

Analyse numérique-1

écriture sous forme normalisée: $0,3410285 \times 10^2$
 Dans un ordinateur, les nombres sont stockés sous forme normalisée.
 Le nb réellement stocké est une approximation du nb fourni par l'utilisateur. car l'ordi dispose d'un nb limité de place.
 virgule flottante en simple précision selon le Standard IEEE-754 $x = (-1)^s \times f \times 2^e$

avec: x : le nombre
 s : le signe ($0 \leq s < 1$ sinon)
 f : Mantisse sur 23 bits
 e : exposant codé sur 8 bits.
 1 bit caché.

0: $exp = 00000000 = \pm 0 \dots 0_2 \times 2^{-127}$
 ∞ : $exp = 11111111$ si $b_i = 0$
 si $\exists b_i \neq 0$ alors NaN.

+ petit nb:

0 0000001 0 ... 0 = mp

+ grand: 0 11111110 111 ... 111 = mg
 $np \approx 1,2 \times 10^{-35}$ $ng \approx 3,4 \times 10^{38}$

nb normalisé: exposant nul et la mantisse a au moins un bit $\neq 0$

La précision eps d'un ordinateur est le + petit nombre positif normalisé tq $1+eps \neq 1$.

Pour la simple précision, on a $eps \approx 2^{-23} \approx 1,192 \times 10^{-7}$
 $f(x) \approx x$ = nb dans l'ordinateur (x vrai).

3 types d'erreur:

- erreur de représentation: $\Delta x = f(x) - x$
- erreur relative de représentation:

$$r(x) = \frac{\Delta x}{f(x)} = \frac{f(x) - x}{f(x)}$$

erreur relative de précision:

$$\eta(x) = \frac{\Delta x}{x} = \frac{f(x) - x}{x} \text{ so } f(x) = x(1 + \eta(x))$$

En règle générale, pour l'analyse des erreurs on utilise la valeur absolue de l'erreur de précision relative, qui pour le standard est:

$$\eta(x) = \begin{cases} \beta^{1-p} & \text{avec } p = \text{nb chiffre mantisse} \\ 0,5 \times \beta^{1-p} & \text{avec } \beta = \text{base de numérotat} \end{cases}$$

Pour calculer l'intervalle de variation $[E_{min}; E_{max}]$ de l'exposant on utilise le biais

$$B = \underbrace{01 \dots 1}_{q-1} = 2^{(q-1)} - 1 \text{ avec } q = \text{nb chiffre expo}$$

i.e. contenu expo = biais $\in \mathbb{Z}_2$

Simple précision	Double précision
nb bit mantisse = 23 + 1 sign	nb bit $f = 52 + 1$
$E_{min} = +127$	nb bit $E = 11$
$E_{max} = 128$	$E_{min} = 1023$

sources d'erreur

- modélisation math du système \mathcal{E} .
- approximato des fonct° analytiques du modèle mathématique
- discrétisation des fonctions obtenues par approximation des fonct° analytiques afin qu'elle soit traitée par l'ordi.
- calcul numérique par l'ordi

$$a \otimes b = fl(fl(a) \otimes fl(b))$$

th erreur des sommes:

$$S = x_1 + \dots + x_n$$

erreur absolue: $\Delta S = \Delta^1 x_1 + \Delta^2 x_2 + \dots + \Delta^m x_m$
 erreur précision entrée: $\eta(S) = \frac{\Delta S}{S} = \frac{x_1}{S} \eta^1(x_1) + \dots + \frac{x_n}{S} \eta^n(x_n)$

On veut borner l'erreur d'entrée:

$$|\eta(S)| \leq \frac{\sum |x_i|}{|\sum x_i|} \times eps$$

$$\eta(a \otimes b) = fl(fl(a) \otimes fl(b))$$

l'erreur est en 2 parties:

- erreur d'entrée: $\eta(a \otimes b)$
- erreur de calcul: $\eta(a \otimes b)$

L'erreur d'entrée dépend de la façon dont sont agencés les calculs.

opération	erreur d'entrée
$a \pm b$	$\eta(a \pm b) = \frac{a}{a \pm b} \eta(a) \pm \frac{b}{a \pm b} \eta(b)$
$a \times b$	$\eta(a \times b) = \eta(a) + \eta(b)$
a / b	$\eta(a / b) = \eta(a) - \eta(b)$
\sqrt{a}	$\eta(a) = 0,5 \times \eta(a)$

On peut la borner:

$$\eta(p) \leq N \times eps \text{ avec } N \text{ nb de facteur}$$

$$\eta(Q) \leq 2 \times eps \text{ avec } Q \text{ quotient}$$

On évalue une fonction grâce au nb conditions

$$k(x) = \left| \frac{\frac{f(x) - f(x')}{f(x)}}{\frac{x - x'}{x}} \right| \approx \left| \frac{f'(x) \times x}{f(x)} \right|$$

si k est petit, la fonction est bien conditionnée

On évalue de la mme manière un algo en décomposant toutes les fonctions.

l'erreur du résultat

$$\Delta y = \begin{bmatrix} \Delta y_1 \\ \Delta y_n \end{bmatrix} \approx \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots \\ \frac{\partial y_n}{\partial x_1} & \dots \end{bmatrix} \Delta x = J[f(x)] \Delta x$$

avec $J[f(x)]$ jacobien

erreur relative:

$$\eta(y) = \begin{bmatrix} \eta(y_1) \\ \eta(y_n) \end{bmatrix} = k(\phi, x) \cdot \eta(x) \text{ avec } k \text{ nombre condition}$$

Algebre lineaire:

norme:

- $\|x\| > 0$ et $\|x\| = 0 \Leftrightarrow x = 0$.
- $\|ax\| = |a| \|x\|$
- $\|x+y\| \leq \|x\| + \|y\|$
- $\|x-y\| \geq \left| \|x\| - \|y\| \right|$

Def norme: $\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}; p \geq 1$.

