



HAL
open science

Contribution à l'étude de la méthode de Lanczos

Marc Lévy

► **To cite this version:**

Marc Lévy. Contribution à l'étude de la méthode de Lanczos. Modélisation et simulation. Université Joseph-Fourier - Grenoble I, 1961. Français. NNT: . tel-00277842

HAL Id: tel-00277842

<https://theses.hal.science/tel-00277842>

Submitted on 7 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.

T H E S E

--

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

pour obtenir

Le titre de Docteur de spécialité
Mathématiques Appliquées

par

Marc LEVY

Ingénieur IEG et IMAG

CONTRIBUTION A L'ETUDE
DE LA METHODE DE LANCZOS

Thèse soutenue le 26 Juin 1961

Devant la commission d'Examen :

MM. KUNTZMANN : Président

GASTINEL

HACQUES

T H E S E

--:--

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

pour obtenir

Le titre de Docteur de spécialité
Mathématiques Appliquées

par

Marc LEVY

Ingénieur IEG et IMAG

CONTRIBUTION A L'ETUDE
DE LA METHODE DE LANCZOS

Thèse soutenue le 26 Juin 1961

Devant la commission d'Examen :

MM. KUNTZMANN : Président

GASTINEL

HACQUES

FACULTE des SCIENCES de L'UNIVERSITE de GRENOBLE

-:~::~:~::~:~::~:~::~:~::~:-

Professeurs titulaires classés d'après leur ancienneté comme titulaires ou professeurs à titre personnel :

Doyens honoraires :

MM. FURETAT P.	1931	Géologie et minéralogie
ROBERT L.	1931	Géologie et Minéralogie

Doyen :

M. WEILL L.	1950	Thermodynamique
MM. DORVILLE A.	1941	Zoologie
HEILMANN R.	1943	Chimie organique
NEEL L.	1945	Physique expérimentale
KRAVTCHEVSKO J.	1946	Mécanique rationnelle
PARDE M.	1946	Botanique
BENOIT J.	1948	Radioélectricité
CHENE L.	1949	Chimie papetière
NOBECOURT P.	1949	Micrographie papetière
FELICI N.	1951	Electrostatique
KUNTSMANN J.	1952	Mathématiques Appliquées
BARBIER	I.10.1953	Géologie appliquée
SANTON L.	I.10.1953	Mécanique des fluides
CHABAUTY C.	I.10.1954	Calcul différentiel et Intégral
OZENEA P.	I.10.1954	Botanique
FALLOT M.	1955	Physique industrielle
GALVANI O.	1957	Mathématiques
TRAYNARD Ph.	I.10.1958	Chimie générale
SOUTIF M.	I.3.1958	Physique générale

MM. CFAYA A.	I.II.1958	Hydrodynamique
REEB G.	I. I.1959	Statistiques Mathématiques
REULOS B.	I.10.1959	Physique générale
WOFLERS F.	I.10.1959	Physique
BESSON J.	I.10.1959	Chimie
MOUSSA A.	I.12.1959	Chimie Nucléaire et Radioactivité
AYANT Y.	I. I.1960	Physique approfondie
M ^{le} LUTZ E.	I.I. 1960	Mathématiques générales
MM. GALLISSOT F.	I.I.1960	Mathématiques pures
BLAMBERT M.	I.10.1960	Mathématiques
BOUCHEZ R.	I. I.1961	Physique Nucléaire
LLI BOUTRY	I. I.1961	Physique
MICHEL	I. I.1961	Géologie et Minéralogie

Professeurs sans chaire

MM. SILBERT R.	1954	Mécanique des Fluides
MOUSSIEGT J.	I. I.1959	Electronique
DESSAUX G.	I.10.1959	Physiologie animale
PILLET E.	I.10.1960	Electrotechnique

Maîtres de conférence

MM. BONNIER E.	I.10.1956	Chimie
BARBIER J.CL.	I.10.1957	Physique
PAUTHENET R.	2. I.1958	Electrotechnique
BUYLE-BODIN M.	I.10.1958	Electronique
M ^{le} NAIM Linda	I.10.1958	Mathématiques
MM. PERRET R.	I.10.1958	Servomécanismes
BREYFUS	I.10.1958	Thermodynamique
VAILLANT F.	I.II.1958	ZOOLOGIE et Hydrobiologie
M ^{me} SOUTIF J.	I.10.1959	Physique
DEBELMAS J.	I.10.1959	Géologie et Minéralogie

Mme KOPLER	I. 10. 1959	Botanique
COHEN J.	I. 10. 1959	Physique
Mme BARBIER	I. 10. 1959	Electrochimie
BEISSONNEAU P.	I. 10. 1959	Physique
ARNAUD P.	I. 10. 1959	Chimie
VAUQUOIS B.	I. 10. 1959	Mathématiques
DEPASSEL P.	I. 11. 1959	Mécanique des Fluides
ROBERT A.	I. 11. 1960	Chimie papetière

Délégués

MM. GERBER	I. 10. 1960	Mathématiques
GIDON P.	I. 10. 1960	Géologie et Minéralogie
DUCROS P.	I. 10. 1960	Minéralogie et Cristallographie
HACQUES J.	I. 10. 1960	Calcul Numérique
PEBAY - PEROLA	I. 10. 1960	Physique
ANGLES D'AURIAC	I. 10. 1960	Mécanique des Fluides
COUMES A.	I. 10. 1960	Electronique
LANCIA R.	I. 10. 1960	Physique automatique
BODU J.	I. 10. 1960	Mécanique des Fluides
BIAREZ	I. 10. 1960	Mécanique physique
GASTINEL } LACAZE }		Chargés d'Enseignement Mathématiques Thermodynamique

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

-:-:-:-:-:-:-:-:-:-

L'idée du travail ci-après m'a été suggérée au cours d'un stage effectué à la compagnie des Machines Bull que je tiens à remercier de son obligeance.

C'est ensuite au laboratoire de Calcul de l'Université de Grenoble que l'essentiel de ces études a été accompli sous la direction de Monsieur le Professeur Kuntzmann, dont les conseils et les directives précieuses m'ont guidé continuellement ; qu'il veuille bien trouver ici l'expression de ma profonde reconnaissance.

Je dois également beaucoup à Monsieur Castinal qui tant en ce qui concerne les études des erreurs que le calcul matriciel m'a apporté un appui constant et précieux.

Enfin je remercie tout le personnel du laboratoire pour son obligeance et son amabilité, et tout particulièrement Monsieur Bolliet.

J'ai voulu, dans ce travail, étudier les erreurs dans la méthode de Lanczos pour le calcul du polynôme caractéristique d'une matrice. Le sujet est très vaste et je suis loin de l'avoir épuisé. Mon travail consiste essentiellement en un dégrossissage et un recensement des principales causes d'erreurs.

Cette étude comporte deux aspects absolument différents. Tout d'abord une étude théorique des erreurs et de leurs propagations. J'ai en particulier, au cours de cette partie, été amené à m'intéresser, aux erreurs sur un produit scalaire en fonction du nombre de dimensions de l'espace. Le résultat obtenu paraît assez intéressant et je regrette beaucoup de n'avoir pu faire que des vérifications incomplètes de ce résultat.

Des vérifications complète. nécessitant une machine beaucoup plus rapide que celle dont je disposais.

Le second aspect, lui, est essentiellement expérimental ; il s'agissait de vérifier certaines hypothèses et de déterminer certaines règles mais de façon purement expérimentale.

Ce second aspect fait l'objet de la dernière partie de l'étude intitulée : "expériences numériques"

Je tiens à remercier Mademoiselle Stankievitch, c'est à ces soins qu'est due la qualité de la présentation matérielle de cette thèse.

Que Monsieur Hacques et Monsieur Gastinel veuillent trouver ici l'expression de ma reconnaissance pour avoir eu l'obligeance d'accepter de faire partie de la commission d'examen.

B I B L I O G R A P H I E

-:-:-:-:-:-:-:-:-:-

a) Etude théorique des valeurs propres

- (I) Wedderburn : Lectures on Matrices
- (II) O. Schreier et E. Sperner : Introduction to Modern Algebra and Matrix Theory
- (III) N. Gastinel : Analyse numérique linéaire (Université de Grenoble)

b) Méthode de Lanczos

- (IV) Lanczos C. An iteration method for the solution of the eigenvalue problem of linear differential and integral operators. *Journal of Research of the NBS* 1950, 45, p 255.
- (V) Rutishauser H. Beiträge zur Kenntnis des Biorthogonalisierungs-Algorithmus von Lanczos *ZAMP*, 1958, 4, p 35
- (VI) F. Genuys : La methode de Lanczos : Séminaire de l'AFCAL. 17 mars 1958
- (VII) Wilkinson : the evaluation of the zeros of ill conditioned polynomials (part II) - *Numerische Mathematik* 1959 p 150

c) Génération de points sur une hypersphère :

- (VIII) Mervin E. Muller. a note on a method for generating points uniformly on N-dimensional spheres. *Communications of the A.C.M.* April 59 p 19

Première. Partie

Remarques Théoriques sur la
Méthode de Lanczos

Chapitre I

Présentation Elementaire de la Méthode

1° - Principe : Cette méthode consiste à trouver une matrice tridiagonale semblable à la matrice donnée, puis à calculer le polynôme caractéristique de cette matrice tridiagonale. Ce dernier calcul est classique et simple. La partie essentielle de la méthode est donc la tridiagonalisation de la matrice.

Je vais d'abord exposer cette méthode dans le cas général où le polynôme caractéristique est identique au polynôme minimal et où le calcul se déroule correctement. Dans une étude ultérieure, j'étudierai les cas particuliers susceptibles de se produire.

Lanczos a montré [(4)] que sa méthode était valable quelle que soit la matrice considérée : singulière ou régulière, défective ou non.

2° - Tridiagonalisation : A partir de maintenant, nous considérons l'espace des vecteurs-colonnes (ou espace des vecteurs à droite de la matrice) désigné par E_{cn} , n étant la dimension de l'espace, ou E_c , si n est égal à la dimension de la matrice, et l'espace des vecteurs lignes (ou espace des vecteurs à gauche de la matrice), désigné par E_{ln} ou E_l .

Un vecteur ligne est orthogonal à un vecteur colonne si la somme des produits des composantes similaires est nul. Soit $b_1 \in E_c$ et $c_1 \in E_l$, deux vecteurs arbitraires.

Nous désignons par $\{a, b, c\}$ l'espace vectoriel engendré par a, b, c ; a, b, c, \dots appartenant à E_c (ou à E_l). Cherchons un vecteur $b_2 \in \{Ab_1; b_1\}$ et orthogonal à c_1 et un vecteur $c_2 \in \{c_1A; c_1\}$ orthogonal à b_1 .

Nous pouvons écrire :

$$b_2 = Ab_1 - \alpha_1 b_1$$

d'où nous tirons $\alpha_1 = \frac{c_1 A b_1}{c_1 b_1}$

et : $c_2 = c_1 A - \alpha_1' c_1$ d'où nous déduisons $\alpha_1 = \alpha_1'$

Nous posons $c_i b_i = \sigma_i$

Nous définirons ensuite :

$$b_3 \in \{Ab_2 ; b_2 ; b_1\} = \{b_1 ; Ab_1 ; A^2 b_1\}; \text{ orthogonal à } c_1 \text{ et } c_2$$

Et $c_3 \in \{c_2 A ; c_2 ; c_1\} = \{c_1 ; c_1 A ; c_1 A^2\}$; orthogonal à b_1 et b_2

soit : $b_3 = Ab_2 - \alpha_2 b_2 - \beta_1 b_1$

d'où $\alpha_2 = \frac{c_2 A b_2}{\sigma_2}$; $\beta_1 = \frac{c_1 A b_2}{\sigma_1}$

et $c_3 = c_2 A - \alpha_2' c_2 - \beta_1' c_1$

d'où nous tirons $\alpha_2 = \alpha_2'$; $\beta_1 = \beta_1'$

De façon générale, nous définirons :

$$b_{k+1} \in \{Ab_k ; b_k ; b_{k-1} ; \dots ; b_1\} \text{ orthogonal à } c_1 ; c_2 ; \dots ; c_k$$

$$c_{k+1} \in \{c_k A ; c_k ; c_{k-1} ; \dots ; c_1\} \text{ orthogonal à } b_1 ; b_2 ; \dots ; b_k$$

Il résulte de ces définitions que :

$$c_i b_j = 0 \quad \text{si } i \neq j \quad (1)$$

$$\text{et que : } c_i A b_j = 0 \quad \text{si } |i - j| > 1 \quad (2)$$

De là il découle que :

$$b_{k+1} = Ab_k - \alpha_k b_k - \beta_{k-1} b_{k-1}$$

$$\text{et } c_{k+1} = c_k A - \alpha_k' c_k - \beta_{k-1}' c_{k-1}$$

avec $\alpha_k = \frac{c_k A b_k}{\sigma_k}$ et $\beta_{k-1} = \frac{c_{k-1} A b_k}{\sigma_{k-1}}$

Remarquons encore que :

$$c_k A b_{k-1} = c_{k-1} A b_k = c_k b_k = \sigma_k$$

$$\text{d'où } \beta_{k-1} = \frac{\sigma_k}{\sigma_{k-1}}$$

On constitue ainsi deux suites de vecteurs respectivement orthogonales. Il résulte de cette orthogonalité que les vecteurs b_{n+1} et c_{n+1} respectivement orthogonaux à n vecteurs indépendants sont nuls.

On a alors :

$$A b_1 = b_2 + \alpha_1 b_1$$

$$A b_2 = b_3 + \alpha_2 b_2 + \beta_1 b_1$$

$$A b_k = b_{k+1} + \alpha_k b_k + \beta_{k-1} b_{k-1}$$

$$A b_n = \alpha_n b_n + \beta_{n-1} b_{n-1}$$

Ceci montre que dans la base $b_1 ; b_2 ; \dots ; b_n$, la matrice A est semblable à la matrice.

$T =$

$$\begin{vmatrix} \alpha_1 & \beta_1 & & & & \\ 1 & \alpha_2 & \beta_2 & & & \\ & 1 & \alpha_3 & \beta_3 & & \\ & & & & \ddots & \\ & & & & & \alpha_n & \beta_{n-1} \\ & & & & & 1 & \alpha_n \end{vmatrix}$$

Obtenir le polynôme caractéristique d'une telle matrice est simple.

Posons :

$$P_1(\lambda) = I$$

$$P_2(\lambda) = \lambda - \alpha_1$$

$$P_3(\lambda) = (\lambda - \alpha_1) P_2(\lambda) - \beta_1 P_1(\lambda)$$

$$P_{k+1}(\lambda) = (\lambda - \alpha_k) P_k(\lambda) - \beta_{k-1} P_{k-1}(\lambda)$$

c'est à dire, encore :

$$\frac{P_{k+1}(\lambda)}{P_k(\lambda)} = \frac{\lambda - \alpha_k - \beta_{k-1}}{\frac{\lambda - \alpha_{k-1} - \beta_{k-2}}{\frac{\lambda - \alpha_{k-2} - \beta_{k-3}}{\dots \frac{\lambda - \alpha_2 - \beta_1}{\lambda - \alpha_1}}}}$$

et $P_{n+1}(\lambda) = (\lambda - \alpha_n) P_n(\lambda) - \beta_{n-1} P_{n-1}(\lambda)$

Alors : $P_{k+1}(\lambda)$ est le polynôme caractéristique de la matrice obtenue en rayant les $(n-k)$ dernières lignes et colonnes de $\|T\|$.

$P_{n+1}(\lambda)$ est donc le polynôme caractéristique de $\|T\|$, c'est à dire celui de A.

3°- Obtention des vecteurs et lignes propres

Supposons que nous ayons calculé les racines $\lambda_1 ; \lambda_2 ; \dots ; \lambda_n$ du polynôme $P_{n+1}(\lambda)$

Soit u_i ($i = 1 ; 2 ; \dots ; n$) les vecteurs propres de la matrice A et v_i ses lignes propres.

Nous pouvons écrire :

$$b_i = \sum_{j=1}^n a_j u_j$$

d'où nous déduisons que :

$$b_i = \sum_{j=1}^n a_j P_i(\lambda_j) u_j$$

En explicitant ainsi tous les vecteurs b_i nous obtenons un système linéaire.

En résolvant celui-ci nous trouverons :

$$u_i = \sum_{j=1}^n a_{ij} b_j$$

Nous pourrions alors écrire :

$$c_k u_i = \sum_{j=1}^n a_{ij} c_k b_j = a_{ik} c_k b_k ;$$

d'où nous tirons $a_{ik} = \frac{c_k u_i}{c_k b_k}$

or : $c_k = \sum_{j=1}^n a_j' P_k(\lambda_j) v_j$

d'où $c_k u_i = P_k(\lambda_i) v_i u_i a_i'$

et $a_{ik} = \frac{P_k(\lambda_i) v_i u_i a_i'}{c_k b_k}$

v_i et u_i étant définis à une homothétie près nous pouvons supposer $a_{ik} = \frac{P_k(\lambda_i)}{\sigma_k}$

c'est à dire :

$$u_i = \sum_{k=1}^n \frac{P_k(\lambda_i)}{\sigma_k} b_k$$

et de la même façon, nous aurons

$$v_i = \sum_{k=1}^n \frac{P_k(\lambda_i)}{\sigma_k} c_k$$

Ce qui permet de déterminer aisément les vecteurs et lignes propres.

Chapitre II

Etude des Itérés d'un vecteur dans le cas
de diviseurs élémentaires non linéaires

J'ai utilisé dans ce chapitre les notations employées dans [I] (chapitres II et III) et dans [II] (§ 22 et 23).

En particulier j'appelle sous-espace invariant d'une matrice un sous-espace qui est transformé globalement en lui-même par la transformation linéaire définie par la matrice. A chaque valeur propre de la matrice peut-être attaché un sous-espace invariant. Si à cette valeur propre correspond un diviseur élémentaire non linéaire d'ordre P on définit une suite de vecteurs invariants.

$$\begin{aligned}
x^1, x^2, \dots, x^P \quad \text{tels que : } Ax^1 &= \lambda x^1 \\
Ax^2 &= \lambda x^2 + x^1 \\
- - - - - \\
Ax^i &= \lambda x^i + x^{i-1}
\end{aligned}$$

x^1 est un vecteur propre et est unique, par contre x^i peut être remplacé par $Z^i = \sum_{j=1}^{j=i} a_{ij}^j x^j$ avec $a_{ii}^i \neq 0$. Une suite de vecteurs invariants : Z^1, Z^2, \dots, Z^P est dite générée par Z^1 et Z^P est appelé générateur de la suite.

