



HAL
open science

Méthode d'estimation à posteriori d'erreurs

Yves Ducrocq

► **To cite this version:**

Yves Ducrocq. Méthode d'estimation à posteriori d'erreurs. Modélisation et simulation. Université Joseph-Fourier - Grenoble I, 1968. Français. NNT: . tel-00281032

HAL Id: tel-00281032

<https://theses.hal.science/tel-00281032v1>

Submitted on 20 May 2008

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

N° d'ordre

T H E S E

présentée à la Faculté des Sciences
de l'Université de Grenoble

pour obtenir
le grade de Docteur de Troisième Cycle
"MATHEMATIQUES APPLIQUEES"

par

Yves DUCROCQ
Licencié ès Sciences

METHODE D'ESTIMATION A POSTERIORI D'ERREURS

Thèse soutenue le 5 février 1968, devant la Commission d'Examen :

Monsieur	J. KUNTZMANN	Président
Messieurs	N. GASTINEL	Examineur
	P.J. LAURENT	Examineur

N° d'ordre

T H E S E

présentée à la Faculté des Sciences
de l'Université de Grenoble

pour obtenir
le grade de Docteur de Troisième Cycle
"MATHEMATIQUES APPLIQUEES"

par

Yves DUCROCQ
Licencié ès Sciences

METHODE D'ESTIMATION A POSTERIORI D'ERREURS

Thèse soutenue le 5 février 1968, devant la Commission d'Examen :

Monsieur	J. KUNTZMANN	Président
Messieurs	N. GASTINEL	Examineur
	P.J. LAURENT	Examineur

Je remercie Monsieur KUNTZMANN, Directeur de l'Institut de Mathématiques Appliquées de Grenoble d'avoir bien voulu présider le Jury de cette thèse,

Je remercie Monsieur GASTINEL, Professeur à la Faculté des Sciences de Grenoble qui m'a suggéré les sujets de cette thèse et qui, par sa patience et ses conseils, a aussi permis son achèvement,

Je remercie aussi Monsieur LAURENT, Maître de Conférences qui m'initia~~x~~ au calcul numérique, d'avoir bien voulu faire partie du Jury de cette thèse,

et enfin, je remercie tous ceux qui m'ont aidé, soutenu, tous ceux qui ont oeuvré si bien à la réalisation matérielle de cette thèse.

FACULTE DES SCIENCES

LISTE DES PROFESSEURS

DOYENS HONORAIRES :

M. MORET

M. WEIL

DOYEN :

M. BONNIER E.

PROFESSEURS TITULAIRES :

MM. NEEL Louis	Chaire de Physique Expérimentale
HEILMANN René	Chaire de Chimie
KRAVTCHEKNO Julien	Chaire de Mécanique Rationnelle
CHABAUTY Claude	Chaire de calcul différentiel et intégral
BENOIT Jean	Chaire de Radioélectricité
CHENE Marcel	Chaire de Chimie Papetière
WEIL Louis	Chaire de Thermodynamique
FELICI Noël	Chaire d'Electrostatique
KUNTZMANN Jean	Chaire de Mathématiques Appliquées
BARBIER Reynold	Chaire de Géologie Appliquée
SANTON Lucien	Chaire de Mécanique des Fluides
OZENDA Paul	Chaire de Botanique
FALLOT Maurice	Chaire de Physique Industrielle
KOSZUL Jean-Louis	Chaire de Mathématiques M.P.C.
GALVANI O.	Mathématiques
MOUSSA André	Chaire de Chimie Nucléaire
TRAYNARD Philippe	Chaire de Chimie Générale

SOUTIF Michel	Chaire de Physique Générale
CRAYA Antoine	Chaire d'Hydrodynamique
REULOS R.	Théorie des Champs
BESSON Jean	Chaire de Chimie
AYANT Yves	Physique Approfondie
GALLISSOT	Mathématiques
Melle LUTZ Elisabeth	Mathématiques
MM. BLAMBERT Maurice	Chaire de Mathématiques
BOUCHEZ Robert	Physique Nucléaire
LLIBOUTRY Louis	Géophysique
MICHEL Robert	Chaire de Minéralogie et Pétrographie
BONNIER Etienne	Chaire d'Electrochimie et d'Electrométallurgie
DESSAUX Georges	Chaire de Physiologie Animale
PILLET E.	Chaire de Physique Industrielle et Electrotechnique
VOCCOZ Jean	Chaire de Physique Nucléaire Théorique
DEBELMAS Jacques	Chaire de Géologie Générale
GERBER R.	Mathématiques
PAUTHENET R.	Electrotechnique
VAUQUOIS B.	Chaire de Calcul Electronique
BARJON R.	Physique Nucléaire
BARBIER Jean-Claude	Chaire de Physique
SILBER R.	Mécanique des Fluides
BUYLE-BODIN Maurice	Chaire d'Electronique
DREVFUS B.	Thermodynamique
KLEIN J.	Mathématiques
VAILLANT F.	Zoologie et Hydrobiologie
ARNOUD Paul	Chaire de Chimie M.P.C.
SENGEL P.	Chaire de Zoologie
BARNOUD F.	Chaire de Biosynthèse de la Cellulose
BRISSONNEAU P.	Physique
GAGNAIRE Didier	Chaire de Chimie Physique

Mme	KÖFLER L.	Botanique
MM.	DEGRANGE Charles	Zoologie
	PEBAY-PEROULA J.C.	Physique
	RASSAT A.	Chaire de Chimie Systématique

PROFESSEURS SANS CHAIRE :

MM.	GIDON P.	Géologie et Minéralogie
	GIRAUD P.	Géologie
	PERRET R.	Servomécanismes
Mme	BARBIER M.J.	Electrochimie
Mme	SOUTIF J.	Physique
MM.	COHEN J.	Electrotechnique
	DEPASSEL R.	Mécanique des Fluides
	GASTINEL N.	Mathématiques Appliquées
	ANGLES-d'AURIAC P.	Mécanique des Fluides
	DUCROS P.	Minéralogie et Cristallographie
	GLENAT R.	Chimie
	LACAZE A.	Thermodynamique
	BARRA J.	Mathématiques Appliquées
	COUMES A.	Electronique
	PERRIAUX J.	Géologie et Minéralogie
	ROBERT A.	Chimie Papetière
	BIAREZ J.P.	Mécanique Physique
	BONNET G.	Electronique
	CAUQUIS G.	Chimie Générale
	BONNETAIN L.	Chimie Minérale
	DEPOMMIER P.	Etude Nucléaire et Génie Atomique
	HACQUES Gérard	Calcul Numérique
	POLOUJADOFF M.	Electrotechnique

MAITRES DE CONFERENCES :

MM. DODU J.	Mécanique des Fluides
LANCIA Roland	Physique Automatique
Mme KAHANE J.	Physique
MM. DEPORTES C.	Chimie
Mme BOUCHE L.	Mathématiques
MM. SARROT-RAYNAUD J.	Géologie Propédeutique
Mme BONNIER M.J.	Chimie
MM. KAHANE A.	Physique Générale
DOLIQUE J.M.	Electronique
BRIERE G.	Physique M.P.C.
DESPRE P.	Chimie S.P.C.N.
LAJZEROWICZ J.	Physique M.P.C.
VALENTIN P.	Physique M.P.C.
BERTRANDIAS J.P.	Mathématiques Appliquées T.M.P.
LAURENT P.	Mathématiques Appliquées T.M.P.
CAUBET J.P.	Mathématiques Pures
PAYAN J.J.	Mathématiques
Mme BERTRANDIAS F.	Mathématiques Pures M.P.C.
MM. LONGEQUEUE J.P.	Physique
NIVAT M.	Mathématiques Appliquées
SOHM J.C.	Electrochimie
ZADWORNY F.	Electronique
DURAND F.	Chimie Physique
CARLIER G.	Biologie Végétale
AUBERT G.	Physique M.P.C.
DELPUECH J.J.	Chimie Organique
PFISTER J.C.	Physique C.P.E.M.
CHIBON P.	Biologie Animale
IDELMAN S.	Physiologie Animale
BLOCH D.	Electrotechnique
BRUGEL L.	I.U.T.
SIBILLE R.	I.U.T.

TABLE DES MATIERES

SOMMAIRE	I
INTRODUCTION GENERALE	II
 <u>PARTIE A - LES MOINDRES CARRES</u>	
INTRODUCTION	A,I,1
ETUDE D'UN CAS PARTICULIER	A,II,1
CAS GENERAL	A,III,1
CALCUL DE D^{-1}	A,III,3
PROCEDURE ALGOL	A,IV,1
EXEMPLES NUMERIQUES	A,IV,7
 <u>PARTIE B - LES POLYNOMES</u>	
INTRODUCTION	B,I,1
CAS PARTICULIER	B,II,1
CAS GENERAL	B,III,1
PROCEDURE ALGOL	B,III,4
EXEMPLES NUMERIQUES	B,IV,1
 <u>PARTIE C - LES VALEURS PROPRES</u>	
INTRODUCTION	C,I,1
CALCUL DE L'ERREUR CONNAISSANT UN VECTEUR PROPRE NUMERIQUE ET LA VALEUR PROPRE ASSOCIEE	C,II,1
CAS GENERAL OU L'ON CONNAIT SEULEMENT UNE VALEUR PROPRE NUMERIQUE	C,III,1
PROCEDURE ALGOL	C,IV,1
EXEMPLES NUMERIQUE	C,IV,4
 <u>PARTIE D - ESPACE DE HILBERT</u>	
INTRODUCTION	D,I,1
RESULTATS PRELIMINAIRES	D,II,1
ETUDE THEORIQUE	D,III,1
APPLICATIONS AUX METHODE DE GALERKIN	D,IV,1
CAS DE L'EQUATION DE FREDHOLM	D,IV,3
CAS DU PROBLEME DE DIRICHLET	D,IV,7
EXEMPLE NUMERIQUE	D,V,1

S O M M A I R E

La partie A de ce travail est consacrée à l'étude à postériori de l'erreur commise dans la résolution d'un système par les moindres carrés. Cette partie est celle qui est susceptible du plus grand nombre d'applications pratiques dans toutes les sciences expérimentales.

La partie B est intéressante au point de vue théorique. Elle parle de l'estimation à postériori de l'erreur sur une ou plusieurs racines numériques d'un polynôme mais me semble beaucoup moins pratique que l'estimation directe.

C traite, elle, du problème des valeurs propres. Nous commençons par estimer à postériori l'erreur quand nous avons une valeur propre et le vecteur propre correspondant puis nous en déduisons l'erreur à postériori sur la valeur propre seule.

Pour ces trois parties, nous donnons une procédure qui réalise effectivement le calcul de cette erreur à postériori. Nous signalons qu'elles sont toutes assez rapides (temps moyen : 20 millièmes d'heures, soit un peu plus d'une minute).

Enfin dans la partie D nous estimons à postériori l'erreur commise dans la résolution d'un système linéaire et ce, dans le cadre d'un espace de Hilbert. Pour cela nous commençons par définir ce que nous appellerons une double norme, puis, à l'aide d'une hypothèse exprimée en utilisant cette double norme, nous étudierons l'erreur à postériori. Nous ferons une application à la méthode de Galerkin des résultats trouvés. Nous particulariserons pour le problème de Dirichlet et l'équation de Fredholm.

La bibliographie de chaque partie se trouve à la fin de celle-ci.

I N T R O D U C T I O N

L'objet de cette thèse étant l'estimation à postériori de l'erreur dans différents cas, nous commencerons par expliquer ce qu'est cette estimation. Nous emprunterons à Jean Gâches [1] ces explications.

Jean Gâches a résolu ce problème dans le cas de systèmes linéaires.

Ayant appliqué un algorithme à un ensemble de données, il importe de connaître la signification du résultat numérique obtenu. C'est-à-dire de préciser dans quelle mesure ce résultat est acceptable, compte-tenu du problème proposé.

La solution exacte ne peut être atteinte qu'exceptionnellement. En effet, on utilise généralement des procédés itératifs qui convergent vers la solution exacte mais on ne peut évidemment aller à la limite, on doit se contenter d'une approximation ; d'autre part il y a aussi les erreurs d'arrondi (dûes au nombre fini et constant de bits par mémoire).

Il y a aussi l'incertitude des données dûes soit à leur nature expérimentale (réalisation non idéale de l'expérience, précision des appareils, qualités de l'observateur) soit, aussi du fait, déjà vu, des troncatures de nombres à leur entrée en machine et de leur transformation en binaire.

Pour caractériser la qualité des résultats obtenus on utilise :

L'analyse directe.

Soit \tilde{x} la solution numérique et x^* la solution exacte d'un problème. C'est le calcul d'erreur habituel et il conduit à des inégalités de la forme $||x^* - \tilde{x}|| < \epsilon$ mais en général si le nombre d'opérations est grand cela conduit à des majorations trop grossières.

L'analyse rétrograde.

Ici on recherche de quels problèmes \tilde{x} est-il solution.

$$(1) \quad f(\tilde{x}) = \varepsilon \quad \varepsilon \text{ est le résidu}$$

$$(2) \quad f'(\tilde{x}) = 0$$

$$||f - f'|| < \varepsilon$$

mais ceci conduit encore à des majorations grossières de l'erreur.

L'analyse à postériori.

C'est un raffinement de la précédente en ce sens que nous chercherons, non pas un problème quelconque vérifié par \tilde{x} , mais le problème le plus "proche" vérifié par \tilde{x} . Nous dirons que f^* est indiscernable de f si $||f^* - f|| < \varepsilon^*$ (quantité donnée à l'avance). Nous dirons alors qu'une solution \tilde{x} est compatible avec un problème f si le problème f^* , vérifié par \tilde{x} , le plus proche est indiscernable de f . La valeur de l'analyse à postériori vient du fait que ε^* est aisément calculé à partir des incertitudes sur les données du problème.

- A -

ETUDE DE L'ERREUR A POSTERIORI

DANS LE CAS

DES MOINDRES CARRES

C H A P I T R E I

INTRODUCTION

On cherche souvent, en physique, à déterminer la valeur numérique d'un groupe de paramètres à l'aide d'une suite nombreuse de mesures. Ces mesures forment un tableau de la forme $A_0 X = B_0$. A_0 et B_0 étant des matrices dont les éléments sont connus avec les marges d'incertitude ϵ_1 et ϵ_2 . X étant le vecteur inconnu à déterminer.

On obtient généralement la valeur des paramètres par la méthode des moindres carrés, si bien qu'on est amené à étudier le système quadratique $A_0^T A_0 X = A_0^T B_0$.

L'objet de la partie A de ce travail est justement de déterminer l'erreur à postériori dans la résolution de ce système.

