



**HAL**  
open science

# 1ère thèse: Matrices du second degré et normes générales en analyse numérique linéaire

Louis Noël Gastinel

## ► To cite this version:

Louis Noël Gastinel. 1ère thèse: Matrices du second degré et normes générales en analyse numérique linéaire: 2ème thèse: Le théorème de Stone - Weirstrass. Modélisation et simulation. Université Joseph-Fourier - Grenoble I, 1960. Français. NNT: . tel-00005173

**HAL Id: tel-00005173**

**<https://theses.hal.science/tel-00005173>**

Submitted on 1 Mar 2004

**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 S  
présentées  
A LA FACULTE DES SCIENCES  
DE L'UNIVERSITE DE GRENOBLE

Pour obtenir  
Le grade de Docteur ès Sciences Mathématiques

Par  
Noël GASTINEL  
Agrégé de l'Université

1° THESE

MATRICES DU SECOND DEGRE ET NORMES GENERALES EN ANALYSE  
NUMERIQUE LINEAIRE

2° THESE

Propositions données par la Faculté

- LE THEOREME DE STONE - WEIRSTRASS

Thèse soutenue le 12 décembre 1960

devant la Commission d'examen

M. CHABAUTY , Président

MM. DE POSSEL  
KUNTZMANN  
HACQUES } examinateurs



Je tiens à assurer de ma respectueuse reconnaissance :

Mon maître, Monsieur le professeur J. KUNZEMANN, directeur du laboratoire de calcul de l'Université de GRENOBLE, sous la direction du quel ce travail a été effectué, pour la bienveillante attention et le soutien constant qu'il m'a accordés, et surtout pour le goût du Calcul qu'il ne cesse de donner à tous ses élèves,

Monsieur le professeur DE POSSEL, de l'Université de PARIS, qui m'a fait l'honneur de se joindre à mon jury ,

Monsieur le professeur CHABAUTY, qui a bien voulu en accepter la présidence,

Monsieur HACQUES, maître de conférences à la faculté des sciences de GRENOBLE.

Je remercie tous mes camarades du Laboratoire de calcul de l'Université de GRENOBLE pour l'aide inappréciable qu'ils m'ont apportée lors de l'élaboration de ce travail.



1° THESE

MATRICES DU 2° DEGRE

&

NORMES GENERALES

EN ANALYSE NUMERIQUE LINEAIRE



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

Le Calcul Numérique présente deux aspects apparemment distincts :  
la partie consacrée à la recherche de méthodes ou de procédés algorithmiques,  
partie algébrique ;  
la partie consacrée au contrôle des calculs et des approximations, partie  
topologique.

Sur le plan théorique, l'Analyse Numérique présuppose donc  
l'algèbre et la topologie. La raison profonde en est la définition même des  
nombres réels en laquelle ce double point de vue se fait déjà sentir.

Ce travail se propose d'utiliser deux idées dirigées dans les  
deux voies dont nous venons de parler. La première est de montrer le rôle  
important joué dans la plupart des processus de calcul linéaire par des  
matrices de polynômes minimaux du second degré, dont l'inverse est, par ce  
fait, immédiatement calculable. La seconde est de montrer comment la  
considération de normes générales peut permettre de définir des notions impor-  
tantes comme les conditionnements numériques et d'obtenir un contrôle plus  
fin des erreurs de calcul.

Les deux premiers chapitres présentent le "matériel" théorique  
utilisé. Le chapitre I contient, après la détermination des matrices du  
2° degré, un certain nombre de propositions prouvant que, par des produits  
à gauche par de telles matrices particulièrement simples, l'on peut trian-  
gulariser, et même, en général, diagonaliser une matrice quelconque.

(algorithmes de GAUSS et de JORDAN). Le chapitre II présente les notions  
de normes de vecteur et de matrice. Après un rappel de résultats classiques,  
dont certains sont dus à OSTROWSKI, nous avons cherché à caractériser les  
normes symétriques de matrices associées à des normes de vecteurs. Il  
résulte des théorèmes (en particulier du théorème VI) que ces normes sont  
les plus "serrées" des normes multiplicatives possibles : donc les plus  
utilisables en calcul. Il est signalé deux caractérisations : l'une  
(théorème V) par l'intermédiaire d'une double association de norme de  
vecteur et de matrice, l'autre (théorème VII) par la propriété, pour ces  
normes, d'être les normes multiplicatives minimales.

Dans le chapitre III, se trouve définie la notion générale de conditionnement numérique. L'idée en est très simple : on veut trouver la solution d'une équation, c'est-à-dire trouver une inconnue  $x$  afin qu'une expression  $T(x)$  soit nulle. Puisque, en calcul numérique, et quel que soit le moyen de calcul utilisé effectivement, toute quantité inférieure à la capacité minimum de ce moyen est considérée comme nulle. La question essentielle est de savoir si le fait qu'une norme de  $T(x)$  soit faible implique la proximité de  $x$  et de la solution théorique. Le conditionnement général est défini comme rapport de deux constantes qui figurent dans le théorème d'équivalence de normes sur un espace vectoriel topologique. Géométriquement, il apparaît comme un coefficient de forme, pour le lieu des  $x$  tels que une norme de  $T(x)$  ait une valeur donnée. On étudie ensuite les conditionnements de système d'équations linéaires. Des comparaisons des différents types de conditionnements entre eux sont établies. Le théorème III de ce chapitre établit l'équivalence de cette définition et de celle de TURING sur les variations relatives de la solution pour des variations des coefficients.

Un conditionnement particulier  $C(A)$  et son complémentaire  $C^0(A)$  sont étudiés au chapitre IV. Leur définition et leurs propriétés géométriques montrent que ces nombres, toujours compris entre zéro et un, ne sont égaux à un que si la matrice  $A$  est orthogonale en lignes. Ils peuvent donc servir à la définition d'un système "bien" ou "mal" conditionné d'une façon plus intrinsèque que les conditionnements généraux. On montre comment les calculer et l'on étudie ensuite les opérations qui les laissent invariants. Des exemples de calcul de conditionnement  $C(A)$  sont donnés : matrice d'HILBERT, matrice d'interpolation d'opérateur différentiel. Dans ce dernier cas, on prouve un résultat bien souvent constaté en calcul : si l'on remplace un problème différentiel aux limites par un système linéaire dû à une interpolation de l'opérateur différentiel, plus l'on augmente la précision en diminuant le pas, plus le système linéaire voit son conditionnement normalisé diminuer.

L'étude des erreurs dans la résolution des systèmes linéaires par élimination est faite dans le chapitre V. Nous établissons la forme de la "matrice d'erreur" en utilisant les matrices du 2<sup>o</sup> degré que l'on a déterminées au chapitre I, aussi bien pour un calcul en virgule fixe qu'en virgule flottante. La matrice d'erreur pour l'inversion en virgule fixe d'une matrice est aussi déterminée. Le théorème VIII indique, en fonction des "matrices d'erreur" et du "second membre d'erreur" la forme générale de la double inégalité à laquelle satisfait une norme de l'erreur absolue, ce qui met en évidence le rôle important joué par le conditionnement de la matrice du premier membre.

Le théorème IX indique une condition suffisante pour choisir la capacité d'un calcul en virgule fixe, afin d'obtenir une précision désirée sur la solution. La forme de la matrice d'erreur montre comment il faut choisir les différents "pivots" dans une résolution en virgule fixe. Des expériences de calcul viennent confirmer cette règle. En virgule flottante, l'on obtient des résultats analogues. Enfin, nous indiquons comment faire en sorte qu'une élimination ne fasse pas diminuer le conditionnement d'un système.

Le chapitre VI traite de la méthode de résolution par orthogonalisation. Une forme de matrice d'erreur est encore établie. Une règle de choix des différentes variantes dans ces calculs est énoncée.

Par l'intermédiaire de la notion de décomposition d'une norme, nous montrons, au chapitre VII, que l'on peut trouver différents procédés de résolution itératifs d'un système linéaire dont la convergence est directement fonction des conditionnements étudiés. L'un est original, les deux autres sont des variations sur les méthodes de gradient, très fréquemment utilisées.

Le chapitre VIII est consacré à l'étude des éléments caractéristiques des matrices. En utilisant toujours des matrices du 2<sup>o</sup> degré, nous obtenons une forme particulièrement simple de transmutation. Deux méthodes de détermination du polynôme caractéristique sont ensuite exposées. La première est basée sur la formule (I) du chapitre I, § 3-d, et permet d'obtenir, en général, une matrice semblable à la matrice donnée de la forme de FROBENIUS. La seconde recherche une matrice semblable, de forme triple-diagonale. Enfin, l'on indique comment obtenir une matrice d'erreur, pour la première méthode, par un raisonnement assez analogue à celui qui est fait au chapitre V. Nous terminons en indiquant comment exploiter cette matrice pour obtenir les variations correspondantes des éléments caractéristiques de la matrice donnée.

---



CHAPITRE I

Matrices du 2° degré

I°) Anneaux de matrices du 2° degré

Définition : On dit qu'une matrice est du 2° degré si son polynôme minimal  $m(u)$ , est du 2° degré; nous écrivons ce polynôme :

$$m(u) = u^2 - pu - q .$$

A partir d'une matrice  $K$  du 2° degré, de type  $(n,n)$ , à éléments appartenant au corps  $C$  des nombres complexes, de polynôme minimal  $m(u)$ , formons l'ensemble  $\mathcal{A}_K$  des matrices  $Z = aI + bK$  avec  $a, b \in C$ ,  $I$  unité de type  $(n,n)$ .

Théorème I :  $\mathcal{A}_K$  est un anneau de matrices. Une condition nécessaire et suffisante pour que  $Z$  soit inversible dans cet anneau est que, pour  $b \neq 0$ ,  $m(\mu) \neq 0$ , si  $-\mu = a/b$ . L'inverse de  $Z = aI + bK$  est alors  $Z' = a'I + b'K$  avec :

$$(1) \begin{cases} a' = \frac{1}{b} & \frac{p-\mu}{m(\mu)} \\ b' = -\frac{1}{b} & \frac{1}{m(\mu)} \end{cases}$$

En effet, si  $Z = aI + bK, Z' = a'I + b'K, Z + Z' = (a+a')I + (b+b')K$  prouve que l'addition dans  $\mathcal{A}_K$  est une loi de groupe commutatif.

$$Z \cdot Z' = (aI + bK) \cdot (a'I + b'K) = aa'I + (ab' + ba')K + bb'K^2 = (aa' + qbb')I + (ab' + ba' + pbb')K$$

puisque  $K^2 = pK + qI$ , qui est bien élément de  $\mathcal{A}_K$ .

La matrice  $Z$  est donc inversible si et seulement si

$$\begin{aligned} aa' + qbb' &= 1 \\ ab' + ba' + pbb' &= 0 \end{aligned}$$

c'est à dire pour  $a, b$  tels que :

$$\text{Det} \begin{vmatrix} a & qb \\ b & a + pb \end{vmatrix} \neq 0$$

ou:  $a^2 + pab - qb^2 \neq 0$

Pour  $a, b, p, q$  donnés,  $a', b'$  se déterminent d'une manière unique par :

$$a' = \frac{a + pb}{a^2 + pab - qb^2}, \quad b' = \frac{-b}{a^2 + pab - qb^2}$$

Pour  $b \neq 0$  ces formules sont les formules (I).

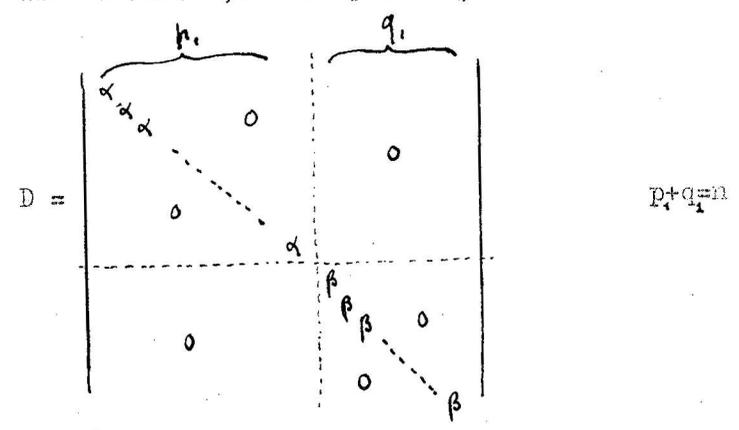
C'est la facilité du calcul de l'inverse d'une matrice  $Z$  qui fait l'importance de telles matrices en analyse numérique. Ces matrices sont construites à partir de matrices  $K$ , nous allons voir comment les déterminer.

2°) Détermination des matrices du 2° degré, de type  $(n, n)$  sur le corps  $C$

La théorie classique de la réduction à la forme de JORDAN [1], nous conduit à décomposer  $m(u) = u^2 - pu - q$  en  $m(u) = (u - \alpha)(u - \beta)$  dans  $C$ .

Deux cas sont possibles:

a)  $\alpha \neq \beta$ . Le polynôme minimal n'ayant que des facteurs irréductibles du premier degré, la matrice  $K$  de polynôme minimal  $m(u)$  n'est autre que la transmuée par une matrice  $T$ , non singulière, d'une matrice diagonale  $D$  de la forme:

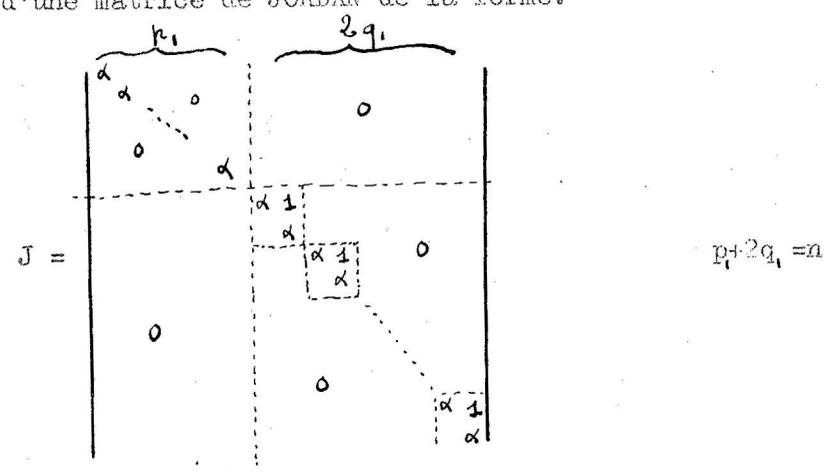


Toute matrice  $K = T^{-1} D T$  convient.

b)  $\alpha = \beta$ . Le polynôme minimal est  $m(u) = (u - \alpha)^2$ . La chaîne des diviseurs élémentaires sera de la forme:

$$\underbrace{(u - \alpha)^2, \dots, (u - \alpha)^2}_{q_1}, \underbrace{(u - \alpha), \dots, (u - \alpha)}_{p_1}$$

Par suite, la matrice  $K$  n'est autre que la transmuée par une matrice  $T$  non singulière d'une matrice de JORDAN de la forme:



Toute matrice  $K = T^{-1} J T$  convient.