Rappelons également que si la nullité de A est zéro, l'espace peut-être décomposé suivant les sous-espaces invariants attachés aux différentes valeurs propres de la matrice et que dans chaque espace invariant une suite de vecteurs invariants constitue une base dudit sous-espace.

Ceci étant j'appellerai sous-espace de profondeur i attachée à la valeur propre λ l'ensemble des vecteurs X tels que $(A - \lambda I)^i X = 0$

Il est clair que cet ensemble de vecteurs constitue bien un sous-espace.

D'autre part ce sous-espace est un sous-espace invariant de la matrice en effet :

$$\begin{aligned} (A - \lambda I)^{i-1} AX &= (A - \lambda I)^{i-1} \lambda X, \\ \text{donc } (A - \lambda I)^i AX &= (A - \lambda I)^i \lambda X = 0 \end{aligned}$$

Ce sous-espace appartient donc au sous-espace invariant attaché à la valeur propre λ .

Les vecteurs invariants x^j tels que $j < i$ appartiennent au sous-espace de profondeur i .

Profondeur d'un vecteur : Un vecteur appartenant à un sous-espace de profondeur i mais pas au sous-espace de profondeur $i+1$ sera dit de profondeur i . C'est le cas du vecteur invariant x^i (lorsque $i \leq p$).

Il est clair qu'un sous-espace de profondeur i contient tous les espaces de profondeurs inférieures.

Profondeur maximale : Si le sous-espace de profondeur $k+1$ est identique au sous-espace de profondeur k , mais que celui-ci soit différent du sous-espace de profondeur $k-1$, alors k sera dit profondeur maximale des sous-espaces attachée à la valeur propre λ .

Ce sous-espace de profondeur maximale est évidemment confondu avec le sous-espace invariant attaché à la valeur propre λ . Nous en déduisons immédiatement que $k=p$, ordre du diviseur élémentaire d'ordre le plus élevé attaché à λ .

2° - Etude des vecteurs transformés d'un vecteur par les puissances successives de la matrice.

Nous pouvons choisir comme base de l'espace une base formées de vecteurs invariants de la matrice, il suffit donc d'étudier les transformés d'un vecteur invariant.

Nous poserons $X_i = 0 ; \forall i \leq 0$

Et les transformés successifs du vecteur X_j seront :

$$AX_j = \lambda X_j + X_{j-1}$$

$$A^2 X_j = A(\lambda X_j + X_{j-1}) = \lambda^2 X_j + 2\lambda X_{j-1} + X_{j-2}$$

$$A^k X_j = \sum_{l=0}^k C_k^l \lambda^{k-l} X_{j-l} ; C_k^l \text{ étant le coefficient du binôme}$$

Théorème : Considérons les vecteurs : $X_i ; AX_i , \dots , A^k X_i$.

Ces vecteurs sont linéairement indépendants si $k < i$ et linéairement dépendants dans le cas contraire.

En effet, supposons que k étant inférieur à i les vecteurs ne soient pas linéairement indépendants, nous pouvons alors écrire une relation de la forme.

$$a_0 X_i + a_1 AX_i + \dots + a_k A^k X_i = 0$$

or $(A - \lambda I) X_i = X_{i-1}$ et $(A - \lambda I)^p X_i = X_{i-p}$

Multiplions la relation précédente par $(A - \lambda I)^{i-1}$

Nous obtenons : $a_0 X_i = 0$ soit $a_0 = 0$

Multiplions là ensuite par $(A - \lambda I)^{i-2}$ nous obtenons

$$a_1 AX_i = 0 , \text{ comme } a_0 = 0 , a_1 AX_i = 0 \text{ et } a_1 = 0$$

Multiplions ensuite par $(A - \lambda I)^{i-3}$, etc... jusqu'à $(A - \lambda I)^{i-k-1}$

Nous montrerons ainsi que a_0 , a_1 , \dots , a_k sont tous nuls ce qui prouve que les vecteurs $X_i , AX_i ; \dots , A^k X_i$ sont bien linéairement indépendants.

D'autre part nous avons :

$(A - \lambda I)^i X_i = 0$ or $(A - \lambda I)^i$ est un polynôme de degré i en A , dont les coefficients ne sont pas tous nuls, ce qui nous donne une relation de dépendance linéaire entre les vecteurs : $X_i, AX_i, \dots, A^i X_i$

Le théorème est donc démontré.

Théorème : Soit une matrice A ayant un polynôme minimal de degré m et un vecteur V ayant pour chaque valeur propre de la matrice une composante non nulle suivant un vecteur de profondeur maximale associé à cette valeur propre ; alors les vecteurs : $V ; AV ; A^2V ; \dots ; A^kV$ seront linéairement indépendants si $k < m$ et linéairement dépendants dans le cas contraire.

En effet soit $V = \sum_{s=1}^r V_s$, V_s composante sur le $s^{\text{ième}}$ espace propre ; une relation de dépendance linéaire s'écrit $p(A)V = 0$; $p(A)$ étant un polynôme en A ; projetons cette relation sur le $s^{\text{ième}}$ espace propre elle s'écrit : $p(A)V_s = 0$

Or si le $s^{\text{ième}}$ espace propre à une profondeur p_s , nous avons $(A - \lambda_s I)^{p_s} V_s = 0$, qui est la relation linéaire de plus faible degré sur le $s^{\text{ième}}$ espace propre. Donc $p(A)$ contient $(A - \lambda_s I)^{p_s}$ en facteur. En faisant ce raisonnement pour chaque espace propre on voit que $p(A)$ est divisible par

$\prod_{i=1}^r (A - \lambda_i I)^{p_i}$ qui est le polynôme minimal de la matrice ce qui démontre le théorème.

Remarque : Supposons que dans un sous-espace V n'ait pas de composante suivant un vecteur de profondeur maximale, (de profondeur p_h) mais seulement suivant un vecteur de profondeur $p_h - f_h$.

Dans ce sous-espace la relation de dépendance linéaire sera alors seulement de degré : $p_h - f_h$ et non plus de degré p_h ; et la suite des vecteurs V, AV, \dots, A^kV sera formée de vecteurs linéairement dépendants à partir de $k = m - f_h$.

Réciproquement : si la suite des vecteurs $V ; AV ; \dots ; A^k V$ est formée de vecteurs linéairement dépendants à partir de $k = m-f$ alors dans au moins un sous-espace, V n'a pas de composante suivant un vecteur de profondeur maximale. De plus si f_s représente dans le $s^{\text{ième}}$ espace propre la différence entre la profondeur de V et la profondeur maximale de l'espace propre ; on a la relation

$$f = \sum_{s=1}^r f_s$$

Chapitre III

Etude de la formation des vecteurs successifs
de la méthode de Lanczos

1° - Retour sur la définition des vecteurs b_i et c_i de la méthode de Lanczos

Définition : Soit un sous-espace E_{c_k} (resp. E_{1k}) appartenant à E_c (resp E_1)

On appelle sous-espace perpendiculaire à E_{c_k} (resp E_{1k}) un sous-espace de E_1 (resp E_c) formé de tous les vecteurs de E_1 (resp E_c) orthogonaux à une base de E_{c_k} (resp E_{1k})

Ce sous-espace perpendiculaire à E_{c_k} a comme dimension $n - k$.

Nous le désignerons par $E_{c_k}^*$ (resp E_{1k}^*) et nous écrirons $E_{c_k} \perp E_{c_k}^*$

Nous noterons : dimension de $E_k = k$ de la façon suivante :

$$D(E_k) = k.$$

$$\text{Posons } J_k = \{b_1; \dots; b_k\} = \{b_1; Ab_1; \dots; A^{k-1}b_1\}$$

$$I_k = \{c_1; \dots; c_k\}$$

b_{k+1} est alors défini de la façon suivante :

$$b_{k+1} \in J_{k+1} \quad \text{et} \quad b_{k+1} \in I_k^*$$

$$\text{donc } b_{k+1} \in (J_{k+1} \cap I_k^*)$$

Cette intersection est au moins de dimension un, pour qu'elle ne soit pas de dimension supérieure il faut et il suffit que

$$J_{k+1} \cdot I_k^* = E_c.$$

Supposons que $\forall j \neq 0; j \leq k$; condition qui est nécessaire pour que l'algorithme se déroule normalement

Nous appellerons \bar{E} le complémentaire d'un espace E

Supposons maintenant que $(\bar{J}_k \cdot I_k^*) \neq \emptyset$. Alors il existe certainement un vecteur $v_0 \in E_1; v_0 \neq 0$ et perpendiculaire à $J_k \cdot I_k^*$.

$V_0 \notin \overline{I}_K$ puisque $V_0 \perp I_K^*$ donc $V_0 \in I_K$

Nous pouvons donc écrire $V_0 = \sum_{i=1}^k a_i c_i$

or V_0 étant perpendiculaire à J_K les produits scalaires $V_0 b_1$; $V_0 b_2 \dots$; $V_0 b_k$ sont tous nuls.

Ce qui peut encore s'écrire :

$$V_0 \cdot b_j = \sum_{i=1}^k a_i c_i \cdot b_j = a_j \sigma_j = 0$$

et ceci pour tout j donc $a_j = 0$; $\forall j$

Et $V_0 = 0$ ce qui contredit l'hypothèse. Donc $\overline{J_K \cdot I_K^*} = \emptyset$;
l'intersection est donc bien de dimension égale à I et b_{k+1} est déterminé de façon unique.

2°- Expression des quantités $\sigma_K, \rho_K, \beta_K$, sous forme de déterminants

Considérons la matrice

$$M_{p+1} = \begin{vmatrix} c_1 b_1 & c_1 A b_1 & \dots & c_1 A^p b_1 \\ c_1 A b_1 & & & \\ \vdots & & & \\ \vdots & & & \\ c_1 A^p b_1 & \dots & \dots & c_1 A^{2p} b_1 \end{vmatrix}$$

Appelons D_{p+1} son déterminant

En effectuant des combinaisons de lignes nous écrirons :

$$D_{p+1} = \begin{vmatrix} c_1 b_1 & c_1 A b_1 & \dots & c_1 A^p b_1 \\ c_1 b_2 & c_1 A b_2 & \dots & c_1 A^p b_2 \\ \dots & \dots & \dots & \dots \\ c_1 b_{p+1} & c_1 A b_{p+1} & \dots & c_1 A^p b_{p+1} \end{vmatrix}$$

En effectuant alors des combinaisons de colonnes il viendra finalement :

$$D_{p+1} = \begin{vmatrix} c_1^{b_1} & c_2^{b_1} & \dots & c_{p+1}^{b_1} \\ c_1^{b_2} & c_2^{b_2} & \dots & c_{p+1}^{b_2} \\ \vdots & \vdots & \ddots & \vdots \\ c_1^{b_{p+1}} & \dots & \dots & c_{p+1}^{b_{p+1}} \end{vmatrix} = \prod_{i=1}^{p+1} \sigma_i$$

Nous en déduisons immédiatement :

$$\sigma_{p+1} = \frac{D_{p+1}}{D_p}$$

Et comme $\Delta_p = \frac{\sigma_{p+1}}{\sigma_p}$

$$\Delta_{p+1} = \frac{D_{p+1} E_{p-1}}{D_p^2}$$

Considérons maintenant la matrice

$$N_p = \begin{vmatrix} c_1 A b_1 & \dots & c_1 A^p b_1 \\ \vdots & \ddots & \vdots \\ c_1 A^p b_1 & \dots & c_1 A^{2p} b_1 \end{vmatrix}$$

Nous pouvons l'écrire :

$$N_p = \begin{vmatrix} c_1 A \\ c_1 A^2 \\ \vdots \\ c_1 A^p \end{vmatrix} \times \begin{vmatrix} b_1 A b_1 & \dots & A^{p-1} b_1 \end{vmatrix}$$

Nous pouvons effectuer des combinaisons de lignes dans la première matrice du produit pour faire apparaître : c_2, c_3, \dots, c_p , et de même des combinaisons de colonnes dans la seconde feront apparaître : b_2, b_3, \dots, b_p

N_p s'écrira alors :

$$N_p = T \begin{vmatrix} c_1^A \\ c_2^A \\ \vdots \\ c_p^A \end{vmatrix} \times \begin{vmatrix} b_1 & b_2 & \dots & b_p \end{vmatrix} T^t ;$$

T étant une matrice triangulaire unitaire.

On voit alors que :

$$N_p = T \begin{vmatrix} \alpha_1 c_1 b_1 & c_2 b_2 & & & \\ \beta_1 c_1 b_1 & \alpha_2 c_2 b_2 & c_3 b_3 & & \\ & & & & \\ & & & & c_p b_p \\ & & & \beta_{p-1} c_{p-1} b_{p-1} & \alpha_p c_p b_p \end{vmatrix} T^t$$

Posons $\Omega_p = \text{Dét } (N_p)$

Alors :

$$\Omega_p = D_p \begin{vmatrix} \alpha_1 & & & & \\ & \alpha_2 & & & \\ & & & & \\ & & & & \\ & & & \beta_{p-1} & \alpha_p \end{vmatrix}$$

A un ordre L , inférieur à n (ordre du polynôme minimal)

a) $\sigma_{K+1} = 0$; $b_{K+1} \neq 0$; $c_{L+1} \neq 0$

b) b_{K+1} (ou c_{L+1}) est nul mais c_{L+1} (ou b_{K+1}) ne l'est pas

c) $b_{K+1} = c_{L+1} = 0$

3°- Etude du cas $\sigma_{K+1} = 0$ avec b_{K+1} et $c_{L+1} \neq 0$

b_{K+1} et c_{L+1} étant tous les deux différents de zéros et σ_{K+1} étant nul cela revient à dire que b_{K+1} est perpendiculaire à c_{L+1}

Dans ce qui suit nous supposerons que b_I et c_I ont des composantes suivants tous les vecteurs de profondeurs maximales respectivement à droite et à gauche de la matrice. Lorsque nous parlerons de vecteurs quelconques cela signifiera à cette restriction près.

Sans restreindre la généralité du raisonnement nous pouvons supposer que à chaque valeur propre n'est associé qu'un seul diviseur élémentaire (celui d'ordre le plus élevé, bien entendu). Nous pouvons, en effet, faire cette supposition puisque la formation de la suite des vecteurs s'arrête à l'ordre du polynôme minimal et que seul le diviseur élémentaire d'ordre le plus élevé intervient dans celui-ci.

Par un choix judicieux de la base, nous pouvons toujours supposer que : $b_i = \sum_{h=1}^r X_h^{ph}$; en appelant p_h la profondeur du sous-espace associé à la $h^{\text{ième}}$ valeur propre λ_h ; le nombre de valeurs propres distinctes étant r .

Nous savons que :

$$\begin{aligned} A^t b_I &= \sum_{h=1}^r A^t X_h^{ph} \\ &= \sum_{h=1}^r \sum_{l=0}^{t-1} C_t^l \lambda_h^{t-l} X_h^{ph-l} \end{aligned}$$

Pour définir le vecteur c_j , nous opérerons de façon différente, nous nous fixerons une base de E_{1n}

$$Y_1^I, \dots, Y_1^{PI}, Y_2^I, \dots, Y_2^{P2}, \dots, Y_r^{Pr};$$

Et nous écrirons alors :

$$c_I = \sum_{h=1}^r \sum_{i=1}^{i=Ph} c_{hi} Y_h^I$$

Effectuons maintenant le produit scalaire :

$$\begin{aligned} \langle c_I, A^t b_I \rangle &= c_j A^t b_j \\ &= \sum_{h=1}^r \sum_{l=0}^t c_t^l h^{t-1} A_{h, Ph-1} \end{aligned}$$

Considérons maintenant le déterminant D_{k+1} (que nous avons défini au paragraphe précédent) il est composé de termes de la forme : $c_I A^t b_I$

Le développement de ce déterminant sera de la forme

$$\sum_{(\text{permutations})} (-1)^q \prod [c_t h^{t-1} a_{h, Ph-1}]$$

Le produit comprenant les différents facteurs de la permutation considérée, dont q est la puissance.

Ce qui nous importe ici c'est que si nous considérons les $c_{h,i}$ (composantes de c_I) comme des inconnues, le développement du déterminant est un polynôme homogène de degré $k+1$.

Si ce polynôme homogène n'est pas identiquement nul ; nous voyons qu'ayant fixé arbitrairement toutes les variables sauf une, nous obtenons un polynôme de degré $k+1$ (au maximum) en la dernière variable et il existe seulement $k+1$ valeurs au maximum qui annulent ce polynôme. Il résulte de cela que si le polynôme homogène n'est pas identiquement nul, il existe une infinité de valeurs de c_I telles que D_{k+1} soit non nul. Or, si D_{k+1} est

différent de zéro, il en est de même de \overline{c}_{k+1} .

Démontrons donc que le polynôme homogène n'est pas identiquement nul. Pour cela nous allons montrer qu'il existe un choix de valeurs telles qu'il n'est pas nul.

Le vecteur b_I étant toujours défini comme ci-dessus nous allons choisir c_I tel que :

$$(I) \quad \left\{ \begin{array}{l} \alpha_I b_I = 0 \\ c_I A b_I = 0 \\ \vdots \\ c_I A^{k-1} b_I = 0 \\ \text{mais} \\ c_I A^k b_I = c \neq 0 \end{array} \right.$$

c_I a n composantes, or nous avons seulement $k+1$ équations nous allons donc imposer comme condition supplémentaire que seules les $k+1$ premières composantes de c_I soient non nulles. Soit : c_I^1 ; c_I^2 ; ... ; c_I^{k+1} ces $k+1$ composantes.