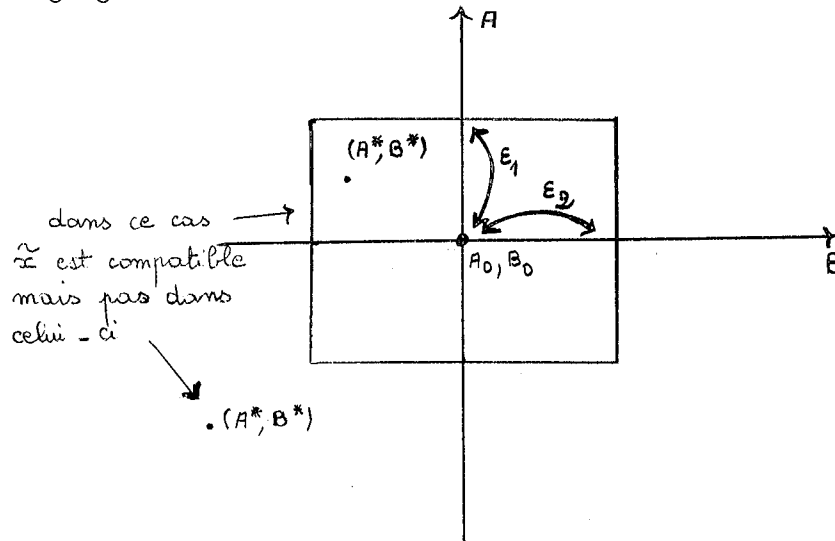
C'est-à-dire, qu'ayant trouvé une solution numérique \hat{x} de $A_0^T A_0 X - A_0^T B_0 = 0$, nous allons examiner l'ensemble des problèmes dont \hat{x} est solution exacte et voir s'il en existe un qui ne soit pas trop "différent" du problème véritable. Pour cela, A_0 étant de dimension $p \times q$ et B_0 de dimension q , on se placera dans l'espace $\mathbb{R}^{q \times (p+1)}$. Dans celui-ci, on cherchera le couple de matrices (A^*, B^*) qui appartient au cône $A^T A X - A^T B = 0$ et qui est à la "distance" minimum de (A_0, B_0) .

- Dans toute la suite nous utiliserons comme norme de matrices ou de vecteurs la racine carrée de la somme des carrés des éléments,

exemple :

$$||A|| = \sqrt{\sum_{i,j} a_{ij}^2}$$

Si (A^*, B^*) est à l'intérieur du domaine dans lequel nous savons que se trouve (A_0, B_0) nous dirons que la solution \tilde{x} est compatible.



Méthode.

Nous utiliserons la méthode des multiplicateurs de Lagrange que je rappelle ici brièvement :

Soit une fonction de plusieurs variables $f(x_1, x_2, \dots, x_n)$ dont nous cherchons un extrémum sachant que x_1, x_2, \dots, x_n ne sont pas quelconques mais soumis à "p" contraintes de la forme $g_i(x_1, \dots, x_n) = 0$ ($i=1, 2, \dots, p$). On montre que l'on trouve cet extrémum en annulant les dérivées partielles par rapport à x_1, \dots, x_n de la fonction auxiliaire $f(x_1, \dots, x_n) + \sum \lambda_i g_i(x_1, \dots, x_n)$. Les λ_i étant des inconnues auxiliaires, nous aurons donc $n + p$ inconnues à déterminer et, pour cela, les n dérivées partielles et les p contraintes.

Dans le cas présent, nous minimiserons d^2 -carré de la distance euclidienne rapportée à l'espace $\mathbb{R}^{p \times q + q}$ du point (A, B) au point (A_0, B_0) avec la seule contrainte pour (A, B) d'appartenir au cône $A^T A \tilde{x} - A^T B = 0$.

Plan Suivi.

Nous commencerons par étudier le cas particulier où A_0 est une matrice à 2 lignes et une colonne, x se réduisant alors à un scalaire, puis nous généraliserons le résultat trouvé au cas d'une matrice colonne de longueur p .

Dans la deuxième partie, nous étudierons le cas le plus général à une restriction près : A_0 sera une matrice ayant plus de lignes que de colonnes. En effet, en pratique, nous aurons bien plus de mesures que d'inconnues. Nous faisons aussi l'hypothèse que $A_0^T A_0$ sera toujours inversible.

Remarque.

Il est évident que si $\tilde{x} \rightarrow x^+$ (solution exacte de $A_0^T A_0 x - A_0^T B_0$) alors $(A^*, B^*) \rightarrow (A_0, B_0)$.

C H A P I T R E II

ETUDE D'UN CAS PARTICULIER

Nous allons étudier le cas particulier suivant :

$$A_o = \begin{pmatrix} a_o \\ a'_o \end{pmatrix}, \quad B_o = \begin{pmatrix} b_o \\ b'_o \end{pmatrix}, \quad X = (x).$$

Appliquons la méthode de Lagrange suivant l'exposé du chapitre I.

Dans notre cas particulier, cela revient à minimiser dans \mathbb{R}^4 la distance euclidienne $d^2 = (a-a_o)^2 + (a'-a'_o)^2 + (b-b_o)^2 + (b'-b'_o)^2$ sachant que $a^{(1)}$ et $b^{(1)}$

vérifient $(a^2+a'^2)_{\tilde{X}} = ab + a'b'$, autrement dit à annuler les dérivées partielles par rapport à a, a', b et b' de la fonction :

$$(a-a_o)^2 + (a'-a'_o)^2 + (b-b_o)^2 + (b'-b'_o)^2 + \lambda \left[(a^2+a'^2)_{\tilde{X}} - (ab+a'b') \right].$$

En ajoutant aux quatre dérivées partielles la contrainte, on obtient le système de 5 équations à 5 inconnues a, a', b, b' et λ :

$\begin{aligned} 2(1+\lambda\tilde{X})a - \lambda b &= 2a_o \\ 2(1+\lambda\tilde{X})a' - \lambda b' &= 2a'_o \\ 2b - \lambda a &= 2b_o \\ 2b' - \lambda a' &= 2b'_o \\ (a^2+a'^2)_{\tilde{X}} &= ab + a'b' \end{aligned}$

En résolvant ce système par rapport à λ , on obtient les quatre valeurs :

$$\begin{array}{l}
 a = \frac{4a_o + 2\lambda b_o}{4(1+\lambda\tilde{x})-\lambda^2} \qquad a' = \frac{4a'_o + 2\lambda b'_o}{4(1+\lambda\tilde{x})-\lambda^2} \\
 b = \frac{2\lambda a_o + 4(1+\lambda\tilde{x})b_o}{4(1+\lambda\tilde{x})-\lambda^2} \qquad b' = \frac{2\lambda a'_o + 4(1+\lambda\tilde{x})b'_o}{4(1+\lambda\tilde{x})-\lambda^2}
 \end{array}$$

On voit que si $\lambda \rightarrow 0$ alors $(a, b, a', b') \rightarrow (a_o, b_o, a'_o, b'_o)$

Nous allons maintenant résoudre l'équation qui nous donnera λ . Nous l'obtiendrons en reportant les 4 valeurs précédentes dans la contrainte.

L'équation ordonnée en λ s'écrit :

$$\begin{aligned}
 &\lambda^2 \left[a_o b_o + a'_o b'_o + (b_o^2 + b'_o{}^2) \tilde{x} \right] + 2\lambda \left[a_o^2 + a'_o{}^2 + b_o^2 + b'_o{}^2 \right] \\
 &- 4 \left[(a_o^2 + a'_o{}^2) \tilde{x} - (a_o b_o + a'_o b'_o) \right] = 0
 \end{aligned}$$

Nous posons :

$$\left\{ \begin{array}{l}
 a_o^2 + a'_o{}^2 = P \\
 b_o^2 + b'_o{}^2 = Q \\
 a_o b_o + a'_o b'_o = R \\
 \text{et } \frac{R}{P} = x^* \quad (\text{c'est la solution exacte})
 \end{array} \right.$$

L'équation précédente s'écrit alors :

$$\lambda^2 (R+Q\tilde{x}) + 2\lambda (P+Q) - 4P(\tilde{x}-x^*) = 0$$

On en déduit les deux solutions.

$$\lambda = \frac{-(P+Q) \pm \sqrt{(P+Q)^2 + 4(R+Q\tilde{x}) P(x-x^*)}}{R + Q\tilde{x}}$$

nous avons vu que si $\lambda \rightarrow 0$ alors (a, b, a', b') tend vers (a_0, b_0, a'_0, b'_0) et donc $\tilde{x} \rightarrow x^*$ ce qui nous impose ici le choix du signe +. En effet dans ce cas lorsque $\tilde{x} \rightarrow x^*$ on a bien $\lambda \rightarrow 0$.

$$\lambda = \frac{-(P+Q) + \sqrt{(P+Q)^2 + 4(R+Q\tilde{x}) P(x-x^*)}}{R + Q\tilde{x}}$$

Remarque.

Le discriminant peut s'écrire :

$$\Delta = (P+2R\tilde{x}-Q)^2 + 4(a_0 b'_0 - a'_0 b_0)^2 (1+\tilde{x}^2)$$

il est donc toujours positif ou nul et donc λ existe toujours.

Calculons maintenant la distance de (A_0, B_0) au cône $A^T A \tilde{x} - A^T B = 0$

On a :

$$a^{(r)} - a_0^{(r)} = \lambda \cdot \frac{(\lambda - 4\tilde{x}) a_0^{(r)} + 2 b_0^{(r)}}{4(1+\lambda\tilde{x}) - \lambda^2}$$

$$b^{(r)} - b_0^{(r)} = \lambda \cdot \frac{2 a_0^{(r)} + \lambda b_0^{(r)}}{4(1+\lambda\tilde{x}) - \lambda^2}$$

On en déduit d^2 qui tous les calculs faits peut s'écrire :

$$d^2 = \frac{\lambda^2}{[4(1+\lambda\tilde{x}) - \lambda^2]^2} \cdot [(\lambda^2 + 4)(P+Q) + 8(P\tilde{x} - R)(2\tilde{x} - \lambda)]$$

Remarque.

d^2 est une somme de carrés donc $d^2 \geq 0$, l'expression de d^2 en fonction de λ (voir page) nous montre que $d^2 \rightarrow 0$ lorsque $\lambda \rightarrow 0$ et que $d^2 \rightarrow \infty$ lorsque $\lambda \rightarrow 2(x \pm \sqrt{1+x^2})$ comme d^2 est continue sur $]2(x-\sqrt{1+x^2}), 2(x+\sqrt{1+x^2})[$ on voit que toutes les valeurs entre 0 et ∞ seront prises on peut donc approcher le vrai problème d'aussi près que l'on voudra.

1^{ère} généralisation.

Nous pouvons aisément généraliser les résultats précédents au cas où A_0 et B_0 sont deux vecteurs colonnes de longueurs quelconques. Il suffit de calculer :

$$P = A_0^T A_0 \quad , \quad Q = B_0^T B_0 \quad \text{et} \quad R = A_0^T B_0.$$

et de les transporter dans les équations précédentes.

C H A P I T R E III

ETUDE DU CAS GENERAL

Nous allons essayer de généraliser les résultats précédents. Nous passerons par l'étape intermédiaire d'une matrice A_0 à 3 lignes et 3 colonnes.

$$A_0 = \begin{pmatrix} a_0 & b_0 & c_0 \\ a'_0 & b'_0 & c'_0 \\ a''_0 & b''_0 & c''_0 \end{pmatrix}, \quad B_0 = \begin{pmatrix} d_0 \\ d'_0 \\ d''_0 \end{pmatrix}, \quad X = \begin{pmatrix} x \\ y \\ z \end{pmatrix}, \quad \lambda = \begin{pmatrix} \lambda \\ \mu \\ \nu \end{pmatrix}$$

écrivons $A^T A X - A^T B$

$$\begin{aligned} & \left[(a^2 + a'^2 + a''^2) \tilde{x} + (ab + a'b' + a''b'') \tilde{y} + (ac + a'c' + a''c'') \tilde{z} - (ad + a'd' + a''d'') \right] \\ & \left[(ab + a'b' + a''b'') \tilde{x} + (b^2 + b'^2 + b''^2) \tilde{y} + (bc + b'c' + b''c'') \tilde{z} - (bd + b'd' + b''d'') \right] \\ & \left[(ac + a'c' + a''c'') \tilde{x} + (bc + b'c' + b''c'') \tilde{y} + (c^2 + c'^2 + c''^2) \tilde{z} - (cd + c'd' + c''d'') \right] \end{aligned}$$

un simple examen nous montre que les termes contenant une des variables $a, a', a'', b, b', \dots, d''$ sont contenues dans une ligne et dans une colonne, cette disposition est générale.

Retour au cas général.

Nous avons à minimiser Trace $\left[(A - A_0)^T (A - A_0) \right] + (B - B_0)^T (B - B_0)$. Nous utiliserons la méthode des multiplicateurs de Lagrange, c'est-à-dire que nous chercherons les extrêmes de

$$\text{TRACE} \left[(A - A_0)^T (A - A_0) \right] + (B - B_0)^T (B - B_0) + \lambda^T (A^T A X - A^T B) \quad (1)$$

Nous avons comme matrices

$$A_0 = \begin{pmatrix} a_{n,1} & \dots & a_{n,m} \\ \vdots & & \vdots \\ a_{n,1} & \dots & a_{n,m} \end{pmatrix}, \quad B_0 = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}, \quad X = \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix}, \quad \lambda = \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_m \end{pmatrix}$$

Si nous dérivons (1) par rapport à a_{ik} , on obtient :

$$2 a_{ik} + \lambda_k (A_{i \cdot} \cdot X) + x_k (A_{i \cdot} \cdot \lambda) - \lambda_k b_i = 2 a_{ik}^0$$

en dérivant par rapport à b_i , on a :

$$2 b_i = 2 b_i^0 + A_{i \cdot} \cdot \lambda. \text{ Nous désignons par } A_{i \cdot} \text{ la } i^{\text{ème}} \text{ ligne de la ma-}$$

trice A.

et, en combinant les deux pour éliminer b_i ,

$$4 a_{ik} + 2 \lambda_k (A_{i \cdot} \cdot X) + 2 x_k (A_{i \cdot} \cdot \lambda) - \lambda_k (A_{i \cdot} \cdot \lambda) = 4 a_{ik}^0 + 2 \lambda_k b_i^0$$

Nous écrivons les unes sous les autres les équations obtenues en faisant varier k . On peut écrire matriciellement :

$$(2) \quad 4 A_{i \cdot}^T + 2 \lambda (A_{i \cdot} \cdot X) + 2 X (A_{i \cdot} \cdot \lambda) - \lambda (A_{i \cdot} \cdot \lambda) = 4 A_{i \cdot}^{\text{OT}} + 2 \lambda b_i^{\text{OT}}$$

Or $A_{i \cdot} \cdot X$ et $A_{i \cdot} \cdot \lambda$ sont des scalaires, on peut donc transposer ces quantités sans changer leur valeur.

$$(2) \text{ s'écrit } 4 A_{i \cdot}^T + 2 \lambda (X^T A_{i \cdot}^T) + 2 X (\lambda^T A_{i \cdot}^T) - \lambda (\lambda^T A_{i \cdot}^T) \\ = [4 I + 2(\lambda X^T + X \lambda^T) - \lambda \lambda^T] A_{i \cdot}^T = 4 A_{i \cdot}^{\text{OT}} + 2 \lambda b_i^{\text{OT}}$$

On pose $D = 4 I + 2(\lambda X^T + X \lambda^T) - \lambda \lambda^T$

On a donc $A_{i_0}^T = D^{-1} (4A_{i_0}^{oT} + 2\lambda b_{i_0}^{oT})$

On en déduit

$$\begin{aligned} A^T &= 2 D^{-1} [2A_{i_0}^T + \lambda B_{i_0}^T] \\ A &= 2 [2A_{i_0} + B_{i_0} \lambda^T] D^{-1} \end{aligned}$$

Or $b_{i_0} = b_{i_0}^o + \frac{1}{2} A_{i_0} \cdot \lambda$ et on en tire donc :

$$B = B_{i_0} + (2A_{i_0} + B_{i_0} \lambda^T) D^{-1} \lambda$$

Nous avons obtenu les valeurs de A et B en fonction de λ , il nous faut déterminer λ . On a l'équation encore inemployée : $A^T A X = A^T B = A^T (B_{i_0} + \frac{1}{2} A \lambda)$ ce qui s'écrit encore $A^T A (2X - \lambda) = 2A^T B_{i_0}$.

