tel que: $\operatorname{cond}(DA) < \operatorname{cond}(A)$.

La détermination des matrices D optimales pour cette propriété porte le nom de problème d'équilibrage. Il s'agit d'un problème pratique et important afin de contrôler les problèmes d'arrondi dans la résolution des systèmes linéaires.

3. Conditionnement d'un problème aux valeurs propres

Soit $N \in \mathbb{N}^*$. Notons $\mathcal{D}_N(\mathbb{C})$ l'ensemble des matrices diagonales de $\mathcal{M}_N(\mathbb{C})$, et considérons une norme $\|\cdot\|$ sur l'espace vectoriel \mathbb{C}^N telle que la norme subordonnée $\|\cdot\|$ à cette norme sur $\mathcal{M}_N(\mathbb{C})$ satisfasse l'identité:

$$\forall D \in \mathcal{D}_N(\mathbb{C}), \|D\| = \max_{1 \leq i \leq N} |D_{ii}|.$$

Exemple: Les normes $\|\cdot\|_1, \|\cdot\|_2, \|\cdot\|_\infty$ et $\|\cdot\|_\infty$ subordonnées aux normes $\|\cdot\|_2, \|\cdot\|_1$ et $\|\cdot\|_\infty$ satisfont cette propriété.

Définition: Soit $A \in \mathcal{M}_N(\mathbb{K})$ une matrice diagonalisable. Le conditionnement de la matrice A relativement au calcul de ses valeurs propres est le nombre strictement positif : $\nu(A) := \inf\{\text{cond}(R), R \in \mathcal{G}_{\mathbb{K}}(N) \text{ t.q. } RAR^{-1} \in \mathcal{D}_N(\mathbb{K})\}$.

Exemples: (i) $\forall D \in \mathcal{D}_N(\mathbb{K}), \nu(D) = 1$.

(ii) Si $\|\cdot\|_2$ désigne la norme subordonnée à la norme hermitienne canonique $\|\cdot\|_2$ sur \mathbb{R}^N , alors :

$$\forall N \in \mathcal{M}_N(\mathbb{K}), \nu_2(N) = 1$$

$$(\text{car: } \forall V \in \mathcal{M}_N(\mathbb{K}), \text{cond}_2(V) = 1).$$

Le conditionnement d'une matrice par rapport au calcul de ses valeurs propres permet d'évaluer les erreurs quant à la détermination des valeurs propres.

Lemme de Bauer - Fike: Soit $A \in \mathcal{M}_N(\mathbb{K})$ une matrice diagonalisable et $\delta A \in \mathcal{M}_N(\mathbb{K})$, alors : $\sigma(A + \delta A) \subset \sigma(A) + \mathcal{D}_\mathbb{K}(0, \nu(A) \|\delta A\|)$.

Preuve:

Soit $\sigma(A) = \{\lambda_1, \dots, \lambda_n\}$ et $\lambda \in \sigma(A + \delta A)$. Quel que soit le nombre ε strictement positif, considérons une matrice inversible $R \in \mathcal{G}_{\mathbb{K}}(N)$ telle que $RAR^{-1} = D \in \mathcal{D}_N(\mathbb{K})$ et :

$$\text{cond}(R) \leq \nu(A) + \varepsilon.$$

Notons que, si $\lambda \in \sigma(A)$, alors, il n'y a rien à démontrer, et supposons donc que $\lambda \notin \sigma(A)$.

Dans ce cas, la matrice $D - \lambda I_N$ est inversible, de sorte que :

$$R^{-1}(A + \delta A - \lambda I_N)R = D - \lambda I_N + R^{-1}(\delta A)R = (D - \lambda I_N) \left[I_N + (D - \lambda I_N)^{-1} R^{-1} \delta A R \right].$$

Comme la matrice $R^{-1}(A + \delta A - \lambda I_N)R$ n'est pas inversible, il est clair que :

$$\| (D - \lambda I_N)^{-1} R^{-1} \delta A R \| \geq 1,$$

De sorte que :

$$\begin{aligned} 1 &\leq \|(0 - \lambda I_n)^{-1}\| \operatorname{cond}(R) \| \delta A \| \\ &\leq \|(0 - \lambda I_n)^{-1}\| (\|A\| + \epsilon) \| \delta A \|. \end{aligned}$$

Par hypothèse, il vient :

$$\|(0 - \lambda I_n)^{-1}\| = \frac{1}{\min_{1 \leq j \leq n} |\lambda_j - \lambda|},$$

De sorte que :

$$\lambda \in \delta(A) + \mathcal{D}_\epsilon(0, (\|A\| + \epsilon) \| \delta A \|).$$

Les théorèmes d'ensuite à la limite $\epsilon \rightarrow 0$.

En conclusion, les erreurs dans la résolution d'un problème aux valeurs propres ne dépendent plus du conditionnement de la matrice désignée ailleurs considérée, mais de celui des matrices de passage associées à ce problème de diagonalisation. En particulier, les erreurs d'arrondi liées au problème aux valeurs propres d'une matrice normale sont peu importantes : ces matrices sont très bien conditionnées pour le problème aux valeurs propres. La détermination des valeurs singulières d'une matrice $A \in \mathbb{R}^{n,n}$, soit des valeurs propres de la matrice A^*A , est donc peu entachée d'erreurs d'arrondi.

III Méthodes directes pour la résolution de systèmes linéaires

Outre les méthodes itératives que nous étudierons dans la suite de ce cours, les méthodes directes de résolution de systèmes linéaires fournissent (au moins au niveau théorique) la (ou les) solution(s) exacte(s) du système considéré. Les possibles erreurs ne sont (en principe) liées qu'aux problèmes d'arrondi (et dépendent donc du conditionnement précédemment introduit).

À noter que ces méthodes ne passent pas par le calcul de l'inverse de la matrice considérée, qui réclame la résolution de N systèmes linéaires (pour chacun des vecteurs $(e_j)_{1 \leq j \leq N}$ de la base canonique de \mathbb{R}^N) pour une matrice carrée de taille $N \times N$.

En particulier, certaines de ces méthodes s'étendent à des systèmes linéaires rectangulaires (même si dans ce cas la solution n'est plus unique).

Cependant, dans le cas où une même matrice $A \in \mathcal{M}_n(\mathbb{K})$ intervient dans la résolution de plusieurs systèmes linéaires, il peut être utile de déterminer son inverse, ou au moins une factorisation de cette matrice, qui facilitera la résolution de ces systèmes.

1. Méthode de Gauss

Soit $N \in \mathbb{N}^*$. Pour $A \in \mathcal{M}_N(\mathbb{K})$ et $b \in \mathbb{K}^N$, considérons la résolution du système linéaire :

$$A(x) = b.$$

La méthode de Gauss pour résoudre ce système repose sur l'observation que la résolution est simple lorsque la matrice A est triangulaire.

Par exemple, si la matrice $A \in \mathcal{M}_N^+(\mathbb{K})$ est triangulaire supérieure, alors, le système $A(x) = b$ se réduit aux N équations :

$$\text{Et } \begin{cases} a_{11}x_1 + \dots + a_{1N}x_N = b_1 \\ \phantom{a_{11}x_1 +} \vdots \\ a_{NN}x_N = b_N \end{cases}$$

dans lequel les nombres $(a_{ij})_{2 \leq j \leq N}$ sont non nuls puisque la matrice A est inversible. Il est ainsi facile de calculer le nombre x_N , qui est égal à :

$$x_N = \frac{b_N}{a_{NN}},$$

puis d'en déduire le nombre :

$$x_{N-1} = \frac{1}{a_{N-1,N-1}} (b_{N-1} - a_{N-1,N}x_N),$$

et d'itérer ce procédé jusqu'à obtenir la valeur de :

$$x_1 = \frac{1}{a_{11}} [b_1 - a_{12}x_2 - \dots - a_{1N}x_N].$$

C'est la méthode de la remontée qui fournit la valeur de la solution

x en $\frac{N(N-1)}{2}$ soustractions, $\frac{N(N-1)}{2}$ multiplications et N divisions, soit

$$\Pi A(x) = \Pi(b),$$

par la méthode de la remontée précédemment introduite.

Sur le plan théorique, la méthode de Gauss conduit à l'énoncé suivant.

Théorème: Soit $A \in \mathcal{M}_N(\mathbb{C})$. Il existe une matrice inversible $\Pi \in \mathcal{GL}_N(\mathbb{C})$ telle que la matrice $\Pi A \in \mathcal{U}_N^+(\mathbb{C})$ soit triangulaire supérieure.