D'autre part nous avons :

$$b_I = \begin{array}{|c|} \hline I \\ \hline 0 \\ \hline \vdots \\ \hline 0 \\ \hline I \\ \hline 0 \\ \hline \vdots \\ \hline 0 \\ \hline I \\ \hline \vdots \\ \hline 0 \\ \hline \vdots \\ \hline \end{array}$$

Les composantes égales à I correspondent aux vecteurs de profondeurs maximales, les composantes nulles aux autres vecteurs.

Nous aurons alors

$$\begin{array}{c}
 \text{Ab}_I = \begin{array}{|c|} \hline \lambda_I \\ \hline I \\ \hline 0 \\ \hline \vdots \\ \hline \lambda_2 \\ \hline I \\ \hline 0 \\ \hline \vdots \\ \hline \end{array} ; A^2 b_I = \begin{array}{|c|} \hline 2 \\ \hline \lambda_I \\ \hline 2\lambda_I \\ \hline I \\ \hline 0 \\ \hline \vdots \\ \hline 0_2 \\ \hline \lambda_2 \\ \hline 2\lambda_2 \\ \hline I \\ \hline 0 \\ \hline \vdots \\ \hline \end{array} ; \dots ; A^i b_I = \begin{array}{|c|} \hline i \\ \hline \lambda_I \\ \hline c_I^i \lambda_I^{i-1} \\ \hline \vdots \\ \hline I \\ \hline 0_i \\ \hline \lambda_2 \\ \hline \vdots \\ \hline \end{array}
 \end{array}$$

Ecrivons maintenant les équations obtenues en effectuant les produits scalaires indiqués :

$$\begin{array}{l}
 c_I^I + 0 + 0 + \dots + 0 + c_I^i + 0 + \dots + 0 + c_I^j + 0 + \dots = 0 \\
 c_I^I \lambda_I + c_I^2 + 0 + \dots + 0 + c_I^i \lambda_2 + c_I^{i+1} + 0 + \dots \\
 \qquad \qquad \qquad + 0 + \dots + c_I^j \lambda_3 + c_I^{j+1} + 0 + \dots = 0 \\
 \\
 c_I^1 \lambda_1^2 + c_I^2 \lambda_1 + c_I^3 + 0 + \dots = 0 \\
 \text{-----} \\
 \text{-----} \\
 c_I^j \lambda_1^j + \dots = c \neq 0
 \end{array}$$

Le déterminant de ce système à la forme suivante :

$$\begin{vmatrix}
 I & 0 & 0 & \dots & I & 0 & \dots\dots\dots \\
 \lambda_1 & I & 0 & & \lambda_2 & I & \\
 \lambda_1^2 & 2\lambda_1 & I & & \lambda_2^2 & 2\lambda_2 & \\
 \vdots & \vdots & \vdots & & \vdots & \vdots & \\
 \lambda_1^k & k\lambda_1^{k-1} & & & \lambda_2^k & k\lambda_2^{k-1} &
 \end{vmatrix}$$

La deuxième colonne est la dérivée de la première par rapport à λ_1 , la troisième la dérivée de la deuxième, etc...

Pour que le système ait une solution il faut et il suffit que le déterminant ci-dessus soit différent de zéro. Il est évident que les différentes colonnes de ce déterminant sont linéairement indépendantes, donc le déterminant est non nul. Ceci nous montre alors que le vecteur c_I satisfaisant le système (I) existe et est unique.

Le déterminant D_{k+1} correspondant à ce vecteur c_I et au vecteur b_I que nous considérons depuis le début de ce paragraphe est triangulaire inférieur. Ses termes diagonaux sont égaux et valent justement $c_I A^k b_I$; donc ici $D_{k+1} = (c_I A^k b_I)^{k+1}$ qui est différent de zéro, par définition.

Ceci nous montre qu'il existe des vecteurs c_I tels que D_{k+1} soit non nul donc le polynôme homogène obtenu en développant D_{k+1} n'est pas identiquement nul ce que nous voulions démontrer.

Nous pouvons effectuer ce raisonnement pour chaque valeur de k inférieure à m , et nous pourrions ainsi énoncé le théorème :

Théorème : Etant donné un vecteur b_I ayant des composantes suivant un vecteur de profondeur maximale correspondant à chacune des valeurs propres. Il existe une infinité de vecteurs c_I tels que aucun des nombres $(c_I A^k b_I)$ ne soit nul pour $k < m$, ordre du polynôme minimal

Autrement dit : ce ne sera que dans des cas particuliers que un vecteur b_K sera perpendiculaire au vecteur c_K correspondant

5°- $\sigma_{K+1} = 0$ mais b_{K+1} ou c_{K+1} nul

Nous allons montrer que si b_I et c_I ont des composantes selon les vecteurs les plus profonds de leurs espaces respectifs, alors il est impossible que b_{K+1} ou c_{K+1} soit nul, si $K < n$ ordre du polynôme minimal.

Dans le premier paragraphe de ce chapitre nous avons donné une définition géométrique des vecteurs b_i et c_i . Reprenons les notations de ce paragraphe : nous y avons montré que :

$$J_K \cdot I_K^* = E_c$$

Si b_I a des composantes suivant tous les vecteurs les plus profonds alors :

$$A^K b_I \not\subset J_K \quad \text{d'ou} \quad Ab_K \not\subset J_K$$

mais alors :

$$J_{K+1} = J_K \cdot Ab_K \text{ a une intersection non nulle}$$

avec I_K^* or :

$$J_{K+1} \cap I_K^* = b_{K+1} \quad \text{ce qui montre que } b_{K+1} \neq 0$$

Par conséquent pour que b_{K+1} soit nul il faut que b_I n'ait pas de composantes suivant tous les vecteurs les plus profonds de son espace.

A *fortiori* pour que b_{K+1} et c_{K+1} soient nuls il faut que pour chacun d'eux, la condition ci-dessus soit remplie.

Finalement nous avons deux sortes d'accidents possibles :

a) b_I ou c_I n'ont pas de composantes suivants tous les vecteurs les plus profonds.

b) b_i et c_i ayant des composantes suivants tous les vecteurs les plus profonds, σ_{K+1}^- est nul pour $K < m$. Dans ce dernier cas b_{L+1} et c_{K+1} sont non nuls. Mais alors nous avons vu qu'il suffisait de modifier un seul des vecteurs de départ pour faire disparaître l'accident en question.

Deuxième Partie

Etude des erreurs sur un pas de
la Méthode de Lanczos

Chapitre I

Erreurs dans un produit scalaire en point décimal flottant

Nous allons maintenant aborder l'étude des erreurs dans la méthode de Lanczos. Cette méthode est théoriquement rigoureuse, mais, du fait du grand nombre d'opérations mises en jeu, on obtient des erreurs de calculs fort importantes.

Ce sont ces erreurs que nous allons étudier.

Les erreurs de calculs dépendant du mode de représentation des nombres, nous supposerons dans ce qui suit que ceux-ci sont écrits en point décimal flottant, c'est-à-dire avec une mantisse de longueur fixe, donnant les chiffres significatifs et un exposant indiquant l'emplacement de la virgule.

Cette représentation est celle qui est le plus souvent utilisée en machine, de plus dans la méthode de Lanczos l'ordre de grandeur des nombres varie dans de larges proportions au cours du calcul et alors le point décimal flottant est à peu près la seule représentation convenable.

Dans ce chapitre, nous commencerons par étudier les erreurs dans les opérations en point décimal flottant (en abrégé P.D.F.) en général. C'est seulement au chapitre suivant que nous aborderons les erreurs dans la méthode de Lanczos proprement dite.

Erreur dans le produit de 2 nombres

Soit à effectuer : $a_1 \cdot a_2$

En point décimal flottant si a_1 et a_2 ont chacun n chiffres significatifs le produit obtenu en aura n également, ceci entraîne

que dans le produit exact les chiffres significatifs ayant les poids les plus faibles ont été tronqués.

Nous représenterons le résultat réellement obtenu par $a_1 \cdot a_2 (I+S)$, S étant un coefficient d'erreur

Ce coefficient est variable pour chaque produit puisqu'il représente la partie tronquée du nombre.

Nous supposerons qu'on peut sans commettre d'erreurs appréciables, le considérer comme une constante.

Erreur dans un produit de plus de deux nombres.

Soit à effectuer $\prod_{i=1}^n a_i$

Nous poserons $a_1 a_2 (I+S) = a_{12}$
nous effectuons ensuite

$$a_{12} \cdot a_3 \rightarrow a_{12} a_3 (I+S) = a_{123} (I+S)^2$$

en appliquant la convention faite ci-dessus

On voit que finalement on obtient

$$(I+S)^{n-1} \prod_{i=1}^n a_i \approx [I + (n-1) S] \prod_{i=1}^n a_i$$

Erreur dans une somme arithmétique.

Nous allons effectuer

$$\sum_{i=1}^n a_i \quad \text{avec } a_i \geq 0$$

La sommation sera effectuée dans l'ordre des indices croissants

c'est à dire que nous aurons l'organigramme suivant

$$\left(\left[(a_1 + a_2) + a_3 \right] + a_4 + \dots \right) + a_n$$

Faisant la même hypothèse que dans le cas d'un produit, nous écrivons que nous obtenons: $(a_1 + a_2) (I+S)$

Cette hypothèse, assez grossière, repose sur le fait que, si le nombre de chiffres significatifs dans la représentation utilisée est n ; dans le cas d'une somme arithmétique, comme dans

celui d'un produit, l'erreur due à la troncature, est, pour un nombre $a : ka10^{-m}$, k étant compris entre 0,1 et 10.

Au résultat obtenu, c'est à dire $(a_1+a_2)(I+S)$ nous ajoutons a_3 et nous obtenons alors: $(a_1+a_2)(I+S)^2 + a_3(I+S)$

Lorsque nous aurons effectué les $n-1$ sommes de notre sommation nous obtiendrons, au lieu du résultat réel, l'expression: $(a_1+a_2)(I+S)^{n-1} + a_3(I+S)^{n-2} + \dots + a_n(I+S)$

Soit, en tenant compte de ce que S est infiniment petit:

$$\sum_{i=1}^n a_i + (n-1)S(a_1+a_2) + (n-2)Sa_3 + \dots + Sa_n$$

Cette expression, pour n assez grand, peut être, sans erreurs importantes remplacée par l'expression approchée :

$$\sum_{i=1}^n a_i + S \sum_{i=1}^n (n-i+1) a_i$$

Une remarque s'impose, ces expressions dépendent des différents termes, et non du résultat global, comme c'était le cas pour le produit. De plus l'erreur sera plus ou moins importante selon l'ordre adopté pour faire la sommation.

Dans le cas où nous ne disposerions que de la somme totale et non des différents termes nous pourrions avoir une expression approchée de l'erreur en supposant tous les termes de la somme égaux. Plaçons nous dans le cas où n est grand ; soit S_n la somme ; l'erreur sera alors donnée par l'expression approchée :

$$\frac{sS_n}{n} \sum_{i=1}^n n-i+1 = \frac{sS_n}{n} \frac{n(n+1)}{2} = \frac{sS_n(n+1)}{2}$$

Erreur dans une somme algébrique :

Lorsque nous effectuons la différence de deux nombres positifs si ces nombres sont très voisins leur différence sera très faible par contre l'erreur de calcul sera du même ordre que celle effectuée sur la somme des 2 nombres.

En effet, posons

$$A = a \cdot 10^{\alpha}$$

$$B = b \cdot 10^{\beta} \quad \text{avec } \alpha \geq \beta \quad A > 0 \quad \text{et} \quad B > 0$$

$$A+B = (a + b \cdot 10^{\beta-\alpha}) \cdot 10^{\alpha}$$

$$A-B = (a - b \cdot 10^{\beta-\alpha}) \cdot 10^{\alpha}$$

L'erreur dans les deux cas est de l'ordre de $k \cdot 10^{\alpha}$, k étant compris entre 0,1 et 10.

Il serait donc complètement faux d'écrire l'erreur s (A-B) comme on serait tenté de le faire par analogie avec la somme arithmétique. Nous aurons une expression donnant une valeur voisine de la valeur exacte de l'erreur si nous prenons comme expression de l'erreur celle correspondant à l'erreur sur la somme des valeurs absolues. Cette expression ne sera qu'approchée mais fournira une évaluation commode de l'erreur effectuée.

Erreur dans un produit scalaire

Dans la méthode de Lanczos l'essentiel des opérations peut-être ramené à une suite de produits scalaires. Il importe donc de connaître l'erreur faite dans une telle opération.

Soit à effectuer le produit du vecteur

$A(a_1, a_2, \dots, a_n)$ par le vecteur $B(b_1, b_2, \dots, b_n)$.

Ce produit scalaire sera

$$\sum_{i=1}^n a_i b_i$$

Cette opération comprend n multiplications suivies d'une sommation algébrique.

Lorsqu'on effectue un produit $a_i b_i$ on obtient au lieu de $a_i b_i$, $(1+S)a_i b_i$, S étant toujours positif car l'erreur était une erreur de troncature à le signe du produit. Nous supposons en effet que les opérations sont faites sous arrondi, comme c'est le cas sur la machine dont nous disposons.

Donc la somme des erreurs faites sur les produits sera

$S \sum_{i=1}^n a_i b_i = s \langle A, B \rangle$, en représentant par $\langle A, B \rangle$ le produit scalaire de A par B.

Il s'agit maintenant d'effectuer la sommation.

Dans la plupart des cas il sera mal commode, voir souvent impossible, de considérer séparément les différents termes de la somme.

Nous utiliserons alors la formule approchée que nous avons obtenue en considérant tous les termes égaux.

Ceci nous donne une erreur sur la somme des termes $a_i b_i$ (indépendamment de l'erreur faite sur les produits)

$$S \frac{n+1}{2} \sum_{i=1}^n |a_i| \cdot |b_i| = S \frac{n+1}{2} \sum_{i=1}^n |a_i| \cdot |b_i| .$$

Appelons \bar{A} le vecteur dont les composantes sont

$|a_1|$, $|a_2|$, ..., $|a_n|$; et \bar{B} le vecteur $(|b_1|, |b_2|, \dots$

$|b_n|$ nous poserons naturellement $|a_i| = \bar{a}_i$
 $|b_i| = \bar{b}_i$

Dès lors:

$$\sum_{i=1}^n |a_i| |b_i| = \sum_{i=1}^n \bar{a}_i \cdot \bar{b}_i = \langle \bar{A}, \bar{B} \rangle = \|\bar{A}\| \cdot \|\bar{B}\| \cos(\bar{A}, \bar{B})$$

$$\text{or } \|\bar{A}\| = \sqrt{\sum_{i=1}^n |a_i|^2} = \sqrt{\sum_{i=1}^n a_i^2} = \|A\|$$

donc ;

$$\sum_{i=1}^n |a_i| |b_i| = \|A\| \cdot \|B\| \cos(\bar{A}, \bar{B}).$$

La valeur de $\cos(\bar{A}, \bar{B})$ est en général inconnue, nous la remplacerons par sa valeur la plus probable. Pour cela nous remarquerons que \bar{A} et \bar{B} étant des vecteurs à composantes toutes positives ou nulles, leurs directions coupent l'hypersphère de rayon 1 en deux points appartenant au premier quadrant (cette notion généralisant la notion identique, en géométrie à 3 dimensions). Ces points de l'hypersphère ayant une probabilité uniforme sur tout le premier quadrant. Une telle probabilité est proportionnelle à l'élément d'aire sphérique.

En définitive chercher la valeur la plus probable de l'angle de \bar{A} et de \bar{B} revient chercher celle de l'angle de deux points X et Y du premier quadrant de l'hypersphère de rayon I. Soit α cet angle

$$I - \cos \alpha = \frac{I X Y I^2}{2}$$

soit x_1, x_2, \dots, x_n les composantes du point X

y_1, y_2, \dots, y_n celles du point Y.