Normes utiles:

- 1 - $\|x\|_1 = \sum |x_i|$
- 2 - $\|x\|_2 = \sqrt{\sum x_i^2}$
- 3 - $\|x\|_\infty = \max |x_i|$

Inegalite de Hölder

$$\sum |x_i y_i| \leq \|x\|_p \cdot \|y\|_q = \left(\sum |x_i|^p \right)^{1/p} \cdot \left(\sum |y_i|^q \right)^{1/q}$$

Pour $p=2$, on retrouve l'inegalite de Cauchy.

Si $\|A \cdot B\| \leq \|A\| \cdot \|B\|$, on parle de norme sous-multiplicative.

• $\text{lub}(A) = \sup_{x \in \mathbb{R}^m} \frac{\|Ax\|}{\|x\|}$ lub = norme subordonnee.

• $\text{Cond}(A) = \|A\| \cdot \|A^{-1}\| = \kappa(A)$

• matrice reguliere = inversible
• matrice singuliere = non inversible et $\kappa(A) = \infty$.

$\|A\|_1 = \max_j \sum_i |a_{ij}|$ (Somme des colonnes)

$\|A\|_2 = \left(\sum_{i,j} |a_{ij}|^2 \right)^{1/2}$

$\|A\|_\infty = \max_i \sum_j |a_{ij}|$ Somme des lignes.

Une norme matricielle est consistante avec une norme vectorielle si $\|Ax\| \leq \|A\| \cdot \|x\|$.

• norme de Frobenius:
 $\|A\|_F = \left(\sum_i \sum_j |a_{ij}|^2 \right)^{1/2} = \sqrt{\text{tr}(A^T A)}$

• Th: Soit $\|\cdot\|$ une norme quelconque, alors le rayon spectral d'une matrice carrée A est sa norme: $\rho(A) \leq \|A\|$.

• A orthogonale $\Leftrightarrow A^T = A^{-1}$

Th de Von Neumann:

Soit A avec $\rho(A) < 1$, $(I-A)$ est reguliere
 $(I-A)^{-1} = \sum_{k=0}^{\infty} A^k$

Th: A matrice carrée, alors

- $\lim_{k \rightarrow \infty} A^k = 0$
- $\forall x \in \mathbb{R}^m, \lim_{k \rightarrow \infty} A^k x = 0$.
- $\rho(A) < 1$
- $\|A\| < 1$ pour au moins une norme subordonnee.

Th: si A reguliere et si $\frac{\|A\|}{\|A\|} \leq \frac{1}{\text{Cond}(A)}$

alors $A+\Delta A$ est aussi reguliere.

• Soit le systeme $Ax=b$
perturbation de b ; b devient $b+\Delta b$.

$$\Leftrightarrow x \rightarrow x+\Delta x \Rightarrow \frac{\|\Delta x\|}{\|x\|} \leq \kappa(A) \cdot \frac{\|\Delta b\|}{\|b\|}$$

si $\kappa(A) \gg 1$, une petite perturbation de b peut provoquer une grande perturbation de la solution.

• Soit le systeme $ax+tb$, les perturbations sur $A \xrightarrow{\Delta} A+\Delta A$

$$\Leftrightarrow \frac{\|\Delta x\|}{\|x\|} \leq \kappa(A) \frac{\|\Delta A\|}{\|A\|}$$

• Si perturbations sur A et b .

$$(A+\Delta A)(x+\Delta x) = b+\Delta b$$

$$\Rightarrow \frac{\|\Delta x\|}{\|x\|} \leq \kappa(A) \cdot \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta b\|}{\|b\|} \right)$$

• Wilkinson: analyse de l'erreur:

$$|x \oplus y - f(x \oplus y)| \leq \epsilon \cdot |f(x \oplus y)|$$

$$\text{le } \kappa \leq \epsilon \cdot \delta_\kappa \text{ avec } \delta_\kappa = \sum_{k=1}^n |x_k| \cdot f'(x_k)$$

• Une matrice A symetrique est definie positive si $x^T A x > 0 \forall x \in \mathbb{R}^n$

• toute matrice reguliere A peut se factoriser de maniere unique selon le produit $A=LU$ avec

$$L = \begin{pmatrix} 1 & & & \\ l_{21} & 1 & & \\ & l_{31} & 1 & \\ & & l_{41} & \dots \end{pmatrix} \text{ et } U = \begin{pmatrix} R_{11} & R_{12} & \dots & \\ & R_{22} & \dots & \\ & & \dots & \\ 0 & & & R_{mm} \end{pmatrix}$$

• une matrice carrée A est diagonalement dominante si $|a_{ij}| \geq \sum_{k=1, k \neq i}^m |a_{ik}|$ et

$$\text{si } \exists i_0 \text{ tq } |a_{i_0 i_0}| > \sum_{\substack{k=1 \\ k \neq i_0}}^m |a_{i_0 k}|$$

une telle matrice est reguliere.