Preuve:

Par la méthode de Gauss précédemment décrite, il est clair que le théorème est vrai lorsque la matrice A est inversible.

Lorsqu'elle n'est plus inversible, il est possible qu'une des colonnes à éliminer soit identiquement nulle. Dans ce cas, cette colonne est déjà sous forme triangulaire supérieure, et il ne reste plus qu'à traiter les autres colonnes.

Sur le plan pratique, le nombre d'opérations élémentaires requises pour l'application de la méthode de Gauss est de l'ordre de N^3 additions ou soustractions, et de N^3 multiplications ou divisions, ce qui demeure acceptable. Plus précisément, chaque étape de la procédure d'élimination nécessite de l'ordre de $(N-k)^2$ additions ou soustractions, et $(N-k)$ multiplications ou divisions, où l'entier k varie entre 1 et $N-1$. Au total, il est donc nécessaire de réaliser de l'ordre de N^3 additions ou soustractions, et N^3 multiplications ou divisions pour cette première étape d'élimination. Les autres étapes, soit le calcul du vecteur $\Pi(b)$, et la méthode de la remontée, ne requièrent que de l'ordre de N^2 additions ou soustractions, et N^2 multiplications ou divisions.

Un autre point pratique important consiste en le choix des pivots. Plus ceux-ci sont grands, moins les erreurs d'arrondis sont importantes. La stratégie du pivot partiel consiste à choisir le pivot à l'étape $k \in \{1, \dots, N-1\}$ comme l'élément

ments $a_{i_0 k}$ de la $k^{\text{ième}}$ colonne tel que :

$$|a_{i_0 k}| = \max_{k \leq i \leq N} |a_{i k}|$$

La stratégie du pivot total consiste quant à elle à choisir le pivot à l'étape $1 \leq k \leq N-1$ comme l'élément $a_{i_0 j_0}$ tel que :

$$|a_{i_0 j_0}| = \max_{k \leq i, j \leq N} |a_{i j}|$$

Il est à noter que cette seconde stratégie nécessite de permuter, si nécessaire, des colonnes de la matrice A , ce qui revient à multiplier celle-ci par une matrice de permutation à droite. L'étape finale de remontée se trouve ainsi compliquée.

Par ailleurs, ces deux méthodes limitent certes les erreurs d'arrondis, mais entraînent un surcoût en opérations élémentaires.

Pour conclure, rappelons que la méthode de Gauss peut se prolonger en méthode de Gauss-Jordan afin d'obtenir in fine une matrice diagonale et qu'elle s'applique aussi dans le cas de matrices rectangulaires. Sur le plan théorique, elle permet ainsi de démentir le résultat suivant.

Lemme : Soit $n \in \mathbb{N}^*$ et $A \in \mathcal{M}_{n, n}(\mathbb{R})$. Il existe deux matrices inversibles $P \in \mathcal{G}_n(\mathbb{R})$ et $Q \in \mathcal{G}_n(\mathbb{R})$ et un entier $r \in \{1, \dots, n\}$ tels que :

$$P A Q = \begin{pmatrix} I_r & 0 \\ 0 & 0 \end{pmatrix}.$$

Notons enfin que la méthode de Gauss fournit aussi l'inverse de la matrice $A \in \mathcal{M}_n(\mathbb{R})$ considérée, ce qui se révèle utile lorsqu'il s'agit de résoudre le système linéaire $A(x) = b$ pour un grand nombre de valeurs du vecteur $b \in \mathbb{R}^n$.

2. Factorisation LU d'une matrice

L'application de la méthode de Gauss conduit au résultat de factorisation

suivant.

Théorème: Soit $A \in \mathcal{GL}_N(\mathbb{C})$. Pour $2 \leq k \leq N$, introduisons les matrices $A_k \in \mathcal{M}_k(\mathbb{C})$ définies par: $\forall 2 \leq i, j \leq k, (A_k)_{i,j} = a_{ij}$. Si les matrices $(A_k)_{2 \leq k \leq N}$ sont inversibles, alors, il existe une matrice triangulaire inférieure $L \in \mathcal{GL}_N(\mathbb{C})$, avec $\forall 2 \leq i \leq N, l_{ii} = 1$, et une matrice triangulaire supérieure $U \in \mathcal{GL}_N(\mathbb{C})$ telles que: $A = LU$. De plus, cette factorisation est unique.

Preuve:

L'existence des matrices L et U résulte de l'application de la méthode de Gauss. Au rang $k=2$, le coefficient a_{22} est non nul, de sorte que l'on peut le choisir comme pivot pour la première colonne. Il existe ainsi une matrice $E_2 =$

$$E_2 = \begin{pmatrix} 1 & & & \\ -\frac{a_{21}}{a_{22}} & 1 & & \\ \frac{a_{31}}{a_{22}} & & 1 & \\ \vdots & & & \ddots \\ -\frac{a_{n1}}{a_{22}} & & & & 1 \end{pmatrix} \text{ telle que: } E_2 A = \begin{pmatrix} a_{11} & * & & \\ 0 & A'_2 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix}$$

alors avoir construit des matrices $(E_j)_{2 \leq j \leq k-2}$ de la forme:

$$\forall 2 \leq j \leq k-2, E_j = \begin{pmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & \ddots \\ & & & & 1 \end{pmatrix}$$

telles que: $E_{k-2} \dots E_2 A =$ $\begin{pmatrix} a_{11} & * & & \\ & a_{22} & * & \\ & & \ddots & \\ & & & 1 \end{pmatrix}$. Les règles usuelles de multiplication par blocs permettent d'écrire cette identité sous la forme:

$$E_{k-2} \dots E_2 A = \begin{pmatrix} E_k & 0 \\ * & I \end{pmatrix} \begin{pmatrix} A_k & * \\ * & * \end{pmatrix} = \begin{pmatrix} a_{11} & * & * \\ & a_{22} & * \\ & & \ddots & \\ & & & 1 \end{pmatrix}, \text{ avec } E_k = \begin{pmatrix} 1 & 0 \\ & \ddots \\ & & 1 \end{pmatrix}$$

De sorte que:

$$E_k A_k = \begin{pmatrix} a_{11} & * \\ & a_{kk} \end{pmatrix}$$

Comme la matrice A_k est inversible, il est clair que le nombre a_{kk} est non nul. Il est donc possible d'introduire une matrice $E_k = \begin{pmatrix} 1 & 0 \\ & \ddots \\ & & 1 \end{pmatrix}$ telle que:

$$E_k \dots E_2 A = \begin{pmatrix} a_{11} & * & & \\ & a_{22} & * & \\ & & \ddots & \\ & & & 1 \end{pmatrix}$$

En définitive, il est permis de trouver une matrice triangulaire inférieure

$$E = E_{N-2} \dots E_2 = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ * & & 1 \end{pmatrix} \text{ telle que:}$$

$$EA \in \mathcal{E}_N^+(\mathbb{K}).$$

Comme l'inverse $L = E^{-1}$ de la matrice E est de même forme, il existe une matrice triangulaire supérieure $U \in \mathcal{E}_N^+(\mathbb{K})$ telle que:

$$A = LU.$$

De plus, s'il existe deux autres matrices $L' \in \mathcal{E}_N^-(\mathbb{K})$, avec $\forall 2 \leq i \leq N, l'_{ii} = 1$, et $U' \in \mathcal{E}_N^+(\mathbb{K})$ telles que:

$$A = L'U',$$

alors:

$$LU = L'U' \Rightarrow (L'^{-1})L = U'(U^{-1}).$$

Comme $(L'^{-1})L \in \mathcal{E}_N^-(\mathbb{K})$ et $U'(U^{-1}) \in \mathcal{E}_N^+(\mathbb{K})$, les matrices $(L'^{-1})L$ et $U'(U^{-1})$ sont diagonales, et dans ce cas:

$$(L'^{-1})L = I_N \Rightarrow L = L' \Rightarrow U' = U,$$

D'où l'unicité de la factorisation LU .

En pratique, l'application de la méthode de Gauss fournit de manière immédiate la valeur de la matrice triangulaire supérieure $U \in \mathcal{E}_N^+(\mathbb{K})$, mais également celle de la matrice triangulaire inférieure $L \in \mathcal{E}_N^-(\mathbb{K})$.