Posons

$$x_1 = \sin \psi_1 \sin \psi_2 \dots \sin \psi_{n-1}$$

$$x_2 = \sin \psi_1 \sin \psi_2 \dots \sin \psi_{n-2} \cos \psi_{n-1}$$

$$x_3 = \sin \psi_1 \sin \psi_2 \dots \sin \psi_{n-3} \cos \psi_{n-2}$$

⋮

$$x_{n-1} = \sin \psi_1 \cos \psi_2$$

$$x_n = \cos \psi_1$$

Et de façon analogue

$$y_1 = \sin \theta_1 \sin \theta_2 \dots \sin \theta_{n-1}$$

$$y_n = \cos \theta_1$$

La distance des deux points s'écrit alors :

$$|XY|^2 = 2-2 \left\{ \cos (\psi_{n-1} - \theta_{n-1}) \prod_{i=1}^{i=n-2} \sin \psi_i \sin \theta_i \right. \\ \left. + \sum_{p=0}^{p=n-3} \cos \psi_{p+1} \cos \theta_{p+1} \prod_{i=1}^{i=p} \sin \psi_i \sin \theta_i \right.$$

d'où :

$$\cos \alpha = \cos (\psi_{n-1} - \theta_{n-1}) \prod_{i=1}^{i=n-2} \sin \psi_i \sin \theta_i \\ + \sum_{p=0}^{p=n-3} \cos \psi_{p+1} \cos \theta_{p+1} \prod_{i=1}^p \sin \psi_i \sin \theta_i$$

L'élément de surface de l'hypersphère, qui représente la probabilité élémentaire d'un point sur cette hypersphère, est : $\sin^{n-2} \varphi_1 \cdot \sin^{n-3} \varphi_2 \dots \sin^2 \varphi_{n-3} \sin \varphi_{n-2} d\varphi_1 d\varphi_2 \dots d\varphi_{n-1}$

La surface du quadrant positif de l'hypersphère est égal à la surface de l'hypersphère : $\frac{2\pi^{n/2}}{\Gamma(\frac{n}{2})}$, divisée par le nombre de

quadrants de cette hypersphère ; soit 2^n ; donc la surface d'un quadrant vaut : $\frac{\pi^{n/2}}{\Gamma(n/2) 2^{n-1}}$

La probabilité cherchée sera alors ;

$$E(\cos \alpha) = \frac{\Gamma(2(n/2)) 2^{2(n-1)}}{\pi^n} \int_0^{\pi/2} \dots \int_0^{\pi/2} [\cos \alpha] \sin^{n-2} \varphi_1 \dots \sin \varphi_{n-2} \sin^{n-2} \theta_1 \dots \sin \theta_{n-2} d\varphi_1 \dots d\theta_{n-1}$$

Le calcul de cette intégrale est tout à fait classique et ne présente pas de difficultés, on trouve finalement :

$$E(\cos \alpha) = \frac{n \Gamma^2(n/2)}{\Gamma^2(\frac{n+1}{2})}$$

Nous allons chercher une expression approchée de $E(\cos \alpha)$, lorsque n est grand, pour cela nous remplacerons les factorielles par leurs expressions par la formule de Stirling soit :

$$E(\cos \alpha) \approx \frac{2 e n}{\pi(n-1)} \left(1 - \frac{1}{n-1}\right)^{n-1}$$

or $\left(1 - \frac{1}{n-1}\right)^{n-1}$ tend vers e^{-1} lorsque n tend vers l'infini

$$\text{Donc } E(\cos \alpha) \approx \frac{2 e n}{\pi(n-1)} \frac{1}{e} \approx \frac{2}{\pi}$$

L'erreur cherchée est alors :

$$\begin{aligned} \frac{n+1}{2} \|A\| \cdot \|B\| E(\cos \alpha) &\approx \frac{n(n+1)}{2} \|A\| \|B\| \frac{2}{\pi} \\ &= \frac{n(n+1)}{\pi} \|A\| \cdot \|B\| \end{aligned}$$

Cette erreur, rappelons-le est celle faite sur la sommation mais, superaveut, nous avons fait une erreur sur les produits deux à deux. L'erreur globale sur un produit scalaire sera alors :

$$\begin{aligned} \bar{\epsilon} (A,B) &= S \left[(A,B) + \|A\| \|B\| \frac{n+1}{2} E(\cos \alpha) \right] \\ &\leq S \|A\| \cdot \|B\| \left[1 + \frac{n+1}{2} E(\cos \alpha) \right] \\ &\neq S \|A\| \cdot \|B\| \left(1 + \frac{n+1}{2} \right) \end{aligned}$$

$\bar{\epsilon} (A,B)$ est l'erreur absolue faite sur le produit scalaire (A,B)
 Nous allons maintenant chercher l'erreur relative faite sur ce produit scalaire soit :

$$\frac{\bar{\epsilon} (A,B)}{(A,B)} = S \left[1 + \frac{\|A\| \cdot \|B\| (n+1) E(\cos \alpha)}{2 (A,B)} \right]$$

or $(A,B) = \|A\| \cdot \|B\| \cos \beta$

en posant β = angle de A et de B

Donc $\frac{\bar{\epsilon} (A,B)}{(A,B)} = S \left[1 + \frac{n+1}{2} \frac{E(\cos \alpha)}{\cos \beta} \right]$

Contrairement à $\cos \alpha$, $\cos \beta$ n'est pas toujours inconnu. Malheureusement, ceci n'est pas le cas général. S'il est inconnu nous chercherons sa valeur moyenne. Pour faire ce calcul nous supposons $\cos \beta$ indépendant de $\cos \alpha$ (ce qui en toute rigueur est faux). Nous remarquerons qu'une des deux directions peut toujours être considérée comme l'axe des pôles et qu'en fait une seule direction varie d'où

$$\begin{aligned} E(\cos \beta) &= \frac{2 \Gamma(n/2)}{2 \pi^{n/2}} \int_0^{\pi/2} \sin^{n-2} \varphi_1 \cos \varphi_1 d\varphi_1 \prod_{i=1}^{n-2} \int_0^{\pi} \sin^{n-i-1} \varphi_i d\varphi_i \\ &\quad \cdot \int_0^{2\pi} d\varphi_{n-1} \\ E(\cos \beta) &= \frac{2 \Gamma(n/2)}{(n-1) \sqrt{\pi} \Gamma(\frac{n-1}{2})} \end{aligned}$$

Appelons \bar{e}_n l'erreur relative

$$\bar{e}_n = S \left[1 + \frac{n+1}{2} \frac{E(\cos \alpha)}{E(\cos \beta)} \right] = S \left[1 + \frac{(n+1) n \Gamma\left(\frac{n}{2}\right)}{2\sqrt{\pi} \Gamma\left(\frac{n+1}{2}\right)} \right]$$

Comme précédemment nous évaluerons les factorielles par la formule de Stirling d'où :

$$\frac{\Gamma\left(\frac{n}{2}\right)}{\Gamma\left(\frac{n+1}{2}\right)} \# \left(1 - \frac{1}{2} \cdot \frac{1}{\frac{n-1}{2}} \right)^{\frac{n-1}{2}} \sqrt{\frac{2}{n-1}} e$$

or $\left(1 - \frac{1}{2} \cdot \frac{1}{\frac{n-1}{2}} \right)^{\frac{n-1}{2}} \rightarrow e^{-1/2}$ quand $n \rightarrow \infty$

Donc $\frac{\Gamma\left(\frac{n}{2}\right)}{\Gamma\left(\frac{n+1}{2}\right)} \# \sqrt{\frac{2}{n-1}}$

Donc : $\bar{e}_n \# S \left[1 + \frac{(n+1) n}{\sqrt{2\pi} (n-1)} \right]$

Et comme nous avons supposé n assez grand nous confondrons $(n+1)$ et $(n-1)$ avec n , ce qui nous donnera :

$$\bar{e}_n \# S \left(1 + \frac{n^{3/2}}{\sqrt{2\pi}} \right)$$

Ce qui nous montre que pour n assez grand l'erreur relative, sur un produit scalaire, varie en moyenne comme $n^{3/2}$

Vérification :

Ce dernier résultat, l'erreur relative sur un produit scalaire proportionnelle en moyenne à $\left(1 + \frac{n^{3/2}}{\sqrt{2\pi}} \right)$ est suffisamment simple pour que nous puissions espérer arriver à le vérifier.

Principe : On génère des couples de points aléatoires sur l'hyper-sphère, de rayon un, à n dimensions. Ces points ayant une densité de répartition uniforme.

Ces points aléatoires sont considérés comme des extrémités de vecteurs ayant pour origine le centre de l'hyper-sphère.

On effectue, pour chaque couple de points, le produit des deux vecteurs. Ce produit est effectué d'une part en simple précision d'autre part en double précision ; la différence entre les deux résultats trouvés représente l'erreur faite en simple précision.

Si on effectue un tel calcul pour un grand nombre de couples de point, la moyenne des erreurs trouvées nous fournira la valeur probable de l'erreur dans l'espace considéré.

Nous savons que l'erreur est de la forme $k n^{3/2}$.

Effectuons le calcul ci-dessus pour un espace à n_1 dimensions, ceci de calculer effectivement k ; connaissant k nous pouvons prévoir la valeur de l'erreur que nous devons avoir dans un espace à n_2 dimensions ; le calcul précédent effectué alors dans l'espace à n_2 dimension nous permettra de vérifier notre théorie.

On peut d'ailleurs, faire mieux et tracer la courbe erreur en fonction de la dimension de l'espace. Cette courbe doit être semblable à la courbe $y = x^{3/2}$

Mise en oeuvre : Le procédé utilisé pour générer des points sur l'hyper-sphère est celui donné par Mervin E. Muller dans les communications de l'A.C.M. d'avril 1959. Ce procédé est basé sur la propriété suivante : Si les x_i , ($i=1, 2, \dots, n$) sont des nombres aléatoires suivant une loi normale. Alors le point y , de l'espace à n dimensions, défini par $y_i = \frac{x_i}{\left(\sum_{i=1}^n x_i^2\right)^{1/2}}, i=1, 2, \dots, n$; et $y = (y_1, y_2, \dots, y_n)$, a une probabilité de répartition uniforme sur la surface de l'hyper-sphère.

Des nombres aléatoires uniformément répartis sont générés par le programme, ces nombres sont transformés à l'aide d'une table de la fonction de Gauss en une répartition normale. Les coordonnées du point y sont ensuite calculées. Un deuxième point de l'hypersphère est ensuite calculé. Il ne reste plus qu'à effectuer le produit en simple et en double précision et à faire la différence de ces deux produits.

Le rapport de cette différence au produit lui-même nous donne l'erreur relative dans le produit scalaire que nous venons d'effectuer.

Le programme écrit effectue ainsi automatiquement 100 produits scalaires successifs, puis imprime la moyenne des erreurs effectuées sur ces 100 produits scalaires.

Le temps de calcul est rigoureusement proportionnel à la dimension de l'espace. Il est malheureusement assez important : le calcul de 100 produits scalaires réclame 10 minutes pour un espace à 5 dimensions. Or ainsi que nous allons le voir la dispersion des résultats oblige à effectuer plusieurs milliers de produits scalaires pour obtenir une moyenne significative.

Si nous effectuons plusieurs séries de produits scalaires, chaque série comprenant 100 produits scalaires, les erreurs moyennes obtenues dans chacune de ces séries constituent une série d'épreuve d'une variable aléatoire. Variable aléatoire, qui du fait du théorème central limite, suit approximativement une loi de Gauss. Nous avons effectué pour chaque point vingt séries d'essais soit 2000 produits scalaires, malheureusement les résultats sont très dispersés.

Ainsi pour un espace à 5 dimensions la moyenne est de $0,417 \cdot 10^{-8}$ et l'écart type (entre les 20 valeurs constituées chacune par le résultat d'une série de 100 essais) est de $0,182 \cdot 10^{-8}$. Pour 15 dimensions nous avons respectivement :

moyenne : $1,797 \cdot 10^{-8}$; écart type $1,967 \cdot 10^{-8}$. Ceci nous montre le peu de précision des résultats obtenus et la nécessité de prendre des échantillons d'autant plus importants que la dimension de l'espace est élevé. Or le temps de calcul d'un produit scalaire étant proportionnel à la dimension de l'espace, les temps de calcul deviennent rapidement prohibitifs.

Ces erreurs étant faites, nous avons tracé sur le même graphique la courbe : $y = 1 + \frac{n^{3/2}}{\sqrt{2\pi}}$; et la courbe des erreurs moyennes, ceci, entre les points 5 et 15.

La courbe d'erreur a été tracée à l'aide des points suivants :

dimension de l'espace	erreur moyenne
5	0,417 10^{-8}
7	0,731 10^{-8}
10	0,890 10^{-8}
13	1,328 10^{-8}
15	1,797 10^{-8}

Les échelles des deux courbes ont été choisies de telle sorte que les deux courbes passent, sur le graphique, par le même point pour $n = 5$.

La coïncidence apparaît comme remarquable, ce qui semble bien confirmer, les études précédentes. De plus, nous voyons que la formule obtenue est valable déjà pour des espaces à cinq dimensions, c'est à dire pour un nombre de dimensions assez réduit. Ceci nous permet de voir que le "n assez grand" dont nous avons parlé plus haut est, en fait, assez faible, et de l'ordre de cinq.

M. Erreur moyenne

- 39 bis -

$20 \rightarrow 1,5 \cdot 10^{-8}$

$15 \rightarrow 1 \cdot 10^{-8}$

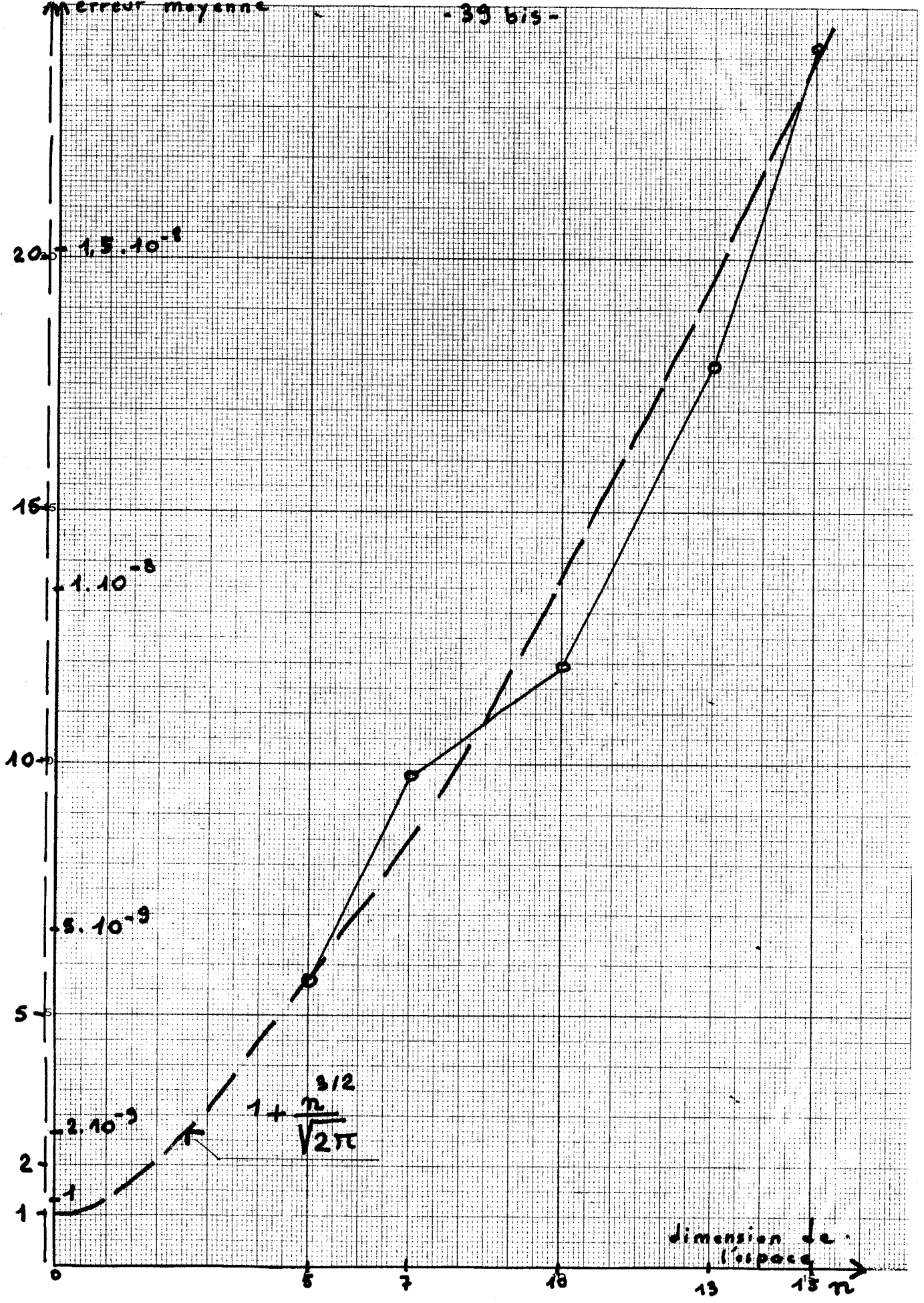
$10 \rightarrow 8 \cdot 10^{-9}$

$5 \rightarrow 2 \cdot 10^{-9}$

1

$$1 + \frac{\pi^{3/2}}{\sqrt{2\pi}}$$

dimension de l'espace \rightarrow
13 15 π



Chapitre II

Erreur à un pas de la méthode de Lanczos

1°- Définition d'un pas :

Nous supposons connu $b_k, c_k, \sigma_{k-1}, \alpha_{k-1}, \beta_{k-2}$ et les quantités analogues d'indices inférieurs. Nous appellerons alors k ième pas l'ensemble des opérations conduisant à la détermination de $\sigma_k, \alpha_k, \beta_{k-1}, b_{k+1}, c_{k+1}$.

Le premier pas est différent des pas suivants : en effet β_0 n'existe pas, nous pouvons le supposer identique en posant

$$\beta_0 = 0$$

2°- Invariance de la méthode par un changement de base

Supposons que nous effectuions un changement de base défini par une matrice H .

La matrice A deviendra la matrice $A^* = H^{-1} A H$ le vecteur c_k deviendra $c_k^* = c_k H$ et le vecteur b_k deviendra le vecteur $b_k^* = H^{-1} b_k$

Effectuons alors le produit de $(c_k H)$ par $(H^{-1} b_k)$

$$c_k H H^{-1} b_k = c_k b_k = \sigma_k$$

de même :

$$(c_k H) \cdot (H^{-1} A H) \cdot (H^{-1} b_k) = c_k A b_k$$

donc $\sigma_k, \alpha_k, \beta_{k-1}$ vont être invariants.

Dans la nouvelle base nous aurons :

$$\begin{aligned} b_{k+1}^* &= A^* b_k^* - \alpha_k b_k^* - \beta_{k-1} b_{k-1}^* \\ &= (H^{-1} A H) (H^{-1} b_k) - \alpha_k (H^{-1} b_k) - \beta_{k-1} (H^{-1} b_{k-1}) \\ &= H^{-1} b_{k+1} \end{aligned}$$

Ceci nous montre bien l'invariance de la méthode. Cette invariance nous sera utile au chapitre suivant pour étudier la propagation des erreurs.

3°- Hypothèses de calcul

a) Nous supposons que les erreurs sont petites ce qui nous permet de les considérer comme indépendantes.

b) Nous négligerons les erreurs dans les produits simples devant celles dans les produits scalaires.

c) Dans ce chapitre nous étudions les erreurs à un pas. En conséquence nous supposons que tous les calculs effectués antérieurement au $k^{\text{ième}}$ pas l'ont été sans erreur.

4°- Recensement des erreurs

Au cours des calculs du $k^{\text{ième}}$ pas nous allons d'abord faire des erreurs sur α_k et β_{k-1} ceci aura pour première conséquence d'altérer la matrice tridiagonale.

De plus les erreurs sur α_k et β_{k-1} fausseront b_{k+1} et c_{k+1} .

En plus de ces erreurs sur α_k et β_{k-1} , il y aura les erreurs propres à c_{k+1} et à b_{k+1} , c'est à dire des erreurs faites sur ces vecteurs et ne provenant pas de l'inexactitude de α_k et β_{k-1} .

Dans ce qui suit, je supposerai toujours que n est assez grand pour qu'il soit possible d'utiliser pour les produits scalaires les formules simplifiées obtenues à l'aide de la formule de Stirling.