Vecteur et valeurs propres

$$Au = \lambda u$$

A matrice carrée et λ vp, u vecteur propre, ma:

- λ^k est valeur propre de A^k .
- $1/\lambda$ est vp de A^{-1}
- $\lambda + R$ est vp de $A + R I_m$
- $c\lambda$ est vp de cA . $\text{tr}(A) = \lambda_1 + \lambda_2 + \dots + \lambda_m$
- si u vecteur propre de A alors $u' = \frac{1}{\|u\|} u$ l'est aussi, de plus $\|u'\| = 1$.
- si B matrice régulière, λ vp de $B^{-1}AB$
- si B mat orthogonale ou unitaire, λ est un vp de $B^T A B$
- si A matrice symétrique déf. positive, toutes les vp sont positives ou nulles.
- si A matrice triangulaire, alors les éléments diagonaux sont des vp.
- Soit u un vp droit et v un vp gauche, alors $\langle u, v \rangle = 0$
- Soit A une matrice carrée diagonalisable et soit U tq $U^{-1}AU = \text{diag}(\lambda_1, \dots, \lambda_m)$.

soit $B = A + \Delta A$ alors $\exists \lambda$ tq

$$|\lambda' - \lambda| \leq k_{\infty}(0) \cdot \|\Delta A\|_{\infty}$$

- $\text{rang}(AB) = \text{rang}(B) - \dim(N(A) \cap R(B))$
- $\text{rg}(AB) \leq \min(\text{rg}(A), \text{rg}(B))$
- $\text{rg}(A) - \text{rg}(B) \leq \text{rg}(A+B) + n$
- $\text{rg}(A^T A) = R(A^T)$ et $R(AA^T) = R(A)$
- $N(A^T A) = N(A)$ et $N(AA^T) = N(A^T)$

Décomposition en valeur singulière:

Toute matrice rectangulaire $A \in \mathbb{R}^{m \times n}$ de rang $R \leq \min(m, n)$ admet une factorisation $A = U \Delta V^T$ avec:

- $U = [u_1, \dots, u_m] \in \mathbb{R}^{m \times m}$ matrice orthogonale dont les colonnes u_i sont les vecteurs propres de la matrice AA^T
- $AA^T = U \Delta V^T (U \Delta V^T)^T U^T = U \Delta^2 U^T$
- $V = [v_1, \dots, v_n] \in \mathbb{R}^{n \times n}$ matrice orthogonale dont les colonnes v_i sont les vecteurs propres de la matrice $A^T A$
- $A^T A = (U \Delta V^T)^T U \Delta V^T = V \Delta^T \Delta V$

$\Delta \in \mathbb{R}^{m \times n}$ la matrice des valeurs singulières $\Delta = \Lambda^{1/2}$ avec Λ étant la matrice diagonale des vp non nuls de la matrice $A^T A$ et aussi des vp non nuls de AA^T .

Soit $A \in \mathbb{R}^{m \times m}$ et la DVS $A = U \Delta V^T$ avec les vs $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m > 0$ alors on a $\|A\| = \sigma_1$, $\|A^{-1}\| = \frac{1}{\sigma_m}$.

Pseudo inverse de $A = U \Delta V^T$ est

$$A^+ = V \Delta^{-1} U^T$$

Th.
Le problème des moindres carrés a au moins une solution a_0 . Si a_1 est une autre solution, alors $X a_0 = X a_1$. Le résidu $r = y - X a_0$ est unique et il est solution de l'équation $X^T r = 0$.

Chaque point a_0 qui réalise $\min \|y - Xa\|$ est aussi solution du système d'éq. minimaux.

$$a_0 = (X^T X)^{-1} X^T y = X^+ y$$

où X^+ est la pseudo inverse.

Une matrice de permutation P est une matrice dont chaque ligne et chaque colonne a un élément = 1 et tous les autres sont nuls.

Théorème de Gershgorin:

Toutes les valeurs propres de la matrice A se trouvent dans l'union de cercles de Gershgorin

$$C_i = \left\{ z \in \mathbb{C} / |z - a_{ii}| \leq \sum_{\substack{k=1 \\ k \neq i}}^m |a_{ik}| \right\}$$

Quotient de Rayleigh:

Si A est une matrice carrée, le coef de Rayleigh est:

$$R(x) = \frac{x^T A x}{x^T x}$$

si A symétrique alors $|R(x)| \leq |\lambda_{\max}|$.