Si nous revenons à la preuve précédente et notons:

$$E_k = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix},$$

pour chaque entier $k \in \{1, \dots, N-1\}$, alors, l'inverse de la matrice E_k est égal à:

$$E_k^{-1} = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix},$$

Et le produit de ces inverses, soit la matrice L , à:

$$L = E_1^{-1} \dots E_{N-1}^{-1} = \begin{pmatrix} 1 & & 0 \\ -l_{21} & \ddots & \\ * & & 1 \end{pmatrix}.$$

L'intérêt de la factorisation LU réside dans la possibilité de résoudre

rapidement de manière répétée des systèmes linéaires de la forme $A(x) = b$ pour différentes valeurs du vecteur $b \in \mathbb{R}^n$. Ces systèmes sont de fait équivalents aux deux systèmes triangulaires suivants :

$$L(y) = b \quad \text{et} \quad V(x) = y,$$

qui peuvent se résoudre par la méthode de la remontée. Le coût en opérations élémentaires est de l'ordre de N^2 au lieu d'un coût de l'ordre de N^3 par la méthode de Gauss, et ce pour chaque résolution du système linéaire $A(x) = b$.

Notons que ce coût est encore réduit lorsque la matrice A est triangulaire. Ceci résulte du théorème suivant.

Théorème: Soit $A = \begin{pmatrix} b_1 & c_1 & 0 \\ a_2 & \dots & c_{n-2} \\ 0 & \dots & b_n \end{pmatrix} \in \mathcal{GL}_n(\mathbb{C})$ une matrice triangulaire. Pour $2 \leq h \leq n$,

introduisons les matrices $A_h \in \mathcal{M}_h(\mathbb{C})$ définies par: $\forall 2 \leq i, j \leq h$,

$$(A_h)_{i,j} = a_{ij}, \quad \text{et notons: } \delta_h = \det(A_h).$$

(i) Les nombres $(\delta_h)_{2 \leq h \leq n}$ sont donnés par les relations de récurrence

$$\delta_2 = b_2 \quad \text{et} \quad \forall 2 \leq h \leq n, \quad \delta_h = b_h \delta_{h-2} - a_h c_{h-2} \delta_{h-2},$$

où nous avons noté $\delta_0 = 1$.

(ii) Si les nombres $(\delta_h)_{2 \leq h \leq n}$ sont tous non nuls, la factorisation LU de la matrice A est donnée par les matrices:

$$L = \begin{pmatrix} 1 & & & 0 \\ a_{21} & \dots & & \\ 0 & \dots & \dots & 1 \\ & & & \dots & \dots & 1 \\ 0 & & & & & \dots & 1 \end{pmatrix} \quad \text{et} \quad V = \begin{pmatrix} b_1 & c_1 & 0 \\ b_2 & & & \\ & \dots & \dots & c_{n-2} \\ & & & \dots & \dots & b_n \\ & & & & & \dots & b_n \end{pmatrix}.$$

Preuve:

(i) La formule de récurrence résulte du développement par rapport à la dernière ligne du déterminant δ_h .

(ii) Si les nombres $(\delta_h)_{2 \leq h \leq n}$ sont tous non nuls, alors, le théorème précédent garantit l'existence et l'unicité de la factorisation LU de la matrice A .

Un calcul direct permet alors d'établir que les matrices L et V ci-dessus

satisfait :

$$[LU]_{1,1} = \frac{\delta_{11}}{\delta_{00}} = b_1 \text{ et } \forall 2 \leq k \leq N, [LU]_{k,k} = \frac{\delta_k + a_{k2}c_{2k} + \dots + a_{k,k-1}c_{k-1,k}}{\delta_{k-1}} = b_k,$$

et, de même,

$$\forall 1 \leq k \leq N-1, [LU]_{k+1,k} = a_{k+1,k} \text{ et } [LU]_{k,k+1} = c_k,$$

De sorte que nous avons bien l'égalité :

$$LU = A$$

Les formules des matrices L et U du théorème permettent alors de calculer les solutions des systèmes triangulaires :

$$L(y) = b \text{ et } U(x) = y,$$

en un nombre d'opérations élémentaires de l'ordre de N^2 , ce qui est un gain considérable par rapport à la méthode de Gauss.

3. Factorisation de Cholesky d'une matrice symétrique définie positive

Comme dans le cas tridiagonal, la factorisation LU d'une matrice A est particulièrement bien adaptée au cas où cette matrice $A \in \mathcal{S}_N^{++}(\mathbb{R})$ est symétrique définie positive. La factorisation de Cholesky, qui repose sur cette remarque se présente sous la forme suivante.

Théorème : Soit $A \in \mathcal{S}_N^{++}(\mathbb{R})$. Il existe une unique matrice triangulaire inférieure $B \in \mathcal{L}_N(\mathbb{R})$ telle que :

$$(i) \quad \forall 1 \leq i \leq N, B_{ii} > 0;$$

$$(ii) \quad A = B \circ B.$$

Preuve :

Pour $1 \leq k \leq N$, introduisons les matrices $A_k \in \mathcal{M}_k(\mathbb{R})$ définies par :

$$\forall 1 \leq i, j \leq k, [A_k]_{i,j} = A_{ij}.$$

Si q désigne la forme quadratique associée à la matrice A , alors, les matrices A_k sont les matrices des formes quadratiques q_k restrictions de q aux

sous-espaces vect (e₁, ..., e_h), où B = (e₁, ..., e_N) désigne la base canonique de l'espace vectoriel E^N. En particulier, les matrices A_h sont symétriques définies positives, de sorte qu'elles sont inversibles. Il est donc permis d'introduire la décomposition (L, U) ∈ GL⁻(R) × GL⁺(R) de la matrice A, avec ∀ 1 ≤ i ≤ N, L_{ii} = 1.

Obtens alors que pour tout entier 1 ≤ h ≤ N:

$$A = L U = \begin{pmatrix} 1 & 0 & 0 \\ * & \ddots & 0 \\ * & * & * \end{pmatrix} \begin{pmatrix} u_{11} & * & * \\ 0 & u_{22} & * \\ 0 & 0 & u_{hh} & * \end{pmatrix},$$

De sorte que:

$$\det(A_h) = \prod_{j=1}^h u_{jj}.$$

Comme tous ces déterminants sont strictement positifs, tous les coefficients (u_{jj})_{1 ≤ j ≤ N} sont également strictement positifs. Il est donc permis d'introduire la matrice diagonale D ∈ D_N(R) définie par:

$$\forall 1 \leq j \leq N, D_{jj} = \sqrt{u_{jj}} > 0.$$

Obtens alors:

$$B = L D \text{ et } C = D^{-1} U.$$

Il vient:

$$A = B C,$$

De sorte que par symétrie de la matrice A:

$${}^t C {}^t B = A \Rightarrow C ({}^t B)^{-1} = B^{-1} ({}^t C).$$

Remarquons alors que ces deux matrices sont de la forme:

$$C ({}^t B)^{-1} = \begin{pmatrix} 1 & * \\ 0 & \ddots \end{pmatrix} \text{ et } B^{-1} ({}^t C) = \begin{pmatrix} 1 & 0 \\ * & \ddots \end{pmatrix}.$$

Il s'ensuit que:

$$C ({}^t B)^{-1} = B^{-1} ({}^t C) = I_N \Rightarrow C = {}^t B.$$

Nous avons donc établi l'existence d'une matrice B ∈ GL⁺(R) telle que:

$$\forall 1 \leq j \leq N, B_{jj} > 0.$$

Pour terminer, notons que si une matrice B satisfait aux conclusions du théorème,

alors, la matrice diagonale D définie par:

$$\forall 1 \leq j \leq N, D_{jj} = b_{jj} > 0,$$

permet d'établir la factorisation suivante de la matrice A ,

$$A = (BD^{-1})(D^t B),$$

qui n'est autre que la factorisation LU de A . Nous avons donc:

$$L = BD^{-1} \text{ et } U = D^t B.$$

En particulier, il vient:

$$\forall 1 \leq j \leq N, U_{jj} = b_{jj}^2 \Rightarrow b_{jj} = \sqrt{U_{jj}},$$

Puis:

$$B = LD,$$

ce qui établit l'unicité de la matrice B .