Définition : On nommera matrice antiscaire une matrice :

$$K = X \cdot Y^T \quad X \text{ et } Y \text{ colonnes non nulles}$$

Théorème II : Toute matrice antiscaire est semblable à l'une des matrices :

$$D = \begin{vmatrix} \lambda & & & \\ & \circ & & \\ & & \ddots & \\ & & & \circ \end{vmatrix}, \lambda \neq 0 ; \quad J = \begin{vmatrix} \circ & 1 & & \\ & \circ & & \\ & & \ddots & \\ & & & \circ \end{vmatrix}$$

et réciproquement.

1°) Pour la propriété directe, on écrira :

$$K = X Y^T \quad K^2 = X(Y^T X) Y^T = \lambda X Y^T$$

puisque  $Y^T X = \lambda$  est un scalaire.

Il en résulte :  $K^2 - \lambda K = 0$

Soit U un vecteur colonne.  $K U = 0$  est vérifié si  $Y^T U = 0$ , donc pour des vecteurs formant un sous-espace de dimension n-1. Il en résulte que 0 est racine d'ordre au moins n-1 pour l'équation caractéristique. Si 0 est racine d'ordre n-1, la forme réduite de JORDAN est D. Si 0 est racine d'ordre n, la matrice est nulle (ce qui exige  $X = 0$  ou  $Y = 0$ , et se trouve donc exclu) ou bien la forme réduite de JORDAN est J.

2°) Réciproquement,

Supposons que K soit semblable à D ou J. Il existe une matrice T inversible telle que

$$K = T^{-1} D T \quad \text{ou} \quad K = T^{-1} J T.$$

Posons  $T = (t_{ij})$  et  $T^{-1} = (\tau_{ij})$  pour  $i, j = 1, \dots, n$

Alors :



Nous allons montrer quelques propriétés importantes des matrices  $X.Y^T$ .

3°) Utilisation de matrices  $X.Y^T$ , antisymétriques.

- a) Décomposition d'une matrice quelconque en une somme de matrices antisymétriques.

Soient  $e_i$  pour  $i=1, \dots, n$ , les vecteurs de la base fondamentale de  $C^n$ , c'est à dire les vecteurs-colonnes

$$e_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \leftarrow i^o \\ \vdots \\ 0 \end{pmatrix}$$

On voit que les matrices unitaires  $E_{ij}$ , base de l'espace à  $n^2$  dimensions des matrices carrées sur  $C$ ,

$$E_{ij} = \begin{pmatrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ i^o \rightarrow & & & 1 & & \\ & & & & & \\ & & & & & \end{pmatrix}$$

ne sont autres que les produits antisymétriques :

$$E_{ij} = e_i \cdot e_j^T$$

Soit alors  $u_k, k=1, \dots, n$ , un système de  $n$  vecteurs linéairement indépendants, puisqu'ils forment une base de  $C^n$ , on peut écrire:

$$e_i = \sum_{k=1}^n \beta_{ki} u_k$$

et, puisque si  $A$  est une matrice quelconque de terme général

$$(a_{ij}), \quad i, j=1, \dots, n, \quad A = \sum_{i,j} a_{ij} E_{ij}$$

$$D'où A = \sum_{i,j} a_{ij} \left( \sum_k \beta_{ki} u_k \right) \left( \sum_{k'} \beta_{k'j} u_{k'}^T \right) = \sum_{m,e} \gamma_{m,e} u_m \cdot u_e^T$$

Proposition I: Etant donné un système de  $n$  vecteurs linéairement indépendants,  $u_k, k=1, \dots, n$ , toute matrice peut être écrite sous la forme

$$A = \sum_{m,e=1, \dots, n} \gamma_{m,e} u_m \cdot u_e^T$$

c'est à dire sous la forme d'une somme de matrices de type  $X.Y^T$ .

Nous verrons ,à la fin de ce chapitre une proposition sur cette décomposition un peu moins générale.

b) L'algorithme d'élimination de GAUSS.

Proposition II: Etant donnée une matrice A quelconque, il existe une suite de matrices du 2° degré,  $G_c$ , en nombre fini, telles que le produit

$$T = G_1 G_2 G_3 \dots G_m A$$

soit une matrice triangulaire supérieure. On peut choisir les  $G_c$  de sorte que  $\text{Det}(A) = \text{Produit des éléments diagonaux de T}$ .

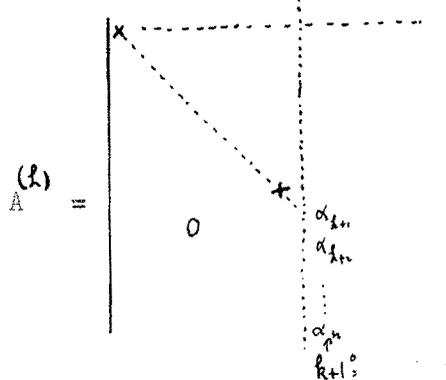
Nous allons utiliser des matrices de type  $Z_{ij} = I + \lambda E_{ij}$  avec  $\lambda \in \mathbb{C}$ ,  $E_{ij}$  antiscaire. Il est clair que se sont bien des matrices du 2° degré, dont l'inverse, pour  $i \neq j$ , n'est autre que  $Z_{ij}^{-1} = I - \lambda E_{ij}$ .

L'effet produit par multiplication à droite ou à gauche d'une matrice quelconque par  $Z$  est bien connu.

$M.Z$  = Matrice M où la j° colonne est remplacée par la somme de celle-ci et de la i°, multipliée par  $\lambda$ .

$Z.M$  = Matrice M où la i° ligne est remplacée par la somme de celle-ci et de la j°, multipliée par  $\lambda$ . [3]

Soit alors une matrice  $A^{(k)}$  de la forme ci-dessous:



c'est à dire ayant des 0 dans ses k premières colonnes et au dessous de la diagonale. Soient  $\alpha_{k+1}, \alpha_{k+2}, \dots, \alpha_n$ , les termes dans la  $k+1$ ° colonne sur la diagonale et au dessous.

a) Supposons  $\alpha_{k+1} \neq 0$ . Il est clair que si l'on considère

$$G_{k+1, k+2} = I - \frac{\alpha_{k+2}}{\alpha_{k+1}} E_{k+2, k+1}$$

$$\vdots$$
$$G_{k+1, n} = I - \frac{\alpha_n}{\alpha_{k+1}} E_{n, k+1}$$

le produit  $G_{k+1, n} \cdot G_{k+1, n-1} \dots G_{k+1, k+2} \cdot A^{(k)}$  est de la forme:  $A^{(k+1)}$

c'est à dire n'a que des 0 dans ses  $k+1$  premières colonnes et au dessous de la diagonale.









Le produit ne diffère de  $A^{(k)}$  que par sa  $k+1^{\circ}$  ligne.

Peut-on choisir les  $\lambda_{k+1,1}, \lambda_{k+1,2}, \dots, \lambda_{k+1,k}$ , de sorte que cette nouvelle  $k+1^{\circ}$  ligne soit orthogonale aux  $k$  premières ? Il suffit d'écrire:

$$h_{k+1} \cdot h_i^T = 0 \quad (i=1, 2, \dots, k)$$

et en tenant compte de l'orthogonalité des  $k$  premières lignes entr'elles, il reste:

$$a_{k+1,i} \cdot h_i^T + \lambda_{k+1,i} \cdot h_i \cdot h_i^T = 0$$

et si l'on suppose (matrices réelles) qu'aucune des  $k$  premières lignes n'est nulle on obtient:

$$\lambda_{k+1,i} = - \frac{a_{k+1,i} \cdot h_i^T}{h_i \cdot h_i^T} \quad (i=1, 2, \dots, k)$$

Soit alors une matrice  $A$  à éléments réels, non singulière,

$$A = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} = (a_{ij}) \quad (i, j = 1, \dots, n)$$

Posons  $A^{(1)} = A$ , et considérons la suite des matrices:

$$A^{(2)} = S_1 \cdot A^{(1)}, \dots, A^{(k)} = S_{k-1} \cdot A^{(k-1)}, \dots, A^{(n)} = S_{n-1} \cdot A^{(n-1)}$$

Il est clair que:

$$A^{(k)} = S_{k-1} \cdot S_{k-2} \cdot \dots \cdot S_1 \cdot A$$

Puisque les déterminants de toutes les matrices  $S_k$  sont égaux à 1 la non-singularité de  $A$  impose la non-singularité des  $A^{(k)}$ , pour tout  $k$  donc quels que soient les  $\lambda_{ij}$ , il est impossible qu'une des  $k$  premières lignes de  $A^{(k)}$  soit nulle, cela permet de voir que pour  $A$  non singulière le procédé suivant, dit de SCHMIDT, permet l'orthogonalisation:

Partant de  $A^{(1)} = A$  on détermine les  $\lambda_{ij}$  par

$$\lambda_{ij} = - \frac{a_i \cdot h_j^T}{h_j \cdot h_j^T}$$

il est clair que  $A^{(k)} = H$  est orthogonale en lignes.

D'où:

Proposition IV

Pour une matrice  $A$ , non singulière il existe une suite,  $S_1, \dots, S_{n-1}$ , de matrices du 2<sup>o</sup> degré, de la forme indiquée, telle que le produit

$$H = S_{n-1} \cdot S_{n-2} \cdot \dots \cdot S_1 \cdot A$$

soit orthogonal en lignes.

d) Transmutations par des matrices du 2° degré .

La plus part des problèmes de détermination numérique de valeurs propres ou de vecteurs propres sont résolus en transmutant la matrice donnée A en A' , avec A' = M . A . M<sup>-1</sup> et M non singulière . si l'on prend pour M, des matrices non simples le travail numérique est considérable et introduit une grosse quantité d'erreurs , d'où l'intéret de chercher des M et M<sup>-1</sup> pour les quelles il n'y a pas d'erreur de calcul.

Or si M = a . I + b . K ∈ Q<sub>n</sub>

M<sup>-1</sup> = a' I + b' K

avec :

a' = 1/b . (p - P) / m(p)

b' = -1/b . 1 / m(p)

Par suite:

A' = (a I + b K) . A . (a' I + b' K) = aa' H + a' b . KA + a b' AK + bb' K.A.K

et enfin :

(I) A' = (p^2 - P.P) / m(p) . A + (P - P) / m(p) . K.A + P / m(p) . A.K - 1 / m(p) . K.A.K

Première application : Matrices - Test :

Afin de vérifier, en calcul numérique, la qualité d'une méthode de détermination des valeurs et vecteurs propres , il est indispensable de traiter par cette méthode des matrices dont on connaisse à l'avance les valeurs propres et les vecteurs propres .

Les matrices antiscales permettent d'écrire

un programme capable de faire sortir , sous forme de données immédiatement ré-utilisables , les éléments d'une matrice dont on se donne à l'avance l'ordre n , la suite λ<sub>1</sub>, λ<sub>2</sub>, ..... λ<sub>n</sub> , des valeurs propres <sup>réelles, mais</sup> pas forcément distinctes , et dont on connaisse les colonnes propres et les lignes propres .

En effet, soit :

M = B . Λ . B<sup>-1</sup>





CHAPITRE I I

Normes de vecteurs et de matrices

I°) Normes de vecteurs. [4]

Etant donné un espace vectoriel  $\mathcal{E}$  sur  $\mathbb{R}$  (ou  $\mathbb{C}$ ), on appelle norme sur cet espace une application  $x \rightarrow p(x)$  de celui-ci dans  $\mathbb{R}_+$  telle que:

- I)  $p(x)=0$  est équivalent à  $x=0$
- II)  $p(\lambda x) = |\lambda|p(x)$  pour tout  $\lambda \in \mathbb{R}$  (ou  $\mathbb{C}$ )
- III)  $p(x+y) \leq p(x) + p(y)$

quels que soient les éléments  $x, y$  de cet espace vectoriel.

Soit donc sur l'espace vectoriel  $\mathcal{E}$  une norme  $p(x)$ . Donnons nous, de plus une application linéaire  $T(x)$  de  $\mathcal{E}$  dans lui-même que nous supposons complet

A  $x \in \mathcal{E}$  faisons correspondre:  $q(x)=p(T(x))$ ; c'est une application de  $\mathcal{E}$  dans  $\mathbb{R}_+$ . Deux questions se posent:

- a) Est-elle une norme sur  $\mathcal{E}$  ?
- b) Est-elle équivalente à  $p(x)$ ?

Proposition I Pour que  $q(x)$  soit une norme sur  $\mathcal{E}$ , il faut et il suffit que  $T(x)$  ne soit pas singulière (c'est à dire qu'il n'y ait pas d'éléments  $\neq 0$  de tels que  $T(x)=0$ ).

En effet à cause de la linéarité de  $T(x)$  on a:

II)  $q(\lambda x) = p[T(\lambda x)] = p[\lambda T(x)] = |\lambda| p[T(x)] = |\lambda| \cdot q(x)$

puis,

III)  $q(x+y) = p[T(x+y)] = p[T(x)+T(y)] \leq p[T(x)] + p[T(y)] \leq q(x) + q(y)$

cela prouve II et III pour  $q(x)$ .

La condition I résulte du fait que  $p(T(x))=0$  est équivalent à  $T(x)=0$ , et que cela ne peut être, à son tour, équivalent à  $x=0$  que si et seulement si  $T(x)$  n'est pas singulière.

Proposition II Pour que la norme  $q(x)$  soit équivalente à  $p(x)$ , il faut et il suffit que  $T(x)$  soit continue dans  $\mathcal{E}$ .

I) La condition est suffisante. Supposons que  $T(x)$  soit continue. Elle l'est pour 0 de  $\mathcal{E}$ . Si l'on se donne la boule  $B_q(\beta)$ ,

ensemble des  $x$  tels que  $q(x)=p(T(x)) \leq \beta$ , il existe  $\alpha$  tel que

$p(x) \leq \alpha \implies p(T(x)) \leq \beta$  ou  $B_q(\beta) \supset B_p(\alpha)$

donc la topologie définie par q n'est pas plus fine que celle définie par p.

Inversement, soit  $B_p(\alpha)$  donnée, je dis qu'il existe un  $\beta$  tel que si  $x \in B_q(\beta)$  cela implique  $x \in B_p(\alpha)$ . Il suffit de prouver un  $a$  positif tel que quelque soit  $x$  :  $p(x) \leq a q(x)$ . Or cela résulte d'un corollaire [5] d'un théorème de BANACH [6] indiquant que toute application linéaire continue et biunivoque d'un espace vectoriel normé et complet est un isomorphisme dans lui-même. Donc, si l'on prend  $\beta = \frac{\alpha}{a}$  On est sûr que  $B_q(\frac{\alpha}{a}) \subset B_p(\alpha)$ , la topologie définie par p est donc moins fine que celle définie par q. D'après cela les deux normes p et q sont bien équivalentes.

2) La condition est nécessaire. Si p et q sont équivalentes c'est qu'il existe deux constantes positives a et b, telles que [7] :

$$a p(x) \leq p(T(x)) \leq b p(x)$$

Si  $\varepsilon$  est donné, il suffit de prendre  $p(x) \leq \frac{\varepsilon}{b}$  pour assurer que  $p(T(x)) \leq \varepsilon$  d'où la continuité en 0, qui implique la continuité en tout point.