Calcul de D^{-1} .

$D = 4 I + 2(\lambda X^T + X \lambda^T) - \lambda \lambda^T$ est de la forme : $aI + bK$. En étudiant K on s'aperçoit qu'il vérifie l'équation

$$(3) K^3 - [2(\lambda^T X + X^T \lambda) - \lambda^T \lambda] K^2 + 4[\lambda^T X X^T \lambda - \lambda^T \lambda X^T X] K = 0$$

or K^2 introduit la matrice $X X^T$ qui n'existe pas dans K. On voit donc que K vérifie un polynôme minimal de degré 3. Or (3) ayant le coefficient de K^3 égal à 1 est donc le polynôme minimal

$$m(u) = u^3 - [2(\lambda^T X + X^T \lambda) - \lambda^T \lambda] u^2 + 4[\lambda^T X X^T \lambda - \lambda^T \lambda X^T X] u$$

K est donc une matrice du 3^o degré.

On pourra donc, sauf cas particulier, inverser D par une matrice de la forme $aI + bK + cK^2$. Procédons par identification

$$(4I+K)(aI+bK+cK^2) = 4 aI + (4b+a)K + (4c+b)K^2 + K^3$$

On remplace K^3 par sa valeur, en fonction de K^2 et K , obtenue à l'aide de (3).

On obtient le système :

$$(4) \quad \begin{cases} a = \frac{1}{4} \\ 4b + c \left\{ 4 \left[(\lambda^T x)^2 - \|\lambda\|^2 \|\times\|^2 \right] \right\} = -\frac{1}{4} \\ b + c \left[4\lambda^T x - \|\lambda\|^2 \right] = 0 \end{cases}$$

- on en tire a , b et c et on obtient tous les calculs faits :

$$D^{-1} = \frac{\left[4(1+\lambda^T x) - \|\lambda\|^2(1+\|\times\|^2) + (\lambda^T x)^2 \right] I - (2+\lambda^T x)(\lambda x^T + x \lambda^T) + (1+\|\times\|^2)\lambda \lambda^T + \|\lambda\|^2 \times \times^T}{16(1+\lambda^T x) - 4\|\lambda\|^2(1+\|\times\|^2) + 4(\lambda^T x)^2}$$

Ceci bien sûr si le système (4) est possible et déterminé, c'est-à-dire le déterminant de la matrice des coefficients est non nul, ou encore si $4(1+\lambda^T x) + (\lambda^T x)^2 \neq \|\lambda\|^2(1+\|\times\|^2)$ la condition sur λ qui en découle est la condition d'inversibilité de D .

Retour au problème.

Nous avons vu que l'équation en λ s'écrivait $A^T A(2x-\lambda) = 2A^T B_0$ donc :

$$4 D^{-1} (2A_0^T + \lambda B_0^T) (2A_0 + B_0 \lambda^T) D^{-1} (2x-\lambda) = 2 D^{-1} (2A_0^T + \lambda B_0^T) B_0$$

ou encore :

$$(5) \quad (2A_0^T + \lambda B_0^T) (2A_0 + B_0 \lambda^T) D^{-1} (2x-\lambda) = (2A_0^T + \lambda B_0^T) B_0$$

Calculons $D^{-1} (2x-\lambda)$. On obtient

$$D^{-1} (2x-\lambda) = \frac{x(2+\lambda^T x) - \lambda(1+||x||^2)}{4(1+\lambda^T x) + (\lambda^T x)^2 - ||\lambda||^2(1+||x||^2)}$$

On pose

$$\begin{aligned} P &= A_o^T A_o \\ Q &= B_o^T B_o \\ R &= A_o^T B_o \end{aligned}$$

Remarque.

Q est un scalaire

On peut encore écrire (5) sous la forme :

$$\begin{aligned} (5') \quad & (4P+2R\lambda^T+2\lambda R^T+Q\lambda\lambda^T) \cdot \left[(2+\lambda^T \tilde{x})\tilde{x} - (1+||\tilde{x}||^2)\lambda \right] \\ & = \left[4(1+\lambda^T \tilde{x}) - ||\lambda||^2(1+||\tilde{x}||^2) + (\lambda^T \tilde{x})^2 \right] (2R+Q\lambda) \end{aligned}$$

Nous poserons que, par définition, x^* est solution exacte du problème posé, c'est-à-dire $Px^* = R$. L'équation (5'), ordonnée en λ s'écrit :

$$\begin{aligned} & \lambda^T \left\{ \left[\tilde{x}\tilde{x}^T - (1+||\tilde{x}||^2)I \right] R - Q\tilde{x} \right\} \lambda \\ & + 2 \left\{ P \left[\tilde{x}\tilde{x}^T - (1+||\tilde{x}||^2)I \right] + (R^T \tilde{x} - Q) - R\tilde{x}^T \right\} \lambda \\ & + 4P(\tilde{x} - x^*) = 0 \end{aligned}$$

Nous ne pouvons résoudre directement cette équation aussi nous utiliserons la méthode de Newton généralisée $\lambda_{i+1} = \lambda_i - \left[F'(\lambda_i) \right]^{-1} F(\lambda_i)$

Nous voulons la racine la plus proche de zéro aussi nous prendrons comme vec-

teur de départ $\lambda_o = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$

ayant cette valeur de λ il suffit de la reporter dans les équations qui donnent $A - A_0$ et $B - B_0$ puis de calculer leurs normes euclidiennes.

C H A P I T R E IV

PROCEDURE ALGOL

PROCEDURE ERPOSYSQUA(CA,CB,X,U,T,IMPOSSIBLE) ;

Reel tableau CA,CB,X ; entier U,T ; ETIQUETTE IMPOSSIBLE ;

COMMENTAIRE cette procédure calcule l'erreur à postériori dans la méthode des moindres carrés. CA est la matrice du premier membre, CB celle du second. CA est de dimension $U \times T$ et CB de $1 \times U$ X est la solution numérique du système quadratique ;

Debut reel tableau GA[1:T,1:T], GB[1:T], A[1:T],

C[1:T], B[1:T,1:T], DELTA[1:T], LANDA[1:T],

MA[1:T], MB[1:U] ; Reel D1,D2,DENOM,PSLX,

NORMX,NORML ; Entier COMPT,K ;

Procédure GRESOLSYSLINE(A,B,X,N,IMPOSSIBLE) ;

reel tableau A,B,X ; Entier N ; étiquette impossible ;

debut triangularisation :

debut entier i,j,k ; reel R ;

pour k := 1 pas 1 jusqua N-1 faire

debut normal :

debut si ABS(A[K,K]) = 0 alors allera

Echange de lignes ;

pour i := k+1 pas 1 jusqua n faire

debut R := A[i,k]/A[k,k] ;

pour j := k+1 pas 1 jusqua N faire

A[i,j] := A[i,j] - R*A[k,j] ;

B[i] := B[i] - R*B[k]

fin

fin ;

```

Allera Retour ;
Echange de lignes :
  debut entier L,M ;
  M := k+1
  pour L := M tant que A[L,k] = 0      A
  L infeg N faire M := M+1 ;
  Si M = N+1 alors allera IMPOSSIBLE ;
  pour j := k pas 1 jusqua N faire
    debut R := A[k,j] ; A[k,j] := A[M,j] ;
    A[M,j] := R ;
    fin ;
  R := B[k] ; B[k] := B[M] ; B[M] := R ;
  Allera NORMAL
  fin ;
Retour
fin
fin triangularisation ;
RESSYSTRI :
  debut entier i,j ; Reel TX ;
  pour i := N pas -1 jusqua 1 faire
    debut TX := 0 ;
    pour j := N pas -1 jusqua i+1 faire
      TX := TX - X[j] * A[i,j] ;
    Si A[i,i] = 0 alors allera impossible ;
    X[i] := (B[i] + TX) / A[i,i]
    fin
  fin Ressystri
  Gresolsysline ;
fin

Q := 0 ;
pour i := 1 pas 1 jusqua U faire
  Q := Q + CB[i] * CB[i] ;
pour i := 1 pas 1 jusqua T faire

```

```

debut GB[i] := 0 ;
pour k := 1 pas 1 jusqu'a U faire
GB[i] := GB[i] + CA[k,i]*CB[k] ;
pour j := 1 pas 1 jusqu'a T faire
  debut GA[i,j] := 0 ;
  pour k := 1 pas 1 jusqu'a U faire
  GA[i,j] := GA[i,j] + CA[k,i]*CA[k,j]
  fin
fin ;
NORMX := 0 ; H := 0 ;
pour i := 1 pas 1 jusqu'a T faire
debut H := H + GB[i]*X[i] ;
C[i] := - 4*GB[i] ;
A[i] := 0 ;
pour j := 1 pas 1 jusqu'a T faire
  debut C[i] := C[i] + 4*GA[i,j]*X[j] ;
  GQ[i,j] := X[i]*X[j] ;
  B[i,j] := - GB[i]*X[j]
  fin
fin ;
pour i := 1 pas 1 jusqu'a T faire
debut GQ[i,i] := GQ[i,i] - (1+NORMX) ;
pour j := 1 pas 1 jusqu'a T faire
  debut A[i] := A[i] + GQ[i,j]*GB[j] ;
  pour k := 1 pas 1 jusqu'a T faire
  B[i,j] := B[i,j] + GA[i,k]*GQ[k,j]
  fin ;
A[i] := A[i] - Q*X[i] ;
B[i,i] := B[i,i] + H-Q
fin ;
pour i := 1 pas 1 jusqu'a T faire
pour j := 1 pas 1 jusqu'a T faire
B[i,j] := 2*B[i,j] ;

```

```

pour i := 1 pas 1 jusqua T faire
LANDA[i] := 0 ; COMPT := 0 ;
ITER : H := P := Q := 0 ;
pour i := 1 pas 1 jusqua T faire
  debut GB[i] := C[i] ;
  GB[i] := GB[i] + H*LANDA[i]
  fin ;
pour i := 1 pas 1 jusqua T faire
pour j := 1 pas 1 jusqua T faire
  debut GA[i,j] := LANDA[i]*A[j] ;
  GA[i,j] := GA[i,j] + B[i,j] ;
  GB[i] := GB[i] + B[i,j]*LANDA[j]
  fin ;
pour i := 1 pas 1 jusqua T faire
GB[i] := - GB[i] ;
GRESOLSYSLINE(GA,GB,DELTA,T,FICHU) ;
pour i := 1 pas 1 jusqua T faire
  debut P := P + ABS(DELTA[i]) ;
  Q := Q + ABS(LANDA[i])
  fin ;
Si COMPT = 0 alors allera Suite ;
Si P/Q infeg 10-3 alors COMPT := COMPT+1 ;
SUITE :
pour i := 1 pas 1 jusqua T faire
LANDA[i] := LANDA[i] + DELTA[i] ;
Si COMPT = 0 alors COMPT := 1 ;
Si COMPT = 5 alors allera CEFINI sinon
allera ITER ;
CEFINI :
PSLX := 0 ; NORML := 0 ;
pour i := 1 pas 1 jusqua T faire
  debut PSLX := PSLX + LANDA[i]*X[i] ;
  NORML := NORML + LANDA[i]*LANDA[i]
  fin ;

```

```

DENOM := 16*(1+PSLX) - 4*NORML(1+NORMX)
        + 4*PSLX*PSLX ;
pour i := 1 pas 1 jusqu'a T faire
pour j := 1 pas 1 jusqu'a T faire
    debut GA[i,j] := NORML*X[i]*X[j] + (1+NORMX)
        *LANDA[i]*LANDA[j] - (2+PSLX)*LANDA[i]
        *X[j] + LANDA[j]*X[i] ;
    GA[i,j] := GA[i,j]/DENOM
    fin ;
pour i := 1 pas 1 jusqu'a T faire
GA[i,i] := GA[i,i] + (4*(1+PSLX) - NORML*
(1+NORMX) + PSLX*PSLX)/DENOM ;
D1 := D2 := 0 ;
pour i := 1 pas 1 jusqu'a U faire
debut pour j := 1 pas 1 jusqu'a T faire
MA[j] := 4*CA[i,j] + 2*CB[i]*LANDA[j] ;
pour j := 1 pas 1 jusqu'a T faire
    debut GB[j] := 0 ;
    pour k := 1 pas 1 jusqu'a T faire
    GB[j] := GB[j] + MA[k]*GA[k,j]
    fin ;
pour j := 1 pas 1 jusqu'a T faire
    debut MA[j] := GB[j] ;
    D1 := D1 + (MA[j] - CA[i,j])*(MA[j]-CA[i,j])
    fin ;
MB[i] := 0 ;
pour j := 1 pas 1 jusqu'a T faire
MB[i] := MB[i] + (MA[j]*LANDA[j])/2 ;
D2 := D2 + MB[i]*MB[i]
fin ;
SAUTLIGNE ;
ECRIRE("DISTANCE_DES_DE","UX_PRØBLEMES") ;
SAUTLIGNE ;
ECRIRE("NØRME_DE_A-AO=", RAC2(D1),,
"NØRME_DE_B-BO=", RAC2(D2)) ;

```

SAUTLIGNE

FICHU : ECRIRE("SYSTEME_ IMPOSSI", "BLE") ;

SAUTLIGNE ; Allera IMPOSSIBLE

FIN ERPOSYSQUA ;

IV - EXEMPLES NUMERIQUES

Soit le système mal conditionné :

$$A_0 = \begin{bmatrix} 3,05 & 4 & 1 & 7 & 6,30 \\ 3,07 & 4,08 & 0,92 & 7,30 & 6,10 \\ 3,01 & 4,01 & 1 & 7,80 & 6,20 \\ 3 & 4,03 & 1,05 & 7,25 & 6,35 \\ 3,05 & 4,08 & 1,02 & 7,50 & 6,24 \\ 3,07 & 4,02 & 1,01 & 7,12 & 6,11 \\ 3,09 & 3,92 & 0,95 & 7,23 & 6,83 \\ 3,12 & 3,84 & 1,05 & 7,14 & 6,37 \\ 3,03 & 4,06 & 0,92 & 7,26 & 6,04 \\ 9,25 & 28,50 & 74,35 & 18,50 & 313,25 \end{bmatrix} \quad B_0 = \begin{bmatrix} 49 \\ 48,98 \\ 49,46 \\ 49,47 \\ 49,16 \\ 48,68 \\ 50,53 \\ 48,90 \\ 48,60 \\ 1165,1 \end{bmatrix}$$

$$x^* = \begin{pmatrix} 2 \\ 4 \\ 1 \\ 1 \\ 3 \end{pmatrix} \quad \tilde{x} = \begin{matrix} 2,146 \\ 3,772 \\ 0,814 \\ 1,033 \\ 3,058 \end{matrix}$$

Gresolsysline

On obtient :

$$\|A^* - A_0\|^2 = 0,96848 \quad 10^{-10}$$

$$\|B^* - B_0\|^2 = 0,15088 \quad 10^{-11}$$

nous voyons donc qu'une infime variation des coefficients de A_0 et B_0 (de l'ordre de 10^{-6}) peut donner des variations de l'ordre du dixième pour la solution.