En pratique, la factorisation de Cholesky d'une matrice $A \in \mathbb{R}^{N \times N}$ symétrique définie positive se calcule par une méthode de coefficients indéterminés. L'identité $A = B^t B$ s'écrit sous la forme:

$$\forall 1 \leq i, j \leq N, A_{ij} = \sum_{k=1}^N b_{ik} b_{jk} = \sum_{k=1}^{\min(i,j)} b_{ik} b_{jk}.$$

Comme la matrice A est symétrique, il suffit de résoudre ces équations pour $i \leq j$. Les calculs des coefficients b_{ij} se font alors par récurrence sur l'ordre $i \in \{1, \dots, N\}$. Pour $i=1$, il vient:

$$\forall 1 \leq j \leq N, A_{1j} = b_{11} b_{j1} \Rightarrow \begin{cases} b_{11} = \sqrt{A_{11}}, \\ \forall 1 \leq j \leq N, b_{j,1} = \frac{A_{1j}}{\sqrt{b_{11}}}, \end{cases}$$

et il est alors possible de résoudre les équations précédentes pour $i=2$, puis $i=3$, et ainsi de suite.

Une fois la factorisation de Cholesky déterminée, la méthode de Cholesky pour la résolution du système linéaire $A(x) = b$ repose sur la résolution des deux systèmes triangulaires:

$$B(y) = b \text{ et } D^t B(x) = y,$$

par la méthode de la remontée comme dans le cas de la factorisation LU .

Il est donc possible de conclure que cette méthode nécessite de l'ordre de la moitié des opérations élémentaires nécessaires à la méthode de Gauss,

sont néanmoins de l'ordre de N^3 opérations élémentaires. Ceci repose sur l'observation que, dans le cas où la matrice A est symétrique, la connaissance (d'un peu plus) de la moitié de ces coefficients suffit à la déterminer entièrement.

Pour terminer ce paragraphe sur les méthodes directes de résolution des systèmes linéaires, rappelons qu'il existe d'autres méthodes (comme la méthode C.R.) qui permettent aussi de résoudre de manière exacte de tels systèmes. Nous renvoyons à l'ouvrage "Introduction à l'analyse numérique matricielle et à l'optimisation" de Ph. Giarlet pour de plus amples détails.

IV Méthodes itératives pour la résolution de systèmes linéaires

Contrairement aux méthodes directes, les méthodes itératives ne cherchent pas à déterminer (au moins au niveau théorique) la (ou les) solution(s) exacte(s) du système à résoudre. Il s'agit au contraire de construire une suite de solutions approchées dont la limite est (l'une des ou) la solution(s) du système analysé.

Le principe de base pour la construction de cette suite est le théorème du point fixe comme nous allons maintenant le voir. Deux questions se posent quant à l'usage de ce théorème: d'une part, la prouva rigoureuse de la convergence de la suite de solutions approchées vers (vers) la solution exacte, d'autre part l'analyse de la vitesse de convergence de cette suite vers sa (possibles) limite. Et notez ici aussi la nécessité de tenir compte en pratique des erreurs d'arrondis (qui dépendent du conditionnement précédemment introduit).

1. Principe général des méthodes itératives

Soit $N \in \mathbb{R}^n$. Considérons une matrice inversible $A \in \mathcal{G}_{n,n}(\mathbb{R})$ et un vecteur $b \in \mathbb{R}^n$ et intéressons-nous à la résolution du système $A(x) = b$. Le principe des méthodes itératives consiste à déterminer une matrice $B \in \mathcal{M}_{n,n}(\mathbb{R})$ telle que la matrice

$I_n - B$ soit inversible, et un vecteur $c \in \mathbb{R}^n$, tels que le système $A(x) = b$ soit équivalent au système linéaire:

$$x = B(x) + c.$$

Une solution $x \in \mathbb{R}^n$ est alors l'unique point fixe de la fonctionnelle F définie par:

$$\forall y \in \mathbb{R}^n, F(y) = B(y) + c.$$

Il est alors possible de construire une suite de solutions approchées $(x_n)_{n \in \mathbb{N}}$ par la méthode du point fixe de Picard. Étant donné un vecteur $x_0 \in \mathbb{R}^n$, il s'agit de définir la suite $(x_n)_{n \in \mathbb{N}}$ par la formule de récurrence:

$$\forall n \in \mathbb{N}, x_{n+1} = F(x_n) = B(x_n) + c.$$

Si la suite $(x_n)_{n \in \mathbb{N}}$ converge vers un vecteur $x \in \mathbb{R}^n$, alors, ce vecteur x sera l'unique solution du système $x = B(x) + c$, soit du système $A(x) = b$.

Nous pouvons établir la convergence de cette méthode sous les hypothèses suivantes.

Théorème: Soit $B \in \mathcal{M}_n(\mathbb{R})$ telle que la matrice $I_n - B$ soit inversible et soit quel que soit le vecteur $x_0 \in \mathbb{R}^n$, considérons la suite $(x_n)_{n \in \mathbb{N}}$ définie par la formule de récurrence:

$$\forall n \in \mathbb{N}, x_{n+1} = B(x_n) + c.$$

(i) Les suites $(x_n)_{n \in \mathbb{N}}$ sont toutes convergentes si le rayon spectral $\rho(B)$ de la matrice B satisfait l'inégalité:

$$\rho(B) < 1.$$

(ii) Dans ce cas, la limite x des suites $(x_n)_{n \in \mathbb{N}}$ est l'unique solution de l'équation:

$$x = B(x) + c.$$

Preuve:

Comme la matrice $I_n - B$ est inversible, il existe un unique vecteur $x \in \mathbb{R}^n$

tel que:

$$x = B(x) + c.$$

Il vient alors:

$$\forall n \in \mathbb{N}, x_{n+2} - x_n = b(x_n - x_n),$$

De sorte que, par récurrence sur l'entier $n \in \mathbb{N}$,

$$\forall n \in \mathbb{N}, x_n - x = b^n(x_0 - x).$$

Si $\rho(b) < 1$, alors,

$$b^n \xrightarrow[n \rightarrow +\infty]{} 0 \text{ dans } \mathcal{M}_N(\mathbb{K}),$$

Et, dans ce cas,

$$x_n \xrightarrow[n \rightarrow +\infty]{} x \text{ dans } \mathbb{K}^N.$$

Réciproquement, si la suite $(x_n)_{n \in \mathbb{N}}$ est convergente (quel que soit le choix du vecteur x_0), alors, sa limite y satisfait l'équation:

$$y = b(y) + c,$$

De sorte que $y = x$. Il s'ensuit que:

$$\forall x_0 \in \mathbb{K}^N, b^n(x_0 - x) \xrightarrow[n \rightarrow +\infty]{} 0 \text{ dans } \mathbb{K}^N.$$

Ceci implique que:

$$\max_{1 \leq i, j \leq N} |(b^n)_{ij}| \xrightarrow[n \rightarrow +\infty]{} 0.$$

Comme cette quantité est une norme sur l'espace vectoriel $\mathcal{M}_N(\mathbb{K})$, il vient:

$$b^n \xrightarrow[n \rightarrow +\infty]{} 0 \text{ dans } \mathcal{M}_N(\mathbb{K}),$$

D'où il découle que $\rho(b) < 1$.

La valeur du rayon spectral $\rho(b)$ de la matrice b détermine donc la convergence de la méthode itérative correspondante. Cette valeur contrôle également la vitesse de convergence comme le souligne l'énoncé suivant.

théorème: Soit $b \in \mathcal{M}_N(\mathbb{K})$ telle que la matrice $I_N - b$ soit inversible et $c \in \mathbb{K}^N$.