Proposition III. Si  $T(x)$  est une application linéaire continue et biunivoque d'un espace vectoriel normé et complet  $\mathcal{E}$  sur lui-même et  $p_1(x)$  et  $p_2(x)$  deux normes équivalentes sur  $\mathcal{E}$ . Il existe alors deux constantes a et b telles que quel que soit x

$$a p_1(x) \leq p_2(T(x)) \leq b p_1(x)$$

Il suffit de remarquer que par la proposition II les deux normes  $p_1(x)$  et  $p_2(T(x))$  sont équivalentes.

Cette propriété est la base théorique de la définition des conditionnements numériques que nous étudierons plus loin.

Après avoir rappelé des propriétés relatives à des espaces vectoriels généraux, nous nous intéresserons dans la suite à des espaces de dimensions finies.

En analyse numérique l'espace  $\mathcal{E}$  est toujours  $R^n$  (ou  $C^n$  à la rigueur) et trois normes de vecteurs sont à envisager:

Si  $x \in R^n$  a pour éléments  $(\xi_i)$  et l'on pose  $\varphi_2(x) = \left( \sum_{i=1}^n |\xi_i|^2 \right)^{\frac{1}{2}}$

(norme de HOLDER pour  $\infty$ )

I) La norme  $\varphi_1(x) = \Phi(x) = \sum_{i=1}^n |\xi_i|$

II) La norme  $\varphi_2(x)$  ou norme "Euclidienne"  $\varphi_2(x) = \|x\| = \left( \sum_{i=1}^n |\xi_i|^2 \right)^{\frac{1}{2}} = (x^T \cdot x)^{\frac{1}{2}}$

III) La norme  $\varphi_\infty(x) = \text{Max}_{i=1, \dots, n} |\xi_i| = M(x)$

Pour ces trois normes il est important de connaître les nombres  $a$  et  $b$  tels que :  $(a \cdot \varphi_i(x) \leq \varphi_j(x) \leq b \cdot \varphi_i(x))$  les voici :

$$I) \quad \|x\| \leq \Phi(x) \leq \sqrt{n} \|x\| \quad , \quad \mathcal{M}(x) \leq \Phi(x) \leq n \cdot \mathcal{M}(x)$$

$$II) \quad \frac{1}{\sqrt{n}} \Phi(x) \leq \|x\| \leq \Phi(x) \quad , \quad \mathcal{M}(x) \leq \|x\| \leq \sqrt{n} \cdot \mathcal{M}(x)$$

$$III) \quad \frac{1}{n} \Phi(x) \leq \mathcal{M}(x) \leq \Phi(x) \quad , \quad \frac{1}{\sqrt{n}} \|x\| \leq \mathcal{M}(x) \leq \|x\|$$

Dans ces inégalités il est toujours un choix pour  $x$  qui assure une égalité.

Enfin, signalons qu'au point de vue géométrique, donner une norme  $p(x)$  de vecteurs sur  $\mathbb{R}^n$  ou sur  $\mathbb{C}^n$ , revient à se donner un corps convexe, la boule  $B_p(a) = \{x; p(x) \leq a\}$ , symétrique par rapport à l'origine et telle que toute demi-droite issue de 0 la coupe en un point distinct de 0 (0 est intérieur à  $B_p(a)$ ). Toutes ces boules, si  $a$  varie, sont homothétiques de l'une  $B_p(1)$  (par exemple), dont la connaissance suffit pour définir la norme [43]

### 28) Normes de matrices

Soit  $\mathcal{M}_n$  l'espace vectoriel à  $n^2$  dimensions des matrices carrées sur  $\mathbb{R}$  (ou  $\mathbb{C}$ ). Une norme générale est, comme on l'a dit plus haut une application  $N(A)$  de  $\mathcal{M}_n$  dans  $\mathbb{R}_+$  telle

$$G_N \quad \begin{cases} I) N(A)=0 \text{ équivaut à } A=0 \\ II) N(\lambda A) = |\lambda| N(A) \\ III) N(A+B) \leq N(A)+N(B) \end{cases}$$

quels que soient  $\lambda \in \mathbb{R}$ ;  $A, B \in \mathcal{M}_n$ .

Une Norme multiplicative, d'après OSTROWSKI [9] est une application  $M(A)$  de  $\mathcal{M}_n$  dans  $\mathbb{R}_+$  telle

$$G_M \quad \begin{cases} I) \text{ Il existe un } A_0 \text{ tel que } M(A_0) \neq 0 \\ II) M(\lambda A) = |\lambda| M(A) \\ III) M(A+B) \leq M(A)+M(B) \\ IV) M(A \cdot B) \leq M(A) \cdot M(B) \end{cases}$$

Notations :  $\mathcal{N}$  désigne l'ensemble des normes sur  $\mathcal{M}_n$   
 $\mathcal{NM}$  désigne l'ensemble des normes multiplicatives sur  $\mathcal{M}_n$   
 $\mathcal{V}$  désigne l'ensemble des normes de vecteurs sur  $\mathbb{R}$  (ou  $\mathbb{C}$ )

Théorème I Toute norme multiplicative est une norme générale dans  $\mathcal{M}_n$ . [8]

- a) Soit I l'unité de  $\mathcal{M}_n$  on a :  $M(A \cdot I) \leq M(A) \cdot M(I)$  (IV)  
 puisque  $M(A) > 0$ , cela implique  $M(I) \geq 1$
- b) Si  $E_{ij}$  sont les matrices unitaires (telles que :  $A = \sum_{i,j} \alpha_{ij} E_{ij}$ )  
 on a  $I = \sum_{i=1}^n E_{ii}$   
 donc  $0 < 1 \leq \sum_i M(E_{ii})$  : il est au moins  $i_0$  tel  $M(E_{i_0 i_0}) > 0$
- c) Mais  $E_{i_0 i_0} = E_{i_0 j} E_{jj} E_{j i_0}$  donc :  $M(E_{i_0 i_0}) \leq M(E_{i_0 j}) \cdot M(E_{jj}) \cdot M(E_{j i_0})$  qq  $j$   
 Par suite tous les  $E_{jj}$  sont tels que :  $M(E_{jj}) > 0$
- d) Mais  $E_{jj} = E_{ji} E_{ii} E_{ij}$  ;  $M(E_{jj}) \leq M(E_{ji}) M(E_{ii}) M(E_{ij})$  donc quelques soient les  $i, j$ ,  $M(E_{ij}) > 0$   
 Or si  $A = (\alpha_{ij})$ , il est clair que :  $\alpha_{ij} E_{jk} = E_{ik} \cdot A \cdot E_{jk}$   
 donc  $|\alpha_{ij}| \cdot M(E_{jk}) \leq M(E_{ik}) \cdot M(A) \cdot M(E_{jk})$

d'où  $M(A) > 0$  dès qu'il existe un  $\alpha_{ij} \neq 0$

Théorème II Toutes les normes sur  $\mathcal{M}_n$  (comme sur tout espace vectoriel de dimension finie) sont équivalentes. C'est à dire que, étant données les deux normes  $N_1(A)$  et  $N_2(A)$ , il est deux constantes positives a et b telles :

$$a \leq \frac{N_2(A)}{N_1(A)} \leq b$$

- a) La démonstration est classique: si  $A = (\alpha_{ij})$ ,  $B = (\beta_{ij})$   
 $|N(A) - N(B)| \leq N(A-B) \leq \sum_{i,j} |\alpha_{ij} - \beta_{ij}| N(E_{ij}) \rightarrow$  continuité de  $N(A)$

- b) de plus si  $N^*(A) = \max_{i,j} |\alpha_{ij}|$  et  $p_1 = \sum_{i,j} N(E_{ij})$   
 $N(A) \leq N^*(A) \cdot p_1$

D'autre part, si les  $(\alpha_{ij})$  sont tels qu'ils varient de sorte que  $N^*(A) = 1$ , le point  $(\alpha_{ij})$  dans l'espace à  $n^2$  dimensions décrit la frontière compacte d'un paralléloèdre, donc il existe pour un choix des  $\alpha_{ij}$  une valeur minimum non nulle pour  $N(A)$ , soit  $p'$ . Mais l'homogénéité de  $\frac{N(A)}{N^*(A)}$  prouve que  $\frac{N(A)}{N^*(A)} = \frac{N(A')}{N^*(A')}$  si  $A' = \frac{1}{N^*(A)} \cdot A$ , telle que  $N^*(A') = 1$  et par suite

$$\frac{N(A)}{N^*(A)} \geq p' \quad , \quad p' \leq \frac{N(A)}{N^*(A)} \leq p_1$$

- c) Mais, de plus  $a_1 \leq \frac{N_1(A)}{N^*(A)} \leq b_1$ ,  $a_2 \leq \frac{N_2(A)}{N^*(A)} \leq b_2$ , on tire :

$$\frac{N_2(A)}{N_1(A)} = \frac{N_2(A)}{N^*(A)} \cdot \frac{N^*(A)}{N_1(A)} \quad , \quad \frac{a_2}{b_2} \leq \frac{N_2(A)}{N_1(A)} \leq \frac{b_2}{a_2} \quad \text{q.e.d.}$$

Il résulte de la démonstration du Théorème II que si  $N(A)$  est donnée, il existe deux constantes  $\alpha, \beta$  positives, telles que :

$$(1) \quad \alpha \leq \frac{N(A)}{N^*(A)} \leq \beta$$

Or, on vérifie immédiatement que  $N^*$  est telle que :

$$N^*(A.B) \leq n N^*(A) \cdot N^*(B) \quad \text{pour tout A et B d'ordre n}$$

donc :

$$N(A.B) \leq \frac{\beta}{\alpha^2} \cdot n \cdot N(A) \cdot N(B)$$

$$\text{Par suite, si } \lambda = \frac{\beta}{\alpha^2} n, \quad \lambda N(A.B) \leq \lambda N(A) \cdot \lambda N(B) \quad (2)$$

puisque  $\lambda N(A) = N_1(A)$  est une norme de matrice (2), s'écrit :

$$N_1(A.B) \leq N_1(A) \cdot N_2(B) \quad \text{donc :}$$

Proposition I : Etant donné une norme générale  $\mathcal{N}$ , il existe toujours un nombre  $\lambda$  tel que  $\lambda \mathcal{N}$  soit une norme multiplicative.

Nous allons voir une autre construction de normes multiplicatives (différente au paragraphe 3°).

Signalons enfin une autre propriété des normes multiplicatives :

Théorème III : Pour toute norme multiplicative  $M(A)$  sur  $\mathcal{E}$ , on a :

$M(A) \geq \lambda_A$ , si  $\lambda_A$  désigne la valeur absolue maximum des valeurs propres de A. [8]

En effet, soit  $x$  un vecteur propre correspondant à la valeur propre  $\lambda$ , je peux écrire la relation  $A.x = \lambda.x$ , où, si  $X$  est la matrice  $(n,n)$  associée à  $x$ , c'est ainsi que je nomme la matrice

$$X = \begin{pmatrix} x_1 & 0 & \dots & 0 \\ x_2 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ x_n & 0 & \dots & 0 \end{pmatrix} \quad \text{si} \quad x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad A.X = \lambda.X$$

Donc :

$$M(A).M(X) \geq \lambda M(X)$$

et, par suite :

$$M(A) \geq \lambda_A$$

### 3°) Normes de matrices associées à des normes de vecteurs.

Dans ce qui va suivre, les vecteurs sont ceux de  $\mathbb{R}^n$  (ou  $\mathbb{C}^n$ ) et les matrices sont toutes de type  $(n,n)$  sur  $\mathbb{R}$  (ou sur  $\mathbb{C}$ ).

Soient  $p(x)$  et  $q(x)$  deux normes ; formons le rapport  $\frac{p(Ax)}{q(x)}$

Posons :

$$S_{pq}(A) = \max_{x \neq 0} \left( \frac{p(Ax)}{q(x)} \right)$$

Théorème IV  $S_{pq}(A)$  est une norme sur  $\mathcal{M}_n$ . Si  $p(x) = q(x)$ , cette norme est multiplicative.

Prouvons d'abord que, quel que soient  $p$  et  $q$ , c'est une norme.

N I) est satisfait car si  $S_{pq}(A) = 0$ , c'est que, quelque soit  $x$ ,  $p(Ax) = 0$  ;  $A = 0$

N II) est satisfait car

$$S_{pq}(\lambda A) = \max_{x \neq 0} \left( \frac{p(\lambda Ax)}{q(x)} \right) = |\lambda| \max_{x \neq 0} \left( \frac{p(Ax)}{q(x)} \right) = \lambda \cdot S_{pq}(A)$$

N(III) est satisfait puisque :

$$\frac{p((A+B)x)}{q(x)} \leq \frac{p(Ax)}{q(x)} + \frac{p(Bx)}{q(x)} \leq S_{pq}(A) + S_{pq}(B)$$

donc  $S_{pq}(A)$  est bien une norme.

Montrons que  $S_{pq}$  n'est pas multiplicative en général.

Si  $p(x)$  est une norme,  $\frac{1}{C} p(x)$  en est une autre, prenons  $q(x) = \frac{1}{C} p(x)$

$$S_{pq}(A \cdot B) \leq S_{pq}(A) \cdot S_{pq}(B)$$

exige :

$$S_{pq}(A \cdot B) \leq C \cdot S_{pq}(A) \cdot S_{pq}(B)$$

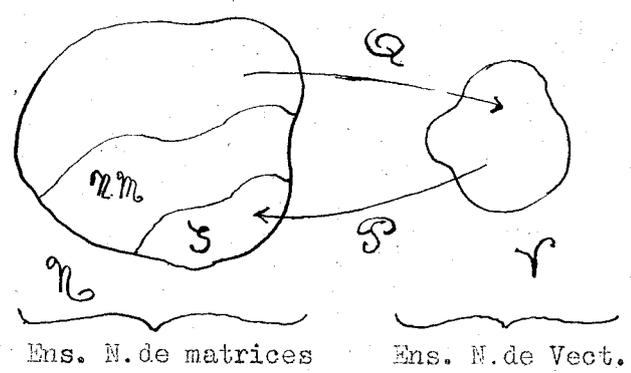
Il y a, certainement, des valeurs de  $C$  pour lesquelles cette inégalité n'est pas vérifiée.

Par contre, de :  $\frac{p(A B x)}{p(x)} = \frac{p(A B x)}{p(B x)} \cdot \frac{p(B x)}{p(x)}$

on tire  $S_{pp}(AB) \leq S_{pp}(A) \cdot S_{pp}(B)$   
 donc  $S_{pp}(A)$  est multiplicative

Pour résumer la situation:

Ayant une norme de vecteur,  $p$ , on peut lui faire correspondre une norme multiplicative  $S_{pp}$  de matrice. Cela définit une application  $\mathcal{P}$  de  $\mathcal{V}$  sur l'ensemble  $\mathcal{S}$  des normes multiplicatives de la forme  $S_{pp}$ .