Factorisation LU:

Pour $j \leftarrow 1$ à $m-1$

Pour $i \leftarrow j+1$ à m

$$A(i,j) = A(i,j) / A(j,j)$$

Pour $k \leftarrow j+1$ à m

$$A(i,k) = A(i,k) - A(i,j) \times A(j,k)$$

Fin Pour

Fin Pour

Fin Pour

Inverser Matrice:

Pour $j = 1$ à m

Pour $i = 1$ à m

Si $(i=j)$

Pour $k = 1$ à m

$$B'(j,k) = B(j,k) / B(j,j)$$

Fin Pour

Si non

Pour $k = 1$ à m

$$B'(i,k) = B(i,k) - B(i,j) \times B'(j,k)$$

Fin Pour

Fin si

Fin P

Fin P

Algo de Cholesky

$$L(1,1) = \sqrt{A(1,1)}$$

Pour $i = 2$ à m

Pour $j = 1$ à $i-1$

$$L(i,j) = \frac{1}{L(j,j)} \left(A(i,j) - \sum_{k=1}^{j-1} L(i,k) \times L(j,k) \right)$$

Fin Pour

$$L(i,i) = \left(A(i,i) - \sum_{k=1}^{i-1} L(i,k)^2 \right)^{1/2}$$

Fin Pour

Méthode itérative

$$\Pi x^{(k+1)} = N x^{(k)} + b \Leftrightarrow x^{(k+1)} = \Pi^{-1} N x^{(k)} + \Pi^{-1} b$$

$$A = \begin{bmatrix} D & F \\ -E & \Omega \end{bmatrix} = \begin{bmatrix} I & R & U \end{bmatrix}$$

Méthode de décomposition $A = \Pi - N$ matrice $\Pi^{-1} N$

$$\text{Jacobi } A = \underbrace{D}_M - \underbrace{(E+F)}_N$$

$$J = \Pi^{-1} N = D^{-1}(E+F) = I - D^{-1}A$$

descript d'1 itérato:

$$D x^{(k+1)} = (E+F) x^{(k)} + b$$

Gauss-Seidel:

$$A = \underbrace{D}_M - \underbrace{(E+F)}_N$$

$$G = \Pi^{-1} N = (D-E)^{-1} F$$

descripto $(D-E) x^{(k+1)} = F x^{(k)} + b$

Richardson:

$$A = \underbrace{I}_M - \underbrace{(I-A)}_N$$

$$R = \Pi^{-1} N = I - A$$

descripto $x^{(k+1)} = (I-A) x^{(k)} + b$

Jacobi

$$A = D - (L+U)$$

Fascicule 1 - Chap 1 - Analyse des erreurs

- forme normalisée: $0, \underbrace{b_1 b_2 \dots b_p}_{\text{mantisse}} \times 10^{\dots}$ forme standardisée: en respectant le nb de places. $110,11,001 \Rightarrow 1,1011001 \times 10^{1001}$
- nbs stockés st des approximat° car sont chiffres infinis en partie décimale.
- nbs flottants norme IEEE-754: s(1) e(8) m(23) en binaire
 - séparer et convertir partie entière et partie décimale. réunir les 2 parties. normaliser
 - on divise par 2 et on lit \uparrow restes \hookrightarrow on multiplie et on lit \downarrow les parties entières
 - exposant $\Rightarrow -1+127$ (on enlève bit caché) • on enlève 1^{er} bit (1) de mantisse
- décodage: $1, \text{mantisse} \times 2^{\text{exposant}-127}$ → binais
- 0 → que des zéros • ⊕ petit nb positif stockable normalisé: $0|0\dots 0|0\dots 0 \Rightarrow 2^{-126}$ pr base 2
- variant° de -e $\Rightarrow E -126$ et 127 • ⊕ grand nb normalisé: $0|1\dots 10|1\dots 1 \Rightarrow 1,111 \times 2^{127}$ pr base 2
- ⊕ petit nb non nul: $2^{-149} \Rightarrow 0|0\dots 0|0\dots 01$
- normalisés: exposant a au ⊕ 1 bit ≠ et au ⊖ 1 = 0 double précision: 127 → 1023
- sous-normalisés: exposant nul, mantisse a au ⊕ 1 bit ≠ 0
- zéro: tout à 0 (1 ⊕ et 1 ⊖)
- NaN: exposant 1...1 et mantisse au ⊖ 1 bit ≠ 0
- infini: exposant 1...1 et mantisse 0...0
- précision eps d'1 pc est le ⊕ petit nb positif normalisé tq $1+eps \neq 1 \Rightarrow 2^{-23}$ pr norme IEEE-754 $\Rightarrow eps = 2^{-p}$ bits mantisse
- 3 types d'erreurs:
 - erreur de représentat°: $\Delta x = f(x) - x$ absolue: $|\Delta x|$
 - erreur relative de représentat°: $r(x) = \frac{\Delta x}{f(x)}$
 - erreur de précision: $\eta(x) = \frac{\Delta x}{x} \Rightarrow f(x) = x(1+\eta(x))$
- on utilise en général val abs de l'erreur de précision relative. pr IEEE-754 $\Rightarrow \eta(x) \leq \begin{cases} \beta^{1-p} & \text{si troncation} \\ 0,5\beta^{1-p} & \text{si arrondi} \end{cases}$
- qd on veut juste convertir, on garde m précision (chiffres significatifs) base 16 → 15 = 1 chiffre

Chap 2 - Standard IEEE-754

- biais = $B = 2^{(q-1)} - 1$ pr base 2, avec q: bits de l'exposant $\Rightarrow E_{min} = -\text{biais} + 1$ et $E_{max} = \text{biais}$.
- binaire: biais = 127 (nb à enlever de l'exposant pr passer du standard au nb)
- $\eta(a \pm b) = \eta^I(a \pm b) + \eta^C(a \pm b) = \frac{a}{a \pm b} \eta^I(a) + \frac{b}{a \pm b} \eta^I(b) + \eta^C(a \pm b)$ I → erreur d'entrée C → erreur de calcul.