Considérons l'unique solution $x \in \mathbb{K}^N$ de l'équation $x = b(x) + c$, et quel que soit le vecteur $x_0 \in \mathbb{K}^N$, la suite $(x_n)_{n \in \mathbb{N}}$ associée par la formule de récurrence:

$$\forall n \in \mathbb{N}, x_{n+2} = b(x_n) + c.$$

Si elle que soit la norme $\|\cdot\|$ sur \mathbb{K}^N , nous avons la convergence:

$$\sup \left\{ \|x_n - x\| \frac{1}{n}, x_0 \in \mathbb{K}^N \text{ t.q. } \|x_0 - x\| = 1 \right\} \xrightarrow[n \rightarrow +\infty]{} \rho(b).$$

Preuve:

Let $\|\cdot\|$ la norme subordonnée à la norme $\|\cdot\|$ sur l'espace vectoriel \mathbb{C}^n . Rappelons que la formule du rayon spectral garantit que:

$$\|B^n\| \xrightarrow[n \rightarrow +\infty]{\frac{1}{n}} \rho(B).$$

Exemple

$$\begin{aligned} \text{Lyp} \{ \|x_{n+1} - x_n\|, x_0 \in \mathbb{C}^n \text{ b.g. } \|x_0 - x\| = 1 \} &= \text{Lyp} \{ \|B^n(x_0 - x)\|, x_0 \in \mathbb{C}^n \text{ b.g. } \\ &\quad \|x_0 - x\| = 1 \} \\ &= \|B^n\|, \end{aligned}$$

il vient:

$$\text{Lyp} \{ \|x_{n+1} - x_n\| \xrightarrow[n \rightarrow +\infty]{\frac{1}{n}} x_0 \in \mathbb{C}^n \text{ b.g. } \|x_0 - x\| = 1 \} \rightarrow \rho(B).$$

Ce second théorème permet en particulier de comparer deux méthodes itératives pour la résolution d'un même système linéaire $A(x) = b$. Plus le rayon de convergence $\rho(B)$ de la matrice B sera faible, et meilleure sera la méthode. Et notez que dans tous les cas la convergence sera géométrique (lorsque $\rho(B) < 1$), ce qui est déjà très rapide.

2. Méthodes de Jacobi et de Gauss-Seidel

Let $N \in \mathbb{N}^*$. Considérons une matrice inversible $A \in \mathcal{G}_N(\mathbb{C})$ et un vecteur $b \in \mathbb{C}^N$, et intéressons-nous à la résolution du système $A(x) = b$.

Définition: Let $A \in \mathcal{G}_N(\mathbb{C})$. Une décomposition régulière de la matrice A est la donnée d'un couple de matrices $(M, N) \in \mathcal{G}_N(\mathbb{C}) \times \mathcal{M}_N(\mathbb{C})$ tel que:

(i) $A = M - N$,

(ii) La matrice $M \in \mathcal{G}_N(\mathbb{C})$ est facile à inverser.

La méthode itérative basée sur cette décomposition consiste alors à choisir un vecteur $x_0 \in \mathbb{C}^N$ et à considérer la suite $(x_n)_{n \geq 0}$ définie par la formule de récurrence:

$$\forall n \in \mathbb{N}, x_{n+1} = M^{-1}N(x_n) + M^{-1}(b).$$

Exemples: (i) Méthode de Jacobi:

Soit $A \in \mathcal{M}_n(\mathbb{C})$ et $D = \begin{pmatrix} a_{11} & & 0 \\ & \ddots & \\ 0 & & a_{nn} \end{pmatrix} \in \mathcal{D}_n(\mathbb{C})$. Si la matrice D

est inversible (c'est-à-dire si tous les coefficients diagonaux de A sont non nuls), alors, la méthode de Jacobi est la méthode itérative associée à la décomposition régulière $A = D - N$, où $D = D$ et $N = -A + D$. Cette méthode consiste donc à choisir un vecteur $x_0 \in \mathbb{C}^n$ et à considérer la suite $(x_n)_{n \in \mathbb{N}}$ définie par la formule de récurrence:

$$\forall n \in \mathbb{N}, x_{n+1} = (I_n - D^{-1}A)x_n + D^{-1}(b).$$

(ii) Méthode de Gauss-Jacobi:

Soit $A \in \mathcal{M}_n(\mathbb{C})$ et $T = \begin{pmatrix} a_{11} & & 0 \\ & \ddots & \\ a_{ij} & & a_{nn} \end{pmatrix} \in \mathcal{M}_n(\mathbb{C})$. Si la matrice

T est inversible (c'est-à-dire si tous les coefficients diagonaux de A sont non nuls), alors, la méthode de Gauss-Jacobi est la méthode itérative associée à la décomposition régulière $A = T - N$, où $T = T$ et $N = -A + T$. Cette méthode consiste donc à choisir un vecteur $x_0 \in \mathbb{C}^n$ et à considérer la suite $(x_n)_{n \in \mathbb{N}}$ définie par la formule de récurrence:

$$\forall n \in \mathbb{N}, x_{n+1} = (I_n - T^{-1}A)x_n + T^{-1}(b).$$

En pratique, lorsqu'elle est définie, la méthode de Jacobi consiste donc à construire la suite $(x_n)_{n \in \mathbb{N}}$ de la façon suivante:

$$\forall n \in \mathbb{N}, \forall 2 \leq i \leq n, (x_{n+1})_i = \frac{1}{a_{ii}} \left[- \sum_{j=2}^n a_{ij} (x_n)_j + b_i \right].$$

En particulier, l'inversion de la matrice D^{-1} est immédiate.

Il en va de même pour la méthode de Gauss-Jacobi qui consiste à construire la suite $(x_n)_{n \in \mathbb{N}}$ de façon à ce que:

$$\forall n \in \mathbb{N}, \forall 2 \leq i \leq n, (x_{n+1})_i = \frac{1}{a_{ii}} \left[- \sum_{j=2}^{i-1} a_{ij} (x_{n+1})_j - \sum_{j=i+1}^n a_{ij} x_n + b_i \right]$$

et note que la complexité de cette seconde méthode est identique à celle de la

méthode de Jacobi, et qu'elle mobilise de plus moins de places en mémoire.

En ce qui concerne la convergence de ces méthodes, nous avons le résultat suivant.

Théorème: Soit $A \in \mathcal{M}_N(\mathbb{K})$. Considérons une décomposition séculaire $(D, N) \in \mathcal{M}_N(\mathbb{K}) \times \mathcal{M}_N(\mathbb{K})$ de la matrice A et la méthode itérative associée, qui repose sur la formule de récurrence:

$$\forall n \in \mathbb{N}, x_{n+1} = D^{-1}N(x_n) + D^{-1}b.$$

La suite $(x_n)_{n \in \mathbb{N}}$ est alors convergente quel que soit le choix du vecteur $x_0 \in \mathbb{K}^N$ si: $\rho(D^{-1}N) < 1$.

Dans ce cas, la limite x de la suite $(x_n)_{n \in \mathbb{N}}$ est l'unique solution de l'équation: $A(x) = b$.

Preuve:

Comme

$$I_N - D^{-1}N = D^{-1}A \in \mathcal{M}_N(\mathbb{K}),$$

l'équivalence découle des résultats précédents, et si x désigne la limite de la suite $(x_n)_{n \in \mathbb{N}}$, alors:

$$x = D^{-1}N(x) + D^{-1}b \Rightarrow A(x) = b.$$

Cet énoncé permet d'établir que la méthode de Jacobi est bien définie et convergente dans le cas des matrices à diagonales dominantes.

Définition: Soit $A \in \mathcal{M}_N(\mathbb{K})$. Une matrice A est à diagonale dominante si:

$$\forall 1 \leq i \leq N, |a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^N |a_{ij}|.$$

Exemple: La matrice identité est à diagonale dominante.

Lemme: Soit $A \in \mathcal{M}_N(\mathbb{K})$. Si la matrice A est à diagonale dominante, alors elle est inversible.

Preuve:

Yait $x \in \mathbb{R}^N$ tel que : $A(x) = 0$, et $i_0 \in \{1, \dots, N\}$ tel que :

$$|x_{i_0}| = \max_{1 \leq i \leq N} |x_i|.$$

Comme $A(x) = 0$, il vient :

$$|a_{i_0 i_0}| |x_{i_0}| = \left| - \sum_{\substack{j=1 \\ j \neq i_0}}^N a_{i_0 j} x_j \right| \leq \sum_{\substack{j=1 \\ j \neq i_0}}^N |a_{i_0 j}| |x_{i_0}|$$

$\Rightarrow |x_{i_0}| = 0$

De sorte que le vecteur x est nul, et la matrice A , inversible.