D'autre part, ayant une norme de matrice  $N(A)$  à  $x \in \mathbb{R}^n$ , faisons correspondre la matrice  $X$  associée et posons  $p_N(x) = N(X)$ . Il est immédiat de voir que  $p_N$  est bien une norme de vecteur. Cela définit une application  $\mathcal{Q}$  de  $\mathcal{N}$  dans  $\mathcal{V}$

Pour caractériser l'ensemble  $\mathcal{S}$  on a le :

Théorème V :  $\mathcal{S}$  est l'ensemble des invariants de l'application composée  $\mathcal{P} \circ \mathcal{Q}$  de  $\mathcal{N}$  dans lui-même.

I) Soit  $S_{pp} \in \mathcal{S}$ , par  $\mathcal{Q}$  il correspond une norme  $p_i(x)$  à cette  $S_{pp} \cdot p_i(x) = S_{pp}(X)$

Or  $Xy = \eta \cdot x$  si  $y = (\eta_i) (i=1, 2, \dots, n)$

donc

$$p_i(x) = S_{pp}(X) = \max_{y \neq 0} \left( \frac{p(\eta \cdot x)}{p(x)} \right) = \left( \max_{y \neq 0} \frac{|\eta_i|}{p(y)} \right) \cdot p(x) = C \cdot p(x)$$

en posant  $C = \max_{y \neq 0} \frac{|\eta_i|}{p(y)}$

Par suite  $S_{p_i p_i}$  qui correspond à  $p_i$  par  $\mathcal{P}$  est telle que

$$S_{p_i p_i}(A) = \max_{y \neq 0} \left( \frac{p_i(Ay)}{p_i(y)} \right) = \max_{y \neq 0} \left( \frac{C \cdot p(Ay)}{C \cdot p(y)} \right)$$

et par cela  $S_{p_i p_i} = S_{pp}$  donc  $\mathcal{P} \circ \mathcal{Q}(S_{pp}) = S_{pp}$

II) Si une norme est invariante par  $\mathcal{P} \circ \mathcal{Q}$  et puisque  $\mathcal{P}$  n'a ses valeurs que dans  $\mathcal{S}$ , elle est du type  $S_{pp}$

Théorème VI (OSTROWSKI): Soit une norme multiplicative  $M(A)$ , si  $\rho_M = \rho(M)$  c'est-à-dire si  $\rho_M(x) = M(X)$  et  $S_{\rho_M \rho_M}$  la norme image de  $M$  par  $\mathcal{P}_0 \mathcal{Q}$ , pour toute matrice  $A$  :

$$S_{\rho_M \rho_M}(A) \leq M(A)$$

Autrement dit, la norme  $S_{\rho_M \rho_M}$  est toujours plus "serrée" que  $M$ .

En effet

$$S_{\rho_M \rho_M}(A) = \max_{y \neq 0} \left( \frac{\rho_M(Ay)}{\rho_M(y)} \right) = \max_{y \neq 0} \left( \frac{M(Y_1)}{M(Y)} \right)$$

si  $Y_1$  est l'associée de  $A.y$

Or, il est clair que  $Y_1 = A.Y$

donc :

$$\frac{M(Ay)}{M(y)} \leq M(A) \cdot \frac{M(y)}{M(y)}$$

et, par suite :

$$S_{\rho_M \rho_M}(A) \leq M(A)$$

D'après la propriété (Prop. I, parag. 2), nous posons :

Définitions : 1) Une norme multiplicative  $M$ , de matrice, est dite "critique" s'il n'y a pas de constante,  $0 < k < 1$ , telle que  $k \cdot M(A)$  soit encore une norme multiplicative.

2) Une norme multiplicative  $M$  de matrice est dite "minimale" s'il n'existe pas de norme multiplicative  $M' \neq M$ , telle que  $M'(A) \leq M(A)$  pour tout  $A$ .

Pour toute norme multiplicative  $M(I) \geq 1$  (Démonstration Th. I).

Donc, toute norme  $M$  telle  $M(I) = 1$  est critique. Les  $S_{pp}$  sont critiques car  $S_{pp}(I) = 1$ . Le théorème VI d'OSTROWSKI prouve qu'il existe toujours une  $S_{pp}$  inférieure à toute norme multiplicative.

Théorème VII : Toute norme  $S_{pp}$  est une norme <sup>multiplicative</sup> minimale.

Nous allons d'abord prouver que si, pour tout  $A$ , on a :

$$S_{qq}(A) \leq S_{pp}(A)$$

c'est que forcément :

$$S_{pp} = S_{qq}$$

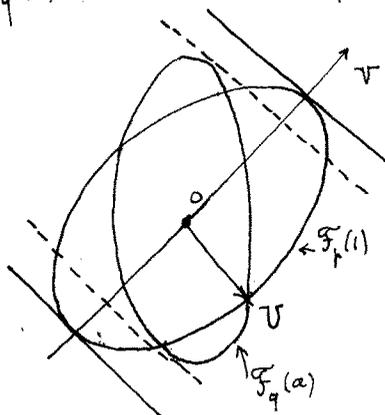
Prenons  $A = U \cdot V^T$ , c'est-à-dire matrice antiscaire où U et V sont deux colonnes quelconques de  $R^n$ .

On a, par définition :

$$\text{Max}_{x \neq 0} \left( \frac{|V^T x|}{|x|} \cdot p(U) \right) \geq \text{Max}_{y \neq 0} \left( \frac{|V^T y|}{|y|} \cdot q(U) \right) \quad (1)$$

puisque :  $U \cdot V^T x = (V^T x) \cdot U$ ,  $p(V^T x \cdot U) = |V^T x| \cdot p(U)$

$\mathcal{F}_p(a)$  désigne la frontière de la boule :  $B_p(a)$  (ensemble des x tels  $p(x) \leq a$ )  
 $\mathcal{F}_q(a)$  " " " " " " :  $B_q(a)$  ( " " " " " "  $q(x) \leq a$ )



Soit  $U \in \mathcal{F}_p(1)$ , donc  $p(U) = 1$   
 alors  $q(U) = a$  et  $U \in \mathcal{F}_q(a)$   
 On peut écrire la condition (1)

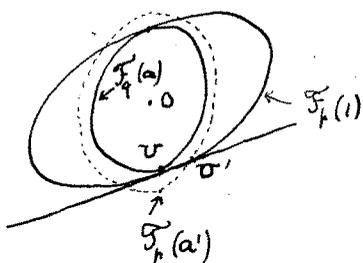
$$(1)': \text{Max}_{x \in \mathcal{F}_p(1)} |V^T x| \geq \text{Max}_{y \in \mathcal{F}_q(a)} |V^T y|$$

puisque

$$\text{Max}_{y \neq 0} \left( \frac{|V^T y|}{q(y)} \cdot q(U) \right) = \text{Max}_{y \in \mathcal{F}_q(a)} |V^T y|$$

(1') se traduit géométriquement en disant que les deux hyper-planes d'appui perpendiculaires à V pour  $B_p(1)$ , contiennent les deux hyper-planes d'appui perpendiculaires à V pour  $B_q(a)$  et cela pour tout V. Donc  $B_q(a)$  est toute dans  $B_p(1)$  et tout hyper-plan d'appui en U pour  $B_p(1)$  l'est pour  $B_q(a)$ . [43]

Or, si a varie, les boules  $B_q(a)$  se déduisent de l'une d'elles par des homothéties. Supposons que pour un point  $U \in \mathcal{F}_p(1)$ ,  $q(U) = a$



Je dis qu'en un autre point  $U' \in \mathcal{F}_p(1)$   
 on a aussi :  $q(U') = a$

En effet, si  $q(U') = a'$ , avec  $a' > a$ , par exemple, pour  $B_q(a')$  il existe des points strictement extérieurs à  $B_p(1)$ , puisqu'il y a des points de  $B_q(a')$  du côté de l'hyper-plan d'appui en U commun non situés du même côté que O. Donc  $B_q(a')$  ne serait pas tout dans  $B_p(1)$ , donc ne peut avoir que  $a = a'$ .

Il en résulte que  $B_p(1)$  est une  $B_q(a)$  pour un  $a$  déterminé.

Donc :  $q(x) = a \cdot p(x)$  ,  $S_{pp} = S_{qq}$  .

D'autre part, si (2) :  $M(A) \leq S_{pp}(A)$  quel que soit  $A$ , on sait qu'il existe  $p_n(x)$  telle que :

$$S_{p_n p_n}(A) \leq M(A) \leq S_{p p}(A)$$

Le résultat ci-dessus prouve que  $S_{p_n p_n} = S_{p p}$

donc (2) implique :  $M(A) = S_{pp}(A)$  .

Les normes  $S_{pp}(A)$  sont les normes multiplicatives minimales.

4°) Groupes de matrices associés à des normes de vecteurs , ou de matrices

L'ensemble  $\mathcal{S}$  des normes  $S_{pp}$  nous semble si important qu'il est intéressant d'essayer de le caractériser d'un autre point de vue.

a) Soit  $p$  une norme de vecteur dans  $R^n$  (ou  $C^n$ ) ,  $\mathcal{G}_p$  désigne l'ensemble des matrices carrées  $Q$  telles :

$$p(Qx) = p(x) \quad \text{pour tout } x \in R^n \text{ (ou } C^n)$$

Il est évident que :

- I)  $I \in \mathcal{G}_p$
- II) si  $Q_1$  et  $Q_2 \in \mathcal{G}_p$ ,  $p(Q_1 Q_2 x) = p(Q_1(Q_2 x)) = p(Q_2 x) = p(x)$   
donc  $Q_1 \cdot Q_2 \in \mathcal{G}_p$
- III) si  $Q \in \mathcal{G}_p$ , elle ne peut être singulière , car si non il existerait  $x_0 \neq 0$  tel  $Qx_0 = 0$  ; par suite :

$$p(x_0) = p(Qx_0) = 0 \quad \text{ce qui est absurde.}$$

Alors , si dans  $p(x) = p(Qx)$  on fait  $y = Qx$

on aura  $p(Q^{-1}y) = p(y)$  donc  $Q^{-1} \in \mathcal{G}_p$  .

Donc  $\mathcal{G}_p$  est un groupe multiplicatif de  $\mathcal{M}_n$  , il est dit associé à la norme de vecteur  $p$  .

b) Soit  $N$  une norme de matrice . On prouve pareillement que l'ensemble  $\mathcal{G}_N$  des matrices carrées  $Q$  telles :

$$N(QA) = N(A) \quad \text{Pour toute } A \in \mathcal{M}_n$$

est un groupe multiplicatif de  $\mathcal{M}_n$  .

( les points I et II sont évidents , pour III, si  $Q \in \mathcal{G}_N$  elle ne peut être singulière , car si non il existe  $x \neq 0$  auquel est associée  $\lambda_0 \neq 0$  telle que  $Q \cdot \lambda_0 = 0$  donc  $N(\lambda_0) = N(Q\lambda_0) = 0$  ce qui est absurde, la suite comme dans la démonstration précédente). Il est dit associé à la norme N.

Théorème VIII Etant donnée la norme N , si l'on forme  $\rho_N = \mathcal{Q}(N)$ , on a:  $\mathcal{G}_N \subset \mathcal{G}_{\rho_N}$   
Le groupe associé à la norme N est un sous-groupe du groupe associé à la norme de vecteur correspondant à N par  $\mathcal{Q}$ .

Car si  $Q \in \mathcal{G}_N$  on a pour tout  $x$  :  $\rho_N(Qx) = N(Qx) = N(x) = \rho_N(x)$

Mais en général si  $Q \in \mathcal{G}_{\rho_N}$  ,  $\rho_N(Qx) = \rho_N(x) \rightarrow N(Qx) = N(x)$   
X n'est pas une matrice quelconque et l'on ne peut affirmer  
 $N(A) = N(QA)$  pour tout A .

Théorème IX Si l'on se donne une norme de matrice  $S_{\rho} \in \mathcal{S}$  , le groupe  $\mathcal{G}_{S_{\rho}}$  est identique au groupe  $\mathcal{G}_{\rho}$ , lui-même identique à  $\mathcal{G}_{\rho}$ .

En effet, je sais que  $\rho_{S_{\rho}}$  est définie par  $\rho_{S_{\rho}}(x) = S_{\rho}(x)$  et j'ai prouvé dans la démonstration du Th.V que  $\rho_{S_{\rho}} = C \cdot \rho = \rho$

donc  $\mathcal{G}_{\rho_{S_{\rho}}} = \mathcal{G}_{\rho}$

Soit donc  $Q \in \mathcal{G}_{\rho} = \mathcal{G}_{\rho_{S_{\rho}}}$

on a:  $S_{\rho}(Q \cdot A) = \max_{y \neq 0} \left( \frac{\rho(QAy)}{\rho(y)} \right) = \max_{y \neq 0} \left( \frac{\rho(Ay)}{\rho(y)} \right) = S_{\rho}(A)$

et par suite:  $Q \in \mathcal{G}_{S_{\rho}}$

Donc  $\mathcal{G}_{\rho_{S_{\rho}}} \subset \mathcal{G}_{S_{\rho}}$  et d'après le théorème précédent comme  $\mathcal{G}_{S_{\rho}} \subset \mathcal{G}_{\rho_{S_{\rho}}}$  on a bien le résultat énoncé.

Théorème X Pour une norme de matrice  $S_{\rho} \in \mathcal{S}$  , toute matrice  $Q \in \mathcal{G}_{\rho}$  est telle que  $S_{\rho}(QA) = S_{\rho}(AQ) = S_{\rho}(A)$

On vient de voir que  $S_{\rho}(QA) = S_{\rho}(A)$

D'autre part,

$$S_{\rho}(AQ) = \max_{x \neq 0} \left( \frac{\rho(AQx)}{\rho(x)} \right) = \max_{x \neq 0} \left( \frac{\rho(AQx)}{\rho(Qx)} \right)$$

puisque Q est inversible Qx prend toute valeur  $y \in \mathbb{R}^n$  (ou  $\mathbb{C}^n$ )

donc:  $S_{\rho}(AQ) = \max_{y \neq 0} \left( \frac{\rho(Ay)}{\rho(y)} \right) = S_{\rho}(A)$

On peut remarquer aussi que  $S_{\rho}(Q) = +1$  pour tout  $Q \in \mathcal{G}_{\rho}$

car  $S_{\rho}(Q) = \max_{x \neq 0} \left( \frac{\rho(Qx)}{\rho(x)} \right) = +1$

5°) Normes de matrices complètement invariantes.

Soit  $N(A)$  une norme générale sur l'ensemble des matrices carrées d'ordre  $n$  et à éléments réels, c'est-à-dire une application de l'ensemble des matrices carrées d'ordre  $n$  sur  $\mathbb{R}$ , dans  $\mathbb{R}$ , et satisfaisant aux axiomes de ce chapitre. On peut se demander s'il existe des normes "complètement invariantes", c'est-à-dire telles

$$N(A) = N(SAS^{-1})$$

pour toute  $A$  et pour toute  $S$  non singulière.

Théorème XII : Il est impossible de trouver pour  $n \geq 2$  une norme complètement invariante.