Chap 3 - Propagat° des erreurs

- erreur d'1 opérat° ⊗ et due aux erreurs de représentat° de a et b et erreur de représentat° du résultat: $a \otimes b = f(f(a) \otimes f(b))$
- ↳ erreur d'entrée = erreur de représentat° du résultat: $\eta^I(a \otimes b)$
- + erreur due au calcul: $\eta^C(a \otimes b)$
- SOMME $S = x_1 + x_2 + \dots \Rightarrow \Delta^I S = \Delta^I x_1 + \Delta^I x_2 + \dots$ et $\eta^I(S) = \frac{x_1}{S} \eta^I(x_1) + \frac{x_2}{S} \eta^I(x_2) + \dots$
- PRODUIT $P = x_1 \cdot x_2 \Rightarrow \Delta^I P = P \cdot \left(\frac{\Delta^I x_1}{x_1} + \frac{\Delta^I x_2}{x_2} \right)$ et $\eta^I(P) = \eta^I(x_1) + \eta^I(x_2)$
- DIVISION $Q = \frac{a}{b} \Rightarrow \Delta^I Q = \frac{b \Delta^I a - a \Delta^I b}{b^2}$ et $\eta^I(Q) = \eta^I(a) + \eta^I(b)$
- RACINE: $\eta^I(x^b) = b \eta^I(x)$
- $\eta^C((a+b)+c) = \eta^C(a+b) + \eta^C(a+b+c)$
- Qualité numérique d'1 fn $f = \text{nb condition}$: $K(x) = \frac{\left| \frac{f(x) - f(x')}{f(x)} \right|}{\left| \frac{x - x'}{x} \right|} \approx \left| \frac{f'(x)}{f(x)} \cdot x \right|$
- K petit \Rightarrow fn bien conditionnée (sinon mal conditionnée)
- K est évalué par stabilité de l'algorithme \Rightarrow si fn utilisées par algo st bien conditionnées.
- erreur du résultat $y = [y_1, \dots, y_m]^T$ d'un algorithme avec c entrées $x = [x_1, \dots, x_n]^T$ est:

$$\Delta y = \begin{bmatrix} \Delta y_1 \\ \vdots \\ \Delta y_m \end{bmatrix} \approx \begin{bmatrix} \frac{\partial \phi_1}{\partial x_1} & \dots & \frac{\partial \phi_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial \phi_m}{\partial x_1} & \dots & \frac{\partial \phi_m}{\partial x_n} \end{bmatrix} \begin{bmatrix} \Delta x_1 \\ \vdots \\ \Delta x_n \end{bmatrix} = J[\phi(x)] \cdot \Delta x$$

• Erreur relative du résultat:

$$\eta(y) = \begin{bmatrix} \eta(y_1) \\ \vdots \\ \eta(y_m) \end{bmatrix} = \begin{bmatrix} K(\phi_1, x_1) & \dots & K(\phi_1, x_n) \\ \vdots & \ddots & \vdots \\ K(\phi_m, x_1) & \dots & K(\phi_m, x_n) \end{bmatrix} \begin{bmatrix} \eta(x_1) \\ \vdots \\ \eta(x_n) \end{bmatrix} = K(\phi, x) \cdot \eta(x)$$

• $f(x) = [f_1(x), \dots, f_m(x)]$. Algorithme: $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^m$ avec $\phi = \phi^{(0)} \circ \dots \circ \phi^{(l)}$

- déroulem^t d'1 algo: $x = x^{(0)} \rightarrow \phi^{(1)}(x^{(0)}) = x^{(1)} \rightarrow \dots \rightarrow \phi^{(l)}(x^{(l-1)}) = x^{(l)} = y$
- $f(x^{(k+1)}) = f(\phi^{(k+1)}(x^{(k)})) \Rightarrow \Delta x^{(k+1)} = f(x^{(k+1)}) - x^{(k+1)}$
- 1 algo A est ⊕ crédible que A' si $Er(A, x) \leq Er(A', x)$.
- Algo numériquement stable \Rightarrow erreurs intermédiaires d'arrondi $H_k x^{(k)}$ st du m ordre de grandeur que ceux inhérents.

Fascicule 2. Algèbre linéaire - Chap 1. Algèbre linéaire et perturbations

Vectorielles:

(2)

$$A = \begin{pmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{pmatrix}$$

$A(n, m)$

(Cauchy avec $p=q=2$)

- norme: $\|x\| \geq 0$ et $\|x\|=0 \Leftrightarrow x=0$, $\|x\| = \|x\| \|1\|$, $\|x+y\| \leq \|x\| + \|y\|$, $\|x-y\| \geq |\|x\| - \|y\||$
- $\|x\|_p = \left(\sum_{i=1}^n |x_i|^p\right)^{1/p} \Rightarrow \|x\|_1 = \sum_{i=1}^n |x_i|$, $\|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2} = \sqrt{x^T x}$, $\|x\|_\infty = \max_{i=1 \dots n} |x_i|$