Théorème: Yait $A \in \mathcal{G}(\mathbb{R}, \mathbb{R})$, une matrice à diagonale dominante. La méthode de Jacobi est bien définie et elle est convergente. Quel que soit le choix de vecteur $x_0 \in \mathbb{R}^N$, la suite $(x_n)_{n \in \mathbb{N}}$ définie par la formule de récurrence :

$$\forall n \in \mathbb{N}, x_{n+2} = (I_N - D^{-1}A)(x_n) + b,$$

où $D = \begin{pmatrix} a_{11} & 0 \\ 0 & \ddots \\ 0 & & a_{nn} \end{pmatrix}$ converge vers l'unique solution x de l'équation :

$$Ax = b.$$

Preuve:

Dans les notations précédentes, les matrices D et N de la décomposition régulière de A sont égales à :

$$D = D \text{ et } N = D - A$$

$$\Rightarrow D^{-1}N = I_N - D^{-1}A = \begin{pmatrix} 0 & -a_{1j} \\ \dots & \dots \\ -a_{ij} & \dots \\ \dots & \dots \\ -a_{ni} & 0 \end{pmatrix}$$

Notons en particulier que les coefficients diagonaux $(a_{ii})_{1 \leq i \leq N}$ sont tous non nuls puisque la matrice A est à diagonale dominante. Les méthode de Jacobi est donc bien définie. De plus, si $\lambda \in \mathbb{R}$ vérifie $|\lambda| \geq 1$, alors, la matrice

$\lambda I_N - D^{-1}N$ est à diagonale dominante, donc inversible. Aussi le nombre

λ n'est-il pas une valeur propre de la matrice $D^{-1}N$. En conclusion,

$$\rho(D^{-1}N) < 1,$$

Et, par le théorème précédent, la méthode de Jacobi est convergente.

Il en va de même pour la méthode de Gauss - Seidel.

Théorème: Soit $A \in \mathcal{G}(\mathbb{R}^n, \mathbb{R})$, une matrice à diagonale dominante. La méthode de Gauss - Seidel est bien définie et elle est convergente. Quel que soit le choix du vecteur $x_0 \in \mathbb{R}^n$, la suite $(x_n)_{n \in \mathbb{N}}$ définie par la formule de récurrence:

$$\forall n \in \mathbb{N}, x_{n+1} = (I_n - T^{-1}A)(x_n) + T^{-1}(b),$$

où $T = \begin{pmatrix} a_{11} & & 0 \\ & \ddots & \\ a_{ij} & & a_{nn} \end{pmatrix} \in \mathcal{G}(\mathbb{R}^n, \mathbb{R})$, converge vers l'unique solution $x \in \mathbb{R}^n$ de

$$l'équation: Ax = b.$$

Preuve:

La méthode de Gauss - Seidel est bien définie puisque les coefficients diagonaux $(a_{ii})_{1 \leq i \leq n}$ sont tous non nuls. Afin de montrer sa convergence, il s'agit de vérifier que:

$$\rho(I_n - T^{-1}A) < 1.$$

Soit donc $\lambda \in \sigma(I_n - T^{-1}A)$ et $x \in \mathbb{R}^n \setminus \{0\}$ tel que:

$$x - T^{-1}Ax = \lambda x.$$

Notons $i, j_0 \in \mathbb{N}$ tel que:

$$|x_{i_0}| = \max_{1 \leq i \leq n} |x_i| \neq 0.$$

Il vient alors:

$$\lambda T(x) = (T - A)(x) \Rightarrow \lambda \sum_{j=0}^{i_0} a_{i_0, j} x_j = - \sum_{j=i_0+1}^n a_{i_0, j} x_j$$

De sorte que:

$$|\lambda| |a_{i_0, i_0}| |x_{i_0}| \leq \left[|\lambda| \sum_{j=0}^{i_0} |a_{i_0, j}| + \sum_{j=i_0+1}^n |a_{i_0, j}| \right] |x_{i_0}|$$

$$\Rightarrow |\lambda| |a_{i_0, i_0}| < \max\{|\lambda|, 1\} |a_{i_0, i_0}|$$

De sorte que $|\lambda| < 1$ et $\rho(I_n - T^{-1}A) < 1$.

La méthode de Gauss - Seidel est également bien définie et convergente lorsque la matrice A est hermitienne définie positive.

Théorème de Ostrowski - Seidel: Soit $A \in \mathcal{H}(\mathbb{R}^n, \mathbb{R})$. La méthode de Gauss - Seidel est bien définie et elle est convergente. Quel que soit le choix du vecteur

$x_0 \in \mathbb{R}^N$, la suite $(x_n)_{n \in \mathbb{N}}$ définie par la formule de récurrence:

$$\forall n \in \mathbb{N}, x_{n+1} = (I_N - T^{-1}A)(x_n) + T^{-1}(b),$$

où $T = \begin{pmatrix} a_{11} & & 0 \\ & \ddots & \\ 0 & & a_{nn} \end{pmatrix} \in \mathcal{L}_\mathbb{R}(\mathbb{R}^N)$, converge vers l'unique solution $x \in \mathbb{R}^N$ de

$$\text{l'équation: } A(x) = b.$$

Preuve:

Comme la matrice A est hermitienne définie positive, ses coefficients diagonaux $(a_{ii})_{1 \leq i \leq N}$ sont tous non nuls de sorte que la méthode de Gauss-Jacobi est bien définie.

Considérons alors la norme $\|\cdot\|_A$ définie sur \mathbb{R}^N par l'expression:

$$\forall x \in \mathbb{R}^N, \|x\|_A = \left[\langle A(x), x \rangle_{\mathbb{R}^N} \right]^{\frac{1}{2}},$$

où $\langle \cdot, \cdot \rangle_{\mathbb{R}^N}$ désigne le produit hermitien canonique de \mathbb{R}^N . Soit $\|\cdot\|_A$ la norme subordonnée à la norme $\|\cdot\|_A$ sur l'espace vectoriel $\mathcal{M}_N(\mathbb{R})$. Soit $v \in \mathbb{R}^N$

tel que $\|v\|_A \leq 1$, alors,

$$\begin{aligned} \|v - T^{-1}A(v)\|_A^2 &= \langle A(v - T^{-1}A(v)), v - T^{-1}A(v) \rangle_{\mathbb{R}^N} \\ &= \|v\|_A^2 - \langle A(T^{-1}A(v)), v \rangle_{\mathbb{R}^N} - \langle A(v), T^{-1}A(v) \rangle_{\mathbb{R}^N} \\ &\quad + \langle A(T^{-1}A(v)), T^{-1}A(v) \rangle_{\mathbb{R}^N} \end{aligned}$$

Or,

$$A(v) = T(T^{-1}A(v)),$$

De sorte que:

$$\begin{aligned} \|v - T^{-1}A(v)\|_A^2 &= \|v\|_A^2 - \langle T^{-1}A(v), T(T^{-1}A(v)) \rangle_{\mathbb{R}^N} + \langle A - T \rangle(T^{-1}A(v)), \\ &\quad T^{-1}A(v) \rangle_{\mathbb{R}^N} \\ &= \|v\|_A^2 - \langle 0(T^{-1}A(v)), T^{-1}A(v) \rangle_{\mathbb{R}^N} \\ &\leq \|v\|_A^2 \leq 1 \end{aligned}$$

En particulier,

$$\|I_N - T^{-1}A\|_A = \max \{ \|v - T^{-1}A(v)\|_A, v \in \mathbb{R}^N \text{ t.q. } \|v\|_A \leq 1 \} \leq 1.$$

Il s'ensuit que $\rho(I_N - T^{-1}A) < 1$, puis que la méthode de Gauss-Jacobi est dans ce cas aussi convergente.

Il est enfin possible de comparer les deux méthodes dans le cas des matrices tridiagonales.

Exemple: Soit $A = \begin{pmatrix} b_{11} & c_{12} & 0 \\ a_{21} & & c_{22} \\ 0 & a_{31} & b_{33} \end{pmatrix} \in \mathcal{M}_3(\mathbb{K})$ une matrice triangulaire.

Considérons la matrice diagonale $D = \begin{pmatrix} b_{11} & 0 \\ & b_{33} \end{pmatrix} \in \mathcal{D}_2(\mathbb{K})$ et la

matrice triangulaire inférieure $T = \begin{pmatrix} b_{11} & 0 \\ a_{21} & \\ 0 & a_{31} & b_{33} \end{pmatrix} \in \mathcal{M}_3(\mathbb{K})$, et supposons que la matrice D est inversible.

Les rayons spectraux des matrices $I_N - D^{-1}A$ et $I_N - T^{-1}A$ satisfont l'identité:

$$\rho(I_N - T^{-1}A) = \rho(I_N - D^{-1}A)^2$$

En particulier, les méthodes de Jacobi et de Gauss-Seidel convergent ou divergent simultanément, et lorsque elles convergent, la méthode de Gauss-Seidel converge plus rapidement.