En effet, considérons la matrice (d'ordre 2)

$$A = \begin{vmatrix} 0 & 1 \\ 0 & 0 \end{vmatrix}$$

et la matrice  $A' = \begin{vmatrix} 0 & \frac{a}{c} \\ 0 & 0 \end{vmatrix}$ , ( $c \neq 0, a \neq 0$ )

Il est clair que :

$$\begin{vmatrix} a & b \\ 0 & c \end{vmatrix} \cdot \begin{vmatrix} 0 & 1 \\ 0 & 0 \end{vmatrix} = \begin{vmatrix} 0 & a \\ 0 & 0 \end{vmatrix}$$

et que :

$$\begin{vmatrix} 0 & \frac{a}{c} \\ 0 & 0 \end{vmatrix} \cdot \begin{vmatrix} a & b \\ 0 & c \end{vmatrix} = \begin{vmatrix} 0 & a \\ 0 & 0 \end{vmatrix}$$

Donc, si  $S = \begin{vmatrix} a & b \\ 0 & c \end{vmatrix}$ ,  $SA = A'S$

$S$  n'étant pas singulière  $A' = S.A.S^{-1}$

Il est clair que si  $N$  est une norme quelconque

$$N(A^c) = \left| \frac{a}{c} \right| \cdot N(A)$$

puisque

$$A^c = \frac{a}{c} \cdot A$$

Donc, on ne peut avoir  $N(A^c) = N(A)$  quel que soit  $A$  et  $S$

Cette démonstration s'étend visiblement à  $n$  quelconque.

Cela prouve qu'il est impossible de trouver une norme de matrice dont la valeur ne dépendrait que de la transformation linéaire qu'elle représente, puisque les matrices  $S.A.S^{-1}$  peuvent s'interpréter comme définissant la même transformation linéaire.

Cela explique les difficultés qu'il y a dans le "dosage" de la proximité des éléments propres d'une matrice (c'est-à-dire attachés à la transformation qu'elle représente) par des mesures de voisinage sur les coefficients de celle-ci.

#### 6°) Normes de matrices usuelles en calcul numérique linéaire.

Il est clair que pour l'ensemble des matrices  $A$  d'ordre  $n$ , à éléments dans  $R$  (ou  $C$ ), en correspondance biunivoque avec  $R^{n^2}$  (ou  $C^{n^2}$ ), on peut prendre une norme de vecteur dans  $R^{n^2}$  (ou  $C^{n^2}$ ).

On obtient ainsi :

$$N_1 = \sum_{i,j} |a_{ij}|, \quad N_2 = \sqrt{\sum_{i,j} |a_{ij}|^2}, \quad N_\infty = \max_{i,j} |a_{ij}|$$

qui correspondent aux trois normes de HOLDER,  $q_1, q_2, q_\infty$

Les normes de matrices obtenues à partir de la construction du paragraphe 3 seront utilisées en détail dans le chapitre suivant.



## CHAPITRE III

### Les conditionnements numériques

#### I) Notion de conditionnement .

Après avoir présenté le matériel théorique dont nous nous servirons nous abordons les problèmes de calcul numérique que nous nous proposons d'étudier. Nous commençons par la notion de conditionnement numérique.

Nous allons considérer, pour commencer, un exemple cité par BODEWIG [9] Soient les trois systèmes de deux équations à deux inconnues du 1<sup>o</sup> degré:

$$\text{I} \quad \begin{cases} 3x + 4y = 7 \\ 3x + 4,00001y = 7,00001 \end{cases} \quad \text{ou} \quad \begin{cases} 3x + 4y - 7 = 0 \\ 3x + 4,00001y - 7,00001 = 0 \end{cases}$$

$$\text{II} \quad \begin{cases} 3x + 4y - 7 = 0 \\ 3x + 3,99999y - 7,00004 = 0 \end{cases}$$

$$\text{III} \quad \begin{cases} 3x + 4y - 7 = 0 \\ 3x + 3,999992y - 7,000042 = 0 \end{cases}$$

Ces systèmes ont pour solution exacte

$$\text{I} \quad x = 1, \quad y = 1$$

$$\text{II} \quad x = 7 + \frac{2}{3}, \quad y = -4$$

$$\text{III} \quad x = 9 + \frac{1}{3}, \quad y = -\left(5 + \frac{1}{4}\right)$$

Il est d'usage de dire que de tels systèmes sont "mal" conditionnés.

Avant de quitter cet exemple, voici un tableau donnant les valeurs des deux formes des équations I, II, III, lorsque l'on y substitue les valeurs solutions de I, II, III.

	Sys I	Sys II	Sys III
Sol I	0	0	0
	0	-0,00005	-0,00005
Sol II	0	0	0
	-0,00005	0	-0,00001
Sol III	0	0	0
	-0,0000625	0,0000105	0

2°) Examen des différentes définitions de conditionnement.

La nécessité d'une définition correcte du conditionnement d'un système linéaire a été ressentie depuis longtemps par les calculateurs.

Remarquons d'abord que la valeur absolue du déterminant d'une matrice n'apprend rien sur le conditionnement du système. Soit, en effet, un système (d'ordre 10 par exemple), de déterminant égal à 1 par exemple ; multiplions par 1/10 toutes les équations, le déterminant du nouveau système est  $10^{-10}$ . Or, la résolution du nouveau système n'est pas numériquement plus délicate que celle du système de départ.

TURING [13] et LONSETH [14] ont le point de vue suivant. Ils cherchent à savoir si, faisant varier de très peu les éléments de A, la solution exacte du système :  $(A + \delta A) x - b = 0$  diffère de peu ou de beaucoup de la solution exacte  $A x - b = 0$ .

Ce point de vue conduit à une définition du conditionnement voisine de la nôtre (cf. Th. III).

TURING [13] considère deux nombres pouvant donner une mesure du conditionnement. (cf. aussi sur ce sujet TODD, [15] )

le premier est :  $T_1(A) = n \cdot \max_{i,j} |a_{ij}| \cdot \max_{i,j} |a'_{ij}|$  ,  $A^{-1} = (a'_{ij})$

le second est :  $T_2(A) = \frac{1}{n} N(A) \cdot N(A^{-1})$  ,  $N(A) = \left( \sum_{i,j} a_{ij}^2 \right)^{\frac{1}{2}}$

Pour VON NEUMANN et GOLDSTINE [16] , la notion de conditionnement vient à l'issue d'un calcul d'erreur dans la résolution numérique des systèmes linéaires à matrice symétrique. Leur nombre de conditionnement égal au rapport de la plus grande à la plus petite valeur propre de la matrice A, est un cas particulier de ceux que nous définissons (à l'inverse près).

### 3°) Définition générale des conditionnements numériques.

Soit  $\mathcal{E}$  un espace vectoriel topologique, normé et complet,  $N_1$  et  $N_2$  deux normes équivalentes sur  $\mathcal{E}$ . Soit  $T$  une application de  $\mathcal{E}$  dans lui-même, satisfaisant en un point  $x_0$  à la propriété:

(P) Il existe un voisinage  $V_0$  de  $x_0$ , tel que pour tout  $x \in V_0$ , deux nombres positifs  $p_1$  et  $p_2$  puissent être trouvés tels que:

$$0 < p_1 \leq \frac{N_1(T(x) - T(x_0))}{N_2(x - x_0)} \leq p_2 \quad \text{pour tout } x \in V_0.$$

Il est clair que si (P) est vérifiée, on ne fait pas de restriction en supposant que  $p_1$  est la Lim des  $p_1$  y satisfaisant

que  $p_2$  est la Lim des  $p_2$  y satisfaisant.

On suppose qu'il en est ainsi dans toute la suite, ~~autre-~~ autrement dit, que les bornes  $p_1$  et  $p_2$  qui figurent dans l'énoncé de (P) sont les "meilleures"

Théorème I : Toute application  $T$  satisfaisant à (P) est continue en  $x_0$  et l'équation

$$T(x) = T(x_0)$$

a dans  $V_0$  une et une seule solution  $x = x_0$ .

a) La continuité est très simple à prouver: si  $\varepsilon$  est donné

$$\text{comme } N_1[T(x) - T(x_0)] \leq p_2 N_2(x - x_0)$$

il suffit de choisir  $x$  de sorte que:  $N_2(x - x_0) \leq \frac{\varepsilon}{p_2}$

mais puisque  $N_1$  et  $N_2$  sont équivalentes il existe  $a$  tel

$N_2(x - x_0) \leq a \cdot N_1(x - x_0)$ : donc il suffit de prendre  $x$  dans le voisinage:

$$V_0 \cap \left\{ x \text{ tel } N_1(x - x_0) \leq \frac{\varepsilon}{p_2 a} \right\}$$

b) Puisque  $p_1 \neq 0$   $p_1 N_2(x_1 - x_0) \leq N_1(T(x_1) - T(x_0))$  implique  $T(x_1) = T(x_0)$

$$N_2(x_1 - x_0) = 0$$

$$\text{donc } x_1 = x_0.$$

Expliquons l'importance en calcul numérique des nombres  $p_1$  et  $p_2$  de la propriété (P):

Tout calcul effectif se conduit avec un certain "moyen de calcul"

Nous entendons par là, soit une machine, soit, dans le cas d'un

calcul "à la main" un certain nombre de règles que doit obligatoirement se donner le calculateur. Dans un cas, comme dans l'autre ce

"moyen" est toujours à capacité limitée, aussi bien supérieurement

qu'inférieurement. Autrement dit, dans tout calcul effectif, on

est obligé de considérer comme nuls des nombres qui ne sont réellement

qu'inférieurs à une certaine valeur  $\varepsilon$ .

Si donc le calcul nous a conduit, pour la solution du problème:

$$T(x) = T(x_0)$$

à un  $x_1$  tel

$$\eta_1 \leq N_1(T(x_1) - T(x_0)) \leq \eta_2$$

tout ce que l'on peut dire relativement à la proximité de  $x_1$  et  $x_0$

est que 
$$N_2(x_1 - x_0) \leq \frac{\eta_2}{p_1} = \varepsilon_{\max}$$

Donc: 
$$\frac{\varepsilon_{\min}}{\varepsilon_{\max}} = \frac{\eta_1}{\eta_2} \cdot \frac{p_1}{p_2} = \frac{N_2(x_1 - x_0) \geq \frac{\eta_1}{p_2} = \varepsilon_{\min}}{N_2(x_1 - x_0) \leq \frac{\eta_2}{p_1} = \varepsilon_{\max}}$$

Propriété I : Le rapport entre les extrêmes des erreurs en norme  $N_2$ , que l'on peut

commettre est égal à 
$$\frac{\eta_1}{\eta_2} \cdot \frac{p_1}{p_2} \quad [10]$$

Si l'ordre de grandeur de  $N_1(T(x_1) - T(x_0))$  est connu  $= \eta = \eta_1 = \eta_2$

on peut dire que ce rapport est égal à : 
$$\frac{p_1}{p_2}$$

Définition I : On appelle conditionnement général d'une équation  $T(x) = T(x_0)$ , avec

T satisfaisant à (P), pour  $x_0$  et relativement aux normes  $N_1$  et  $N_2$

le rapport 
$$\frac{p_1}{p_2}$$
.

#### 4°) Cas d'applications linéaires

D'après les propriétés exposées au chapitre II, si T est linéaire et non singulière la propriété (P) est toujours satisfaite, les nombres  $p_1$  et  $p_2$  ne dépendent pas de  $x_0$ . La définition I est donc toujours applicable. Il est cependant obligatoire de particulariser pour les applications numériques afin de rendre utilisable en calcul cette définition.

Soit alors A une matrice carrée non singulière définissant une application linéaire de  $R^n$  dans  $R^n$ ; formons les 9 rapports

$$r(\varphi_i, \varphi_j, A, x) = \frac{\varphi_i(Ax)}{\varphi_j(x)} \quad (j, i = 1, 2, \dots, n)$$

On sait que  $\text{Max}_x r(\varphi_i, \varphi_j, A, x)$  est une norme de matrice sur  $M_n$

c'est ce qu'au chapitre II on a posé  $= S_{\varphi_i, \varphi_j}(A) = S_{i,j}(A)$

(elle n'est multiplicative que si  $i=j$ )

Donc le  $p_2$  de la propriété (P) n'est autre que  $S_{i,j}(A)$

D'autre part, écrivons

$$r(\varphi_i, \varphi_j, A, x) = \frac{\varphi_i(x_2)}{\varphi_j(A^{-1}x_2)} \quad , \quad x_2 = Ax \quad , \quad r(\varphi_i, \varphi_j, A, x) = \frac{1}{r(\varphi_j, \varphi_i, A^{-1}, x_2)}$$

Puisque A n'est pas singulière

Donc le  $p_1$  de la propriété (P) n'est autre que  $(S_{j,i}(A^{-1}))^{-1}$

D'où le théorème important :

Théorème II : Le conditionnement général d'une matrice A non singulière pour les normes

$\varphi_i$  et  $\varphi_j$  est :

$$\gamma_{ij} = 1 / S_{ij}(A) \cdot S_{ji}(A^{-1})$$

Définition : Si l'on se donne un système d'équations linéaires dans  $\mathbb{R}^n$  :

$$A x - b = 0$$

ayant une solution unique  $x_0$ , son conditionnement général n'est autre que celui de sa matrice de premier membre, puisque le système s'écrit :

$$A x = A x_0$$

Remarquons que le théorème II peut, très généralement, conduire à poser pour définition d'un nombre de conditionnement

$$\Gamma_{12} = 1 / N_1(A) \cdot N_2(A^{-1})$$

$N_1$  et  $N_2$  étant deux normes générales.

C'est ce qu'a fait TURING (à l'inverse près).

Les différentes interprétations de ce nombre sont plus difficiles.

Nous nous bornerons à l'étude de normes de matrices du type  $S_{ij}$ .

Enfin, pour rapprocher les points de vue, on a :

Théorème III : Soit  $x_0$  la solution exacte d'un système linéaire d'ordre n

$$A x - b = 0$$

Soit  $x_1$  la solution exacte du système varié

$$(A + \delta A)x - b = 0$$

$\varphi_i$  et  $\varphi_j$  deux normes de vecteurs sur  $\mathbb{R}^n$  (ou  $\mathbb{C}^n$ ).

(1) Si les variations  $\delta A$  sont prises de sorte que

soit toujours d'un même ordre de grandeur  $\eta$  ;

alors, le rapport des variations relatives extrêmes de la solution,

en normes  $\varphi_i$  et  $\varphi_j$ , n'est autre que le conditionnement général  $\gamma_{ij}$ .