Le même calcul fait avec les coefficients multipliés par 10 puis par 100 donne :

$$x_{10}^2 = \begin{cases} 2,140 \\ 3,776 \\ 0,811 \\ 1,034 \\ 3,058 \end{cases}$$

$$x_{100}^2 = \begin{cases} 2,144 \\ 3,772 \\ 0,806 \\ 1,034 \\ 3,060 \end{cases}$$

	X1	X10	X100
$\ A^* - A_0\ ^2 =$	$0,96848 \cdot 10^{-10}$	$0,674624 \cdot 10^{-8}$	$0,62370 \cdot 10^{-6}$
$b^* - b_0^2 =$	$0,15088 \cdot 10^{-10}$	$0,148445 \cdot 10^{-9}$	$0,97964 \cdot 10^{-8}$

nous voyons que l'erreur a posteriori dépend un peu de la façon dont nous ren-
trons les nombres puisque l'erreur suivant le cas n'est pas multipliée par
100 ou 10000 mais reste inférieure à cette estimation. Cela prouve que, d'après
la première remarque, en entrant un nombre en mémoire, nous commettons une erreur
d'arrondi qui peut provoquer une erreur sur le résultat assez grave. L'étude
du troisième cas où le nombre s'écrit en binaire de façon simple et où il n'y a
donc pas d'erreur d'arrondi nous montre bien que c'est la méthode qui est défi-
ciente.

Et voici un exemple mieux conditionné :

$$A_0 = \begin{bmatrix} 1 & 12 & 24 & 110 \\ 3 & 4 & 2 & 160 \\ 7 & 1 & 31 & 200 \\ 12 & 7 & 50 & 520 \\ 18 & 9 & 110 & 1040 \\ 34 & 11 & 500 & 2700 \\ 36 & 15 & 23 & 512 \\ 41 & 18 & 4 & 43 \\ 50 & 30 & 1 & 28 \end{bmatrix}$$

$$B_0 = \begin{bmatrix} 712 \\ 246 \\ 844 \\ 1614 \\ 3366 \\ 12878 \\ 1194 \\ 385 \\ 448 \end{bmatrix}$$

$$x^* = \begin{vmatrix} 2 \\ 10 \\ 20 \\ 1 \end{vmatrix}$$

$$\tilde{x} = \begin{vmatrix} 2,000011 \\ 9,999897 \\ 19,99989 \\ 1,000019 \end{vmatrix}$$

$$\|A^x - A_0\|^2 = 0,5436398 \cdot 10^{-8}$$

$$\|B^x - B_0\|^2 = 0,3213597 \cdot 10^{-12}$$

On voit que l'erreur sur x et la distance au problème vérifié le plus voisin sont du même ordre (10^{-6}).

BIBLIOGRAPHIE

PARTIE - A

On trouvera une étude des multiplicateurs de Lagrange dans :

J. BASS - COURS DE MATHEMATIQUES
Masson Editeur 1961

Une introduction plus détaillée à l'analyse à postériori des erreurs peut être trouvée dans :

[1] J. GACHES - COMPATIBILITE D'UNE SOLUTION CALCULEE AVEC LES DONNEES D'UN SYSTEME LINEAIRE A COEFFICIENTS INCERTAINS.
Thèse Besançon 1966

W. DETTLI - W. PRAGER - COMPATIBILITY OF APPROXIMATE SOLUTIONS OF LINEAR EQUATIONS WITH GIVEN ERROR BOUNDS FOR COEFFICIENTS AND RIGHT-HAND SIDES
Numerische Mathematik 6, (405-409) 1964

W. DETTLI - W. PRAGER - J.H. WILKINSON - ADMISSIBLE SOLUTIONS OF LINEAR SYSTEMS WITH NOT SHARPLY DEFINED COEFFICIENTS
J. Siam Numer. Anal. Ser. V. Vol. 2, N° 2 - 1965

On trouvera tous renseignements sur du second degré, polynômes minimaux et la méthode d'inversion de matrice utilisée dans :

N. GASTINEL - MATRICES DU SECOND DEGRE ET NORMES GENERALES EN ANALYSE NUMERIQUE LINEAIRE
Publications scientifique et techniques du Ministère de l'Air - 1962.

ETUDE DE L'ERREUR A POSTERIORI

DANS LE CAS

DU POLYNOME



CHAPITRE I

INTRODUCTION

Considérons une équation de la forme

$x^n + a_1^0 x^{n-1} + a_2^0 x^{n-2} + \dots + a_{n-1}^0 x + a_n^0 = 0$. Nous avons trouvé, par hypothèse, p solutions numériques de ce polynôme : $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_p$. Nous allons chercher s'il existe des jeux de coefficients a_1, a_2, \dots, a_n qui soient tels que $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_p$ en soient des solutions exactes, et si, parmi ces jeux, il s'en trouve un qui soit plus "proche" au sens de la norme euclidienne dans \mathbb{R}^n du jeu (a_1^0, \dots, a_n^0) . Nous utiliserons comme en A les multiplicateurs de Lagrange.

Nous avons ici " p " contraintes pour les a_i .

A savoir que $\sum_{i=1}^n a_i \tilde{x}_k^{n-i} + a_n = 0$ pour $k = 1, \dots, p$.

Suivant un plan similaire de celui suivi dans A nous étudierons d'abord le cas particulier du polynôme du second degré dont on connaît une puis deux solutions numériques.

C H A P I T R E II

ETUDE D'UN CAS PARTICULIER1^{er} Cas.

Nous connaissons une racine \tilde{x} d'une équation du second degré.

Nous allons donc chercher le doublet (b^*, c^*) de \mathbb{R}^2 sur la droite d'équation $\tilde{x}^2 + b\tilde{x} + c = 0$ qui se trouve à la distance minimale de (b_0, c_0) .

On appelle ε_i la valeur du polynôme pour $x = \tilde{x}_i$. Appliquons la méthode de Lagrange. On aura donc à dériver $(b-b^0)^2 + (c-c^0)^2 + \lambda(\tilde{x}^2 + b\tilde{x} + c)$ ce qui nous donne :

$$\begin{aligned} 2(b^* - b^0) + \lambda\tilde{x} &= 0 \\ 2(c^* - c^0) + \lambda &= 0 \end{aligned}$$

d'où l'on tire les valeurs de b^* et c^* en fonction de λ :

$$\begin{aligned} b^* &= b_0 - \frac{\lambda}{2} \tilde{x} \\ c^* &= c_0 - \frac{\lambda}{2} \end{aligned}$$

Reportons ces valeurs dans la contrainte :

$$\tilde{x}^2 + (b_0 - \frac{\lambda \tilde{x}}{2}) \tilde{x} + c_0 - \frac{\lambda}{2} = 0$$

$$2(\tilde{x}^2 + b_0 \tilde{x} + c_0) = \lambda(1 + \tilde{x}^2)$$

$$2 \varepsilon = \lambda(1 + \tilde{x}^2)$$

d'où :

$$\lambda = \frac{2 \varepsilon}{1 + \tilde{x}^2}$$

On en déduit d^2

$$\begin{aligned} d^2 &= (b^* - b_0)^2 + (c^* - c_0)^2 \\ &= \frac{\lambda^2}{4} \tilde{x}^2 + \frac{\lambda^2}{4} = \frac{\lambda^2}{4} (1 + \tilde{x}^2) \\ &= \frac{\varepsilon^2}{(1 + \tilde{x}^2)^2} (1 + \tilde{x}^2) = \frac{\varepsilon^2}{1 + \tilde{x}^2} \end{aligned}$$

2^{ème} Cas.

Nous connaissons les deux racines \tilde{x} et \tilde{y} .

Nous avons à dériver cette fois :

$$(b - b_0)^2 + (c - c_0)^2 + \lambda(\tilde{x}^2 + b\tilde{x} + c) + \mu(\tilde{y}^2 + b\tilde{y} + c).$$

Nous obtenons :

$$\begin{aligned} 2b^* &= 2b_0 - \lambda \tilde{x} - \mu \tilde{y} \\ 2c^* &= 2c_0 - \lambda - \mu \end{aligned}$$

d'autre part :

$$\begin{cases} \tilde{x}^2 + b^* \tilde{x} + c^* = 0 \\ \tilde{y}^2 + b^* \tilde{y} + c^* = 0 \end{cases}$$

Si on reporte les valeurs de b^* et c^* dans les deux contraintes on obtient :

$$\begin{aligned} \lambda(1+\tilde{x}^2) + \mu(1+\tilde{x}\tilde{y}) &= 2 \varepsilon_1 \\ \lambda(1+\tilde{x}\tilde{y}) + \mu(1+\tilde{y}^2) &= 2 \varepsilon_2 \end{aligned}$$

d'où l'on tire les valeurs de λ et μ :

$$\begin{aligned} \lambda &= 2 \frac{(1+\tilde{y}^2)\varepsilon_1 - (1+\tilde{x}\tilde{y})\varepsilon_2}{(1+\tilde{y}^2)(1+\tilde{x}^2) - (1+\tilde{x}\tilde{y})^2} \\ \mu &= 2 \frac{(1+\tilde{x}^2)\varepsilon_2 - (1+\tilde{x}\tilde{y})\varepsilon_1}{(1+\tilde{y}^2)(1+\tilde{x}^2) - (1+\tilde{x}\tilde{y})^2} \end{aligned}$$

CHAPITRE III

ETUDE DU CAS GENERAL

Nous avons un polynôme :

$x^n + a_1^o x^{n-1} + a_2^o x^{n-2} + \dots + a_{n-1}^o x + a_n^o = 0$ et p solutions numériques de ce polynôme : $\tilde{x}_1, \tilde{x}_2, \tilde{x}_3, \dots, \tilde{x}_p$ ($p < n$). En utilisant, comme au chapitre II, la méthode des multiplicateurs de Lagrange, nous aurons à dériver

$$\sum_{i=1}^n (a_i - a_i^o)^2 + \sum_{j=1}^p \lambda_j \left(\sum_{k=1}^n a_k \tilde{x}_j^{n-k} + x_j^n \right).$$

En dérivant par rapport à a_i on obtient :

$$2(a_i^* - a_i^o) + \sum_{j=1}^p \lambda_j \tilde{x}_j^{n-i} = 0$$

$$\text{d'où : } a_i^* = a_i^o - \frac{1}{2} \sum_{j=1}^p \lambda_j \tilde{x}_j^{n-i}$$

ou encore avec :

$$A^* = \begin{pmatrix} a_1^* \\ a_2^* \\ \vdots \\ a_n^* \end{pmatrix}$$

$$B = \begin{pmatrix} 1 & 1 & \dots & 1 \\ \tilde{x}_1^2 & \tilde{x}_2^2 & \dots & \tilde{x}_p^2 \\ \tilde{x}_1^2 & \tilde{x}_2^2 & \dots & \tilde{x}_p^2 \\ \vdots & \vdots & \dots & \vdots \\ \tilde{x}_1^{2n} & \tilde{x}_2^{2n} & \dots & \tilde{x}_p^{2n} \end{pmatrix}$$

$$A_0 = \begin{vmatrix} a_1^0 \\ a_2^0 \\ \vdots \\ \vdots \\ a_n^0 \end{vmatrix} \qquad \lambda = \begin{vmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \vdots \\ \lambda_p \end{vmatrix}$$

Notre résultat se présente sous la forme matricielle suivante :

$$A^* = A_0 - \frac{1}{2} B \cdot \lambda$$

il nous faut maintenant calculer les λ_j .

Nous savons que :
$$\sum_{i=1}^n a_i^* x_k^{n-i} + x_k^n = 0$$

et que :
$$\sum_{i=1}^n a_i^0 x_k^{n-i} + x_k^n = \epsilon_k$$

et par soustraction de ces deux égalités

$$\sum_{i=1}^n (a_i^* - a_i^0) x_k^{n-i} = -\epsilon_k$$

or nous avons vu que :

$$a_i^* - a_i^0 = -\frac{1}{2} \cdot \sum_{j=1}^p \lambda_j x_j^{n-i}$$

on en déduit :

$$\sum_{i=1}^n \left(\sum_{j=1}^p \lambda_j x_j^{n-i} \right) x_k^{n-i} = 2 \epsilon_k$$

Nous avons affaire à des sommes finies, on peut donc inverser les signes sommes et écrire :

$$\sum_{j=1}^p \left[\sum_{i=1}^n (x_j x_k)^{n-i} \right] \lambda_j = 2 \epsilon_k \quad k = 1, 2, \dots, p$$

ou encore matriciellement

$$B^T (A^* - A_0) = -\epsilon$$

Or $A^* - A_0 = -\frac{1}{2} B\lambda$

donc $B^T B\lambda = 2\epsilon$

et $\lambda = 2(B^T B)^{-1} \epsilon$

Nous avons obtenu un système linéaire de p équations à p inconnues que nous traiterons soit par la méthode de Gauss soit par une autre et on a :

$$d^2 = \left\| (B^T B)^{-1} \epsilon \right\|^2$$

PROCEDURE ERPOPOL(AO,X,N,P,IMPOSSIBLE) ;

REEL TABLEAU AO,X ; ENTIER N,P ; ETIQUETTE IMPOSSIBLE ;

COMMENTAIRE cette procédure calcule l'erreur à postériori commise dans le cas où nous avons trouvé P solutions numériques réunies dans le tableau X du polynôme du N^{ième} degré dont les coefficients sont rangés dans le tableau AO de telle façon que la somme de l'indice et de la puissance soit toujours égale à N ;

DEBUT ENTIER I,J,K ; REEL TABLEAU GX [1:N,1:P], GA [1:P,1:P],

EPS [1:P], LANDA [1:P] ;

REEL D, DELTA ; 50-10-10

PROCEDURE GRESOLSYSLINE(A,B,X,N,IMPOSSIBLE) ;

REEL TABLEAU A,B,X ; ENTIER N ; ETIQUETTE IMPOSSIBLE ;

CODE ;

DEBUT TRIANGULARISATION :

DEBUT ENTIER I,J,K ; REEL R ;

POUR K := 1 PAS 1 JUSQUA N-1 FAIRE

DEBUT NORMAL :

DEBUT SI ABS(A[K,K]) = 0 ALORS ALLERA ECHANGE DE LIGNES ;

POUR I := K+1 PAS 1 JUSQUA N FAIRE

DEBUT R := A[I,K]/A[K,K] ;

POUR J := K+1 PAS 1 JUSQUA N FAIRE

A[I,J] := A[I,J] - R*A[K,J] ;

B[I] := B[I] - R*B[K]

FIN

FIN ;

ALLERA RETOUR ;

ECHANGE DE LIGNES :

DEBUT ENTIER L,M ; M := K+1 ;

POUR L := M TANTQUE A[L,K] = 0 ET L INFEG N

FAIRE M := M+1 ;

SI M = N+1 ALORS ALLERA IMPOSSIBLE ;

POUR J := K PAS 1 JUSQUA N FAIRE
DEBUT R := A[K,J] ; A[K,J] := A[M,J] ;
 A[M,J] := R
FIN ;
 R := B[K] ; B[K] := B[M] ; B[M] := R ;
ALLERA NORMAL

FIN ; RETOUR :

FIN

FIN TRIANGULARISATION ;

RESSYSTRI :

DEBUT ENTIER I,J ; REEL TX ;

POUR I := N PAS -1 JUSQUA 1 FAIRE

DEBUT TX := 0 ;

POUR J := N PAS -1 JUSQUA I+1 FAIRE

TX := TX - X[J] × A[I,J] ;

SI A[I,I] = 0 ALORS ALLERA IMPOSSIBLE ;

X[I] := (B[I] + TX) / A[I,I]

FIN

FIN RESSYSTRI ;

FIN GRESOLSYSLINE ;

POUR I := 1 PAS 1 JUSQUA P FAIRE

GX[1,I] := 1.0 ;

POUR I := 2 PAS 1 JUSQUA N FAIRE

POUR J := 1 PAS 1 JUSQUA P FAIRE

GX[I,J] := GX[I-1,J] × X[J] ;

POUR I := 1 PAS 1 JUSQUA P FAIRE

DEBUT EPS[I] := 2 × GX[N,I] × X[I] ;

POUR J := 1 PAS 1 JUSQUA N FAIRE

EPS[I] := EPS[I] + 2 A0[J] × GX[N+1-J,I]

FIN ;

POUR I := 1 PAS 1 JUSQUA P FAIRE

POUR J := 1 PAS 1 JUSQUA P FAIRE

POUR K := 1 PAS 1 JUSQUA N FAIRE

504040
GA[I,J] := GA[I,J] + GX[K,I]×GX[K,J] ;
GRESOLSYSLINE(GA,EPS,LANDA,P,IMPOSSIBLE) ;
POUR I := 0 PAS 1 JUSQUA N-1 FAIRE
DEBUT DELTA := 0 ;
 POUR J := 1 PAS 1 JUSQUA P FAIRE
 DELTA := DELTA + GX[N-I,J]×LANDA[J] ;
 D := D + 0.25 DELTA DELTA
FIN ;
SAUTLIGNE ; ECRIRE (D) ;
FIN ERPOPOL ;

Exemples numériques.