Preuve:

Soit

$$\forall \lambda \in \mathbb{C}, \begin{cases} P_J(\lambda) = \det(\lambda D + A - D), \\ \text{et } P_{GS}(\lambda) = \det(\lambda T + A - T). \end{cases}$$

$$\text{et } P_{GS}(\lambda) = \det(\lambda T + A - T).$$

Comme les matrices D et T sont inversibles, nous avons les équivalences:

$$\begin{cases} \forall \lambda \in \mathbb{C}, \lambda \in \sigma(I_N - D^{-1}A) \Leftrightarrow P_J(\lambda) = 0, \\ \text{et } \forall \lambda \in \mathbb{C}, \lambda \in \sigma(I_N - T^{-1}A) \Leftrightarrow P_{GS}(\lambda) = 0. \end{cases}$$

De plus,

$$\begin{aligned} \forall \lambda \in \mathbb{C}, P_{GS}(\lambda^2) &= \det(\lambda^2 T + A - T) \\ &= \det \begin{pmatrix} \lambda^2 b_{11} & c_{12} & 0 \\ \lambda^2 a_{21} & & c_{22} \\ 0 & \lambda^2 a_{31} & \lambda^2 b_{33} \end{pmatrix}, \end{aligned}$$

ou,

$$\begin{pmatrix} \lambda^2 b_{11} & c_{12} & 0 \\ \lambda^2 a_{21} & & c_{22} \\ 0 & \lambda^2 a_{31} & \lambda^2 b_{33} \end{pmatrix} = \begin{pmatrix} \lambda & 0 \\ & \lambda \end{pmatrix} \begin{pmatrix} \lambda^2 b_{11} \lambda & \lambda c_{12} & 0 \\ \lambda a_{21} & & \lambda c_{22} \\ 0 & \lambda a_{31} & \lambda b_{33} \end{pmatrix} \begin{pmatrix} \lambda^{-2} & 0 \\ & \lambda^{-2} \end{pmatrix}$$

De sorte que:

$$\forall \lambda \in \mathbb{C}, P_{GS}(\lambda^2) = \det \begin{pmatrix} \lambda^2 b_{11} \lambda & \lambda c_{12} & 0 \\ \lambda a_{21} & & \lambda c_{22} \\ 0 & \lambda a_{31} & \lambda b_{33} \end{pmatrix} = \lambda^N \det(\lambda D + A - D).$$

En conclusion, il vient:

$\lambda \in \sigma(I_N - T^{-1}A) \setminus \{0\}$ ou $\lambda \in \sigma(I_N - U^{-1}A) \setminus \{0\}$,

De sorte que :

$$\rho(I_N - T^{-1}A) = \rho(I_N - U^{-1}A)^2.$$

Les autres conclusions du théorème résultent alors des énoncés précédents.

Au moins dans le cas des matrices tridiagonales, la méthode de Gauss-Seidel est meilleure que celle de Jacobi. Et noter pour conclure que la méthode de Gauss-Seidel fait partie des méthodes dites de relaxation pour lesquelles il est encore permis d'étendre ce résultat de comparaison de convergence. Nous renvoyons à l'ouvrage "Introduction à l'analyse numérique matricielle et à l'optimisation" de Ph. Giarlet pour de plus amples détails.

V Méthodes de calculs des valeurs propres et des vecteurs propres

Le problème de la recherche des éléments propres d'une matrice $A \in M_N(\mathbb{C})$ est délicat. Quel que soit le problème d'aujourd'hui que nous avons déjà évoqué, il n'existe pas de méthode directe de recherche des valeurs propres au delà de la dimension $N=5$. En effet, il s'avérerait en effet possible de calculer les racines d'un polynôme de degré supérieur ou égal à 5.

Nous allons donc présenter deux méthodes itératives pour déterminer certaines des valeurs propres de la matrice $A \in M_N(\mathbb{C})$: les méthodes de la puissance itérée, et de la puissance inverse. Il existe bien sûr d'autres méthodes itératives pour déterminer les éléments propres d'une matrice, par exemple, les méthodes de Jacobi, de Givens-Householder, ou la méthode QR. Nous renvoyons à l'ouvrage "Introduction à l'analyse numérique matricielle et à l'optimisation" de Ph. Giarlet pour de plus amples détails.

2. Méthode de la puissance itérée

Soit $N \in \mathbb{N}^*$. Considérons une matrice $A \in \mathcal{M}_N(\mathbb{C})$ telle que A est diagonalisable, et notons $\sigma(A) = \{\lambda_1, \dots, \lambda_N\}$ son spectre avec multiplicité, et $(e_i)_{1 \leq i \leq N}$ une base de vecteurs propres de A telle que:

$$\forall 1 \leq i \leq N, A(e_i) = \lambda_i e_i$$

Quel vecteur $x \in \mathbb{C}^N$ se décompose sous la forme:

$$x = \sum_{i=1}^N x_i e_i,$$

On note que:

$$\forall n \in \mathbb{N}, A^n(x) = \sum_{i=1}^N x_i A^n(e_i) = \sum_{i=1}^N x_i \lambda_i^n e_i.$$

En particulier, s'il existe une valeur propre λ_1 de plus grande module, ce s'est - à - dire telle que:

$$\forall 2 \leq i \leq N, |\lambda_i| < |\lambda_1|,$$

et, si x est un vecteur de \mathbb{C}^N tel que $x_1 \neq 0$, alors,

$$\|A^n(x)\|_{\infty} \underset{n \rightarrow +\infty}{\sim} |x_1| |\lambda_1|^n \|e_1\|,$$

où $\|\cdot\|$ désigne une norme quelconque de \mathbb{C}^N . La norme des puissances itérées $A^n(x)$ est contrôlée par la plus grande valeur propre de A , ce qui permet donc d'estimer cette valeur propre. Cette remarque conduit à la méthode de la puissance itérée suivante.

Algorithme: Soit $A \in \mathcal{M}_N(\mathbb{C})$ une matrice diagonalisable. Notons $\sigma(A) = \{\lambda_1, \dots, \lambda_N\}$ son spectre avec multiplicité, et $(e_i)_{1 \leq i \leq N}$ une base de vecteurs propres de A telle que:

$$\forall 1 \leq i \leq N, \|e_i\|_2 = 1 \text{ et } A(e_i) = \lambda_i e_i,$$

où $\|\cdot\|_2$ désigne la norme hermitienne canonique de \mathbb{C}^N . Choisissons un vecteur $x_0 \in \mathbb{C}^N$ tel que $\|x_0\|_2 = 1$ et définissons la suite

$(x_n)_{n \in \mathbb{N}}$ par la formule de récurrence:

$$\forall n \in \mathbb{N}, x_{n+1} = \frac{A(x_n)}{\|A(x_n)\|_2}.$$

Posons enfin:

$$\forall n \in \mathbb{N}, \mu_n = \langle A(x_n), x_n \rangle_2.$$

Supposons que:

$$(i) \forall \lambda \leq \lambda_2 < \lambda_1$$

(ii) Le vecteur $x_0 = \sum_{i=2}^n (x_0)_i e_i$ vérifie la condition: $(x_0)_2 \neq 0$.

Alors, les suites $(x_n)_{n \in \mathbb{N}}$ et $(\mu_n)_{n \in \mathbb{N}}$ sont bien définies et vérifient:

$$(i) \mu_n \xrightarrow[n \rightarrow +\infty]{} \lambda_2$$

$$(ii) \exists c \in E_{\lambda_2}(A) \text{ t.q. } \frac{\|x_n\|}{\lambda_2^n} \xrightarrow[n \rightarrow +\infty]{} c.$$

Preuve:

Comme $\lambda_2 > 0$ et $(x_0)_2 \neq 0$, le vecteur $A(x_0)$ n'est pas nul, de sorte que le vecteur x_1 est bien défini et vaut:

$$x_1 = \frac{A(x_0)}{\|A(x_0)\|_2}$$

Supposons alors que la suite $(x_n)_{n \in \mathbb{N}}$ est bien définie jusqu'au rang n et que:

$$\forall 0 \leq k \leq n, x_k = \frac{A^k(x_0)}{\|A^k(x_0)\|_2}$$

En particulier, il existe un nombre strictement positif λ tel que:

$$x_n = \lambda A^n(x_0) = \lambda \sum_{i=2}^n (x_0)_i \lambda_i^n e_i$$

$$\Rightarrow x_n \notin \text{Ker}(A).$$

Le vecteur x_{n+1} est donc bien défini et:

$$x_{n+1} = \frac{A^{n+1}(x_0)}{\|A^{n+1}(x_0)\|_2} = \frac{A^{n+1}(x_0)}{\|A^n(x_0)\|_2} \times \frac{\|A^n(x_0)\|_2}{\|A^{n+1}(x_0)\|_2} = \frac{A^{n+1}(x_0)}{\|A^n(x_0)\|_2}$$