- En effet, de  $A x_0 = b$  et  $(A + \delta A) x_1 - b = 0$

on tire :

$$A(x_1 - x_0) + \delta A x_1 = 0$$

donc :

$$\varphi_i(A(x_1 - x_0)) = \varphi_i(\delta A x_1)$$

$$\frac{\varphi_i(\delta A x_1)}{\varphi_i(x_1)}$$

$$\frac{\varphi_i(A(x_1 - x_0))}{\varphi_j(x_1 - x_0)} = \frac{\varphi_i(\delta A x_1)}{\varphi_i(x_1)} \cdot \frac{\varphi_i(x_1)}{\varphi_j(x_1 - x_0)} \neq \eta \cdot \frac{\varphi_i(x_1)}{\varphi_j(x_1 - x_0)}$$

si  $p_1, p_2$  sont les barres du rapport du 1er terme.

$$E_{z, \text{Max}} = \text{Max}_{x_1} \left( \frac{\varphi_j(x_1 - x_0)}{\varphi_i(x_1)} \right) \neq \eta \cdot \frac{1}{p_1}$$

$$E_{z, \text{min}} = \text{Min}_{x_1} \left( \frac{\varphi_j(x_1 - x_0)}{\varphi_i(x_1)} \right) \neq \eta \cdot \frac{1}{p_2}$$

$$\frac{E_{z, \text{min}}}{E_{z, \text{Max}}} \neq \frac{p_1}{p_2} = \gamma_{ij}$$

En ce sens, le  $\gamma_{ij}$  apparaît comme un indice de forme.

Nota : La condition (1) est toujours satisfaite en ordre de grandeur, ce qui est en calcul la chose la plus importante.

5°) Comparaison des différents conditionnements généraux d'une matrice.

Soient, pour une matrice A donnée non singulière, les deux rapports :

$$r_{ij} = \frac{\varphi_i(Ax)}{\varphi_j(x)} \quad , \quad r_{k\ell} = \frac{\varphi_k(Ax)}{\varphi_\ell(x)} \quad , \quad (i, j, k, \ell = 1, 2, \infty)$$

on sait que :

$$(1) \quad a_{ki} \varphi_k(x) \leq \varphi_i(x) \leq b_{ki} \varphi_k(x)$$

donc ,

$$a_{ki} \varphi_k(Ax) \leq \varphi_i(Ax) \leq b_{ki} \varphi_k(Ax)$$

Mais

$$a_{ej} \varphi_e(x) \leq \varphi_j(x) \leq b_{ej} \varphi_e(x)$$

d'où

$$a_{ki} \frac{\varphi_k(Ax)}{\varphi_j(x)} \leq \frac{\varphi_i(Ax)}{\varphi_j(x)} \leq b_{ki} \frac{\varphi_k(Ax)}{\varphi_j(x)}$$

puis

$$\frac{a_{ki}}{b_{ej}} \frac{\varphi_k(Ax)}{\varphi_e(x)} \leq \frac{\varphi_i(Ax)}{\varphi_j(x)} \leq \frac{b_{ki}}{a_{ej}} \frac{\varphi_k(Ax)}{\varphi_e(x)} ; \quad \frac{a_{ki}}{b_{ej}} r_{k\ell} \leq r_{ij} \leq r_{kj} \leq \frac{b_{ki}}{a_{ej}} r_{k\ell}$$

enfin, en posant  $\gamma_{ij}$  et  $\gamma_{k\ell}$  pour les conditionnements relatifs aux couples de normes  $\varphi_i, \varphi_j$  et  $\varphi_k, \varphi_\ell$  , on peut énoncer :

Théorème IV Les conditionnements généraux d'une même matrice A, par rapport à deux couples de normes  $\varphi_i, \varphi_j$  et  $\varphi_x, \varphi_e$  sont toujours tels: (les constantes  $a_{pr}$  et  $b_{pr}$  sont définies par (1))

$$\gamma_{ij} \cdot \frac{a_{je}}{b_{je}} \cdot \frac{a_{il}}{b_{il}} \leq \gamma_{le} \leq \gamma_{ij} \cdot \frac{b_{ej}}{a_{ej}} \cdot \frac{b_{li}}{a_{li}}$$

Nous allons écrire ces relations pour les cas envisagés:

Voici les tableaux relatifs aux valeurs des  $a_{ij}$  et des  $b_{ij}$ .

$i \rightarrow a_{ij}$	$a_{ij}$	1	2	$\infty$
1	1	1	$\frac{1}{\sqrt{n}}$	$\frac{1}{n}$
2	1	1	1	$\frac{1}{\sqrt{n}}$
$\infty$	1	1	1	1

$i \rightarrow b_{ij}$	$b_{ij}$	1	2	$\infty$
1	1	1	1	1
2	$\sqrt{n}$	1	1	1
$\infty$	n	$\sqrt{n}$	1	1

D'après ces tableaux on voit que la double inégalité la moins "serrée" est obtenue pour:  $k, l = 1, i = j = \infty$   
 elle s'écrit:  $\gamma_{\infty} \cdot \frac{1}{n^2} \leq \gamma_{11} \leq \gamma_{\infty} \cdot n^2$   
 cela nous permet de conclure:

Propriété II Tous les conditionnements généraux que nous venons de définir sont équivalents.

Théorème V Le conditionnement  $\gamma_{22}$  relatif aux deux normes  $\varphi_x$  et  $\varphi_e$ , identiques à la norme euclidienne sur  $R^n$ , d'une matrice A non singulière n'est autre que le rapport  $\frac{\sqrt{\lambda_1}}{\sqrt{\lambda_n}}$ , où  $\lambda_1$  et  $\lambda_n$  désignent la plus petite et la plus grande des valeurs propres de la matrice  $A^T A$ .

En effet dans ce cas on peut écrire:  $\gamma_{22} = \frac{\|Ax\|}{\|x\|}$

ou si X désigne la colonne des composantes de x

$$\gamma_{22}^2 = \frac{X^T A^T A X}{X^T X}$$

Or, la matrice  $A^T.A$  est définie positive et symétrique, et non singulière, donc a des valeurs propres  $0 < \lambda_1 < \lambda_2 < \dots < \lambda_n$  strictement positives.

D'autre part, il existe une matrice  $Q$ , orthonormale ( $Q^T.Q = I$ ) telle que :

$$Q^T A^T . A Q = \Lambda$$

ou

$$\Lambda = \begin{vmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & & \\ \vdots & & \ddots & \\ 0 & & & \lambda_n \end{vmatrix}$$

Or, si l'on pose  $X \equiv Q Y$ , on a :

$$r_{22}^2 = \frac{Y^T . Q^T A^T . A . Q . Y}{Y^T Q^T Q Y} = \frac{Y^T \Lambda Y}{Y^T Y} = \frac{\sum_{i=1}^n \lambda_i y_i^2}{\sum_{i=1}^n y_i^2}$$

sous cette forme, on voit que :

$$r_{1,22}^2 = \lambda_1, \quad r_{n,22}^2 = \lambda_n$$

donc :

$$r_{22} = \sqrt{\frac{\lambda_1}{\lambda_n}}$$

Si  $A$  est donnée définie positive et symétrique si  $0 < u_1 < \dots < u_n$  sont les valeurs propres de  $A$ , elle-même,  $u_i^2 = \lambda_i$

donc, dans ce cas :

$$r_{22} = \frac{u_1}{u_n}$$

C'est le nombre de VON NEUMANN. Il est assez laborieux à calculer.

Les nombres de TURING le sont un peu moins.

Il paraît normal de considérer un système comme le mieux conditionné possible si le système de "plans" dont l'intersection est la solution du système linéaire est un système orthogonal. Le chapitre qui va suivre expose les propriétés de nombres ayant les qualités d'être assez faciles à construire, et d'indiquer la condition d'orthogonalité simplement .

CHAPITRE IV

Le conditionnement normalisé C(A)

1°) Notations : Dans ce qui va suivre, je vais particulièrement étudier le rapport

$$r_{12} = r(\varphi_1, \varphi_2, A, x) = \frac{\varphi_1(Ax)}{\varphi_2(x)}$$

On posera :

$$r_{12}(x) = \frac{\Phi(Ax)}{\|x\|}, \quad m_A = \min_{x \neq 0} (r_{12}(x)), \quad M_A = \max_{x \neq 0} (r_{12}(x)),$$
$$\gamma_{12} = \gamma = \frac{m_A}{M_A}.$$

L'intérêt de ces normes dissymétriques est d'avoir des propriétés géométriques simples comme on va le voir.

Afin de mieux utiliser cette interprétation géométrique, nous commencerons par l'étude d'un cas particulier :

2°) Etude géométrique du rapport  $\frac{m_A}{M_A}$ , pour A "normée en lignes"

Nous disons que A est "normée en lignes", si tout vecteur-ligne  $A_i$  est de module 1 :

$$\sum_{k=1}^{k=n} a_{ik}^2 = 1 \quad (i=1, 2, \dots, n)$$

On suppose pour la suite A non singulière.

$$\text{Ici } \Phi(x) = \sum_{i=1}^{i=n} \left| \sum_{k=1}^{k=n} a_{ik} x_k \right| = \sum_{i=1}^{i=n} |f_i(x)|$$

Pour trouver  $m_A$  et  $M_A$  il suffit de déplacer le point (x) sur le lieu des points (x) tels que  $\Phi(x) = 1$

Or ce lieu

est la frontière d'un polyèdre  $\mathcal{P}$  convexe ayant l'origine pour centre de symétrie, cette frontière:  $\Pi$  est constituée par des "faces" planes qui sont des parties des plans dont les équations sont :

$$\sum_{i=1}^{i=n} \varepsilon_i f_i(x) = 1 \quad (\varepsilon_i = \pm 1)$$

La condition de normalisation des lignes de A entraîne que

$\Pi$  est le lieu des points dont la somme des distances géométriques aux plans  $P_i$  d'équation  $f_i(x) = 0$  est égale à 1

Par suite:

a)  $m_A = \min_{(x) \in \mathcal{P}} \left( \frac{1}{\|x\|} \right) = \frac{1}{d_M}$ , si  $d_M$  est la distance maximum de l'origine à un point de  $\Pi$ , et puis que  $\Pi$  limite un polyèdre convexe,  $d_M$  est la distance de 0, au sommet

le plus éloigné de 0;

- b)  $M_A = \text{Max}_{(x) \in \mathcal{P}} \left( \frac{1}{\|x\|} \right) = \frac{1}{h_m}$ , si  $h_m$  est la distance minimum de 0 à un point de  $\mathcal{P}$ . Toujours à cause de la convexité de  $\mathcal{P}$ ,  $h_m$  est la distance de 0 à la face la plus proche de 0.

### 3) Etude analytique:

- a) Soient  $D_i, (i=1, 2, \dots, n)$  les  $n$  droites:

$$D_i = \bigcap_j P_j \quad (j=1, 2, \dots, i-1, i+1, \dots, n)$$

les plans  $P_i$ , d'équations  $f_i(x) = 0$ .

Le fait que  $A$  soit non singulière assure que les droites  $D_i$  peuvent être prises comme support d'un système de vecteurs  $(\vec{e}_1, \vec{e}_2, \dots, \vec{e}_n)$  formant une base pour  $\mathbb{R}^n$ . Nous choisissons un tel système de sorte que  $\|\vec{e}_i\|^2 = 1$ .

- b) Détermination des faces du polyèdre  $\mathcal{P}$  dans le système d'axes  $(0; \vec{e}_1, \vec{e}_2, \dots, \vec{e}_n)$ .

Un plan coupant les droites  $D_i$  en des points  $A_i$  d'abscisse  $\alpha_i$ ,

a pour équation:

$$\sum_{i=1}^n \frac{x_i}{\alpha_i} = 1, \quad \text{ou si } V \text{ est la colonne: } \begin{pmatrix} \frac{1}{\alpha_1} \\ \vdots \\ \frac{1}{\alpha_n} \end{pmatrix}, \quad x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

cela s'écrira: (1)  $x^T \cdot V = 1$

Or le produit scalaire de deux vecteurs de composantes  $(x_i)$  et  $(y_i)$

s'écrit:  $\mu = (x_1 \vec{e}_1 + \dots + x_n \vec{e}_n) \cdot (y_1 \vec{e}_1 + \dots + y_n \vec{e}_n) = x^T \cdot G \cdot y$

où  $G = (g_{ij})$  avec  $g_{ij} = \vec{e}_i \cdot \vec{e}_j \quad (i=1, \dots, n; j=1, \dots, n)$

( $G$ : matrice de "métrique")

Si  $\vec{U}$  est un vecteur de module 1, on aura entre ses composantes

$$(2) \quad u^T G u = 1$$

Considérons le plan orthogonal en  $H$ , de coordonnées  $h \cdot u$  ( $h > 0$ ) au

vecteur  $\vec{U}$ , son équation est:  $(x^T - h \cdot u^T) \cdot G \cdot u = 0$

ou  $x^T G u = h \quad (3)$

En identifiant (3) et (1) cela donne:  $h \cdot V = G u$

ou  $u = h G^{-1} V$

et portant dans (2), puisque  $G^{-1}$  est symétrique:  $h^2 V^T G^{-1} V = 1$

Donc:

$$(4) \quad \frac{1}{h^2} = V^T G^{-1} V$$

Mais les points  $A_i$  sont tels que le plan  $\frac{x_i}{\alpha_i} = 1$  soit à la distance 1 de 0, car le polyèdre  $\mathcal{P}$  n'est autre que le lieu des points dont la somme des distances géométriques aux plans de notre nouveau système de coordonnées est inférieure ou égale à 1.

Donc si

$$V_i = \begin{pmatrix} 0 \\ \vdots \\ \frac{1}{\alpha_i} \leftarrow i^e \\ \vdots \\ 0 \end{pmatrix}$$

Il résulte de (4) que :

$$1 = V_i^T G^{-1} V_i$$

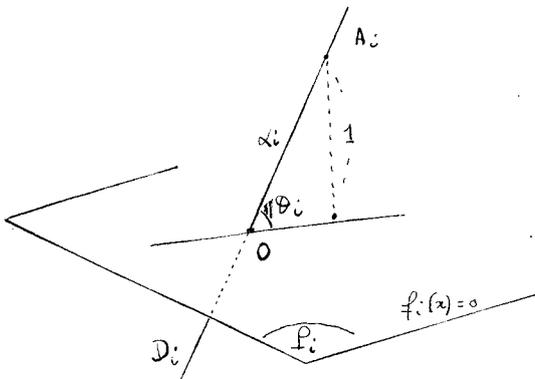
Donc :  $\frac{g^{ii}}{\alpha_i^2} = 1$ ,  $\alpha_i = \varepsilon_i \sqrt{g^{ii}}$ , si  $G^{-1} = (g^{ii})$

On peut remarquer que : 1°)  $g^{ii} \geq 1$

Car si  $\theta_i$  est l'angle aigu de  $D_i$  avec  $P_i$  d'équation  $f_i(x) = 0$

on a  $|\alpha_i| = \frac{1}{\sin \theta_i}$

2°) L'indétermination de signe provient du fait que sur  $D_i$ , le point  $A_i$  et son symétrique par rapport à 0 sont sommets de  $\mathcal{P}$ . Les faces de  $\mathcal{P}$  sont les simplexes de sommets : n points  $A_i$  pris de sorte que deux d'entre-eux ne soient pas symétriques par rapport à 0



c) Calcul des  $d_n$  et  $h_n$ .