Nous avons appliqué les résultats précédents à l'équation :

$$x^4 - 3,58 x^3 + 4,8059 x^2 - 2,867222 x + 0,6414408$$

dont les solutions exactes sont : 0,88 ; 0,89 ; 0,90 ; 0,91.

$$p = 3 \quad \tilde{x}_1 = 0,8775 ; \tilde{x}_2 = 0,9012 ; \tilde{x}_3 = 0,89$$

$$d^2 = 0,515038 \quad 10^{-9}$$

$$p = 2 \quad \tilde{x}_1 = 0,8775 ; \tilde{x}_2 = 0,9012$$

$$d^2 = 0,460637 \quad 10^{-16}$$

$$p = 1 \quad \tilde{x}_1 = 0,88$$

$$d^2 = 0,39201 \quad 10^{-16}$$

$$p = 1 \quad \tilde{x}_1 = 0,89989$$

$$d^2 = 0,1158634 \quad 10^{-16}$$

nous voyons qu'une variation de l'ordre de 10^{-8} des coefficients entraîne une variation de 10^{-3} sur les solutions. Cette équation est très mal conditionnée car les racines sont très proches. Les derniers chiffres montrent que le seul fait d'introduire un nombre en machine donne, par suite des troncatures dans la transformation de décimal en binaire, une erreur en 10^{-8} qui pourra influencer sur les racines de façon sensible.

Equation mieux conditionnée.

$$x^3 - 106 x^2 + 605 x - 500$$

les racines sont : 1, 100, 5 ;

$$p = 2 \quad \tilde{x}_1 = 99,5 ; \tilde{x}_2 = 4,5 ;$$

$$d^2 = 0,821305 \cdot 10^{+2}$$

$$p = 1 \quad \tilde{x}_1 = 100,000002$$

$$d^2 = 0,5959271 \cdot 10^{-19}$$

nous voyons qu'il faut une variation de l'ordre de 10 pour faire varier les racines de 0,5.

BIBLIOGRAPHIE

PARTIE - B

Aucune nouvelle notion n'est introduite.

On trouvera cependant une étude du même problème mais traitée par des méthodes topologiques dans :

J. GACHES - COMPATIBILITE D'UNE SOLUTION CALCULEE AVEC LES DONNEES D'UN SYSTEME LINEAIRE A COEFFICIENTS INCERTAINS
Thèse Besançon 1966.

NGUYEN HUU VINH - SEMI NORME DUALE GENERALISEE ET APPROXIMATION D'UN VECTEUR
C.R. Acad. Sc. Paris t. 262 - pp. 1456-59 (1966)

- C -

ETUDE DE L'ERREUR A POSTERIORI

DANS LE CAS

DES VALEURS PROPRES

INTRODUCTION

L'analyse à postériori des erreurs dans le cas des valeurs propres n'est pas possible par la méthode que nous avons employée jusqu'à maintenant. En effet la contrainte : la valeur propre numérique $\tilde{\lambda}$ trouvée doit être solution exacte du polynôme caractéristique d'une certaine matrice A , n'est pas du tout aisément manipulable. Nous allons donc utiliser un biais. Au lieu de considérer que nous avons seulement une valeur propre $\tilde{\lambda}$, nous admettrons par hypothèse que nous connaissons aussi un vecteur propre \tilde{x} .

Ce qui nous permettra de résoudre, au passage, le problème de l'erreur à postériori dans ce cas. En effet, la contrainte dans ce cas se réduit à $A\tilde{x} = \tilde{\lambda}\tilde{x}$ qui est linéaire. Nous remarquons aussi que \tilde{x} apparaît dans les deux membres, nous pourrions donc considérer dorénavant uniquement les vecteurs \tilde{x} de norme unité.

Ayant résolu ce problème, nous reviendrons au problème initial. Pour cela nous essaierons de minimiser, en faisant varier \tilde{x} sur la boule unité, la norme de l'erreur à postériori dépendante de \tilde{x} que nous avait donnée la solution du problème précédent.

C H A P I T R E II

ERREUR A POSTERIORI QUAND ON CONNAIT UNE VALEUR PROPRE NUMERIQUE $\tilde{\lambda}$ ET LE VECTEUR PROPRE NUMERIQUE \tilde{x} CORRESPONDANT

Ce cas est redevable de la méthode qui nous est maintenant habituelle.

Posons le problème : nous avons une matrice carrée A_0 . Nous cherchons une valeur propre et le vecteur propre correspondant par une méthode numérique quelconque (Givens, Jacobi, Rutishauser). Nous trouvons $\tilde{\lambda}$ et \tilde{x} . Nous allons chercher s'il existe des matrices A qui admettent $\tilde{\lambda}$ et \tilde{x} comme valeur et vecteur propre exacts, et nous verrons si, parmi ces matrices A , il y en a une qui soit plus "proche", toujours au sens de la norme euclidienne sur \mathbb{R}^{n^2} , de A_0 . Pour cela nous utiliserons encore la méthode de Lagrange. Nous dériverons donc :

$$\sum_i \sum_j (a_{ij} - a_{ij}^0)^2 + \mu^T (A\tilde{x} - \tilde{\lambda}\tilde{x}) \quad \text{on a } \mu = \begin{matrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{matrix}$$

On obtient :

$$2(a_{ij}^* - a_{ij}^0) + \mu_i x_j^* = 0$$

ce qui s'écrit matriciellement :

$$\boxed{2(A^* - A_0) + \mu \tilde{x}^T = 0}$$

il nous faut maintenant calculer la valeur de μ . Nous disposons pour cela de la contrainte :

$$A^* \tilde{x} = \tilde{\lambda} \tilde{x} \quad \text{d'où}$$

$$(A_0 - \frac{\mu}{2} \tilde{x} \tilde{x}^T) \tilde{x} = \tilde{\lambda} \tilde{x}$$

d'où encore :

$$A_0 \tilde{x} - \tilde{\lambda} \tilde{x} = \frac{\mu}{2} \tilde{x}$$

nous voyons donc que μ est le double du résidu η on peut donc encore écrire :

$$A^* - A_0 = \eta \tilde{x} \tilde{x}^T$$

et
$$\|A^* - A_0\|^2 = \|(A_0 - \lambda I) \tilde{x} \tilde{x}^T\|^2$$

c'est le résultat cherché.

Remarque.

La matrice, admettant $\tilde{\lambda}$ et \tilde{x} comme valeur et vecteur propre, qui était la plus proche de A_0 est donc $A^* = A_0 - \eta \tilde{x} \tilde{x}^T$ or elle se déduit très facilement de l'équation matricielle $A_0 \tilde{x} - \tilde{\lambda} \tilde{x} = \eta \tilde{x}$ puisque $\|\tilde{x}\| = 1$.

C H A P I T R E III

ERREUR A POSTERIORI DANS LE CALCUL DES VALEURS PROPRES $\tilde{\lambda}$

Nous allons utiliser dans ce chapitre les résultats du chapitre II.

Nous avons vu que l'erreur dans le cas où nous connaissions un vecteur propre numérique \tilde{x} en plus de la valeur propre $\tilde{\lambda}$ était égale à $|(A_0 - \tilde{\lambda}I) \tilde{x} \tilde{x}^T|$ d'autre part \tilde{x} est de norme 1.

Pour résoudre le problème posé dans ce chapitre nous allons essayer d'éliminer \tilde{x} pour cela nous chercherons le minimum de $|(A_0 - \tilde{\lambda}I) \tilde{x} \tilde{x}^T|$ lorsque \tilde{x} se déplace sur la boule unité. Pour ce faire, nous utiliserons encore une fois les multiplicateurs de Lagrange. Calculons d'abord la valeur des éléments de $(A_0 - \tilde{\lambda}I) \tilde{x} \tilde{x}^T$ sachant que nous appelons dorénavant B la quantité $A_0 - \tilde{\lambda}I$.

On a donc :

$$B \tilde{x} = \left[\sum_i b_{ji} \tilde{x}_i \right]_{j=1 \dots n}$$

$$B \tilde{x} \tilde{x}^T = \left[\sum_i b_{ji} \tilde{x}_i \tilde{x}_k \right]_{\substack{j=1 \dots n \\ k=1 \dots n}}$$

et
$$||B \tilde{x} \tilde{x}^T||^2 = \sum_{j=1}^n \sum_{k=1}^n \left(\sum_{i=1}^n b_{ji} \tilde{x}_i \tilde{x}_k \right)^2 = F$$

Si nous dérivons ce carré de norme par rapport à x_ℓ on obtient :

$$\frac{\partial F}{\partial x_\ell} = 2 \sum_{j=1}^n \left[\left(\sum_{i=1}^n b_{ji} x_i \right) \left\{ b_{j\ell} \left(\sum_{k=1}^n x_k^2 \right) + x_\ell \left(\sum_{i=1}^n b_{ji} x_i \right) \right\} \right]$$

Ce que l'on peut encore écrire :

$$2 \sum_j \left[b_{j\ell} \left(\sum_i b_{ji} x_i \right) \left(\sum_k x_k^2 \right) + x_\ell \left(\sum_i b_{ji} x_i \right)^2 \right]$$

et matriciellement :

$$2 \left[||x||^2 B_\ell^T \cdot BX + x_\ell X^T B^T BX \right]$$

nous utilisons la méthode des multiplicateurs de Lagrange. Nous avons ici une seule contrainte $||x||^2 = 1$. Nous devons ajouter $\mu \frac{\partial (||x||^2 - 1)}{\partial x_\ell}$ au résultat précédent. Nous obtiendrons, en écrivant toutes les dérivées partielles ($\ell=1,2,\dots,n$) les unes sous les autres, un système qui, matriciellement, s'exprimera par :

$$\frac{\partial \varphi}{\partial x_\ell} = X^T X B^T B X + X X^T B^T B X = \mu X$$

$B X X^T B^T$ est un scalaire, on a donc

$$\frac{\partial \varphi}{\partial x} = \left[||x||^2 B^T B + (||BX||^2 - \mu) I \right] X$$

ou encore :

$$\frac{\partial \varphi}{\partial x} = \left[B^T B + \left(\frac{\|Bx\|^2 - \mu}{\|x\|^2} \right) I \right] x = 0$$

et

$$\left[B^T B - kI \right] x = 0$$

pour avoir une solution non nulle de ce problème il faut que

$\det [B^T B - kI] = 0$ donc que k soit valeur propre. $B^T B$ est symétrique, donc,

d'après le théorème de SCHUR, il y a n valeurs propres réelles. D'autre part, x est un des vecteurs propres. Comme μ ne sert à rien d'autre qu'à assurer la normalisation de x les vecteurs \tilde{x} cherchés seront donc les n vecteurs propres normalisés de $B^T B$.

Or :

$$\|B\tilde{x}_i\tilde{x}_i^T\|^2 = \text{trace} (\tilde{x}_i\tilde{x}_i^T B^T B \tilde{x}_i\tilde{x}_i^T)$$

Comme

$$B^T B \tilde{x}_i = \lambda_i \tilde{x}_i$$

donc

$$\begin{aligned} \|B\tilde{x}_i\tilde{x}_j^T\|^2 &= \text{trace} (\lambda_i \underbrace{x_i x_i^T x_i x_i^T}_{=1}) \\ &= \text{trace} (\lambda_i x_i x_i^T) \end{aligned}$$

$$\text{comme} \quad \sum_{j=1}^n \sum_{i=1}^n (x_i x_j)^2 = \sum_{i=1}^n x_i^2 \sum_{j=1}^n x_j^2 = 1$$

$$\text{on a} \quad \|B\tilde{x}_i\tilde{x}_j^T\|^2 = \lambda_i \sum_{j=1}^n x_j^2 = \lambda_i$$

la solution qui nous intéresse est donc celle qui correspond à λ_i minimum.

C,IV,1
CHAPITRE-IV

PROCEDURE ALGOL

Procédure ERPOSVALPRO (AoR, LANDAR,N,ERREUR) ;
réel tableau AoR ; Réel LANDAR, ERREUR ; Entier N ;
Commentaire Cette procédure calcule l'erreur a postériori dans le calcul des valeurs
propres d'une matrice $N \times N$ appelée Aor. La valeur propre approchée s'appelle LANDAR ;
début entier i, j, k ; réel tableau MUDEU [1 : N], $DI [1:N], D [1:N]$;
VR [1 : N, 1 : N] , VI [1:N, 1:N] ;
réel EPS ; $SO2050$ DI, VI ;
procédure JACOBI (N, AR, LAMBDA, VR, EPSILON) ;
entier N ; réel EPSILON ; réel tableau AR, LAMBDA, VR ; CODE ;
début entier i, j, p, q, k ; réel X, Y, WR, E, F, D1, D2, D, MUR, R, CR, G, H,
UR, ZR, S, T ;
procédure TRANSFO (R, TR, UR, WR, ZR) ;
réel R, TR, UR, WR, ZR ;
début ZR := R * UR + TR * WR ;
si ABS(ZR) $\leq 10^{-15}$ alors ZR := 0 ;
fin
pour i := 1 pas 1 jusqua N faire
pour j := 1 pas 1 jusqua N faire
VR[I,j] := si i = J alors 1 sinon 0 ;
itération : pour p := 1 pas 1 jusqua N-1 faire
pour Q := P+1 pas 1 jusqua N faire
si ABS (AR[P,Q]) $\neq 0$ alors
début X := AR[P,P] ; Y := AR[Q,Q] ; WR := AR[P,Q] ;
E := X-Y ; F = RAC3(E \uparrow 2 + 4 * WR \uparrow 2) ;
D1 := E+F ; D2 := E-F ;
D := si ABS (D1) > ABS (D2) alors D1 sinon D2 ;
MUR := 2 * WR/D ;
R := 1/RAC2 (1+MUR \uparrow 2) ;
CR := R * MUR ; G := 2 * (CR * AR[P,Q]) * R ;
H := X * R \uparrow 2 + Y * CR \uparrow 2 ;
T := Y * R \uparrow 2 + X * CR \uparrow 2 ;
AR[P,P] := G + H ; AR[Q,Q] := T-G ;
AR[P,Q] := 0 ;
si P ≥ 2 alors pour K := 1 pas 1 jusqua P-1 faire