D'où le fait que les suites $(x_n)_{n \in \mathbb{N}}$ et $(\mu_n)_{n \in \mathbb{N}}$ sont bien définies et que:

$$\forall n \in \mathbb{N}, x_n = \frac{A^n(x_0)}{\|A^n(x_0)\|_2}$$

Comme

$$\forall n \in \mathbb{N}, A^n(x_0) = \sum_{i=2}^n (x_0)_i \lambda_i^n e_i,$$

il vient:

$$A^n(x_0) \underset{n \rightarrow +\infty}{\sim} (x_0)_2 \lambda_2^n e_2 \text{ et } \|A^n(x_0)\|_2 \underset{n \rightarrow +\infty}{\sim} |(x_0)_2| |\lambda_2|^n,$$

D'où on a que:

$$\mu_n \underset{n \rightarrow +\infty}{\sim} \frac{|\lambda_2|^{2n} \lambda_2 (x_0)_2^2}{|\lambda_2|^{2n} (x_0)_2^2} \xrightarrow[n \rightarrow +\infty]{} \lambda_2,$$

et

$$\frac{\|x_n\|}{\lambda_2^n} \xrightarrow[n \rightarrow +\infty]{} c = \frac{(x_0)_2}{|(x_0)_2|} e_2 \in E_{\lambda_2}(A).$$

L'hypothèse sur le vecteur x_0 n'est pas une difficulté en pratique. Des fait des exercices d'exercices, il est probable que l'un des vecteurs $(x_n)_{n \in \mathbb{N}}$ aura une

composante non nulle suivant le vecteur e_1 .

De plus, la décomposition en blocs de Jordan d'une matrice quelconque de $M_n(\mathbb{C})$ permet d'étendre la méthode aux matrices $A \in M_n(\mathbb{C})$ telles que $\exists \lambda \in \sigma(A)$ t.q. $\forall \mu \in \sigma(A) \setminus \{\lambda\}, |\mu| < |\lambda|$.

Pour contre, cette méthode ne fonctionne plus nécessairement lorsque la matrice A a deux valeurs propres différentes de même module.

Afin de déterminer les autres valeurs et vecteurs propres de la matrice, il est possible d'utiliser une méthode de déflation. Sous les hypothèses du théorème précédent, la transposée ${}^t A$ de la matrice A est elle-même diagonalisable et son spectre est identique à celui de A . Il existe donc une base $(f_i)_{1 \leq i \leq n}$ de vecteurs propres de ${}^t A$ tels que:

$$\forall 1 \leq i \leq n, {}^t A(f_i) = \lambda_i f_i \text{ et } \|f_i\|_2 = 1.$$

En appliquant la méthode de la puissance itérée aux matrices A et ${}^t A$, nous pouvons déterminer des valeurs approchées de λ_2 et des vecteurs propres e_2 et f_2 . Supposons alors que:

$${}^t f_2 e_2 \neq 0.$$

Dans ce cas, il est possible de définir la matrice:

$$B = A - \lambda_2 \frac{e_2 e_2^t f_2}{{}^t f_2 e_2}.$$

Il vient alors:

$$B(e_2) = 0$$

et,

$$\forall 2 \leq j \leq n, {}^t f_2 A e_j = \lambda_j {}^t f_2 e_j = \lambda_2 {}^t f_2 e_j \Rightarrow {}^t f_2 e_j = 0,$$

De sorte que:

$$\forall 2 \leq j \leq n, B(e_j) = \lambda_j e_j.$$

Le spectre de la matrice B est donc égal à:

$$\sigma(B) = \{0, \lambda_0, \dots, \lambda_n\}$$

Si cette matrice vérifie les hypothèses du théorème précédent, alors, il est possible de calculer sa plus grande valeur propre par la méthode de la puissance itérée, et d'en déduire une seconde valeur propre de A , et une

vecteur propre associé.

Il faut cependant que cette méthode de déflation se borne au problème des auto-
diagonal en raison desquels les identités précédentes ne sont pas exactes. En pratique,
il est préférable d'utiliser la méthode de la puissance inverse pour déterminer
les autres éléments propres de la matrice A .

2. Méthode de la puissance inverse

Soit $M \in \mathbb{R}^n$. Afin de déterminer les valeurs propres d'une matrice $A \in \mathbb{R}^{n \times n}$,
la méthode de la puissance inverse repose sur un argument de localisation
des valeurs propres de A , puis sur la méthode de la puissance itérée afin
de déterminer leurs valeurs approchées.

Plus précisément, la première étape consiste à déterminer des valeurs approchées
des valeurs propres de la matrice A . Elle repose sur des arguments de localisation
tels que le théorème suivant.

Théorème de Gershgorin - Hadamard : Soit $A \in \mathbb{R}^{n \times n}$. Si $\lambda \in \sigma(A)$, alors,
il existe $i \in \{1, \dots, n\}$ t.q. $|\lambda - a_{ii}| \leq \sum_{j \neq i} |a_{ij}|$

Preuve :

Par l'absurde, si cette conclusion n'est pas vérifiée, alors, la matrice
 $\lambda I_n - A$ est à diagonale dominante, donc inversible, ce qui est contradictoire.

Une fois localisée grossièrement une valeur propre λ_0 de A à travers une valeur
approchée λ_0 , deux cas de figures peuvent se produire :

- si $\lambda_0 \in \sigma(A)$, alors, le problème est résolu et il est possible de chercher
les autres valeurs propres de A ;
- sinon, la matrice $A - \lambda_0 I_n$ est inversible et le spectre de son inverse
est donné par les nombres $\frac{1}{\mu - \lambda_0}$ où $\mu \in \sigma(A)$. En particulier, si λ_0 est

suffisamment proche de λ , alors, la valeur propre $\frac{1}{\lambda - \lambda_0}$ sera la valeur propre dominante de la matrice $(A - \lambda_0 I_N)^{-1}$ et il sera possible de la calculer par la méthode de la puissance itérée. D'un point de vue théorique, cette discussion est résumée par l'énoncé suivant.

Théorème: Soit $A \in \mathcal{M}_N(\mathbb{C})$ une matrice diagonalisable. Soient $\sigma(A) = \{\lambda_1, \dots, \lambda_N\}$ son spectre avec multiplicité, et $(e_i)_{1 \leq i \leq N}$ une base de vecteurs propres de A telle que:

$$\forall 1 \leq i \leq N, \|e_i\|_2 = 1 \text{ et } A(e_i) = \lambda_i e_i,$$

où $\|\cdot\|_2$ désigne la norme hermitienne canonique de \mathbb{C}^N .

Supposons alors donné un nombre complexe μ tel que:

$$\mu \neq \lambda_i \text{ et } \forall 1 \leq j \leq N, |\mu - \lambda_j| > |\mu - \lambda_i|.$$

Choisissons alors un vecteur $x_0 = \sum_{i=1}^N (x_0)_i e_i \in \mathbb{C}^N$ tel que $\|x_0\|_2 = 1$ et $(x_0)_i \neq 0$ et définissons la suite $(x_n)_{n \in \mathbb{N}}$ par la formule de récurrence:

$$\forall n \in \mathbb{N}, x_{n+1} = \frac{(A - \mu I_N)^{-1} x_n}{\|(A - \mu I_N)^{-1} x_n\|_2}.$$

La suite $(x_n)_{n \in \mathbb{N}}$ est alors bien définie et elle satisfait:

$$(i) \langle (A - \mu I_N)^{-1} x_n, x_n \rangle_2 \rightarrow \frac{1}{\lambda_i - \mu} \text{ as } n \rightarrow +\infty$$

$$(ii) \exists c \in E_{\lambda_i}(A) \text{ s.t. } \frac{(x_n)_i}{\|x_n\|_2} \rightarrow c \text{ as } n \rightarrow +\infty$$

Preuve:

Évidente.

En sus des difficultés inhérentes à la méthode de la puissance itérée, la principale difficulté tient ici au choix de la valeur approchée μ de la valeur propre λ_i . La convergence de la méthode de la puissance inverse est d'autant meilleure que μ est proche de λ_i , mais alors, la matrice $(A - \mu I_N)^{-1}$ est bien plus singulière. En pratique, il s'agit de trouver un équilibre entre ces deux propriétés ce qui n'est pas évident.