D'autre part soit une quantité a , et a^* la valeur approchée de cette quantité, je désignerai par $d(a)$ l'erreur faite sur cette quantité

Ce qui donnera

$$a^* = a + d(a)$$

5° - Erreurs sur α_k et β_{k-1}

Nous avons deux façons d'envisager l'étude de ces erreurs : ou bien nous utiliserons les erreurs absolues sur les produits scalaires, ou bien les erreurs relatives.

Ces deux procédés ne nous donneront pas exactement le même résultat. En effet pour passer du stade des erreurs absolues à celui des erreurs relatives nous avons fait une étude statistique supplémentaire ce qui enlève de la précision aux résultats obtenus. Par contre l'utilisation des erreurs absolues pour α_k et β_{k-1} nous conduira à introduire des majorations que nous ne ferions pas si nous utilisions les erreurs relatives. Ajoutons à cela que les formules obtenues à l'aide des erreurs relatives sont plus maniables que celles obtenues en utilisant les erreurs absolues.

Il ne semble pas possible de comparer la précision des résultats obtenus par les deux méthodes. Dans la suite des calculs il sera plus intéressant, suivant le cas, d'utiliser tantôt l'une tantôt l'autre évaluation.

a) Utilisation des erreurs absolues

Etant donné une matrice A, réelle, quelconque, d'éléments a_{ij} ; j'appelle norme de A, la quantité :

$$\left[\sum_{i=1}^n \sum_{j=1}^n a_{ij}^2 \right]^{1/2},$$

Je désignerai cette quantité par N (A).

J'appelle norme de la i^{ème} ligne de A la quantité :

$$\left[\sum_{j=1}^n a_{ij}^2 \right]^{1/2} \quad \text{et je la désigne par } \|A_i\|$$

J'ai la relation

$$N(A) = \left[\sum_{i=1}^n \|A_i\|^2 \right]^{1/2}$$

Ces définitions étant posées ; considérons les calculs d'un pas de la méthode de Lanczos.

J'ai :

$$\alpha_k = \frac{c_k A b_k}{\sigma_k} \quad \text{avec} \quad \sigma_k = c_k b_k$$

Pour calculer $c_k A b_k$: je multiplie b_k par A puis j'effectue le produit scalaire du vecteur obtenu par c_k

Donc :

$$d(\alpha_k) = \frac{c_k d(Ab_k) + d[c_k (Ab_k)]}{\sigma_k} - \alpha_k \frac{d(\sigma_k)}{\sigma_k}$$

Considérons Ab_k : Ab_k est un vecteur

La $i^{\text{ème}}$ composante de ce vecteur est obtenue par le produit scalaire de la $i^{\text{ème}}$ ligne de A par b_k . Appelons $(Ab_k)_i$ cette $i^{\text{ème}}$ composante j'ai :

$$d[(Ab_k)_i] = \|A_i\| \cdot \|b_k\| S \frac{(\pi + n + 1)}{\pi}$$

puisque $\|A_i\|$ est la norme de la $i^{\text{ème}}$ ligne de A considéré comme un vecteur.

Soit maintenant le produit scalaire :

$c_k d(Ab_k)$. Il est majoré par :

$$\begin{aligned} \text{or } \|c_k\| \cdot \|d(Ab_k)\| &= \left[\sum_{i=1}^n d[(Ab_k)_i]^2 \right]^{1/2} \\ &= \left\{ \sum_{i=1}^n \left[S \|A_i\| \|b_k\| \frac{(\pi + n + 1)}{\pi} \right]^2 \right\}^{1/2} \\ &= S \|b_k\| \frac{(\pi + n + 1)}{\pi} \left\{ \sum_{i=1}^n \|A_i\|^2 \right\}^{1/2} \\ &= S \frac{(\pi + n + 1)}{\pi} \|b_k\| \cdot N(A) \end{aligned}$$

Etudions maintenant $d(c_k (Ab_k))$.

Il s'agit d'un produit scalaire et l'erreur sera :

$$S \|c_k\| \cdot \|Ab_k\| \frac{(\overline{\sigma}_{k-1} + \overline{\sigma}_k)}{\overline{\sigma}_k}$$

Cette erreur est majorée par

$$s \frac{(\overline{\sigma}_{k-1} + \overline{\sigma}_k)}{\overline{\sigma}_k} \|c_k\| \cdot \|b_k\| N(A)$$

Enfin

$$d(\overline{\sigma}_k) = S \frac{(\overline{\sigma}_{k-1} + \overline{\sigma}_k)}{\overline{\sigma}_k} \|c_k\| \cdot \|b_k\|.$$

En définitive l'erreur sur d_k sera majorée par

$$|d(d_k)| \leq \frac{S (\overline{\sigma}_{k-1} + \overline{\sigma}_k)}{\overline{\sigma}_k} \|c_k\| \cdot \|b_k\| [2 N(A) + |d_k|]$$

Etudions maintenant les erreurs sur β_{k-1} . β_{k-1} a comme expression :

$$\beta_{k-1} = \frac{\overline{\sigma}_k}{\overline{\sigma}_{k-1}}$$

Donc

$$d(\beta_{k-1}) = \frac{d(\overline{\sigma}_k)}{\overline{\sigma}_{k-1}} - \beta_{k-1} \frac{d(\overline{\sigma}_{k-1})}{\overline{\sigma}_{k-1}}$$

Ici se pose une difficulté, théoriquement $d(\overline{\sigma}_{k-1})$ est nul, mais par contre $d\overline{\sigma}_k$ interviendra dans $d(\beta_k)$ à la même place que $d(\overline{\sigma}_{k-1})$ dans $d(\beta_{k-1})$. C'est pourquoi je raisonnerai pour $d(\beta_{k-1})$ comme si $d(\overline{\sigma}_{k-1})$ était différent de zéro. D'ailleurs dans le calcul réel, $d(\overline{\sigma}_{k-1})$ est différent de zéro, et indépendamment de l'erreur transmise par β_{k-2} , fait apparaître une erreur sur β_{k-1} , au $k^{\text{ième}}$ pas.

$$d(\beta_{k-1}) = \frac{1}{\overline{\sigma}_{k-1}} (d\overline{\sigma}_k - \beta_{k-1} d\overline{\sigma}_{k-1}).$$

d'où :

$$|d\beta_{k-1}| \leq \left| \frac{1}{\sigma_{k-1}} \frac{S (n+n+1)}{\pi} (\|c_k\| \cdot \|b_k\| + |\beta_{k-1}| \cdot \|c_{k-1}\| \cdot \|b_{k-1}\|) \right|$$

b) Utilisation des erreurs relatives :

Nous nous plaçons dans les mêmes hypothèses qu'en a).

Nous avons alors :

$$d(Ab_k) = S \left(1 + \frac{n^{3/2}}{\sqrt{2}\pi} \right) Ab_k \text{ puisque chaque composante est}$$

un produit scalaire.

$c_k \cdot (Ab_k)$ est également un produit scalaire et l'erreur totale sera :

$$\frac{|d(c_k \cdot Ab_k)|}{|c_k \cdot Ab_k|} = 2 S \left(1 + \frac{n^{3/2}}{\sqrt{2}\pi} \right)$$

D'autre part

$$\frac{|d(\sigma_k)|}{|\sigma_k|} = S \left(1 + \frac{n^{3/2}}{\sqrt{2}\pi} \right)$$

Donc, comme les erreurs relatives s'ajoutent dans un quotient :

$$\frac{|d(\sigma_k)|}{|\sigma_k|} = S \left(1 + \frac{n^{3/2}}{\sqrt{2}\pi} \right)$$

et

$$\frac{|d(\beta_{k-1})|}{|\beta_{k-1}|} = 2 S \left(1 + \frac{n^{3/2}}{\sqrt{2}\pi} \right)$$

6° - Erreurs sur les vecteurs b_K et c_K

Compte-tenu des hypothèses que nous avons formulées, nous pouvons écrire :

$$d(b_{K+I}) = d(Ab_K) - d(\alpha_K) b_K - d(\beta_{K-I}) b_{K-I}$$

Nous voyons qu'en plus des erreurs dues à α_K et β_{K-I} il y aura le vecteur $d(Ab_K)$. Selon l'hypothèse a) ou b) nous aurons pour $d(Ab_K)$ les expressions suivantes :

a) $d(Ab_K)$ est un vecteur dont les composantes sont :

$$[d(Ab_K)]_i = \|A_i\| \cdot \|b_K\| \frac{S(n+n+I)}{\pi}$$

$$b) \quad d(Ab_K) = \frac{S(I+n^{3/2})}{\sqrt{2\pi}} Ab_K$$

je vais maintenant chercher quelle est, dans le cas a) la valeur de $\|d(b_{K+I})\|$

$$\|d(b_{K+I})\| \leq \|d(Ab_K)\| + |d(\alpha_K)| \|b_K\| + |d(\beta_{K-I})| \|b_{K-I}\|$$

cherchons d'abord

$$\|d(Ab_K)\| = N(A) \|b_K\| \frac{S(n+n+I)}{\pi}$$

Nous l'avons vu plus haut ,

donc

$$\|d(b_{K+I})\| < \frac{S(n+n+I)}{\pi} \left\{ \frac{I}{|\sigma_K|} \|c_K\| \|b_K\|^2 [2N(A) + |\alpha_K|] \right. \\ \left. + N(A) \|b_K\| + \frac{I}{|\sigma_{K-I}|} [\|c_K\| \cdot \|b_K\| + |\beta_{K-I}| \cdot \|c_{K-I}\| \cdot \|b_{K-I}\|] \|b_{K-I}\| \right\}$$

Ou encore en groupant différemment les termes :

$$\|d(b_{k+I})\| < \frac{S(\pi+n+I)}{\pi} \left\{ \|c_k\| \cdot \|b_k\| \left[\frac{\|b_k\|}{|\sigma_k|} (2N(A) + |\alpha_k|) + \frac{\|b_{k-I}\|}{|\sigma_{k-I}|} \right] \right. \\ \left. + \frac{|\beta_{k-I}|}{|\sigma_{k-I}|} \|c_{k-I}\| \cdot \|b_{k-I}\|^2 + N(A) \|b_k\| \right\}$$

7°- Importance de σ_k :

Reprenons l'expression $d(\alpha_{k-I}) < \frac{S(\pi+n+I)}{\pi} [2N(A) + |\alpha_k|]$

$$\cdot \frac{\|c_k\| \cdot \|b_k\|}{|\sigma_k|} \cdot \frac{\|c_k\| \cdot \|b_k\|}{|\sigma_k|} = \frac{I}{\cos(c_k, b_k)}$$

cette dernière expression peu varier de I à $+\infty$ alors que

$$\frac{S(\pi+n+I)}{\pi} [2N(A) + |\alpha_k|]$$

varie assez peu, puisqué, seul

$|\alpha_k|$ varie.

$\left(\frac{I}{\cos(c_k, b_k)} \right)$ intervient également dans $d(\beta_{k-I})$ puisque :

$$|d(\beta_{k-I})| \leq \frac{I}{|\sigma_k|} \frac{S(\pi+n+I)}{\pi} \|c_{k+I}\| \|b_{k+I}\| + \frac{S(\pi+n+I)}{\pi} \left(\frac{I}{\cos(c_k, k)} \right)$$

Ici son influence est plus faible, mais elle n'en demeure pas moins importante.

Ceci nous montre que dès que l'angle de c_k avec b_k s'approche de $\frac{\pi}{2}$ les erreurs vont croître dans des proportions astronomiques il suffit que $\cos(c_k, b_k)$ soit divisé par 10 pour que l'erreur soit multipliée approximativement par la même quantité.

Calculer à chaque pas $\cos(c_K, b_K)$ (et nous verrons ultérieurement que cela est très facile) permettra donc de contrôler que l'erreur garde une valeur normale. Car il faut remarquer que, si l'angle (c_K, b_K) devient très voisin de $\frac{\pi}{2}$, ceci est lié au choix des vecteurs de départ. Nous pouvons donc dire que nous avons d'une part une erreur "normale" qu'on ne peut éviter, c'est celle que nous avons étudiée de façon probabiliste dans le chapitre précédent et dans celui-ci, d'autre part une erreur exceptionnelle due au mauvais choix de vecteurs de départ ; erreur qui se traduit par un angle (c_K, b_K) voisin de $\frac{\pi}{2}$. Cette erreur-ci, qui peut devenir beaucoup plus importante que l'erreur normale, peut-être évitée en modifiant le choix des vecteurs de départ.

Remarquons que s'il n'y avait pas d'erreur de calcul $(c_K, b_K) = \frac{\pi}{2}$ nous donnerait un cas singulier. Avec les erreurs de calcul ce n'est plus seulement la direction perpendiculaire à c_K que b_K ne doit pas avoir, mais c'est dans un certain cône, autour de cette direction, que b_K ne doit pas être.

Chapitre III

Etude de la propagation d'une erreur sur α_k ou β_k

Notations de ce chapitre : Considérons une certaine quantité a nous appellerons a' l'erreur faite sur cette quantité.

1° - Propagation des erreurs sur un pas

Supposons que des erreurs aient été commises au $k^{\text{ième}}$ pas et que aux pas suivants les seules erreurs proviennent de la propagation de ces erreurs du $k^{\text{ième}}$ pas.

Nous avons :

$$c'_{k+1} = c'_k \Delta - \alpha'_k c'_k - \alpha'_k c'_k - \beta'_{k-1} c'_{k-1} - \beta'_{k-1} c'_{k-1} \quad (1)$$

d'ou nous déduisons que :

$$c'_{k+1} b_{k+p} = c'_k \Delta b_{k+p} - \alpha'_k c'_k b_{k+p} - \beta'_{k-1} c'_{k-1} b_{k+p} \quad (2)$$

ceci pour $p \geq 1$.

remplaçons Δb_{k+p} par : $b_{k+p+1} + \alpha'_{k+p} b_{k+p} + \beta'_{k+p-1} b_{k+p-1}$

Nous obtenons :

$$c'_{k+1} b_{k+p} = c'_k b_{k+p+1} - \alpha'_{k+p} c'_k b_{k+p} - \beta'_{k+p-1} c'_k b_{k+p-1} - \alpha'_k c'_k b_{k+p} - \beta'_{k-1} c'_{k-1} b_{k+p} \quad (3)$$

L'étude faite dans ce chapitre sera basée sur ces trois formules.

2° - Propagation d'une erreur sur α_s

Supposons qu'aucune erreur n'ait été faite avant le s ième pas et que à ce pas la seule erreur faite le soit sur α_s : soit α'_s .

Nous allons commencer par étudier l'influence de cette erreur sur les produits $c_j b_i$

Nous avons par hypothèse : $c'_s = 0$

donc $c''_s b_i = 0$ pour tout i

L'équation nous permet d'écrire

$$c'_{s+1} = -\alpha'_s c_s$$

$$c'_{s+1} b_i = -\alpha'_s c_s b_s$$

donc :

$$c'_{s+1} b_i = 0 \quad \forall i > s$$

L'équation (I) nous donne ensuite :

$$c''_{s+2} = c'_{s+1} A - \alpha'_{s+1} c'_{s+1} - \alpha'_{s+1} c'_{s+1}$$

$$c''_{s+2} b_{s+2+P} = c'_{s+1} A b_{s+1+P} - \cancel{\alpha'_{s+1} c'_{s+1} b_{s+2+P}} - \alpha'_{s+1} c'_{s+1} b_{s+2+P}$$

Ceci pour tout $P \geq 0$

en remplaçant c'_{s+1} par sa valeur soit : $-\alpha'_s c_s$

$$c''_{s+2} b_{s+2+P} = -\alpha'_s \cancel{c_s A b_{s+1+P}} - \alpha'_{s+1} \cancel{\alpha'_s c_s b_{s+2+P}}$$

$$\text{donc } c''_{s+2} b_{s+2+P} = 0 \quad \forall P \geq 0$$

Supposons maintenant que :

$$c'_j b_i = 0 ; \quad \forall j \leq k ; \quad \forall i \geq j$$

et démontrons que : $c'_{l+1} b_{k+P} = 0 ; \quad \forall P \geq 1$

(3) nous permet d'écrire :

$$c'_{k+1} b_{k+p} = c'_k b_{k+p+1} - \alpha_{k+p} c'_k b_{k+p} + \beta_{k+p-1} c'_k b_{k+p-1} - \alpha_k c'_k b_{k+p} - \beta_{k-1} c'_{k-1} b_{k+p}$$

en vertu de la supposition faite nous avons :

$$c'_k b_{k+p+1} = 0 ; c'_k b_{k+p} = 0 ; c'_k b_{k+p-1} = 0 ; c'_k b_{k+p} = 0 ;$$

$$c'_{k-1} b_{k+p} = 0$$

puisque $P \geq 1$.

$$\text{donc : } c'_{k+1} b_{k+p} = 0$$

Ceci nous montre que si le théorème est vrai jusqu'à l'ordre k , il est vrai jusqu'à l'ordre $k+1$. Or nous l'avons démontré directement jusqu'à l'ordre 3.