C,IV,2

```

début UR := AR[P,K] ; WR := AR[K,Q] ;
    TRANSFO (R,-CR, UR, WR, ZR) ;
    AR[K,P] := ZR ;
    TRANSFO (R, CR, WR, UR, ZR) ;
    AR[K,Q] := ZR
fin ;
si Q > P + 2 alors pour K := P+1 pas 1 jusqu Q-1 faire
début VR := AR[P,K] ; WR := AR[K,Q] ;
    TRANSFO (R, -CR, WR, UR, ZR) ;
    AR[K,Q] := ZR ;
    TRANSFO (R, CR, UR, WR, ZR) ;
    AR[P,K] := ZR
fin ;
si Q < N-1 alors pour K := Q+1 pas 1 jusqu N faire
début UR := AR[P,K] ; WR := AR[Q,K] ;
    TRANSFO (R, CR, UR, WR, ZR) ;
    AR[P,K] := ZR ;
    TRANSFO (R, -CR, WR, UR, ZR) ;
    AR[Q,K] := ZR
fin ;
pour i := 1 pas 1 jusqu N faire
début UR := VR [I,P] ; WR := VR [I,Q] ;
    TRANSFO (R, CR, UR, WR, ZR)
    VR[I,P] := ZR ;
    TRANSFO (R, -CR, WR, UR, ZR) ;
    VR[I,Q] := ZR
fin ;
fin itération ;
S := T := 0 ;
Pour i := 1 pas 1 jusqu N-1 faire
début pour j := i+1 pas 1 jusqu N faire
    S := S + ABS (AR[I,j]) ; T := T + ABS (AR[I,I]) ;
fin ;

```

C,IV,3

```

si T = 0 alors allera FINI ;
si S/T > EPSILON alors allera ITERATION ;
FINI : pour J := 1 pas 1 jusqua N faire
      LAMBDA[J] := AR [J,J]

```

```

fin JACOBI ;

```

```

pour i := 1 pas 1 jusqua N faire
  AoR[I,I] := AoR[I,I] - LANDAR ;
pour i := 1 pas 1 jusqua N faire
pour j := 1 pas 1 jusqua N faire
  début VR [i,j] := 0 ; VR [I,I] := 0 ; VR [I,I] := 0 ;
    pour K := 1 pas 1 jusqua N faire
      VR[i,j] := VR[i,j] + AoR[I,K] * AoR[J,K] ;

```

```

fin ;

```

```

EPS := 10-8 ;
502050
JACOBI (N, VR, MUDEU, AOR, EPS) ;

```

```

ERREUR := MUDEU [1] ;

```

```

pour i := 2 pas 1 jusqua N faire
  si MUDEU [i] < ERREUR alors ERREUR := MUDEU [i] ;

```

```

fin ERPOSVALPRO ;

```

Exemples numériques

Soient la matrice

$$\begin{bmatrix} 9 & 1 & -2 & 1 \\ 1 & 8 & -3 & -2 \\ 2 & -3 & 7 & -1 \\ 1 & -2 & -1 & 6 \end{bmatrix}$$

Cette matrice a comme valeur propre exacte 6.

Si on prend comme valeur propre approchée $\tilde{\lambda} = 6,0004$ on trouve comme carré de l'erreur $d^2 = 0,10425 \cdot 10^{-7}$ si $\tilde{\lambda}$ vaut 6,75 alors $d^2 = 0,56249$

Considérons maintenant la matrice

$$\begin{bmatrix} 33 & -3 & 0 & -4 & 0 & 8 & 0 & -4 \\ 3 & 33 & 4 & 0 & -8 & 0 & 4 & 0 \\ 0 & 4 & 29 & 1 & -12 & -2 & -8 & -2 \\ 4 & 0 & 1 & 29 & -2 & -12 & -2 & -8 \\ 0 & -8 & -12 & -2 & 25 & 1 & -4 & -2 \\ 8 & 0 & -2 & -12 & 1 & 25 & -2 & -4 \\ 0 & 4 & -8 & -2 & -4 & -2 & 21 & 1 \\ 4 & 0 & -2 & -8 & -2 & -4 & 1 & 21 \end{bmatrix}$$

qui admet 6 et 48 comme valeur propre exacte

$$\tilde{\lambda} = 6,08 \text{ donne } d^2 = 0,63906 \cdot 10^{-2}$$

$$\text{et } \tilde{\lambda} = 48,005 \text{ donne } d^2 = 0,23329 \cdot 10^{-4}$$

Le programme est rapide : 22 millièmes d'heures pour 6 calculs d'erreurs appartenant à des matrices 8×8 .

BIBLIOGRAPHIE

PARTIE - C

On trouvera des procédés d'analyse directe d'erreur dans le calcul des valeurs propres dans les livres de :

A.M. TURING - ROUNDING-OFF ERRORS IN MATRIX PROCESSES
Quart. J. MECH. APPL. MATH. pp. 287-308

[2] J.H. WILKINSON - ROUNDING ERRORS IN ALGEBRAIC PROCESSES
Her majesty's stationery office 1963

On trouvera d'ailleurs dans ce dernier livre une étude du problème de l'estimation à postériori mais qui n'est résolue que pour une matrice symétrique.
Les méthodes de calcul des vecteurs et valeurs propres sont tirées du livre.

N. GASTINEL - ANALYSE NUMERIQUE LINEAIRE
Hermann Editeur.

- D -

ERREURS A POSTERIORI

DANS LE CAS

DES ESPACES DE HILBERT

C H A P I T R E I

INTRODUCTION

Dans cette partie, nous étudierons l'erreur à postériori dans un espace de Hilbert pour les problèmes linéaires. Nous nous placerons dans un espace de Hilbert H séparable sur \mathbb{R} . Le produit scalaire en sera noté $(,)$ et la norme correspondante $\| \cdot \|$.

Soit $\mathcal{L}(H,H)$ l'ensemble des endomorphismes continus de H dans H . On définira la norme de l'un deux A , par $\|A\| = \sup_{\|x\|=1} \|Ax\|$.

Nous définirons ensuite ce que nous appellerons un double produit scalaire et, avec la norme correspondante, on définira une "distance entre problèmes".

Nous pourrons ainsi, à l'aide d'une adaptation de la méthode de Lagrange, estimer le minimum de la distance entre le problème résolu et le problème effectif.

C H A P I T R E II

RESULTATS PRELIMINAIRES

Soit donc H un espace de Hilbert séparable sur \mathbb{R} $A \in \mathcal{L}(H, H)$ et soient $\{x_p\}, \{y_q\}$ ($p, q=1, 2, \dots$) deux systèmes orthonormés complets quelconques de H . C'est-à-dire tels que :

$$(x_i, x_j) = (y_i, y_j) = \delta_{ij} \quad i, j = 1, 2, \dots$$

Définissons la quantité :

$$N^2(A, \{x_p\}, \{y_q\}) = \sum_{p,q} |(Ax_p, y_q)|^2 \quad (1)$$

Celle-ci va nous servir à caractériser une partie α de $\mathcal{L}(H, H)$ qui sera l'ensemble des endomorphismes A tels que $N^2(A, \{x_p\}, \{y_q\})$ soit fini.

Nous allons d'abord montrer que α ne dépend pas de $\{x_p\}$ ou de $\{y_q\}$. Il suffit de vérifier que cela est vrai pour un A quelconque $\varepsilon \alpha$ $\sum_{p,q} |(Ax_p, y_q)|$ étant une série à termes positifs nous n'avons pas à nous préoccuper de l'ordre des sommations et nous pouvons donc écrire :

$$\sum_{p,q} |(Ax_p, y_q)|^2 = \sum_{p=1}^{\infty} \left(\sum_{q=1}^{\infty} |(Ax_p, y_q)|^2 \right)$$

mais (Ax_p, y_q) sont les coefficients de Fourier de Ax_p sur la base $\{y_q\}$. On a donc l'égalité de Bessel-Parseval :

$$\sum_{q=1}^{\infty} |(Ax_p, y_q)|^2 = \|Ax_p\|^2$$

et (1) s'écrit donc :

$$\sum_{p,q} |(Ax_p, y_q)|^2 = \sum_{p=1}^{\infty} \|Ax_p\|^2$$

or $(Ax_p, y_q) = (x_p, A^* y_q) = (A^* y_q, x_p)$

donc $\sum_{p,q} |(Ax_p, y_q)|^2 = \sum_{q=1}^{\infty} \|A^* y_q\|^2$

Comme A^* ne dépend que de A et que les deux systèmes $\{x_p\}$ et $\{y_q\}$ sont indépendants, on voit donc que la quantité N^2 ne dépend que de A (et aussi que $N(A) = N(A^*)$).

Montrons que \mathfrak{a} est un sous espace vectoriel de $\mathcal{L}(H, H)$

$\forall \lambda \in \mathbb{R}$ -on a $N(\lambda A) = |\lambda| \cdot N(A)$ donc $\lambda A \in \mathfrak{a}$ si A y appartient

Montrons que si A et $B \in \mathfrak{a}$ alors $A + B$ aussi

En effet, soit $\{x_p\}$ un système quelconque orthonormé. On a :

$$\|Ax_p + Bx_p\|^2 = \|Ax_p\|^2 + \|Bx_p\|^2 + 2(Ax_p, Bx_p)$$

$$\|Ax_p - Bx_p\|^2 = \|Ax_p\|^2 + \|Bx_p\|^2 - 2(Ax_p, Bx_p)$$

et en sommant ces deux identités

$$\|Ax_p + Bx_p\|^2 \leq \|Ax_p + Bx_p\|^2 + \|Ax_p - Bx_p\|^2 = 2(\|Ax_p\|^2 + \|Bx_p\|^2)$$

ce qui montre bien que $(A+B) \in \mathfrak{a}$ si A et B y appartiennent.

D,II,3

Montrons maintenant que N est une norme $N(\lambda A) = |\lambda| \cdot N(A)$ est trivial.

Montrons que $N(A+B) \leq N(A) + N(B)$

Pour cela considérons l'espace $\ell_2 \subset (\mathbb{R} \times \mathbb{R} \times \mathbb{R} \times \dots)$. Un élément $x = \{x_1, x_2, \dots\}$ est dit appartenir à ℓ_2 si $\sum x_i^2$ existe. C'est d'ailleurs le carré de la norme habituellement attachée à ℓ_2 .

$$\begin{aligned} \left\| \sum_{p=1}^{\infty} \|Ax_p + Bx_p\|_H \right\|_{\ell_2} &\leq \left\| \sum_{p=1}^{\infty} (\|Ax_p\|_H + \|Bx_p\|_H) \right\|_{\ell_2} \\ &\leq \left\| \sum_{p=1}^{\infty} \|Ax_p\|_H \right\|_{\ell_2} + \left\| \sum_{p=1}^{\infty} \|Bx_p\|_H \right\|_{\ell_2} \end{aligned}$$

donc $N(A+B) \leq N(A) + N(B)$

et $N(A) = 0$ entraînant $Ax_p = 0 \quad \forall x_p \in$ système orthonormé quelconque, on a donc

$A = 0$

N est bien une norme.

Remarque.

Si $\|x\| = 1$ on peut prendre $x = \{x_1, 0, 0, \dots\}$

et donc $N^2(A) \geq \|Ax_1\|^2 = \|Ax\|^2$

donc $\|A\| \leq N(A)$.

La norme N dérive du double produit scalaire $((A,B)) = \sum_{p=1}^{\infty} (Ax_p, Bx_p)$

(double ne signifie rien en lui-même, ce n'est qu'une notation) on voit que $((A,A)) = N^2(A)$ montrons que si A et B $\in \mathcal{A}$ et si $\{x_p\}$ est un système orthonormé alors $((A,B))$ converge absolument et, de plus, $((A,B))$ ne dépend pas de $\{x_p\}$.

L'inégalité de Schwarz donne immédiatement le résultat :

$$|(Ax_p, Bx_p)| \leq \|Ax_p\| \cdot \|Bx_p\|$$

or $(\|Ax_p\| - \|Bx_p\|)^2 \geq 0$

donne $\|Ax_p\|^2 + \|Bx_p\|^2 \geq 2\|Ax_p\| \cdot \|Bx_p\|$

donc $|(Ax_p, Bx_p)| \leq \frac{1}{2} (\|Ax_p\|^2 + \|Bx_p\|^2)$

or A et B $\in \mathcal{A}$, on a $\sum_p \|Ax_p\|^2$ et $\sum_p \|Bx_p\|^2$ qui convergent.

D'autre part, soit $\{y_p\}$ un autre système orthonormé nous fixons p.

L'évaluation du produit scalaire en fonction des coefficients de Fourier dans le système $\{y_q\}$ donne :

$$(Ax_p, Bx_p) = \sum_{q=1}^{\infty} (Ax_p, y_q) \cdot (Bx_p, y_q)$$

On peut donc écrire (les séries étant convergentes on ne se préoccupera pas de l'ordre de sommation

$$\begin{aligned}
((A,B)) &= \sum_{p,q}^{\infty} (Ax_p, y_q) \cdot (Bx_p, y_q) \\
&= \sum_{p,q}^{\infty} (x_p, A^* y_q) \cdot (x_p, B^* y_q) \\
&= \sum_{q=1}^{\infty} (A^* y_q, B^* y_q)
\end{aligned}$$

Cette expression ne dépend pas de $\{x_p\}$. (voir début de II).

$((,))$ est un produit scalaire. (La linéarité est évidente, la commutativité aussi) \mathcal{A} est donc un espace pré-hilbertien.

Montrons que \mathcal{A} est complet.

Soit donc $\{A_n\}$ une suite de Cauchy d'endomorphisme de \mathcal{A} .

$\forall \epsilon > 0 \rightarrow \exists \delta$ tel que m et $n \geq \delta$ entraîne $N(A_n - A_m) \leq \epsilon$

or $N(A) \geq ||A||$ entraîne $||A_n - A_m|| \leq \epsilon$ si m et $n \geq \delta$

Soit $\hat{x} \in H$.

On a $||A_n \hat{x} - A_m \hat{x}|| \leq ||A_n - A_m|| \cdot ||\hat{x}|| \leq \epsilon \cdot ||\hat{x}||$

donc $\{A_n \hat{x}\}$ est une suite de Cauchy dans H et H étant complet. On montre alors (cf yoshida page 69) que $A_n \hat{x}$ tend vers $T \hat{x}$ et ce quel que soit \hat{x} . T est un opérateur linéaire et borné de $\mathcal{L}(H, H)$.

Montrons que $T \in \mathcal{A}$

$N(A_n) \leq N(A_{m_0} - A_n) + N(A_{m_0})$ si n et $m_0 \geq \delta_0$

alors on a :

$$N(A_n) \leq \varepsilon + N(A_{m_0}) = K$$

$$N(A_n) \leq K \text{ quel que soit } n$$

Soit donc $\{x_p\}$ une suite orthonormale de H fixons P

$$\sum_{p=1}^P \|A_n x_p\|^2 < N^2(A_n) \leq K^2$$

nous voyons que la sommation de gauche est bornée par la même quantité quel que soit n et P

$$\text{donc } \sum_{p=1}^{\infty} \|T x_p\|^2 \leq K^2 \text{ donc } N(T) < +\infty \text{ et par conséquent } T \in \mathfrak{a}.$$

\mathfrak{a} est donc complet et donc hilbertien.

De plus, \mathfrak{a} est un idéal bilatère pour l'algèbre $\mathcal{L}(H, H)$.