1°) Il est clair que :  $d_n = \text{Max } |\alpha_i| = \text{Max } \sqrt{g^{ii}}$

2°)  $h_n = \min(h)$ ,  $h$ , précédemment trouvé fromule (4)

Or celle-ci, s'écrit en remarquant que  $V = V_1 + V_2 + \dots + V_n$

$$\frac{1}{h^2} = (V_1 + \dots + V_n)^T G^{-1} (V_1 + \dots + V_n)$$

Ou 
$$\frac{1}{h^2} = n + 2 \sum_{\substack{i < j \\ (1, \dots, n)}} \frac{g^{ij}}{\alpha_i \alpha_j} = n + 2 \sum_{\substack{i < j \\ (1, \dots, n)}} \varepsilon_i \varepsilon_j \frac{g^{ij}}{\sqrt{g^{ii}} \sqrt{g^{jj}}} \quad (6)$$

Etudions les diverses valeurs de  $h$  par (6)

Posons  $Z_{ij} = \frac{g^{ij}}{\sqrt{g^{ii}} \sqrt{g^{jj}}}$   
et

$$L = \sum_{i < j} \varepsilon_i \varepsilon_j Z_{ij}, \quad \frac{1}{h^2} = n + 2.L$$

D'autre part, écrivons :

$$L = \varepsilon_1 [\varepsilon_2 Z_{12} + \varepsilon_3 Z_{13} + \dots + \varepsilon_n Z_{1n}] + \varepsilon_2 [\varepsilon_3 Z_{23} + \dots + \varepsilon_n Z_{2n}] + \dots + \varepsilon_{n-2} [\varepsilon_{n-1} Z_{n-2, n-1} + \varepsilon_n Z_{n-2, n}] + \varepsilon_{n-1} \varepsilon_n Z_{n-1, n}$$

Je regarde le signe de  $Z_{n-1, n}$  et je choisis les signes  $\varepsilon_{n-1}, \varepsilon_n$  de sorte que le dernier terme soit positif, puis je regarde le signe du crochet du terme précédent (une fois que  $\varepsilon_{n-1}, \varepsilon_n$  sont fixés ce crochet a un signe bien déterminé) je choisis  $\varepsilon_{n-2}$  de sorte que ce terme soit positif, etc ..... Il est clair que de proche en proche je peux choisir les  $\varepsilon_i$  de sorte que  $L \geq 0$

Donc on peut affirmer, la plus courte distance de l'origine à une face de  $\mathcal{P}$  est obtenue certainement pour un choix  $\varepsilon_i$  tels que  $L \geq 0$

Donc : 
$$\frac{1}{h_m^2} \geq n$$

Par suite on a toujours :

$$\frac{m_A^2}{M_A^2} = \frac{h_m^2}{d_m^2} = \frac{1}{\text{Max } |d_i|^2} \cdot h_m^2 \leq \frac{1}{n}$$

Enfin si  $\frac{m_A^2}{M_A^2} = \frac{1}{n}$  c'est que  $\text{Max } |d_i|^2$  qui est, comme on l'a vu toujours  $\geq 1$  est effectivement tel = 1 (puisque  $h_m^2 \leq \frac{1}{n}$ )

Mais alors si  $|d_i|^2 = \text{Max } |d_i|^2 = 1$

donc :  $|d_i| = 1$  pour tout  $i$

et puisque  $d_i = \frac{1}{m_i \theta_i}$

cela entraîne

$$\theta_i = \frac{1}{2}$$

D'où :

### Théorème I

Pour une matrice  $A$ , à termes réels, normée en lignes, non singulière

on a 
$$\gamma = \frac{m_A}{M_A} \leq \frac{1}{\sqrt{n}}$$

et une condition nécessaire et suffisante pour qu'elle soit orthonormale

est que 
$$\frac{m_A}{M_A} = \frac{1}{\sqrt{n}} \quad [10]$$

4°) Etude du rapport  $\frac{M_n}{M_n}$  pour A à termes réels, non singulière.

Nous allons généraliser la démonstration précédente au cas où A n'est pas normalisée en lignes. Prenons encore le même système d'axes de coordonnées  $(0; \vec{e}_1, \vec{e}_2, \dots, \vec{e}_n)$ .

Pour le plan,  $\sum_{i=1}^{i=n} \frac{x_i}{\alpha_i} = 1$

la distance à l'origine est toujours donnée par :

où  $V = \begin{pmatrix} \frac{1}{\alpha_1} \\ \vdots \\ \frac{1}{\alpha_i} \\ \vdots \\ \frac{1}{\alpha_n} \end{pmatrix}$   $\frac{1}{h^2} = V^T \cdot G^{-1} \cdot V$

Le sommet  $A_i$  sur la droite  $D_i$ , est tel que les composantes de  $\vec{OA}_i$

sont:  $\vec{OA}_i = \begin{cases} 0 \\ \vdots \\ \alpha_i \leftarrow i \\ \vdots \\ 0 \end{cases}$

Il est aussi tel que  $f_j(A_i) = 0$  pour  $j \neq i$  ( $j = 1, \dots, i-1, i+1, \dots, n$ )

donc  $\Phi(A_i) = |f_i(A_i)| = d_i \cdot p_i$

si l'on pose  $p_i = \|A_i\| = \sqrt{a_{i1}^2 + a_{i2}^2 + \dots + a_{in}^2}$

Pour un sommet du polyèdre  $\mathcal{P}$ :  $\Phi(x) = 1$ , on aura  $d_i = \frac{1}{p_i}$ .

D'autre part, le plan  $\frac{x_i}{\alpha_i} = 1$  (parallèle à  $f_i(x) = 0$  et passant par  $A_i$ ) est à la distance  $d_i$  de O.

Par suite :

$\frac{1}{d_i^2} = V_i^T G^{-1} V_i$  si:  $V_i = \begin{pmatrix} 0 \\ \vdots \\ \frac{1}{\alpha_i} \\ \vdots \\ 0 \end{pmatrix} \leftarrow i$

donc:  $\frac{g^{ii}}{\alpha_i^2} = \frac{1}{d_i^2} = p_i^2$  et  $\alpha_i = \varepsilon_i \frac{\sqrt{g^{ii}}}{p_i}$

On aura encore :

1°)  $d_m = \text{Max}_i |d_i| = \text{Max}_i \frac{\sqrt{g^{ii}}}{p_i}$

2°)  $\frac{1}{h^2} = \sum_{i=1}^{i=n} V_i^T G^{-1} V_i + 2 \sum_{\substack{i < j \\ (i, \dots, n)}} V_i^T G^{-1} V_j$

$\frac{1}{h^2} = \sum_{i=1}^{i=n} p_i^2 + 2 \sum_{i < j} \varepsilon_i \varepsilon_j \frac{g^{ij}}{\sqrt{g^{ii}} \cdot \sqrt{g^{jj}}} p_i \cdot p_j$

Le raisonnement précédemment fait prouve donc que  $h_m$  (puisque'on peut faire un choix des  $\varepsilon_i$  de sorte que le 2° terme soit  $> 0$ )

$\frac{1}{h_m^2} \geq \sum_{i=1}^{i=n} p_i^2 = N^2(A)$

si

$N(A) = \sqrt{\sum_{i,j} a_{ij}^2}$

Donc  $\frac{m_A^2}{M_A^2} = \frac{p_m^2}{\max_i |x_i|^2} \leq \frac{1}{N(A)} \cdot \frac{1}{\max_i |d_i|^2}$

Or posons  $g^{kk} = \max_i g^{ii} \geq 1$  et  $p_0 = \min_i (p_i)$

on aura certainement:  $\frac{m_A^2}{M_A^2} \leq \frac{p_0^2}{N(A)} \cdot \frac{1}{g^{kk}} \leq \frac{p_0^2}{N(A)}$

Si  $\frac{m_A^2}{M_A^2} = \frac{p_0^2}{N(A)}$  c'est que  $g^{kk} = 1$  donc tous les  $g^{ii} = 1 \rightarrow \theta_i = \frac{\pi}{2}$   
 donc A sera orthogonale en lignes d'où:

Théorème II Pour une matrice A, à termes réels, non singulière, on a  $\gamma = \frac{m_A}{M_A} \leq \frac{p_0}{N(A)}$

et une condition nécessaire et suffisante pour que A soit orthogonale en

lignes est que:  $\frac{m_A}{M_A} = \frac{p_0}{N(A)}$

où  $p_0 = \min_i \|A_i\|$  et  $N(A) = \left( \sum_{ij} a_{ij}^2 \right)^{\frac{1}{2}}$ . (10)

Ce théorème va nous fournir la possibilité de définir un conditionnement normalisé.

5° Le conditionnement normalisé C(A)

Soit une matrice A quelconque, à termes réels

Posons  $\begin{cases} C(A) = 0 & \text{si A est singulière} \\ C(A) = \frac{m_A}{M_A} / \frac{p_0}{N(A)} & \text{si A est non singulière} \end{cases}$

Théorème III On a toujours:  $0 \leq C(A) \leq 1$ .

Pour que A soit singulière il faut et il suffit  $C(A) = 0$ ,

pour que A soit orthogonale en lignes il faut et il suffit  $C(A) = 1$ ,

si A est telle que ses éléments varient de sorte que  $A \rightarrow A_1$  et  $A_1$  singulière, sans ligne de longueur = 0,  $C(A) \rightarrow 0$

Les trois premières parties de l'énoncé du théorème II et de la définition de C(A). Pour démontrer la dernière partie il suffit de remarquer

que  $C^2(A) \leq \frac{1}{g^{kk}} = \sin^2 \theta_k$

d'après les relations vues plus haut.

Or si  $A \rightarrow A_1$  singulière <sup>sans ligne nulle</sup>, l'un des  $\theta_i \rightarrow 0$ , donc  $\theta_k \rightarrow 0$

$(\min \theta_k = \min_i |\sin \theta_i|)$

Définition:

le nombre C(A) sera appelé "conditionnement normalisé"

ou, conditionnement (sans plus) si aucune confusion n'est à craindre.

C'est l'étude de cette fonction de matrice que nous allons faire dans ce qui va suivre.

Pour le calcul effectif de  $C(A)$ , on utilisera les déterminations suivantes de  $m_A$  et  $M_A$ .

$$\text{Soit } \Phi(Ax) = \sum_{i=1}^n |f_i(x)| \text{ et formons } r_{12} = \frac{\Phi(Ax)}{\|x\|}$$

si je pose  $Ax = X$  ( $A$  non singulière)

$$\text{on aura: } r_{12} = \frac{\Phi(X)}{\|A^{-1}X\|}$$

Je peux pour chercher le maximum de ce rapport me placer dans l'espace de point courant  $X = (X_i)$  ( $i=1, \dots, n$ )

Soit  $\mathcal{E}$  l'hyper-ellipsoïde d'équation  $\|A^{-1}X\|^2 = 1$ ,  $X^T (A^{-1})^T A^{-1} X = 1$ . Si  $X_0$  se déplace sur celui-ci, chercher le maximum de  $r_{12}$ , revient à chercher le maximum de  $|X_1| + |X_2| + \dots + |X_n| = \Phi(X)$ . Or les polyèdres  $\Phi(X) = \Phi(X_0)$ , varient tous en restant homothétiques et leur convexité assure que le maximum cherché est certainement obtenu l'orsque l'une des faces de ce polyèdre est tangente à  $\mathcal{E}$ .

Or les faces de ces polyèdres, sont perpendiculaires aux vecteurs

$$E = \begin{cases} \varepsilon_1 = \pm 1 \\ \varepsilon_2 = \pm 1 \\ \vdots \\ \varepsilon_n = \pm 1 \end{cases}$$

Dont les composantes sont des  $\pm 1$  ( $2^n$  tels vecteurs)

D'autre part, au point  $X_0$  le plan tangent à  $\mathcal{E}$  a pour équation:

$$X^T (A^{-1})^T A^{-1} X_0 = 1$$

Ce plan est perpendiculaire à  $E$ , s'il existe un scalaire  $k$  tel:

$$(A^{-1})^T A^{-1} X_0 = k \cdot E$$

donc  $X_0 = k \cdot A A^T E$ .

Pour déterminer  $k$ , il suffit d'écrire que  $X_0 \in \mathcal{E}$

$$\text{ce qui donne: } k^2 (A A^T E)^T (A^{-1})^T A^{-1} (A A^T E) = 1$$

$$\text{d'où } k^2 \cdot E^T A A^T E = 1, \quad k^2 = \frac{1}{\|A^T E\|^2}$$

$$\text{D'autre part, le plan } k \cdot X^T E = 1 \quad \text{ou} \quad X_1 \varepsilon_1 + \dots + X_n \varepsilon_n = \frac{1}{k}$$

coupe l'un des axes du système de référence en un point d'abscisse sur  $e_i$

ce point est tel que pour lui,  $\Phi(p_i) = \frac{1}{|k|}$

par suite le maximum cherché, n'est autre que le maximum de  $1/|k|$

d'où:

Théorème IV Le maximum  $M_A$  de  $\gamma$  n'est autre que  $\text{Max}_E \|A^T E\|$  où  $E$  parcourt l'ensemble des vecteurs dont les composantes ont pour valeur absolue des 1.

II)  $m_A$

Pour déterminer  $m_A$  je considère toujours  $e_i$  mais je suppose déplacer le

point  $X_1$  sur le polyèdre  $\Phi(X) = 1$

la valeur de  $\|A^T X_1\|$  varie et, toujours à cause de la convexité le

maximum de  $\|A^T X_1\|$  donc le minimum de  $\gamma$  ne peut être atteint qu'en un sommet

de  $\Phi(X)=1$ . Or ces sommets sont les points

$$e_i = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \leftarrow i \\ 0 \\ 0 \end{pmatrix}$$

et la valeur de

$$m_A = \frac{1}{\text{Max}_i \|A^T e_i\|}$$

d'où :

Théorème V Le minimum  $m_A$  de  $\gamma$  n'est autre que  $1/\text{Max}_i \|A^T e_i\|$ , c'est à dire, l'inverse de la longueur de la colonne de plus grande longueur de la matrice inverse de  $A$ .

Par suite, on aura

$$(I) \quad C(A) = \frac{N(A)}{p_0(A)} \cdot \frac{1}{\text{Max}_i \|A^T e_i\|} \cdot \frac{1}{\text{Max}_E \|A^T E\|}$$

C'est cette formule qui en définitive nous servira à calculer des  $C(A)$ ,

dès que l'inverse d'une matrice est calculable il est en effet assez commode de l'utiliser, bien que  $\text{Max}_E \|A^T E\|$  soit, en général, assez difficile à évaluer.

Nous étudierons dans la suite des propriétés générales de  $C(A)$ , soit au cours de transformations de  $A$ , soit au cours d'opérations matricielles sur  $A$ .

7°) Le conditionnement normalisé complémentaire C'(A) .

Nous allons rapidement voir comment on peut définir un autre nombre de conditionnement C'(A) qui a des propriétés analogues à celles de C(A) précédemment défini. Cette fois partons du rapport  $r_{\infty 2}(x)$

on posera :  $r_{\infty 2} = \frac{M_b(Ax)}{\|x\|}$  ,  $m'_A = \min_{x \neq 0} (r_{\infty 2}(x))$  ,  $M'_A = \max_{x \neq 0} (r_{\infty 2}(x))$   
 $\gamma_{\infty 2} = \gamma' = \frac{m'_A}{M'_A}$

Comme précédemment étudions avec les mêmes notations le rapport  $\gamma'$ .