$$\text{donc : } \boxed{c'_j b_i = 0 \quad \forall i \geq j ; \quad \forall j}$$

Par symétrie nous en déduisons que : $c_j b'_i = 0 ; \forall i ; \forall j \geq i$

$$\text{Il en résulte que } \boxed{\sigma'_j = 0 \quad \forall j}$$

en effet :

$$\sigma'_j = c'_j b_j + c_j b'_j = 0$$

de même, comme $\beta_j = \frac{\sigma'_{j+1}}{\sigma_j}$

$$\boxed{\beta'_j = 0 ; \forall j}$$

Etudions maintenant les erreurs sur α_{s+p}

$$\alpha'_{s+1} = \frac{c'_{s+1} A b_{s+1} + c_{s+1} A b'_{s+1}}{\sigma_{s+1}} \quad \text{puisque } \sigma'_{s+1} = 0$$

$$\alpha'_{s+1} = -\alpha'_s \frac{c_s A b_{s+1} + c_{s+1} A b_s}{\sigma_{s+1}} = -2\alpha'_s \frac{\sigma'_{s+1}}{\sigma_{s+1}}$$

$$\text{done } \boxed{\alpha'_{s+1} = -2\alpha'_s}$$

$$\alpha'_{s+2} = \frac{c'_{s+2} A b_{s+2} + c_{s+2} A b'_{s+2}}{\sigma_{s+2}}$$

Etudions le terme $c'_{s+2} A b_{s+2}$

$$\begin{aligned} c'_{s+2} A b_{s+2} &= (c'_{s+1} A - \alpha'_{s+1} c_{s+1} - \alpha_{s+1} c'_{s+1}) (b_{s+3} + \alpha_{s+2} b_{s+2} + \beta_{s+1} b_{s+1}) \\ &= c'_{s+1} A (b_{s+3} + \alpha_{s+2} b_{s+2} + \beta_{s+1} b_{s+1}) - \alpha'_{s+1} c_{s+1} \beta_{s+1} b_{s+1} \\ &= c'_{s+1} \left[b_{s+4} + \alpha_{s+3} b_{s+3} + \beta_{s+2} b_{s+2} + \alpha_{s+2} (b_{s+3} + \alpha_{s+2} b_{s+2} \right. \\ &\quad \left. + \beta_{s+1} b_{s+1}) + \beta_{s+1} (b_{s+2} + \alpha_{s+1} b_{s+1} + \beta_s b_s) \right] \\ &\quad - \alpha'_{s+1} \beta_{s+1} \sigma_{s+1} \\ &= c'_{s+1} \beta_{s+1} \beta_s b_s - \alpha'_{s+1} \beta_{s+1} \sigma_{s+1} \\ &= c'_{s+1} b_s \frac{\sigma_{s+2}}{\sigma_s} - \alpha'_{s+1} \sigma_{s+2} \\ &= -\alpha'_s c_s b_s \frac{\sigma_{s+2}}{\sigma_s} - \alpha'_{s+1} \sigma_{s+2} \\ &= -(\alpha'_s + \alpha'_{s+1}) \sigma_{s+2} \end{aligned}$$

par symétrie :

$$c_{s+2} A b_{s+2} = -(\alpha'_s + \alpha'_{s+1}) \sigma_{s+2}$$

$$\text{d'où } : \alpha'_{s+2} = -2(\alpha'_s + \alpha'_{s+1}) = -2(\alpha'_s - 2\alpha'_s) = 2\alpha'_s$$

$$\boxed{\alpha'_{s+2} = 2\alpha'_s}$$

Etudions maintenant le produit scalaire :

$$\begin{aligned}
 f'_{s+p} &= c'_{s+p} \Lambda b_{s+p} \\
 f'_{s+p} &= (c'_{s+p-1} \Lambda - \alpha'_{s+p-1} c'_{s+p-1} - \chi_{s+p-1} c'_{s+p-1} - \beta_{s+p-2} c'_{s+p-2}) \Lambda b_{s+p} \\
 &= (c'_{s+p-1} \Lambda - \alpha'_{s+p-1} c'_{s+p-1} - \beta_{s+p-2} c'_{s+p-2}) (b_{s+p+1} + \chi_{s+p} b_{s+p} \\
 &\quad + \beta_{s+p-1} b_{s+p-1}) - c'_{s+p-1} \sigma_{s+p} \\
 &= c'_{s+p-1} \Lambda (b_{s+p+1} + \chi_{s+p} b_{s+p} + \beta_{s+p-1} b_{s+p-1}) - \alpha'_{s+p-1} \sigma_{s+p} \\
 &= c'_{s+p-1} [b_{s+p+2} + \alpha_{s+p+1} b_{s+p+1} + \beta_{s+p} b_{s+p} + \alpha_{s+p} (b_{s+p+1} \\
 &\quad + \chi_{s+p} b_{s+p} + \beta_{s+p-1} b_{s+p-1}) + \beta_{s+p-1} (b_{s+p} + \alpha_{s+p-1} b_{s+p-1} \\
 &\quad + \beta_{s+p-2} b_{s+p-2})] - \alpha'_{s+p-1} \sigma_{s+p} \\
 &= c'_{s+p-1} \beta_{s+p-1} / \beta_{s+p-2} b_{s+p-2} - \alpha'_{s+p-1} \sigma_{s+p} \\
 &= \frac{\sigma_{s+p}}{\sigma_{s+p-2}} c'_{s+p-1} b_{s+p-2} - \alpha'_{s+p-1} \sigma_{s+p} \\
 &= \sigma_{s+p} \left[\frac{(c'_{s+p-2} \Lambda - \alpha'_{s+p-2} c'_{s+p-2} - \alpha_{s+p-2} c'_{s+p-2})}{\sigma_{s+p-2}} \right. \\
 &\quad \left. - \beta_{s+p-3} c'_{s+p-3} \right] b_{s+p-2} - \alpha'_{s+p-1} \sigma_{s+p} \\
 &= \sigma_{s+p} \left[\frac{1}{\sigma_{s+p-2}} (c'_{s+p-2} \Lambda b_{s+p-2} - \alpha'_{s+p-2} \sigma_{s+p-2}) - \alpha'_{s+p-1} \right] \\
 &= \sigma_{s+p} \left[\frac{c'_{s+p-2} \Lambda b_{s+p-2}}{\sigma_{s+p-2}} - \alpha'_{s+p-1} - \alpha'_{s+p-2} \right]
 \end{aligned}$$

Il en résulte que :

$$\frac{f'_{s+p}}{\sigma_{s+p}} = \frac{f'_{s+p-2}}{\sigma_{s+p-2}} - \alpha'_{s+p-1} - \alpha'_{s+p-2}$$

posons :

$$\frac{f'_{s+p}}{J_{s+p}} = g'_{s+p}$$

il vient alors :

$$g'_{s+p} = g'_{s+p-2} - \alpha'_{s+p-1} - \alpha'_{s+p-2} \quad (4)$$

Calculons maintenant :

$$\begin{aligned} \alpha'_{s+p} &= \frac{c'_{s+p} Ab_{s+p} + c_{s+p} Ab'_{s+p}}{J_s} \\ &= \frac{c'_{s+p-2} Ab_{s+p-2} + c_{s+p-2} Ab'_{s+p-2}}{J_{s+p-2}} - 2 [\alpha'_{s+p-1} + \alpha'_{s+p-2}] \\ &= \alpha'_{s+p-2} - 2 [\alpha'_{s+p-1} + \alpha'_{s+p-2}] \end{aligned}$$

$$\alpha'_{s+p} = - (2 \alpha'_{s+p-1} + \alpha'_{s+p-2})$$

Nous avons donc une formule de récurrence qui nous donne α'_{s+p} .

supposons que pour :

$$1 \leq j < p \quad . \quad \alpha'_{s+j} = -\alpha'_{s+j-1}$$

$$\text{alors } \alpha'_{s+p-1} = -\alpha'_{s+p-2} \quad \text{d'où : } \alpha'_{s+p} = \alpha'_{s+p-2} = -\alpha'_{s+p-1}$$

Donc le théorème est encore vrai pour p

$$\text{Or : } \alpha'_{s+1} = -2 \alpha'_s$$

$$\alpha'_{s+2} = 2 \alpha'_s = -\alpha'_{s+1}$$

donc en définitive

$$\left. \begin{aligned} \alpha'_{s+2r} &= 2 \alpha'_s \\ \alpha'_{s+2r-1} &= -2 \alpha'_s \end{aligned} \right\} r > 0$$

3° - Propagation d'une erreur sur β_s

Supposons qu'aucune erreur n'ait été faite avant la sième pas et qu'à ce pas la seule erreur soit sur β_s .

donc $c'_s = 0$ $\alpha'_s = 0$ et $\beta_s \neq 0$

$$c_{s+1} = c_s A - \alpha_s c_s - \beta_{s-1} c_{s-1}$$

β_s n'intervenant pas nous aurons :

$$c'_{s+1} = 0 \quad \text{donc } \sigma'_{s+1} = 0 \quad \text{et } \alpha'_{s+1} = 0$$

$$c'_{s+2} = -\beta'_s c_s$$

$$\sigma'_{s+2} = c'_{s+2} b_{s+2} + c_{s+2} b'_{s+2} = -\beta'_s (c_s b_{s+2} + c_{s+2} b_s) = 0$$

et puisque $\sigma'_{s+2} = 0$, $\beta'_{s+1} = 0$

$$\alpha'_{s+2} = \frac{c'_{s+2} A b_{s+2}}{\sigma_{s+2}} + \frac{c_{s+2} A b'_{s+2}}{\sigma_{s+2}} = 0$$

donc $\sigma'_{s+2} = \alpha'_{s+2} = 0$;

considérons maintenant c'_{s+3}

$$\begin{aligned} c'_{s+3} &= c'_{s+2} A - \alpha'_{s+2} c'_{s+2} \\ &= -\beta'_s [c_s A - \alpha'_{s+2} c_s] \end{aligned}$$

c'_{s+3} est une combinaison linéaire de $c_s A$ et de c_s ; c'_{s+3} s'exprime donc linéairement en fonction de c_s et de c_{s+1}

Supposons maintenant que c'_{s+k} s'exprime linéairement en fonction de c_s à c_{s+k-2} et ceci pour tout k compris entre 1 et p inclus.

$$\text{Soit } c'_{s+k} = \sum_{i=1}^{i=k-2} e_{s+i} c_{s+i} \quad \forall 1 \leq k \leq p$$

Démontrons que c'_{s+p+1} s'exprime linéairement en fonction de c_s à c_{s+p-1}

$$c'_{s+k} b_{s+k} = \left(\sum_{i=1}^{i=k-2} e_{s+i} c_{s+i} \right) b_{s+k} = 0$$

$$c'_{s+k} A b_{s+k} = \left(\sum_{i=1}^{i=k-2} e_{s+i} c_{s+i} \right) A b_{s+k} = 0$$

donc $\forall 1 \leq k \leq p \quad \alpha'_{s+k} = 0$ et $\beta'_{s+k} = 0$

$\forall 1 < k \leq p \quad \beta'_{s+k} = 0$

On a alors :

$$\begin{aligned} c'_{s+p+1} &= c'_{s+p} A - \alpha'_{s+p} c'_{s+p} - \beta'_{s+p-1} c'_{s+p-1} \\ &= \sum_{i=1}^{p-2} e_i c_i A - \alpha'_{s+p} \sum_{i=1}^{p-2} e_i c_i - \beta'_{p-1} \sum_{i=1}^{p-3} f_i c_i \\ &= \sum_{i=1}^{p-1} h_i c_i \quad \text{c.q.f.d.} \end{aligned}$$

Or $c'_{s+1} = 0$

$c'_{s+2} = -\beta'_s c_s$, s'exprime linéairement en fonction de c_s .

c'_{s+3} s'exprime linéairement en fonction de c_s et c_{s+1} donc c'_{s+p} s'exprime linéairement en fonction de c_s à c_{s+p-2} pour p quelconque

Mais alors : $\left| \alpha'_{s+p} = \beta'_{s+p} = 0 \quad \forall p > 0 ; \alpha'_s = 0 \right|$

Donc une erreur sur β'_s ne donne aucune erreur sur les α et les β ultérieurs.

En ce qui concerne les altérations des coefficients β_j j'ai obtenu les résultats suivants :

Tout d'abord les coefficients α ne sont absolument pas modifiés. Quand aux coefficients β , voici ce qu'ils deviennent :

Si je multiplie β_3 par 1,001 j'obtiens $\beta_3 = 9,787025$;
 $\beta_4 = 5,841128$; $\beta_5 = 5,716857$

Alors que je devais avoir pour β_4 et β_5 les valeurs suivantes :
 $\beta_4 = 5,841113$; $\beta_5 = 5,716838$

Donc il y a une erreur parasite sur β_4 et β_5 mais cette erreur est très faible elle vaut, en effet, environ $\beta_3' \times 10^{-3}$

Si je multiplie β_2 par 1,01 j'obtiens :

$\beta_2 = 9,088119$; $\beta_3 = 9,778148$; $\beta_4 = 5,842839$; $\beta_5 = 5,722064$

L'erreur parasite ici est de l'ordre $\beta_2' \times 10^{-2}$

Ceci nous montre donc que si les erreurs transmises sur les coefficients α suivent rigoureusement la règle il se produit une erreur parasite sur les coefficients β particulièrement lorsque l'erreur transmise provient d'une erreur initiale sur un coefficient β

mais cette erreur parasite est faible, très faible même pour les erreurs provenant d'une erreur initiale sur un coefficient α .

De plus il faut remarquer que l'erreur de 1/100 faite dans notre dernier essai est d'un ordre très supérieur à celui des erreurs de calculs courantes.

Or ainsi que nous le montre l'essai précédent l'erreur parasite est d'autant plus réduite (en valeur relative par rapport à l'erreur initiale) que l'erreur initiale est faible. Donc on peut dire que la propagation des erreurs sur α et β suit très étroitement la règle trouvée.

4°- Vérifications :

Pour vérifier l'exactitude de ces résultats j'ai utilisé un programme accessoire qui multiplie au choix un des coefficients α ou β par 1,01 ou 1,001. N'importe lequel des coefficients α ou β peut ainsi être modifié. j'ai effectué de nombreux essais, tant avec 1,01 qu'avec 1,001 ; et toujours les erreurs trouvées se sont révélées parfaitement conformes aux formules ci-dessus.

Ainsi pour une matrice 6 x 6 dont la forme tridiagonale exacte est :

$$\begin{pmatrix} -0,5 & 14,91666 & & & & \\ I & -0,5 & 3,998137 & & & \\ & I & -0,5 & 9,777248 & & \\ & & I & -0,5 & 5,341113 & \\ & & & I & -0,5 & 5,716333 \\ & & & & I & -0,5 \end{pmatrix}$$

J'ai obtenu comme matrice tridiagonale, après multiplication de α par 1,01, la matrice ci-dessous :

$$\begin{pmatrix} -0,5 & 14,91666 & & & & \\ I & -0,5 & 3,998137 & & & \\ & I & -0,505 & 9,777273 & & \\ & & I & -0,4900001 & 5,341147 & \\ & & & I & -0,5099998 & 5,716339 \\ & & & & I & -0,4900005 \end{pmatrix}$$

Si j'appelle e l'erreur prévue par l'étude ci-dessus on voit que sur α l'erreur est environ de $e(1+10^{-5})$

Sur β l'erreur parasite est environ 10 fois plus grande que sur α , (soit $2 \cdot 10^{-4}$, e étant l'erreur sur α , puisque'il n'est pas prévu d'erreur sur β)

Donc on peut dire que la formule est parfaitement vérifiée.

Chapitre IV

Propagation d'une erreur sur b_K

Une erreur faite sur b_K va commencer par altérer α_K et $c_K A b_K$ donc α_K et β_{K-1} . Ce n'est pas tout : indépendamment des erreurs faites sur α_K et β_{K-1} nous aurons une erreur sur b_{K+1} due à ce que b_K est erroné et dans l'expression :
$$b_{K+1} = (A - \alpha_K I) b_K - \beta_{K-1} b_{K-1}$$
, il faudra remplacer b_K par $(b_K + b'_K)$.

En définitive les quantités qui nous intéressent sont les coefficients de la matrice tridiagonale. Or nous connaissons la propagation d'une erreur sur un de ces coefficients.

Nous allons en conséquence procéder de la façon suivante pour faire notre étude :

Nous supposons une erreur sur b_K , cette erreur fausse α_K et β_{K-1} de α'_K et β'_{K-1} . Mais une fois calculés α'_K et β'_{K-1} nous ne nous en préoccuperont plus, car nous connaissons leurs propagations b'_K fausse également b_{K+1} de b'_{K+1} puis en chaîne tous les b_{K+p} . Ces erreurs sur les b_i donnent des erreurs sur les coefficients de la matrice tridiagonale mais nous ne nous en soucions pas, car nous connaissons la propagation de telles erreurs.

Par conséquent, nous allons commencer par étudier la propagation de b'_K par les vecteurs b_{K+p} sans tenir compte des erreurs commises sur les α_i et β_i (Ceci est possible en égard à l'hypothèse de linéarité des erreurs). Ensuite nous étudierons l'erreur engendrée sur α_K et β_{K-1}

Notations : nous prenons les mêmes qu'au chapitre précédent mais nous faisons en plus l'hypothèse que b'_k peut être décomposé suivant les vecteurs propres u_i

$$\text{Nous posons alors } b'_k = \sum_{i=1}^n (b'_k)_i u_i$$

I° - Propagation d'une erreur par les vecteurs b_{k+p}

Nous supposons, ce que nous pouvons toujours faire, que par une transmutation l'espace est rapporté aux vecteurs propres.

$$b_{k+1} = (A - \alpha_k I) b_k - \beta_{k-1} b_{k-1}$$

Donc, si b_k est remplacé par $b'_k + b_k$

$$\begin{aligned} b'_{k+1} &= (A - \alpha_k I) b'_k = \sum_{i=1}^n (A - \alpha_k I) (b'_k)_i u_i \\ &= \sum_{i=1}^n (\lambda_i - \alpha_k) (b'_k)_i u_i = \sum_{i=1}^n (b'_{k+1})_i u_i \end{aligned}$$

en appelant λ_i la valeur propre associée au vecteur propre u_i

$$b_{k+2} = (A - \alpha_{k+1} I) b_{k+1} - \beta_k b_k \quad \text{Donc}$$

$$b'_{k+2} = (A - \alpha_{k+1} I) b'_{k+1} - \beta_k b'_k$$

$$b'_{k+2} = \sum_{i=1}^n \left[(\lambda_i - \alpha_{k+1}) (b'_{k+1})_i - \beta_k (b'_k)_i \right] u_i$$

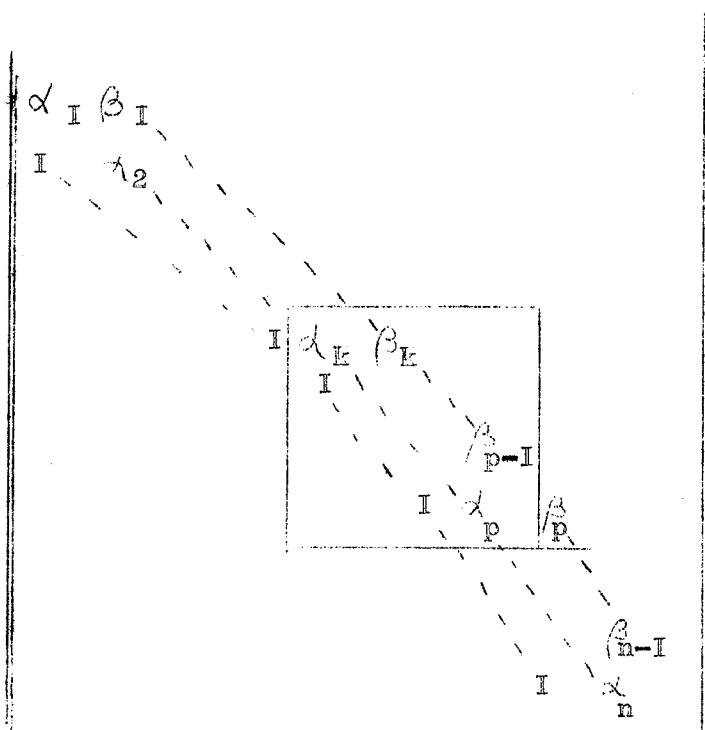
et de façon générale

$$b'_{k+p+1} = (A - \alpha_{k+p} I) b'_{k+p} - \beta_{k+p-1} b'_{k+p-1}$$

$$= \sum_{i=1}^n \left[(\lambda_i - \alpha_{k+p}) (b'_{k+p})_i - \beta_{k+p-1} (b'_{k+p-1})_i \right] u_i$$

Nous voyons immédiatement apparaître l'analogie avec la formation du polynôme caractéristique de la matrice tridiagonale.