En effet soit $\{x_p\}$ un système orthonormé $A \in \mathfrak{a}$ et $B \in \mathcal{L}(H, H)$ alors

$$\sum_{p=1}^P \|B \cdot A x_p\|^2 \leq \|B\|^2 \cdot \sum_{p=1}^P \|A x_p\|^2 \leq \|B\|^2 \cdot N^2(A)$$

Or $\mathcal{L}(H, H)$ est l'espace des endomorphismes continus donc bornés, par conséquent :

$N(B \cdot A) < +\infty$ donc $B \cdot A \in \mathfrak{a}$ d'autre part si $B \in \mathcal{L}(H, H)$, B^* aussi si $A^* \in \mathfrak{a}$ alors $B^* \cdot A^* \in \mathfrak{a}$ donc $(B^* \cdot A^*)^*$ aussi et finalement $A \cdot B \in \mathfrak{a}$.

Nous allons étudier un endomorphisme particulier car les résultats obtenus nous serviront pour la suite.

Soient u et v deux éléments fixés $\in]-[$

On appelle trans-vection et on note uv^T l'endomorphisme qui à $x \in H$ fait correspondre le vecteur $(v, x)u$

Montrons d'abord que $uv^T \in \mathfrak{a}$

En effet, soit encore $\{x_p\}$ un système orthonormé

$$uv^T(x_p) = (v, x_p) \cdot u \quad \text{et par conséquent}$$

$$\|uv^T(x_p)\|^2 = |(v, x_p)|^2 \cdot \|u\|^2$$

$$\begin{aligned} \text{donc} \quad N^2(uv^T) &= \sum_{p=1}^{\infty} \|uv^T(x_p)\|^2 = \|u\|^2 \cdot \sum |(v, x_p)|^2 \\ &= \|u\|^2 \cdot \|v\|^2 \end{aligned}$$

d'après l'égalité de Bessel-Parseval.

$$N(uv^T) = \|u\| \cdot \|v\|$$

Etudions maintenant le double produit scalaire $((A, uv^T))$

$$\begin{aligned} ((A, uv^T)) &= \sum_{p=1}^{\infty} (Ax_p, uv^T x_p) \\ &= \sum_{p=1}^{\infty} (Ax_p, (v, x_p)u) \\ &= \sum_{p=1}^{\infty} (v, x_p) (A^*u, x_p) \\ &= (v, A^*u) = (u, Av) \end{aligned}$$

$$((A, uv^T)) = (u, Av)$$

C H A P I T R E III

ETUDE THEORIQUE

Soit $A_0 \in \mathcal{L}(H,H)$ et $b_0 \in H$

Considérons l'équation $A_0 x = b_0$

Nous supposons avoir trouvé une solution numérique \tilde{x} de ce problème.

Suivant la démarche habituelle, nous sommes conduit à considérer des couples (A,b)

$A \in \mathcal{L}(H,H)$ et $b \in H$ tels que $A\tilde{x} = b$

Nous appellerons par définition :

distance d'un problème (A_0, b_0) à un autre (A,b) la quantité :

$$\phi(A,b) = N^2(A-A_0) + \|b - b_0\|^2$$

Pour que cette mesure ait un sens nous nous restreindrons par convention à étudier l'erreur à postériori uniquement pour les A tels que $A - A_0 \in \mathcal{A}$.

Suivant l'habitude, nous conviendrons que \tilde{x} est "compatible" si

$$\phi(A,b) \leq \eta$$

η étant un seuil fixé d'avance.

Nous cherchons donc le minimum de $\phi(A,b)$ avec la contrainte suivante :

$$A\tilde{x} = b.$$

D,III,2

Nous utiliserons encore les multiplicateurs de Lagrange mais sous une forme adaptée.

Pour trouver le minimum de $\phi(A,b)$ sachant que $A\tilde{x} = b$ nous étudierons :

$$\phi(A,b) = N^2(A-A_0) + \|b - b_0\|^2 + 2(z, A\tilde{x} - b)$$

$z \in H$ étant une espèce de multiplicateur de Lagrange.

Faisons varier A et b de δA et δb

$$\begin{aligned} \phi(A+\delta A, b+\delta b) &= ((A-A_0+\delta A, A-A_0+\delta A)) \\ &+ (b-b_0+\delta b, b-b_0+\delta b) + 2(z, A\tilde{x}-b+\delta A\tilde{x}-\delta b) \\ &= N^2(A-A_0) + N^2(\delta A) + 2((\delta A, A-A_0)) + \|b-b_0\|^2 \\ &+ \|\delta b\|^2 + 2(\delta b, b-b_0) + 2(z, A\tilde{x}-b) + 2(z, \delta A\tilde{x}-\delta b) \\ &= \phi(A,b) + N^2(\delta A) + \|\delta b\|^2 + 2(z, \delta A\tilde{x}) + 2(\delta b, b-b_0-z) \end{aligned}$$

Or nous avons vu que $(z, \delta A\tilde{x}) = ((\delta A, z\tilde{x}^T))$

$$\begin{aligned} \text{donc } \phi(A+\delta A, b+\delta b) &= \phi(A,b) + N^2(\delta A) + \|\delta b\|^2 \\ &+ 2((\delta A, A-A_0+z\tilde{x}^T)) \\ &+ 2(\delta b, b-b_0-z) \end{aligned}$$

nous savons que nous aurons un extrémum de ϕ lorsque la variation de ϕ en fonction de $\pm\delta A$ et $\pm\delta b$ sera minimum. Dans notre cas cela revient à annuler

$$((\delta A, A-A_0+z\tilde{x}^T))$$

$$\text{et } (\delta b, b-b_0-z)$$

Nous ne disposons que de A, b et z

Nous aurons donc les trois équations à trois inconnues

$$\left\{ \begin{array}{l} A - A_0 + z\tilde{x}^T = 0 \quad (1) \\ b - b_0 - z = 0 \quad (2) \\ A\tilde{x} - b = 0 \quad (3) \end{array} \right.$$

Réolvons par rapport à z

$$A^* = A_0 - z\tilde{x}^T$$

$$b^* = b_0 + z$$

Reportons dans (3)

$$(A_0 - z\tilde{x}^T)\tilde{x} - b_0 - z = 0$$

posons $A_0\tilde{x} - b_0 = \eta$ alors $z = \frac{\eta}{1 + \|\tilde{x}\|^2}$

$$A^* - A_0 = \frac{\eta}{1 + \|\tilde{x}\|^2} \tilde{x} \quad b^* - b_0 = \frac{\eta}{1 + \|\tilde{x}\|^2}$$

On en déduit

$$\phi(A^*, b^*) = \frac{\|\eta\|^2}{1 + \|\tilde{x}\|^2}$$

C H A P I T R E IV

APPLICATIONS AUX METHODES DE GALERKIN

Soit une forme linéaire $a(u,v)$ qui applique $H \times H$ sur \mathbb{R} telle que :

$$\begin{cases} |a(u,v)| \leq \gamma \cdot \|u\| \cdot \|v\| & \gamma > 0 & \textcircled{1} \\ a(u,u) \geq \delta \|u\|^2 & \delta > 0 & \textcircled{2} \end{cases}$$

Soit d'autre part $L(v)$ une forme linéaire qui applique H sur \mathbb{R} . L est continue.

Nous voulons résoudre le problème : trouver u tel que $a(u,v) = L(v)$ quel que soit $v \in H$.

Remarquons d'abore que les hypothèses $\textcircled{1}$ et $\textcircled{2}$ associées au lemme de LAX-MILGRAM (Yoshida page 92) nous indiquent qu'il existe un isomorphisme de H sur lui-même. (Endomorphisme bi-univoque et bi-continu) tel que :

$$a(u,v) = (Su,v) \quad \forall u,v \in H$$

De même, la continuité de L entraîne l'existence de $\varphi \in H$ tel que : $L(v) = (\varphi,v) \quad \forall v$ le problème admet donc la forme variationnelle trouver u tel que $(Su,v) = (\varphi,v)$ ou encore $Su = \varphi$.

Appliquons les résultats trouvés dans III à l'équation $S_0 u = \varphi_0$.

Nous avons la solution numérique \tilde{u} . Elle est solution exacte de

$$S^* u = \varphi^*$$

$$s^* = s_0 \cdot \frac{1}{1 + \|\tilde{u}\|^2} (s_0 \tilde{u} - \varphi_0) \cdot \tilde{u}$$

$$\varphi^* = \varphi_0 + \frac{(s_0 \tilde{u} - \varphi_0)}{1 + \|\tilde{u}\|^2}$$

$$\text{et } \phi(s^*, \varphi^*) = \frac{\|s_0 \tilde{u} - \varphi_0\|^2}{1 + \|\tilde{u}\|^2}$$

si on appelle n la quantité $s_0 \tilde{u} - \varphi_0$

alors

$$\phi(s^*, \varphi^*) = \frac{\|n\|^2}{1 + \|\tilde{u}\|^2}$$

Supposons que pour résoudre le problème nous ayons utilisé la méthode de Galerkin. Pour cela, on dispose d'un système orthonormé $\{\varphi_1, \varphi_2, \dots\}$ d'éléments de H .

Nous cherchons une solution de la forme

$$u_n = \sum_{i=1}^n a_i \varphi_i$$

telle que :

$$a(u_n, \varphi_i) = L(\varphi_i) \quad (i=1, 2, \dots, n)$$

Nous ne connaissons pas n mais nous savons que

$$\begin{aligned} (n, \varphi_j) &= a(u_n, \varphi_j) - (\varphi, \varphi_j) \\ &= a(u_n, \varphi_j) - L(\varphi_j) \end{aligned}$$

D, IV, 3

On voit donc que $n = \sum_{i=n+1}^{\infty} [a(u_n, \varphi_i) - L(\varphi_i)] \varphi_i$

et par conséquent :

$$\phi(S^*, *) = \frac{\sum_{i=n+1}^{\infty} [a(u_n, \varphi_i) - L(\varphi_i)]^2}{1 + \|u_n\|^2}$$

On ne peut, bien sûr, avoir seulement qu'une estimation grossière de cette erreur car nous ne pouvons calculer une série infinie

ESTIMATION A POSTERIORI DE L'ERREUR
DANS LE CAS DE L'EQUATION DE FREDHOLM

Soit à résoudre l'équation de Fredholm de seconde espèce : trouver $u \in H$ tel que :

$$Au(p) = u(p) - \int_{\Omega} K(p, Q) u(Q) d\Omega_Q = f(p)$$

sous les hypothèses classiques :

1° K est tel que $\int_{\Omega_p} \int_{\Omega_Q} K^2(p, Q) d\Omega_p d\Omega_Q < +\infty$

2° $\int_{\Omega} f^2 d\Omega$ existe

3° il existe au plus une solution $\|u\| < \infty$

la formulation variationnelle de ce problème est classique : il s'agit d'égaliser les coefficients de Fourier de Au et de f .

c'est-à-dire qu'il faut que :

$$(Au, v) = (f, v) \text{ quel que soit } v \in H$$

on notera aussi (Au, v) par $a(u, v)$

$$a(u, v) = (Au, v) = \int_{\Omega_p} \left[u - \int_{\Omega_Q} K(p, Q) u(Q) d\Omega_Q \right] v(p) d\Omega_p$$

Montrons que cet opérateur a est borné (ou continu) c'est-à-dire que

$$|a(u, v)| \leq \gamma \cdot \|u\| \cdot \|v\| \quad \gamma > 0$$

$$|a(u, v)| \leq \left| \int_{\Omega_p} u \cdot v \cdot d\Omega_p \right| + \left| \int_{\Omega_p} v(p) d\Omega_p \cdot \int_{\Omega_Q} K(p, Q) u(Q) d\Omega_Q \right|$$

or $\left| \int_{\Omega_p} u \cdot v \cdot d\Omega_p \right| \leq \|u\| \cdot \|v\|$ inégalité de Schwartz.

Etudions l'autre partie :

$$\begin{aligned} \left| \int_{\Omega_p} v(p) d\Omega_p \int_{\Omega_Q} K(p, Q) u(Q) d\Omega_Q \right| &\leq \left| \int_{\Omega_p} |v| d\Omega_p \right| \left| \int_{\Omega_Q} K u d\Omega_Q \right| \\ &\leq \left| \int_{\Omega_p} \|u\| \cdot v(p) \left[\int_{\Omega_Q} K^2(p, Q) d\Omega_Q \right]^{\frac{1}{2}} d\Omega_p \right| \end{aligned}$$

à cause de l'inégalité de Schwartz

$$\leq \|u\| \cdot \int_{\Omega_p} v(p) \left(\int_{\Omega_Q} K^2(p, Q) d\Omega_Q \right)^{\frac{1}{2}} d\Omega_p$$

d'après la même inégalité.

$$\leq ||u|| \cdot ||v|| \cdot \int_{\Omega_p} \int_{\Omega_Q} K^2(p,Q) \leq \alpha ||u|| \cdot ||v||$$

donc $|a(u,v)| \leq (1+\alpha) ||u|| \cdot ||v||$

ce qui montre bien que a est bornée et continue.

Cette démonstration est aussi valable pour l'équation de Fredholm de première espèce :

trouver u tel que $\int_{\Omega} K(p,Q) u(Q) d\Omega_Q = f(p)$

Montrons que l'opérateur A est coercitif c'est-à-dire que (Au,u) ou $a(u,u) > \delta \cdot ||u||^2$ d'après la première démonstration

$$a(u,u) = ||u||^2 - \int_{\Omega_p} u(p) d\Omega_p \int_{\Omega_Q} K(p,Q) u(Q) d\Omega_Q$$

$$> ||u||^2 - \alpha ||u||^2$$

il faut que $a(u,u) > \delta ||u||^2 \quad \delta > 0$

donc que $1 - \alpha > \delta$

$\alpha < 1 - \delta$;

c'est-à-dire que : $\int_{\Omega_p} \int_{\Omega_Q} K^2(p,Q) d\Omega_p d\Omega_Q < 1 - \delta$

alors le formalisme de IV est applicable.

Soit une base orthonormée $\{\varphi_1, \varphi_2, \dots, \varphi_n, \dots\}$

$$\begin{aligned}
 a(\varphi_j, \varphi_i) &= (A\varphi_j, \varphi_i) = \int_{\Omega} (A\varphi_j(p)) \varphi_i(p) d\Omega_p \\
 &= \int_{\Omega} \left[\varphi_j(p) - \int_{\Omega} K(p, Q) \varphi_j(Q) d\Omega_Q \right] \varphi_i(p) d\Omega_p \\
 &= \int_{\Omega} \varphi_j(p) \cdot \varphi_i(p) d\Omega_p - \int_{\Omega} \varphi_i(p) d\Omega_p \int_{\Omega} K(p, Q) \varphi_j(Q) d\Omega_Q \\
 &= \delta_{ij} - \gamma_{ij}
 \end{aligned}$$

Nous cherchons une solution de la forme

$$u_n = \sum_{i=1}^n a_i \varphi_i$$

avec, comme hypothèse : $(Au_n, \varphi_i) = (f, \varphi_i)$

pour $i = 1, 2, \dots, n$.

$$(Au_n, \varphi_i) = a_i - \sum_{j=1}^n \gamma_{ij} a_j = (f, \varphi_i)$$

alors
$$r = \sum_{i=n+1}^{\infty} \left[(Au_n, \varphi_i) - (f, \varphi_i) \right] \varphi_i$$

et
$$\Phi(S_1, \varphi_i) = \frac{\sum_{i=1}^{n+1} \left[(Au_n, \varphi_i) - (f, \varphi_i) \right]^2}{1 + \sum a_i^2}$$

CAS DU PROBLEME DE DIRICHLET

Soient $a_{ij}(x)$ et $a_0(x)$ des fonctions mesurables bornées de Ω (ouvert borné de \mathbb{R}^n) et $f(x) \in L_2(\Omega)$. En prenant les dérivées faibles (dérivées au sens des distributions) le problème de Dirichlet s'énonce ainsi :

trouver $u \in H_0^1(\Omega)$ tel que :

$$-\sum_{i,j=1}^n \frac{\partial}{\partial x_i} \left(a_{ij}(x) \frac{\partial u}{\partial x_j} \right) + a_0(x) \cdot u = f(x)$$

Ce problème admet la forme variationnelle suivante :

$$a(u,v) = \sum_{i,j} \int_{\Omega} a_{ij}(x) \frac{\partial u}{\partial x_i} \cdot \frac{\partial v}{\partial x_j} dx + \int_{\Omega} a_0(x) u \cdot v dx = \int_{\Omega} f \cdot v dx$$

il faut trouver u de façon qu'on ait bien l'égalité quel que soit $v \in H_0^1(\Omega)$

le produit scalaire dans H_0^1 sera défini par :

$$(u,v)_{H_0^1} = \int_{\Omega} \left(u \cdot v + \sum_{i,j=1}^n \frac{\partial u}{\partial x_i} \cdot \frac{\partial v}{\partial x_j} \right) dx$$

De plus nous avons les hypothèses suivantes :

$$\textcircled{a} \quad \sum_{i,j} a_{ij}(x) \xi_i \xi_j \geq \alpha (|\xi_1|^2 + \dots + |\xi_n|^2)$$

$$\alpha > 0$$

$$\textcircled{b} \quad a_0(x) \geq \alpha_0 > 0$$

Montrons que nous pouvons appliquer la méthode de Galerkin.