- inégalité de Holder:

$$\sum_{i=1}^n |x_i y_i| \leq \|x\|_p \|y\|_q \Rightarrow \sum_{i=1}^n |x_i y_i| \leq \sqrt{\sum_{i=1}^n |x_i|^2} \cdot \sqrt{\sum_{i=1}^n |y_i|^2}$$

Matricielles:

- n définitions avec $x \rightarrow Ax$ sauf dernière. Si $\|AB\| \leq \|A\| \|B\| \Rightarrow$ norme sous-multiplicative

- $\|A\| = \sup_{x \in \mathbb{R}^n, \|x\|=1} \|Ax\| = \sup_{x \in \mathbb{R}^n, x \neq 0} \frac{\|Ax\|}{\|x\|} \Rightarrow$ norme sous-ordonnée: $\text{lub}(A) = \sup_{x \in \mathbb{R}^n, x \neq 0} \frac{\|Ax\|}{\|x\|}$

- $\|A\|_1 = \max_{j=1 \dots n} \sum_{i=1}^n |a_{ij}|$ (somme colonnes) $\|A\|_2 = \sqrt{\sum_{i,j=1}^n |a_{ij}|^2}$ (diag) $\|A\|_\infty = \max_{i=1 \dots n} \sum_{j=1}^n |a_{ij}|$ (somme lignes)

- Norme Frobenius: $\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\text{tr}(A^t A)}$ $A \in \mathbb{R}^{m \times n}$

- Un pb est bien conditionné si 1 petite variat° des entrées (données) n'entraîne pas 1 grande variat° des sorties (résultats).

- Scalaires: $\eta(x) = \frac{\Delta x}{x} = \frac{m(x) - x}{x}$. Vecteurs: $\eta(x) = \frac{m(\|x\|_\infty) - \|x\|_\infty}{\|x\|_\infty}$. Conditionnement: $K(A) = \|A\| \cdot \|A^{-1}\|$

- A matrice carrée. $\lim_{k \rightarrow \infty} A^k = 0$ et THÉO VON NEUMANN: matrice $\|x\|_\infty$ $I - A$ régulière $\Leftrightarrow (I - A)^{-1} = \sum_{k=0}^{\infty} A^k$

- Système linéaire: Si A et b ont perturbat°, $Ax = b \Rightarrow (A + \Delta A)(x + \Delta x) = b + \Delta b$

nb machine: $m(x) = x(1 + \eta) = x + \eta(x)$ $|\eta| < \epsilon$

- erreur en retard pr 1 colonne de C: $m(C_j) = (A + \Delta A)b_j$ avec $\|\Delta A\| \leq \gamma_n \|A\|$ pour $C = A \cdot B$

complexité: $C = A \cdot B \Rightarrow c_{ij} = \sum_{k=1}^n a_{ik} b_{kj} \Rightarrow n^3$ multiplicat° et $n^2(n-1)$ additions \Rightarrow complexité $O(n^3)$.

Définitions:

- appli linéaire $\Rightarrow f(\alpha u + \beta u') = \alpha f(u) + \beta f(u')$ appli lin. de U dan V = forme linéaire
- A régulière $\Leftrightarrow A^{-1}$ existe (sinon singulière)
- A (carrée) orthogonale $\Leftrightarrow A^t A = I \Leftrightarrow A^{-1} = A^t$
- A symétrique $\Leftrightarrow A = A^t$. définie positive $\Leftrightarrow x^t A x > 0 \forall x \in \mathbb{R}^n$
- Toute matrice régulière A peut se factoriser de manière unique $A = LU$ avec:

$$L = \begin{pmatrix} 1 & 0 & \dots & 0 \\ l_{21} & 1 & & \\ \vdots & \vdots & \ddots & \vdots \\ l_{n1} & l_{n2} & \dots & 1 \end{pmatrix} \quad R = \begin{pmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ 0 & r_{22} & & \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & r_{nn} \end{pmatrix}$$

L unique si on met des 1 sur diagonale

- rayon spectral: $\rho(A) = \max \{ |\lambda|, \lambda \text{ val propre de } A \}$
- $A(m, n) \Rightarrow n = \dim \text{Im}(A) + \dim \text{Ker}(A)$

Chap 2. Méthodes de résolution de systèmes linéaires $Ax = b$ $x = ?$

- diagonalisat°: $A = PDP^{-1}$

- factorisat° LU: $A = LU \Leftrightarrow$ on résout à la place 2 systèmes $Ly = b$ puis $Ux = y$.