Considérons cette matrice



Nous avons isolé une sous-matrice, également tridiagonale.
 Nous appellerons une telle sous-matrice B_k^p

C'est une matrice carrée d'ordre $(p-k+1)$ matrice ayant la même diagonale de la matrice tridiagonale entre les termes de rang k et de rang p

Il est clair que l'erreur sur le vecteur b_{p+1} aura pour composante suivant le vecteur propre u_i

d'où

$$b_{p+1}'' = \sum_{i=1}^n \left[\text{Dét} \left| B_k^p - \lambda_i I \right| (b_k'')_i u_i \right]$$

2° - Erreur sur les termes λ_k et β_k due à une erreur b_k

Au lieu de calculer \mathcal{J}_k nous calculons \mathcal{J}_k^* qui vaut

$$\mathcal{J}_k^* = b_k \cdot c_k + b_k' c_k = \mathcal{J}_k + \mathcal{J}_k'$$

De même au lieu de calculer $c_k Ab_k$ nous calculons

$$(c_k A)(b_k + b_k') = c_k Ab_k + c_k Ab_k'$$

Donc au lieu de λ_k nous aurons:

$$\lambda_k^* = \frac{c_k Ab_k + c_k Ab_k'}{c_k b_k + c_k b_k'} \neq \frac{c_k Ab_k + c_k Ab_k'}{c_k b_k} \left(1 - \frac{c_k b_k'}{c_k b_k} \right)$$

$$\text{d'où } \alpha'_k = \frac{c_k A b'_k}{c_k b_k} - \frac{\alpha_k c_k b'_k}{\sigma_k}$$

$$\alpha'_k = \frac{1}{\sigma_k} \left[c_k A - \alpha_k c_k \right] b'_k$$

$$\alpha'_k = \frac{1}{\sigma_k} \left[c_{k+1} + \beta_{k-1} c_{k-1} \right] b'_k$$

Pour $\beta_{k-1} = \frac{\sigma_k}{\sigma_{k-1}}$ nous obtenons de même

$$\beta_{k-1}^* = \frac{\sigma_k + \sigma_k}{\sigma_{k-1}} \text{ ce qui nous donne une erreur}$$

$$\beta'_{k-1} = \frac{\sigma_k}{\sigma_{k-1}} = \frac{c_k b'_k}{\sigma_{k-1}} = \frac{\beta_{k-1} c_k b'_k}{\sigma_k}$$

En résumé nous obtiendrons comme erreurs :

$$\alpha'_k = \frac{c_{k+1} b'_k}{\sigma_k} + \frac{c_{k-1} b'_k}{\sigma_{k-1}}$$

$$\beta'_{k-1} = \frac{c_k b'_k}{\sigma_{k-1}}$$

Nous pouvons majorer ces erreurs en remplaçant chaque produit scalaire par le produit des normes de vecteurs, nous obtenons alors :

$$|\alpha'_k| \leq \left(\frac{\|c_{k+1}\|}{|\sigma_k|} + \frac{\|c_{k-1}\|}{|\sigma_{k-1}|} \right) \|b'_k\|$$

$$\text{et } |\beta'_{k-1}| \leq \frac{\|c_k\| \cdot \|b'_k\|}{|\sigma_{k-1}|}$$

Tous les α et β postérieurs seront aussi entachés d'erreurs, les formules seront analogues il suffit de remplacer b'_k par b'_{k+p} pour avoir l'erreur sur α_{k+p} et β_{k+p-1}

Comme b'_{k+p} est une erreur due à la propagation de b'_k . Il serait intéressant d'avoir α'_{k+p} et β'_{k+p} en fonction de b'_k

$$\text{Or } b'_{k+p} = \sum_{i=1}^n \left[\text{Det} \left| B_L^{k+p-1} - \lambda_i I \right| (b'_k)_i u_i \right]$$

Nous pouvons facilement trouver un majorant de α'_k et de β'_k

En effet $|(b'_{k+p})_i| \leq \text{Max}_i | \text{Det} [B_E^{k+p-I} - \lambda_i I] | |(b')_i|$

On en déduit que $\|b'_{k+p}\| \leq \text{Max}_i | \text{Det} (B_E^{k+p-I} - \lambda_i I) | \|b'_k\|$

D'où les formules

$$|\alpha'_{k+p}| \leq \text{Max}_i | \text{Det} (B_E^{k+p-I} - \lambda_i I) | |\alpha'_k|$$

$$|\beta'_{k+p-I}| \leq \text{Max}_i | \text{Det} B_E^{k+p-I} - \lambda_i I | |\beta'_{k-I}|$$

Remarque : Sur α'_{k+p} il y a outre cette erreur les erreurs propagées à partir des erreurs sur les α'_{k+j} antérieures ($j < p$), erreurs qui sont dues aux produits scalaires faux. Sur β'_{k+p-I} , au contraire il n'y a pas d'erreurs propagées.

Nous allons maintenant chercher l'erreur totale causée sur un terme α'_{k+p} .

Nous avons : $|\alpha'_{k+j}| \leq \text{Max}_i | \text{Det} (B_E^{k+j-I} - \lambda_i I) | |\alpha'_k|$

Cette erreur se transmet sur les termes suivants de la façon suivante :

$$\begin{aligned} \alpha'_{k+j+2i} &= 2 \alpha'_{k+j} \\ \alpha'_{k+j+2i-1} &= -2 \alpha'_{k+j} \end{aligned}$$

Donc $|\alpha'_{k+j+i}| = 2 |\alpha'_{k+j}|$

Je vais désigner par α'_{k+p} l'erreur totale sur α'_{k+p} , on aura :

$$|\alpha'_{k+p}| < |\alpha'_k| + 2 \sum_{j=1}^{p-1} \left\{ \text{Max}_i | \text{Det} (B_E^{k+j-I} - \lambda_i I) | \right\} |\alpha'_k| + \text{Max}_i | \text{Det} (B_E^{k+p-I} - \lambda_i I) | |\alpha'_k|$$

donc

$$|\alpha'_{k+p}| < |\alpha'_k| \left\{ 1 + \text{Max}_i | \text{Det} (B_E^{k+p-I} - \lambda_i I) | + 2 \sum_{j=1}^{p-1} \text{Max}_i | \text{Det} (B_E^{k+j-I} - \lambda_i I) | \right\}$$

et en remplaçant $|\alpha'_k|$ par sa valeur

$$|\alpha'_{k+p}| \leq \|b'_k\| \left(\frac{\|c_{k+1}\|}{\sqrt{k!}} + \frac{\|c_{k-1}\|}{\sqrt{(k-1)!}} \right) \left\{ 1 + \text{Max}_i | \text{Det} (B_E^{k+p-I} - \lambda_i I) | + 2 \sum_{j=1}^{p-1} \text{Max}_i | \text{Det} (B_E^{k+j-I} - \lambda_i I) | \right\}$$

Tandis que

$$|\beta'_{k+p-I}| \leq \text{Max}_i | \text{Det} (B_E^{k+p-I} - \lambda_i I) | \frac{\|c_k\| \cdot \|b'_k\|}{\sqrt{k-1}}$$

Or nous avons vu au chapitre II de cette partie une majoration de $\|b_k^i\|$

$$\|b_k^i\| < \frac{S(n+n+1)}{n} \left\{ \|c_{k-1}\| \cdot \|b_{k-1}\| \left[\frac{\|b_{k-1}\|}{\sqrt{k-1}} (2N(A) + |\alpha_{k-1}|) + \frac{\|b_{k-2}\|}{\sqrt{k-2}} \right] \right. \\ \left. + \frac{|\beta_{k-2}|}{\sqrt{k-2}} \|c_{k-2}\| \cdot \|b_{k-2}\|^2 + N(A) \|b_{k-1}\| \right\}$$

Ceci nous permet donc d'écrire une majoration de $|\alpha'_{k+p}|$ et de $|\beta'_{k+p-1}|$ fonction uniquement de quantités soit connues, soit calculables :

Nous aurons alors :

$$|\alpha'_{k+p}| < \frac{S(n+n+1)}{n} \left\{ \|c_{k-1}\| \cdot \|b_{k-1}\| \left[\frac{\|b_{k-1}\|}{\sqrt{k-1}} (2N(A) + |\alpha_{k-1}|) + \frac{\|b_{k-2}\|}{\sqrt{k-2}} \right] \right. \\ \left. + \frac{|\beta_{k-2}|}{\sqrt{k-2}} \|c_{k-2}\| \cdot \|b_{k-2}\|^2 + N(A) \|b_{k-1}\| \right\} \left\{ 1 + \right. \\ \left. \max_i |\text{Det}(B_k^{k+p-1} - \lambda_i I)| + 2 \sum_{j=1}^{j=p-1} \max_i |\text{Det}(B_k^{k+j-1} - \lambda_i I)| \right\} \\ \cdot \left(\frac{\|c_{k+1}\|}{\sqrt{k}} + \frac{\|c_{k-1}\|}{\sqrt{k}} \right)$$

et :

$$|\beta'_{k+p-1}| < \max_i |\text{Det}(B_k^{k+p-1} - \lambda_i I)| \frac{S(n+n+1)}{n} \left\{ \|c_{k-1}\| \cdot \|b_{k-1}\| \right. \\ \left[\frac{\|b_{k-1}\|}{\sqrt{k-1}} (2N(A) + |\alpha_{k-1}|) + \frac{\|b_{k-2}\|}{\sqrt{k-2}} \right] + \frac{|\beta_{k-2}|}{\sqrt{k-2}} \|c_{k-2}\| \cdot \|b_{k-2}\|^2 \\ \left. + N(A) \|b_{k-1}\| \right\} \frac{\|c_k\|}{\sqrt{k-1}}$$

3°- Vérifications :

Le principe consiste à ajouter au vecteur b_k un vecteur b_k^i notablement supérieur au vecteur d'erreur due au calcul.

Malheureusement il est à peu près impossible de calculer effectivement $|\alpha'_{k+p}|$. Ce qui fait que pratiquement seule l'erreur faite au premier pas pourra être vérifiée.

J'ai les formules

$$\alpha'_k = \frac{c_{k+1} b_k^j}{\sqrt{k}} + \frac{c_{k-1} b_k^j}{\sqrt{k-1}}$$

et

$$\beta'_{k-1} = \frac{c_k b_k^j}{\sqrt{k-1}}$$

Les produits scalaires $\sqrt{\quad}$ sont calculés au cours de la méthode de Lanczos. Un programme accessoire m'a permis de les extraire et de faire les produits scalaires : $\langle c_{k+1}, b_k^j \rangle$; $\langle c_k, b_k^j \rangle$; $\langle c_{k-1}, b_k^j \rangle$;

Avec une matrice 8x8 dont j'avais altéré le vecteur b_3 en y ajoutant le vecteur $b_3'' = (0,1 ; 0,1 ; \dots 0,1)$

Vecteur erreur extrêmement important puisque les composantes de b_3 étaient comprises entre 0,5 et 2,5 (en valeurs absolues). J'ai obtenu les résultats suivants :

$$\begin{aligned} \alpha_3 &= -0,5 \\ \alpha_3 \text{ modifié} &= -0,512 \text{ II}6 \text{ 47} \\ d' \text{ ou } d_3^j &= -0,012 \text{ II}6 \text{ 47} \neq 0,012 \end{aligned}$$

d'autre part :

$$\begin{aligned} c_4 b_3^j &= -93,511 \text{ 485} \\ c_2 b_3^j &= 3,149 \text{ 999 94} \\ \tau_2 &= 200 \\ \sqrt{3} &= 3379,699 \text{ 92} \end{aligned}$$

donc on doit avoir

$$\begin{aligned} \alpha'_3 &= \frac{-93,511 \text{ 485}}{3379,699 \text{ 92}} + \frac{3,149 \text{ 999 94}}{200} \\ &= -0,011 \text{ 913 576} \neq 0,012 \end{aligned}$$

Pour β_2 j'obtiens

$$\begin{aligned} \beta_2 &= 16,673 \text{ 266} \\ \beta_2 \text{ modifié} &= 16,731 \text{ 187} \end{aligned}$$

$$d' \text{ ou } \beta_2' = 0,057 \text{ 921} \neq 0,058$$

$$\text{or } c_3 b_3'' = 11,7$$

$$\text{et } \sqrt{2} = 200$$

$$d' \text{ ou } \beta_2' = \frac{11,7}{200} = 0,0585 \neq 0,058$$

Nous voyons donc que l'erreur ne diffère que de très peu de la valeur prévue. Il faut noter d'ailleurs que nous sommes dans un cas qui ne se produira jamais en réalité, en effet on n'obtiendra jamais d'erreur aussi importante que celle introduite ici. Or l'étude faite suppose la linéarité des erreurs, celle-ci sera vérifiée bien plus étroitement dans les cas réels qu'ici ; ce qui permet d'estimer que la formule d'erreur trouvée est satisfaisante.

Troisième Partie

Expériences Numériques

Chapitre I

Détail du programme.

Un programme destiné à mettre en oeuvre la méthode de Lanczos a été écrit sur le calculateur du laboratoire de calcul de l'Université de Grenoble, un Gamma ET Bull. Je n'entrerai pas ici dans le détail de ce programme il me suffira d'indiquer qu'il utilise des nombres écrits en point décimal flottant (P.D.F.) en simple précision, c'est à dire, avec une mantisse de 9 chiffres et un exposant pouvant varier de -50 à + 49. Dans les opérations les troncatures sont effectuées sans arrondis.

Le programme écrit permet de normaliser les vecteurs obtenus et de réorthogonaliser les vecteurs b_k (resp c_k) aux vecteurs c_j (resp b_j) ; $j < k$. Ces deux compléments ont été proposés par Rutishauser et repris depuis par les différents auteurs.

La réorthogonalisation a pour but de maintenir la propriété d'orthogonalité, propriété essentielle, puisque c'est sur elle que repose toute la méthode.

La normalisation, elle, ne vise pas à diminuer les erreurs, mais comme il a été constaté, qu'en général, les modules des vecteurs avaient tendance à diminuer il s'agit d'éviter les pertes de précision dues aux exposants trop faible.

Voici plus précisément de quoi il s'agit :

Réorthogonalisation :

A chaque pas on remplace le vecteur b_k calculé par b_k^* défini par

$$b_k^* = b_k - \sum_{j=1}^{k-1} \gamma_{k,j} b_j \quad \text{avec} \quad \gamma_{kj} = \frac{c_j b_k}{\sigma_j}$$

Chapitre II

Essais

Nous avons utilisé le programme précédent pour étudier expérimentalement trois problèmes différents :

- L'utilité de la réorthogonalisation
- Le choix des vecteurs de départ
- L'influence des valeurs caractéristiques de la matrice étudiée.

1°- Matériel expérimental.

Pour faire ces expériences j'ai utilisé un programme de génération de matrices ayant des vecteurs et des valeurs propres connus, que je dois à l'obligeance de Monsieur Gastinel.

Ce programme fabrique des matrices ayant les vecteurs propres suivants :

$$\begin{array}{|c|c|c|c|c|c|}
 \hline
 \beta & I & I & & I & I \\
 \hline
 I & \beta & I & & I & I \\
 \hline
 I & I & \beta & & I & I \\
 \hline
 I & I & I & & \vdots & \vdots \\
 \hline
 \vdots & \vdots & \vdots & & \vdots & \vdots \\
 \hline
 \vdots & \vdots & \vdots & & I & I \\
 \hline
 \vdots & \vdots & \vdots & & \beta & I \\
 \hline
 I & I & I & & I & \beta \\
 \hline
 \end{array}$$

et des valeurs propres quelconques.

Le coefficient β est à la disposition de l'utilisateur, si on pose

$$B = \begin{array}{|c|}
 \hline
 \beta I I \dots \dots \dots I \\
 \hline
 I \beta I \dots \dots \dots I \\
 \hline
 \vdots & & I \\
 \hline
 \vdots & & \vdots \\
 \hline
 \vdots & & I \\
 \hline
 \vdots & & \vdots \\
 \hline
 \vdots & & \vdots \\
 \hline
 \vdots & & I \\
 \hline
 I \dots \dots \dots I \beta \\
 \hline
 \end{array}$$

B est la matrice des vecteurs propres et alors on a :

$$B^{-1} = \begin{vmatrix} \alpha' + \beta' & \beta' & \beta' & \dots & \beta' \\ \beta' & \alpha' + \beta' & \beta' & \dots & \beta' \\ \vdots & \beta' & \alpha' + \beta' & \dots & \beta' \\ \vdots & \beta' & \beta' & \dots & \alpha' + \beta' \\ \beta' & \beta' & \beta' & \dots & \beta' \end{vmatrix}$$

avec $\alpha' = \frac{1}{\beta - 1}$ et $\beta' = - \frac{1}{(\beta - 1)(\beta - 1 + n)}$

n étant, bien entendu l'ordre de la matrice.