Pour cela, montrons d'abord que $a(u,u)$ est coercitive c'est-à-dire que :

$$a(u,u) \geq \delta ||u||^2$$

$$a(u,u) = \int_{\Sigma} \sum_{i,j} a_{ij} \left(\frac{\partial u}{\partial x_i} \right) \left(\frac{\partial u}{\partial x_j} \right) + a_0 u^2$$

d'après (a) on a :

$$\begin{aligned} a(u,u) &\geq \int \alpha \sum \left(\frac{\partial u}{\partial x_i} \right)^2 + \alpha_0 u^2 \\ &\geq \min(\alpha, \alpha_0) \int \sum \left(\frac{\partial u}{\partial x_i} \right)^2 + u^2 \\ &= \min(\alpha, \alpha_0) ||u||^2 \end{aligned}$$

Montrons maintenant que $a(u,v)$ est bornée ou continue, c'est-à-dire qu'il existe γ tel que :

$$|a(u,v)| \leq \gamma \cdot ||u|| \cdot ||v|| \quad \gamma > 0$$

$$\begin{aligned} |a(u,v)| &= \left| \int_{\Omega} \sum_{i,j} a_{ij} \left(\frac{\partial u}{\partial x_i} \right) \left(\frac{\partial v}{\partial x_j} \right) d\Omega + \int_{\Omega} a_0 u \cdot v d\Omega \right| \\ &\leq \int \sum |a_{ij}| \cdot \left| \frac{\partial u}{\partial x_i} \right| \cdot \left| \frac{\partial v}{\partial x_j} \right| d\Omega + \int |a_0| \cdot |u| \cdot |v| d\Omega \end{aligned}$$

prenons $M = \text{Max} (a_{ij}(x) \text{ et } a_0(x))$ quels que soient i, j et x .

$$\text{alors } |a(u,v)| \leq M \int \left(\sum_{i,j} \left| \frac{\partial u}{\partial x_i} \right| \cdot \left| \frac{\partial v}{\partial x_j} \right| + |u| \cdot |v| \right) d\Omega$$

$$\leq M \int \left(\left(\sum_i \left| \frac{\partial u}{\partial x_i} \right| \right) \left(\sum_j \left| \frac{\partial v}{\partial x_j} \right| \right) + |u| \cdot |v| \right) d\Omega$$

$$\leq M \int \left(\sum_{i=1}^n \left| \frac{\partial u}{\partial x_i} \right| + |u| \right) \left(\sum_j \left| \frac{\partial v}{\partial x_j} \right| + |v| \right) d\Omega$$

Or
$$\sum_{i=1}^n \frac{|x_i|}{n} \leq \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

donc
$$|a(u, v)| \leq M \int (n+1)^2 \left(\sum \left(\frac{\partial u}{\partial x_i} \right)^2 + u^2 \right)^{\frac{1}{2}} \left(\sum \left(\frac{\partial v}{\partial x_j} \right)^2 + v^2 \right)^{\frac{1}{2}} d\Omega$$

et d'après l'inégalité de Schwartz

$$\begin{aligned} &\leq M(n+1)^2 \left[\int \sum \left(\frac{\partial u}{\partial x_i} \right)^2 + u^2 \right]^{\frac{1}{2}} \left[\int \sum \left(\frac{\partial v}{\partial x_j} \right)^2 + v^2 \right]^{\frac{1}{2}} \\ &\leq M(n+1)^2 \|u\| \cdot \|v\| \end{aligned}$$

Posons $M(n^2+1) = \gamma$ alors

$$|a(u, v)| \leq \gamma \|u\| \cdot \|u\|$$

c'est le résultat cherché.

Soit donc une base $\{\varphi_1, \dots, \varphi_n\}$ de $H_0^1(\Omega)$.

Nous chercherons une solution de la forme

$$u_n = \sum_{i=1}^n a_i \varphi_i$$

où les a_i sont tels que :

$$a(u_n, \varphi_i) = \int f, \varphi_i d\Omega$$

cherchons les composantes de $\eta = Su_n - \varphi$

Ce sont :

$$\begin{aligned} (Su_n - \varphi, \varphi_j) &= a(u_n, \varphi_j) - (\varphi, \varphi_j) \\ &= a(u_n, \varphi_j) - L(\varphi_j) \end{aligned}$$

donc :

$$\eta = \sum_{i=n+1}^{\infty} (a(u_n, \varphi_i) - L(\varphi_i)) \varphi_i$$

et

$$\phi = \frac{\sum_{i=n+1}^{\infty} [a(u_n, \varphi_i) - L(\varphi_i)]^2}{1 + \sum a_i^2}$$

EXEMPLE NUMERIQUE

Soit à résoudre l'équation de Fredholm de seconde espèce

$$u(x) - \int_{-1}^{+1} \frac{3}{16} (x^2+y^2) u(y) dy = x^2 \quad (1)$$

$$\int_{-1}^{+1} dx \int_{-1}^{+1} \frac{9}{256} (x^2+y^2)^2 dy = \frac{1008}{11520}$$

Nous voyons que cette quantité est plus petite que 1 donc on aura bien continuité et coercivité.

Nous pouvons donc utiliser la méthode de Galerkin. Nous utiliserons comme système orthonormé l'ensemble des polynômes de Legendre : p_i ($i=1,2,\dots$). Ceux-ci sont bien orthogonaux sur $[-1,+1]$, ils seront orthonormés si nous les multiplions par un facteur convenable. C'est-à-dire pour le polynôme p_n par la quantité $\sqrt{\frac{2}{2n+1}}$.

Nous cherchons une solution numérique de la forme $u_n = \sum_{i=0}^n a_i p_i$.

Pour cela nous égalons les n premiers coefficients de Fourier des 2 membres de (1) où nous remplacerons u par u_n .

Cela donne :

$$\int_{-1}^{+1} \left[\sum_{i=0}^n q_i p_i(x) - \frac{3}{16} \int_{-1}^{+1} \sum_{i=0}^n a_i p_i(x) (x^2+y^2) dy \right] p_k(x) dx = \int_{-1}^{+1} x^2 p_k(x) dx$$

pour k variant de 0 à n

nous savons que $p_i(x) p_k(x) = 0$ si $i \neq k$

donc :

$$\begin{aligned} \int_{-1}^{+1} x^2 p_k(x) dx &= a_k - \frac{3}{16} \int_{-1}^{+1} \sum a_i p_i(x) p_k(x) dx \int_{-1}^{+1} y^2 dy \\ &\quad - \frac{3}{16} \int_{-1}^{+1} \sum a_i p_i(x) p_k x^2 dx \cdot \int_{-1}^{+1} dy \\ &= a_k - \frac{1}{8} a_k - \frac{3}{8} \int_{-1}^{+1} \sum a_i p_i(x) p_k(x) x^2 dx \end{aligned}$$

étudions la dernière intégrale.

Nous savons qu'un polynôme de Legendre est orthogonal à tout polynôme de degré inférieur sur $[-1, +1]$. On en déduit que les seuls termes qui ne s'annuleront pas sont ceux où la différence entre i et k sera ≤ 2 . Calculons toutes les intégrales $\int p_i(x) p_k(x) x^2 dx$ qui vérifient cette condition. (Nous nous limiterons à 4).

$$p_0 = \sqrt{2}$$

$$p_1 = \sqrt{\frac{2}{3}} \cdot x$$

$$p_2 = \sqrt{\frac{2}{5}} \left(\frac{3}{2} x^2 - \frac{1}{2} \right)$$

$$p_3 = \sqrt{\frac{2}{7}} \left(\frac{5}{2} x^3 - \frac{3}{2} x \right)$$

$$p_4 = \sqrt{\frac{2}{9}} \left(\frac{35}{8} x^4 - \frac{15}{4} x^2 + \frac{3}{8} \right)$$

D,V,3

$$\int p_0 p_0 x^2 = 1,3333$$

$$\int p_0 p_1 x^2 = 0 \text{ nous remarquons que lorsque } i + k \text{ est impair}$$

$$\int p_i p_k x^2 \text{ est nul sur } [-1, +1]$$

$$\int p_0 p_2 x^2 = 0,2385$$

$$\int p_0 p_3 x^2 = \int p_0 p_4 x^2 = 0$$

$$\int p_1 p_1 x^2 = 0,2666$$

$$\int p_1 p_2 x^2 = 0$$

$$\int p_1 p_3 x^2 = 0,0498$$

$$\int p_1 p_4 x^2 = 0$$

$$\int p_2 p_2 x^2 = 0,0838$$

$$\int p_2 p_3 x^2 = 0$$

$$\int p_2 p_4 x^2 = 0,2271$$

$$\int p_3 p_4 x^2 = 0$$

Nous cherons une solution de la forme $\sum_{i=0}^2 a_i p_i$

On a donc le système :

$$\frac{7}{8} a_0 + 1,3333 a_0 + 0,2385 a_2 = \int_{-1}^{+1} x^2 p_0(x) dx$$

D,V,4

$$\frac{7}{8} a_1 + 0,2666 a_1 = \int_{-1}^{+1} x^2 p_1(x) dx$$

$$\frac{7}{8} a_2 + 0,2385 a_0 + 0,0838 a_2 = \int_{-1}^{+1} x^2 p_2(x) dx$$

$$\int_{-1}^{+1} x^2 p_0(x) dx = 0,9428$$

$$\int_{-1}^{+1} x^2 p_1(x) dx = 0$$

$$\int_{-1}^{+1} x^2 p_2(x) dx = 0,1696$$

$$2,2083 a_0 + 0,2385 a_2 = 0,9428$$

$$1,1416 a_1 = 0$$

$$0,9588 a_2 + 0,2385 a_0 = 0,1696$$

On en déduit $a_1 = 0$

$$a_0 = 0,4190$$

$$a_2 = 3,8792$$

D,V,5

$$\text{d'où } u_2 = 3,8792 \cdot \sqrt{\frac{2}{7}} \left(\frac{3}{2} x^2 - \frac{1}{2} \right) + 0,4190 \times \sqrt{2}$$

$$u_2(x) = 3,0879 x^2 - 0,4368$$

Entamons maintenant le calcul à postériori de l'erreur.

Calculons $a(u_2, P_3) - L(P_3)$, il vaut :

$$\sqrt{\frac{2}{7}} \cdot \int_{-1}^{+1} \left[3,0879 x^2 - 0,4368 - \frac{3}{16} \int_{-1}^{+1} (x^2 + y^2) (3,0879 x^2 - 0,4368) dy \right] \cdot$$

$$\left(\frac{5}{3} x^3 - \frac{3}{2} x \right) dx$$

Le terme dans le crochet est de degré pair, l'intégrale en x portera donc sur des termes de degré impair elle sera donc nulle.

Passons au terme suivant :

$$\sqrt{\frac{2}{9}} \int_{-1}^{+1} \left[3,0879 x^2 - 0,4368 - \frac{3}{16} \int_{-1}^{+1} (x^2 + y^2) (3,0879 x^2 - 0,4368) dy \right]$$

$$\cdot \frac{35 x^4 - 30 x^2 + 3}{8} dx =$$

D,V,6

$$\sqrt{\frac{2}{9}} \int_{-1}^{+1} \left[3,0879 x^2 - 0,4368 - \frac{3}{16} (3,0879 x^2 - 0,4368) \left(\frac{5}{3}\right) \right] \times \frac{35 x^4 - 30 x^2 + 3}{8} dx$$

$$\times \sqrt{\frac{2}{9}} \int_{-1}^{+1} (2,1230 x^2 - 0,3004) \frac{35 x^4 - 30 x^2 + 3}{8} dx$$

$$\frac{\sqrt{2}}{8\sqrt{9}} \int_{-1}^{+1} (74,3050 x^6 - 169,1 x^4 + 15,408 x^2 - 0,9012) dx$$

$$\frac{\sqrt{2}}{24} \times 37,04 = 2,1855$$

Nous savons que l'erreur vaut :

$$\frac{\sum_{i=3}^{\infty} (a(u_2, p_i) - L(p_i))^2}{1 + \sum a_i^2}$$

Nous avons calculé ici :

$$\frac{\sum_{i=3}^4 (a(u_2, p_i) - L(p_i))^2}{1 + (3,08)^2 + (0,43)^2}$$

qui sera une minoration de l'erreur, et on trouve d'ailleurs pour cette erreur tronquée la valeur 0,20.

Ce qui prouve qu'il était inutile de poursuivre les calculs au delà de la deuxième décimale dans le cas où l'on s'arrête à $n = 2$.

BIBLIOGRAPHIE

PARTIE - D

Pour l'étude des espaces de Hilbert, on pourra se référer à :

J. DIEUDONNE - FOUNDATION OF MODERN ANALYSIS
Academic Press 1960

L'utilisation de ces espaces en analyse numérique est traitée dans :

P.J. LAURENT - THEORIE DE L'APPROXIMATION
Cours de l'Université de Grenoble - 1964

On pourra se référer pour des précisions sur les opérateurs à :

K. YOSIDA - FUNCTIONAL ANALYSIS
Springer Verlag - Berlin 1965

N. DUNFORD & J.T. SCWHARTZ - LINEAR OPERATOR
Intersciences Publishers Inc. 1958

On aura de même tous les renseignements sur la méthode de Galerkin, l'équation de Fredholm, le problème de Dirichlet dans les livres suivants :

S.G. MIKHLIN - INTEGRAL EQUATIONS
Pergamon Press. New York 1957

J.L. LIONS - METHODES D'APPROXIMATION NUMERIQUE DES PROBLEMES AUX LIMITES DE
LA PHYSIQUE MATHEMATIQUE
Laboratoire de Calcul Numérique de l'Institut Blaise Pascal.

Ce dernier livre donne d'ailleurs beaucoup de renseignements sur les dérivées faibles.

VU

Grenoble, le

Le Président de la Thèse

VU

Grenoble, le

Le Doyen de la Faculté des Sciences

VU, et permis d'imprimer,

Le Recteur de l'Académie de GRENOBLE