- pour trouver L et U, factorisat° de GAUSS

Chap 3. Méthodes itératives linéaires

$Ax=b$

(3)

on veut décomposer A en $A=M.N$.

méthodes non relaxées: A peut s'écrire $A=D.E.F=D+L+U$ avec $A = \begin{bmatrix} D & F \\ -E & \end{bmatrix} = \begin{bmatrix} D & U \\ L & \end{bmatrix}$

pour $A=D.E.F$

méthode	Décomp. $A=M.N$	matrice $M^{-1}N$	Descript° d'une itération
Jacobi	$A=D-(E+F)$	$J=M^{-1}N=D^{-1}(E+F)$ $=I-D^{-1}A$	$Dx^{(k+1)}=(E+F)x^{(k)}+b$
Gauss-Seidel	$A=(D-E)-F$	$G=M^{-1}N=(D-E)^{-1}F$	$(D-E)x^{(k+1)}=Fx^{(k)}+b$
Richardson	$A=I-(I-A)$	$R=M^{-1}N=I-A$	$x^{(k+1)}=(I-A)x^{(k)}+b$
Jacobi	$A=D-(L+U)$	$J=M^{-1}N=-D^{-1}(L+U)$	$Dx^{(k+1)}=-(L+U)x^{(k)}+b$
Gauss-Seidel	$A=(D+L).U$	$G=M^{-1}N=(D+L)^{-1}U$	$(D+L)x^{(k+1)}=Ux^{(k)}+b$

↑ matrice d'itérat°

Fascicule 3. Valeurs et vecteurs propres - Chap 1 - Valeurs et vecteurs propres

- Le pb algébrique des valeurs propres consiste en la sol° du système de $n+1$ eq° linéaires $(A-\lambda B)u=0 \Rightarrow Au=\lambda Bu$.

- Si $B=I$, on obtient le pb des val propres: trouver $\lambda \in \mathbb{R}$ tq $Au=\lambda u$ avec $A \in \mathbb{R}^{n \times n}$, $u \in \mathbb{R}^n$, $u \neq 0$.

- $A(u)=Au$. P matrice de passage de $A \rightarrow B$ et $T: B \rightarrow A$. $B=T^{-1}AP$. Si $E=F \Rightarrow B=P^{-1}AP$.

- λ val propre de $A \Rightarrow \lambda^k$ val propre de A^k . Si A sym. déf. positive, toutes les val. propres et réelles positives.

- Si A nilpotente ($\exists k > 0$ tq $A^k=0$) $\Rightarrow \text{tr}(A)=0$. $\text{tr}(A)=\sum \lambda_i$.

- diagonalisat°: $D=P^{-1}AP$ est semblable à A, $A=PD P^{-1}$. Si val propres distinctes \Rightarrow diagonalisable.

Si on regarde dim des sous-espaces propres.

- Théo de Schur: (triagonalisat°) il \exists une matrice unitaire U et 1 matrice triangulaire sup T tq $A=U^t T U$

- $A \mapsto B=U^t A U$ une transp. unitaire de similitude et A normale $\Rightarrow U^t A U = \text{Diag}(\lambda_1, \dots, \lambda_n)$.

- Soit A une matrice carrée d'ordre n. On associe à chaque ligne i de la matrice le disque de Gerschgorin, D_i dont le centre est l'élément diagonal a_{ii} et le rayon la somme des modules des autres

éléments de la ligne: $D_i^{(e)} = \{z \in \mathbb{C} / |z - a_{ii}| \leq \sum_{j \neq i} |a_{ij}|\}$ où $R_i^{(e)} = \sum_{j \neq i} |a_{ij}|$ est le rayon de $D_i^{(e)}$

- Toutes les val. propres appartiennent à au \ominus un des $\bigcup_{i=1}^n$ disques.

- u vect. pr. de A et λ val. propre de A $\Leftrightarrow Au=\lambda u$.

- estimat° des erreurs en cas de perturbat°:

- on reprend u ($Au=\lambda u$) et v tq $v^t A = \mu v^t$. les val. propres λ et μ st identiques. On a: $v^t u = 0$ (orthogonalité).

- matrice perturbée: $A(E)=A+\Delta A$. Si λ val propre de A $\Rightarrow |\lambda(E)-\lambda| \leq 0$ ($E^{1/m}$) où m =multiplicité de λ .

si λ est simple $\rightarrow \lambda = \lim_{\epsilon \rightarrow 0} \frac{A(E)-\lambda}{\epsilon} = \frac{v^t \Delta A u}{v^t u}$

Chap 2. Décomposition en valeurs singulières

- $R \rightarrow$ image et $N \rightarrow$ noyau

- $\text{Ker}(A) = \{0\}$ ssi $\text{rg} A = n$

- E un espace. F, G \subseteq E st complémentaires $\Leftrightarrow E=G+F$ et $G \cap F = \{0\} \Leftrightarrow E=G \oplus F$

- $E=G \oplus F$. Si $x=y+z$ avec $x \in E$, $y \in G$ et $z \in F$, alors y est la project° de x dans G et z celle de x dans F.

- project° orthogonale de x sur $\text{Vect}(u) = \frac{u}{\|u\|} \left(\frac{u}{\|u\|} \right)^t x$

projecteur orthogonal P_u ds $\text{Vect}(u) \Rightarrow P_u = \frac{u u^t}{u^t u}$ dans $\text{Vect}(u^i) \Rightarrow P_{u^i} = I - P_u$.