B^{-1} est la matrice des lignes propres.

Les vecteurs propres ainsi formés sont régulièrement répartis, ils forment un parapluie dont β caractérise l'ouverture.

si $\beta = 1$, ils sont tous confondus ; si $\beta = \infty$ ils sont tous orthogonaux.

Donc, dans notre matrice nous aurons à notre disposition les valeurs propres et le parapluie des vecteurs propres.

J'ai effectué mes essais sur des matrices 15x15

Interprétation des résultats.

On pourrait résoudre le polynôme caractéristique obtenue, malheureusement ceci est long et la résolution introduit des erreurs importantes.

J'ai préféré la méthode qui consiste à mesurer la précision par le nombre de chiffres exacts sur les coefficients du polynôme caractéristique. Cette méthode est critiquable : en effet, les études sur la résolution des polynômes montrent qu'à de faibles erreurs sur les coefficients ne correspondent pas obligatoirement de faibles erreurs sur les racines. De plus tous les coefficients ne sont pas également sensible aux erreurs.

A cela il faut ajouter ici les faits suivants : le nombre de chiffres exacts n'est pas absolument constant d'un coefficient à l'autre. D'autre part le problème se pose de connaître les coefficients

exacts. Or si les coefficients extrêmes peuvent être assez facilement calculés exactement à partir des racines il n'en est pas de même pour les coefficients centraux.

En pratique j'ai après quelques tâtonnements mis au point le procédé suivants : je calcule directement les coefficients extrêmes. De quelques résolutions faites dans les cas les plus favorables et de ces calculs exacts je tire une table des coefficients du polynôme avec, par exemple, 6 chiffres exacts. (Ce sont les chiffres communs aux différents calculs) Comme tous les autres calculs seront fait dans des cas moins favorables, cette table est en pratique suffisante et suffisamment sûre.

2°- Influence de la réorthogonalisation :

Pour réorthogonaliser le vecteur b_k on est amené à calculer $k-1$ produits scalaires, à faire $k-1$ divisions et à additionner $k-1$ vecteurs.

Ceci va demander beaucoup de temps. D'autre part le nombre d'opérations effectuées est grand et on commet des erreurs sur ces opérations. On peut alors se demander s'il y aura réellement amélioration.

On peut d'autant plus se le demander que la situation est beaucoup plus grave qu'il ne paraît à première vue.

En effet b_k devait être orthogonal à c_j et s'il n'y est pas, il diffère peu d'un vecteur orthogonal donc $(c_j, b_k) \neq \frac{0}{2}$.

Nous sommes là justement dans le cas où les erreurs de calculs sont très importantes.

Cependant, il faut remarquer que plus le calcul direct manquera de précision, meilleure sera la correction. En effet, plus b_k s'écartera de la normale à c_j plus la correction sera efficace. On peut donc prévoir que la correction sera en quelque sorte "auto-compensatrice" et qu'elle sera d'autant plus efficace que le vecteur à corriger sera moins bon.

Le but des calculs effectués dans ce cadre devait nous permettre de vérifier ces hypothèses :

Tout d'abord j'ai constaté que la réorthogonalisation double approximativement le temps de calcul, son coût est donc loin d'être négligeable.

En ce qui concerne son efficacité, je n'ai pas constaté de cas où elle n'apporte pas une amélioration. Toutefois, dans les cas que je qualifierai, de favorables (il s'agissait dans ces cas de matrices ayant leurs valeurs propres bien distinctes, et d'une étude menée avec des vecteurs de départ bien choisis) l'amélioration était très faible.

Voici par exemple, le résultat obtenu avec une matrice ayant comme valeurs propres : 1 ; -2 ; 3 ; -4 ; ... ; -14 ; 15, et un $\beta=2$. Les vecteurs de départ étant choisis pour donner de bons résultats. Le polynôme caractéristique trouvé en utilisant la réorthogonalisation est :

$$\begin{array}{r}
 x^{15} - 8,000\ 000\ 24\ x^{14} \\
 - 587,999\ 997\ x^{13} + 4255,999\ 99\ x^{12} \\
 + 132\ 901, 995\ x^{11} - 853\ 103, 935\ x^{10} \\
 - 14\ 661\ 123, 2\ x^9 + 31\ 223\ 125, 3\ x^8 \\
 + 323\ 730\ 733\ x^7 - 3\ 211\ 351\ 630\ x^6 \\
 - 23\ 033\ 632\ 600\ x^5 + 32\ 303\ 254\ 300\ x^4 \\
 + 273\ 343\ 303\ 000\ x^3 - 629\ 324\ 753\ 000\ x^2 \\
 - 943\ 550\ 132\ 000\ x + 1\ 307\ 674\ 363\ 000
 \end{array}$$

Les erreurs relatives sur les coefficients sont, au maximum de $5 \cdot 10^{-3}$, ainsi, le deuxième coefficient est exactement : -8 ; le troisième : - 533 ; le dernier 1 307 674 363 000

Le même essai, mais sans réorthogonalisation, nous a donné les résultats suivants :

$$\begin{array}{r}
 x^{15} - 8,000\ 000\ 57\ x^{14} \\
 - 0,587\ 999\ 973 \cdot 10^3\ x^{13} + 0,425\ 600\ 013 \cdot 10^4\ x^{12} \\
 + 0,132\ 901\ 994 \cdot 10^6\ x^{11} - 0,853\ 104\ 036 \cdot 10^6\ x^{10} \\
 - 0,146\ 611\ 232 \cdot 10^8\ x^9 + 0,312\ 231\ 322 \cdot 10^8\ x^8 \\
 + 0,323\ 730\ 739 \cdot 10^9\ x^7 - 0,331\ 135\ 210 \cdot 10^{10}\ x^6 \\
 - 0,230\ 336\ 326 \cdot 10^{11}\ x^5 + 0,323\ 032\ 635 \cdot 10^{11}\ x^4 \\
 + 0,273\ 343\ 310 \cdot 10^{12}\ x^3 - 0,629\ 324\ 960 \cdot 10^{12}\ x^2 \\
 - 0,943\ 550\ 113 \cdot 10^{12}\ x + 0,130\ 767\ 507 \cdot 10^{13}
 \end{array}$$

Comme on peut le voir le nombre de chiffres exacts passe en gros de 7 à 6.

Si on essaye la même matrice, mais avec un vecteur b_I très voisin du vecteur propre associé à la valeur propre de plus grand module et un vecteur c_I voisin de la ligne propre associée à la même valeur propre ; ce qui, comme nous le verrons ci-après, est un très mauvais choix de vecteurs de départ. Alors ; avec la réorthogonalisation, la précision des résultats n'est pas modifiée.

Par contre sans réorthogonalisation, on obtient de très mauvais résultats :

$- 0,588\ 273\ 272 \cdot 10^3$ $+ 0,138\ 037\ 808 \cdot 10^6$ $- 0,146\ 853\ 319 \cdot 10^8$ $+ 0,830\ 765\ 925 \cdot 10^{10}$ $- 0,281\ 143\ 558 \cdot 10^{11}$ $+ 0,274\ 430\ 059 \cdot 10^{12}$ $- 0,952\ 024\ 403 \cdot 10^{12}$	$x^{15} - 0,795\ 209\ 316 \cdot 10^1 x^{14}$ $x^{13} + 0,422\ 775\ 314 \cdot 10^4 x^{12}$ $x^{11} - 0,846\ 744\ 653 \cdot 10^6 x^{10}$ $x^9 + 0,805\ 377\ 937 \cdot 10^8 x^8$ $x^7 - 0,377\ 440\ 640 \cdot 10^{10} x^6$ $x^5 + 0,318\ 619\ 460 \cdot 10^{11} x^4$ $x^3 - 0,630\ 371\ 731 \cdot 10^{12} x^2$ $x + 0,129\ 150\ 227 \cdot 10^{13}$
--	--

Ainsi qu'on peut s'en rendre compte il n'y a plus, en moyenne, que deux chiffres exacts.

De nombreux essais ont tous donné le même résultat, à savoir : influence très faible de la réorthogonalisation dans les cas favorables, mais par contre amélioration substantielle dans les autres cas.

3°- Influence du choix du vecteur de départ (sans réorthogonalisation)

Il s'agissait surtout ici de voir s'il était intéressant de partir d'un vecteur voisin du vecteur propre associé à la valeur propre de plus grand module. En effet, un tel vecteur peut être obtenu facilement en multipliant un vecteur arbitraire par les puissances de la matrice.

Si nous supposons un vecteur décomposé suivant les vecteurs propres u_i de la matrice, nous aurons :

$$b = \sum_{i=1}^n a_i u_i \quad a_i \text{ étant la composante suivant le vecteur } u_i$$

$$A b = \sum_{i=1}^n a_i \lambda_i u_i$$

Appelons u_1 le vecteur propre associé à la valeur propre de plus grand module

$$\text{Si } |a_i| \ll |a_1| \quad \forall i \neq 1$$

$$\text{Comme } |\lambda_1| > |\lambda_i| \quad \forall i \neq 1 \quad \frac{|a_i| |\lambda_i|}{|a_1| |\lambda_1|} < \frac{|a_i|}{|a_1|}$$

Alors Ab est intérieur au cône centré sur u_1 et de génératrice b . La méthode de Lanczos va avoir à former une base avec des vecteurs tous intérieurs à ce cône, donc tous très voisins on conçoit qu'alors les erreurs vont être importantes.

L'expérience devait vérifier cette hypothèse. Dans la matrice considérée dans le 2^{ème} § le vecteur propre associé à la valeur propre de plus grand module

$$\text{était } \begin{pmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 2 \end{pmatrix}$$

et la ligne propre correspondante $(1, 1, \dots, -15)$.

$$\text{en partant avec un vecteur } b_1 = \begin{pmatrix} 1 \\ 1 \\ \cdot \\ \cdot \\ 2, 2 \end{pmatrix}$$

et un vecteur $a_1 = (1, 1, \dots, -14)$ On obtenait le polynôme indiqué à la fin du 2^{ème} § on le voit les résultats sont très mauvais.

Par contre en partant avec les vecteurs:

$$b_1 = \begin{pmatrix} 1 \\ 1 \\ \cdot \\ 1 \\ 3 \end{pmatrix} \quad C_1 = (1, 1, \dots, 1, -12)$$

qui sont relativement éloigné la précision était du même ordre qu'avec les vecteurs

$$\begin{array}{l} I \text{ et } (I, I, \dots, I) \\ | \\ I \\ \cdot \\ \cdot \\ | \\ I \end{array}$$

vecteurs qui, égaux à la somme des vecteurs propres se sont révélés donner de très bons résultats.

J'ai alors essayé de partir près d'un autre vecteur propre, en particulier de celui associé à la valeur propre de plus petit module. Je n'ai constaté alors aucune influence défavorable.

J'ai également essayé de partir avec un vecteur différant peu d'un vecteur orthogonal ou vecteur propre associé à la valeur propre de plus grand module (et d'un vecteur ligne équivalent). Un tel vecteur est perpendiculaire à un vecteur $A^k b$. Les résultats ont été absolument catastrophiques; j'ai essayé ceci pour plusieurs matrices, pour toutes j'ai sorti des polynômes caractéristiques ayant des coefficients sans aucun chiffre exact.

Tout au plus l'ordre de grandeur de l'exposant était-il conservé. Ce résultat peut s'expliquer par le fait que la composante sur le vecteur propre associé à la valeur propre de plus grand module est voisine de zero, or nous avons vu que si une composante suivant un vecteur propre est nulle, il s'agit d'un cas particulier et que l'algorithme ne peut s'achever.

Nous sommes donc très près d'un cas particulier c'est ce qui explique le résultat.

Un fait ressort de ces expériences: c'est l'influence très grande des composantes correspondantes au vecteur et à la ligne propres associés à la valeur propre de plus grand module.

Le résultat pratique étant qu'il ne faut partir ni suivant leurs direction, ni perpendiculairement.

4^{me} : Influence des valeurs propres voisines.

Soit une matrice ayant deux valeurs propres voisines: λ et $\lambda+d$.
 λ est associée au vecteur propre u_1 et $\lambda+d$ au vecteur propre u_2 .

Je pose $u = u_1 + u_2$ et j'appelle (u, u^I) une base de l'espace $\{u_1 \times u_2\}$ dont un des vecteurs de la base est justement $u = u_1 + u_2$.

J'appelle E_{n-1} l'espace

$\{u \times u_3 \times u_4 \times \dots \times u_n\}$
 et

En l'espace $\{u_1 \times u_2 \times u_3 \times \dots \times u_n\} = E_{n-1} \times u^I$

Je pose $b_1 = \sum_{i=1}^n u_i$ Ce que j'ai toujours le droit de faire

Alors

$$Ab_1 = \left(\sum_{i=1}^n \right) u_i + \lambda u + u_2 d \lambda$$

Or $u_2 d \lambda$ est un vecteur infiniment petit.

Et je peux écrire $u_2 d \lambda = k_1 u d \lambda + k_2 u^I d \lambda$.

Seul $k_2 u^I d \lambda$ n'appartient pas à E_{n-1}

Donc la partie n'appartenant pas à E_{n-1} est un infiniment petit

je peux écrire:

$$b_2 = b_2^* + db_2$$

avec $b_2^* \in E_{n-1}$ et $d \notin E_{n-1}$ (u^I posons $b_2^* = \sum_{i=1}^n (b_2^*)_i u_i$

$$Ab_2 = \sum_{i=1}^n \lambda_i (b_2^*)_i u_i + \lambda \left[(b_2^*)_i + (b_2^*)_2 \right] u + \text{infiniments petits}$$

De même à chaque pas: la composante selon u^I sera infiniment petite.

Au pas $n-1$ nous aurons déterminé dans E_{n-1} $n-1$ vecteurs indépendants: la composante de b_n appartenant à E_{n-1} sera nulle, seule celle selon u^I sera non nulle.

Or cette composante est très petite par rapport aux vecteurs des pas précédents. Les erreurs faites aux pas précédents ont été déterminées par le module des vecteurs à ces pas. Ces erreurs sont

Ceci apparait parfaitement aux essais :

Les matrices essayées avaient comme valeurs propres

1; -2; 3; -4; 5; 5,01; -6; 7; -8; 9; -10; 11; -12; 13; -14.

Alors que les différents coefficients σ_k et γ_k sont de l'ordre de quelques unités. γ_{n-1} et σ_n sont de l'ordre du centième; donc la factorisation apparait nettement.

Quand au polynôme caractéristique alors qu'il devrait être

$$\begin{aligned} & x^{15} + 0,199\ 000 \cdot 10 \ x^{14} - 0,513\ 070 \cdot 10^3 \ x^{13} \\ - & 0,569\ 169 \cdot 10^3 \ x^{12} + 0,103\ 042 \cdot 10^6 \ x^{11} + 0,266\ 852 \cdot 10^5 \ x^{10} \\ - & 0,998\ 679 \cdot 10^7 \ x^9 + 0,487\ 846 \cdot 10^7 \ x^8 + 0,494\ 955 \cdot 10^9 \ x^7 \\ - & 0,539\ 470 \cdot 10^9 \ x^6 - 0,120\ 333 \cdot 10^{11} x^5 + 0,177\ 317 \cdot 10^{11} x^4 \\ + & 0,124\ 454 \cdot 10^{12} x^3 - 0,191\ 990 \cdot 10^{12} x^2 - 0,374\ 876 \cdot 10^{12} x \\ + & 0,436\ 763\ 239 \cdot 10^{12} \end{aligned}$$

(Je n'ai indiqué les coefficients qu'avec 6 chiffres car je n'ai pas pu obtenir une précision supérieure, sauf pour les coefficients extrêmes).

J'obtenais avec les vecteurs (1, 1,, 1) et

I
:
:
:
I

$$\begin{aligned} & x^{15} + 0,198\ 782 \cdot 10 \ x^{14} - 0,513\ 036 \cdot 10^3 \ x^{13} \\ - & 0,568\ 112 \cdot 10^3 \ x^{12} + 0,103\ 049 \cdot 10^6 \ x^{11} - 0,264\ 913 \cdot 10^5 \ x^{10} \\ - & 0,998\ 735 \cdot 10^7 \ x^9 + 0,489\ 539 \cdot 10^7 \ x^8 + 0,495\ 020 \cdot 10^9 \ x^7 \\ - & 0,540\ 193 \cdot 10^9 \ x^6 - 0,120\ 358 \cdot 10^{11} x^5 + 0,177\ 461 \cdot 10^{11} x^4 \\ + & 0,124\ 487 \cdot 10^{12} x^3 - 0,192\ 101 \cdot 10^{12} x^2 - 0,375\ 002 \cdot 10^{12} x \\ + & 0,436\ 957 \cdot 10^{12} \end{aligned}$$

Comme on peut le constater il n'y a plus, en moyenne, que 3 chiffres exacts par coefficient.

Avec la réorthogonalisation la précision monte à 6 chiffres (au lieu de 7, lorsque les valeurs propres sont toutes bien

séparées). Donc on voit, ainsi que je l'ai dit dans le paragraphe 2, que la réorthogonalisation est très utile dans les cas où la précision des résultats a tendance à diminuer.

S O M M A I R E

--:--:--:--

	Page :
Introduction	I
Bibliographie	3
Première Partie : Remarques Théoriques	4
Chapitre I : Présentation élémentaire de la méthode.	5
Chapitre II: Etude des itérés d'un vecteur dans le cas de diviseurs élémentaires non linéaires.	10
Chapitre III : Etude de la formation des vecteurs successifs de la méthode de Lanczos	15
Deuxième Partie : Etude des erreurs sur un pas	27
Chapitre I : Erreurs dans I produit scalaire en point décimal flottant	28
Chapitre II : Erreur à un pas de la méthode de Lanczos	40
Chapitre III : Etude de la propagation d'une erreur sur χ_k ou β_k	49
Chapitre IV : Propagation d'une erreur sur b_k	59
Troisième Partie : Expériences Numériques	67
Chapitre I : Détail du programme	68
Chapitre II : Essais	70