- A inversible $\Rightarrow \det A \neq 0$ (\Rightarrow) A est une base possible de l'espace

exemple de transformate:

$f_1 = e_2$ $f_2 = e_1 + 2e_2$ $f_3 = -e_1 - 3e_2 + e_3$

$B = \begin{pmatrix} 0 & 1 & -1 \\ 1 & 2 & -3 \\ 0 & 0 & 1 \end{pmatrix}$ $P = B^t = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 2 & -3 \\ -1 & -3 & 1 \end{pmatrix}$ $P^{-1} = \begin{pmatrix} -2 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}$ et $A' = P^{-1}AP$

- Décomposit° en val. singulières de A: on peut factoriser $A = U\Delta V^t$ avec: (DVS)

- U vecteurs propres de AA^t
- V vecteurs propres de A^tA
- Δ matrice des val. singulières, $\Delta = \sqrt{\Lambda} = \sqrt{\text{diag}(\lambda_1, \dots, \lambda_n)}$ val. propres de A^tA ou AA^t non nulles

- $B = \begin{pmatrix} 0 & A \\ A^t & 0 \end{pmatrix} \Rightarrow$ val. sing. de A st $\sigma_1 \dots \sigma_r$ (\Rightarrow) val. sing. non nulles de B st $\sigma_1, \dots, \sigma_r, -\sigma_1, \dots, -\sigma_r$.

- Soit $A = U\Delta V^t$ avec val. sing. $\sigma_1 \geq \dots \geq \sigma_r > 0 \Rightarrow \|A\| = \sigma_1$. Si $m=n$ et A régulière, $\|A^{-1}\| = \frac{1}{\sigma_n}$.

$\hookrightarrow K(A) = \frac{\sigma_1}{\sigma_n}$

Chap 3. Pseudoinverse et moindres carrés

- Pseudoinverse de $A = U\Delta V^t \Rightarrow A^+ = V\Delta^{-1}U^t$ Si Δ diagonale $\Rightarrow \Delta^{-1} = \begin{pmatrix} \frac{1}{\sigma_1} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\sigma_n} \end{pmatrix}$

- $K^+(A) = \|A\| \|A^+\|$ (nb condit° de A pour inversion gauche)

- Soit $Ax = b$ avec perturbat° $A(x + \Delta x) = b + \Delta b$. L'erreur relative est bornée par:

$\frac{\|\Delta x\|}{\|x\|} \leq K^+(A) \frac{\|\Delta b\|}{\|b\|}$

Partiel année dernière:

- ⊕ petits et ⊕ grand nb normalisés représentables: $B = 2^p - 1 \Rightarrow x_{\max} = 0|1\dots 10|1\dots 1 = 1,11\dots 1 \times \beta^{ep-B}$
 $x_{\min} = 0|0\dots 01|0\dots 0$

- ⊕ grde erreur: $|r(x)| = \beta^{1-p}$ avec $\beta =$ base et $p =$ bits mantisse si troncature, si arrondi $\Rightarrow 0,5\beta^{1-p}$

- eps est le ⊕ petit positif tq $1+eps \neq 1$: $1 \rightarrow 0|01\dots 1|0\dots 0 = 1,0\dots 0 \times 2^{ep-B}$
 $1+eps \rightarrow 0|01\dots 1|0\dots 01 = 1,0\dots 01 \times 2^{ep-B}$
 $\Rightarrow eps = 0,00\dots 01_2$

- calcul par ordi a 1 précision finie. De dès que le terme dépasse la capacité de représentat° de la machine, il est converti à 0 et son ajout à la somme ne change pas la valeur de la somme. Il ya de convergence numérique.

$A = \begin{pmatrix} x & x & & 0 \\ x & \ddots & \ddots & \\ 0 & \ddots & \ddots & x \\ & & x & x \end{pmatrix}$ $L = \begin{pmatrix} x & & & \\ x & \ddots & & \\ 0 & \ddots & \ddots & \\ 0 & & x & 1 \end{pmatrix}$ $U = \begin{pmatrix} x & x & & 0 \\ & \ddots & \ddots & \\ 0 & \ddots & x & \\ & & x & x \end{pmatrix}$

Leat triang. inf. avec des 1 sur la diagonale et ne comporte qu'il sous-diag. non nulle car hérite de la structure bande de A. U est triang. sup et ne comporte qu'une sur diag. non nulle car idem.

- On peut alors calculer quelques u_{ij} et l_{ij} . On peut démontrer des résultats ($u_{ii} = \frac{i+1}{i}$) par récurrence, d'ab $\det(A) = \det(LU) = \det(L) \det(U) = \prod_{i=1}^n \frac{i+1}{i} = n+1$.

$J = \frac{1}{2} \begin{pmatrix} 0 & 1 & 0 \\ 1 & \ddots & \ddots \\ 0 & \ddots & 0 \end{pmatrix}$

- ⊕ petite val propre λ^i de A: on calcule $Au^{(i)} = \lambda^i u^{(i)}$

- conditionnement de A: $\text{cond}_2(A) = \frac{\|A\|_2}{\|A^{-1}\|_2}$

- méthode converge? ai car consistante et $\rho(J) < 1$

- $\rho(J)^2 = \rho(G)$ car A tridiagonale. $\rho(G) < \rho(J)$ donc Gauss-Seidel est méthode ⊕ rapide.

- $A = \begin{pmatrix} J & 1 & 1 \\ 0 & 6 & -J \\ 1 & 0 & -J \end{pmatrix}$ on cherche disques et on trouve 1 val. pr. ds le disque $(-5, 1)$ et 2 ds 2 disques qui se recoupent $(5, 1)$ et $(6, 1)$

- u vectem propre si $Au = \lambda u$, dc vectem propre doit être multiple.