Pour déterminer  $m'_A$  et  $M'_A$  il suffit encore, pour A non singulière de déplacer le point (x) sur le lieu des points (x) tels que  $M_b(Ax) = 1$ . Ce lieu est la frontière  $\Pi'$  d'un paralléloétope  $\mathcal{P}'$ , ayant l'origine pour centre de symétrie, les faces planes de celui-ci sont des parties des plans d'équation:

$$f_i(x) = \pm 1$$

Donc encore  $m'_A = \frac{1}{d'_n}$  et  $M'_A = \frac{1}{h'_m}$

mais ici,

1°)  $h'_m$  est très facile à déterminer: c'est  $\min_i \left( \frac{1}{p_i} \right) = \frac{1}{\max_i (p_i)} = \frac{1}{p^0}$

2°) Pour pouvoir calculer  $d'_n$  il suffit de remarquer que les sommets de  $\mathcal{P}'$  sont les points dont les coordonnées sont les nombres  $\alpha_i$  du § 3°) dans  $(0; \vec{e}_1, \vec{e}_2, \dots, \vec{e}_n)$

Si  $V'_i = \begin{pmatrix} 0 \\ \vdots \\ \alpha_i \leftarrow i^o \\ \vdots \\ 0 \end{pmatrix}$

Le vecteur  $U' = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix}$  joint 0 à un sommet de  $\mathcal{P}'$ .

Le carré de la longueur est fourni par:  $l^2 = U'^T G U'$

ou:  $l^2 = \sum_{i=1}^{i=n} V_i'^T G V_i' + 2 \sum_{\substack{i < j \\ (1, \dots, n)}} V_i'^T G V_j'$

D'après cela,

$$l^2 = \sum_{i=1}^{i=n} \frac{g_{ii}}{p_i^2} + 2 \sum_{i < j} \varepsilon_i \varepsilon_j \frac{\sqrt{g_{ii}} \cdot \sqrt{g_{jj}}}{p_i p_j} g_{ij}$$

puisque ( d'après le même raisonnement que celui fait plus haut)

1° on peut en déduire:  $l_n^2 \geq \sum_{i=1}^{i=n} \frac{1}{\sin^2 \theta_i \cdot p_i^2}$

$$d_n^{1/2} \geq \sum_{i=1}^{i=n} \frac{1}{p_i^2} = 1/N^{1/2}(A) \quad (*)$$

(\*) La notation  $N^{1/2}(A)$  n'est utilisée que par analogie;  $N(A)$  n'est pas une norme de matrice.

Par suite :  $\frac{m_A'^2}{M_A'^2} = \frac{h_m'^2}{d_M'^2} = \frac{1}{p^{\circ 2}} \cdot \frac{1}{d_M'^2} \leq \frac{1}{(p^{\circ})^2} \cdot \frac{1}{\sum_i (\frac{1}{p_i})^2}$

d'où :  $\frac{m_A'}{M_A'} \leq \frac{1}{p^{\circ}} \cdot \left( \frac{1}{\sum_i (\frac{1}{p_i})^2} \right)^{\frac{1}{2}} = \frac{N'(A)}{p^{\circ}}$

Inversement, supposons que A soit telle que :  $\frac{m_A'}{M_A'} = \frac{N'(A)}{p^{\circ}}$

puisque  $M_A' = p^{\circ}$

et que toujours  $m_A'^2 = \frac{1}{d_m'^2} \leq 1 / \sum_{i=1}^{i=n} \frac{1}{h_m'^2 \theta_i \cdot p_i^2} \leq N'^2(A)$

C'est que  $\sin^2 \theta_i = 1$  pour tout  $\theta_i$ .

Donc, on peut aussi, en posant :

$$\begin{cases} C'(A) = 0 & \text{si A est singulière} \\ C'(A) = \frac{m_A'}{M_A'} \frac{p^{\circ}}{N'(A)} & \text{si A est non singulière} \\ & = \frac{m_A'}{N'(A)} \end{cases}$$

énoncer :

Théorème VI : On a toujours  $0 \leq C^{\circ}(A) \leq 1$

Pour que A soit singulière, il faut et il suffit que  $C^{\circ}(A) = 0$ .

Pour que A soit orthogonale en lignes, il faut et il suffit que  $C^{\circ}(A) = 1$ .

Si A tend vers une matrice singulière  $A_0$ , sans ligne nulle, alors  $C^{\circ}(A) \rightarrow 0$ .

On peut appeler  $C^{\circ}(A)$  le conditionnement normalisé complémentaire de  $C(A)$

Étudions  $m_A'$  d'une autre façon :

$$m_A' = \min_x \left( \frac{\mathcal{M}_0(Ax)}{\|x\|} \right) = \min_x \left( \frac{\mathcal{M}_0(x)}{\|A^{-1}x\|} \right)$$

Si je déplace le point X sur la frontière de  $\mathcal{M}_0(x) \leq 1$ , le maximum de  $\|A^{-1}x\|$  ne peut arriver que en un sommet de ce polyèdre, dont les coordonnées sont les nombres

$$E = \begin{pmatrix} \pm 1 \\ \pm 1 \\ \vdots \\ \pm 1 \end{pmatrix} \quad \text{donc} \quad m_A' = \frac{1}{\max_E \|A^{-1}E\|}$$

et  $C'(A) = 1 / N'(A) \cdot \max_E \|A^{-1}E\|$

formule à rapprocher de celle du théorème V.

### 8°) Relations entre ces deux conditionnements.

Les rapports  $\gamma$  et  $\gamma'$  sont liés par : (théorème IV du chapitre III)

$$\frac{1}{n} \gamma' \leq \gamma \leq n \gamma'$$

ou

$$\frac{1}{n} \gamma \leq \gamma' \leq n \gamma$$

Donc, on peut écrire :

$$\frac{1}{n} \gamma' \cdot \frac{p_0}{N'(A)} \leq \gamma \cdot \frac{N(A)}{p_0} \cdot \frac{p_0 \cdot p_0}{N(A) \cdot N'(A)} \leq n \gamma' \cdot \frac{p_0}{N'(A)}$$

et, par suite,

$$\frac{1}{n} C'(A) \leq C(A) \cdot \frac{p_0 \cdot p_0}{N(A) \cdot N'(A)} \leq n C'(A)$$

enfin :

$$\frac{1}{n} C'(A) \cdot \frac{N(A) \cdot N'(A)}{p_0 \cdot p_0} \leq C(A) \leq n C'(A) \cdot \frac{N(A) \cdot N'(A)}{p_0 \cdot p_0}$$

D'où les relations :

$$\frac{1}{n} C(A) \cdot \frac{p_0 \cdot p_0}{N(A) \cdot N'(A)} \leq C'(A) \leq n C(A) \cdot \frac{p_0 \cdot p_0}{N(A) \cdot N'(A)}$$

On peut remarquer que :

$$\frac{p_0}{p_0} \leq \frac{p_0 \cdot p_0}{N(A) \cdot N'(A)} \leq \frac{p_0}{p_0}$$

Il résulte de ces relations que les deux conditionnements  $C(A)$  et  $C'(A)$  sont pratiquement d'utilisation équivalente avec, pour le calcul, une supériorité pour  $C(A)$  qui, d'après ce qui précède, se calcule plus simplement. C'est la raison pour laquelle cette fonction sera surtout utilisée dans la suite.

### 9°) Relations entre $C(A)$ et les nombres de VON NEUMANN et de TURING.

I) Pour  $A$ , symétrique, définie positive, posons :

$$P(A) = \frac{\lambda_M}{\lambda_m}, \quad \text{nombre de VON NEUMANN,}$$

( $\lambda_M, \lambda_m$  : plus grande et plus petite des valeurs propres de  $A$ )

Comparons le à  $\gamma$  et  $C(A)$ .

- a) Puisque : 
$$\Phi(Ax) = \sum_{i=1}^{i=n} |f_i(x)|$$

on a (cf. Chapitre II § I) 
$$\|Ax\| \leq \Phi(Ax) \leq \sqrt{n} \|Ax\|$$

donc : 
$$\frac{\|Ax\|^2}{\|x\|^2} \leq \frac{\Phi^2(Ax)}{\|x\|^2} \leq n \frac{\|Ax\|^2}{\|x\|^2}$$

par suite :  $\frac{x^T A^T A x}{x^T x} \leq m_A^2 \leq \frac{\Phi^2(Ax)}{\|x\|^2} \leq M_A^2 \leq n \cdot \frac{x^T A^T A x}{x^T x}$

et, dans le cas où A est symétrique :  $\lambda_m^2 \leq m_A^2 \leq M_A^2 \leq n \lambda_M^2$

Donc  $\frac{1}{n} \frac{\lambda_m^2}{\lambda_M^2} \leq \frac{m_A^2}{M_A^2}$  ou  $\boxed{\frac{1}{n} \leq P \cdot \gamma}$  (I)

l'égalité est réalisée pour  $A = I$

et, par suite  $\frac{1}{P(A)} \leq \sqrt{n} \cdot C(A) \cdot \frac{P_0}{N(A)}$

- b) Je dis que :  $P \cdot \gamma \leq \sqrt{n}$

Soit x le vecteur propre correspondant à la plus petite valeur propre, de composantes  $x_i$  (i=1, ..., n)

y " " " " à la plus grande valeur propre, de composantes  $y_i$  (i=1, ..., n)

et ces deux vecteurs de module euclidien = 1

$$m_A \leq \lambda_m \cdot \sum_i |x_i|, \quad M_A \geq \lambda_M \sum_i |y_i|$$

$$\frac{m_A}{M_A} \leq \frac{\lambda_m}{\lambda_M} \cdot \frac{\sum |x_i|}{\sum |y_i|}$$

Il est évident que :

$$1 \leq \sum_i |x_i| \leq \sqrt{n}, \quad \text{si } \|x\| = 1$$

donc :

$$\frac{m_A}{M_A} \leq \frac{\lambda_m}{\lambda_M} \cdot \sqrt{n}, \quad \boxed{P \cdot \gamma \leq \sqrt{n}} \quad \text{(II)}$$

On aurait pu obtenir (I) et (II) des inégalités générales établies au chapitre III, Th. IV.

Si l'on considère :

$$A = \begin{pmatrix} \lambda & 0 & 0 & \dots & 0 \\ 0 & (\lambda + \varepsilon - 1)\eta & -\eta & \dots & -\eta \\ 0 & -\eta & (\lambda + \varepsilon - 1)\eta & & \\ \vdots & & & \ddots & \\ \vdots & & & & -\eta \\ 0 & -\eta & -\eta & \dots & (\lambda + \varepsilon - 1)\eta \end{pmatrix}$$

on peut choisir  $\varepsilon$  et  $\eta$  de sorte que  $P(A) \cdot \gamma(A)$  soit aussi près que l'on veut de  $\sqrt{n}$

II) Soit le premier nombre de TURING

$$T_1(A) = n \cdot \text{Max}_{ij} |a_{ij}| \cdot \text{Max}_{ij} |\alpha_{ij}|$$

$\alpha_{ij}$  terme général de  $A^{-1}$

Il est à remarquer d'abord que  $T_1(A) \geq 1$

car  $\sum_k a_{ik} \cdot \alpha_{ki} = 1$  implique  $1 \leq n \cdot \text{Max}_{ik} |a_{ik}| \cdot \text{Max}_{ki} |\alpha_{ki}|$

et à plus forte raison,  $T_1(A) \geq 1$

- a)  $\frac{1}{\gamma} = \text{Max}_i \|A^{-1} e_i\| \cdot \text{Max}_j \|A^T E_j\|$   
 $\leq \sqrt{n} \text{Max}_{ij} |\alpha_{ij}| \sqrt{n} \cdot n \text{Max}_{ij} |a_{ij}| \leq n T_1$

(III)  $\frac{1}{n} \leq T_1 \cdot \gamma$

- b) d'autre part

$$\text{Max}_i \|A^{-1} e_i\| \geq \text{Max}_{ij} |\alpha_{ij}|$$
$$\text{Max}_j \|A^T E_j\| \geq \text{Max}_{ij} |a_{ij}|$$

donc :

$\frac{1}{\gamma} \geq \frac{T_1}{n}$  (IV)  $T_1 \cdot \gamma \leq n$

Remarquons que des exemples simples prouvent que les bornes sont atteintes.

III) Soit enfin le 2ème nombre de TURING

$$T_2(A) = \frac{1}{n} N(A) \cdot N(A^{-1})$$

- a) Puisque

$$\text{Max}_i \|A^{-1} e_i\| \geq \frac{1}{\sqrt{n}} N(A^{-1})$$
$$\text{Max}_j \|A^T E_j\| \geq N(A)$$

on a :

(V)  $T_2 \cdot \gamma \leq \frac{1}{\sqrt{n}}$

- b) Enfin

$$\text{Max}_i \|A^{-1} e_i\| \leq N(A^{-1})$$
$$\text{Max}_j \|A^T E_j\| \leq \sqrt{n} \cdot N(A)$$

donc :

(VI)  $T_2 \gamma \geq \frac{1}{n\sqrt{n}}$

En résumé, voici les inégalités obtenues :

$\frac{1}{n} \leq P \cdot \gamma \leq \sqrt{n}$	(I), (II)
$\frac{1}{n} \leq T_1 \cdot \gamma \leq n$	(III), (IV)
$\frac{1}{n\sqrt{n}} \leq T_2 \cdot \gamma \leq \frac{1}{\sqrt{n}}$	(V), (VI)

les inégalités pouvant être atteintes.

On obtient les relations avec  $C(A)$  en multipliant par  $\frac{N(A)}{p_0}$   
puisque  $C(A) = \gamma(A) \cdot \frac{N(A)}{p_0}$

10°) OPérations laissant C(A) invariant .

Théorème VII Le conditionnement général  $\gamma_{ij}$ , relatif aux deux normes  $\varphi_i$  et  $\varphi_j$ , d'une matrice A ne change pas si A est remplacée par  $\lambda.A$ , pour  $\lambda$  scalaire non nul. Il en est de même pour C(A).

En effet si  $m = \min_{x \neq 0} \left( \frac{\varphi_i(Ax)}{\varphi_j(x)} \right)$  et  $M = \max_{x \neq 0} \left( \frac{\varphi_i(Ax)}{\varphi_j(x)} \right)$   
 on aura:  $m' = \min_{x \neq 0} \left( \frac{\varphi_i(\lambda Ax)}{\varphi_j(x)} \right) = |\lambda| m$  ,  $M' = |\lambda| M$

donc  $\gamma'_{ij} = \frac{m'}{M'} = \frac{m}{M} = \gamma_{ij}$

De plus, dans le cas du conditionnement normalisé C(A)

donc:  $\frac{\rho_0(A)}{N(A)} = \frac{\rho_0(\lambda A)}{N(\lambda A)} = \frac{|\lambda| \rho_0(A)}{|\lambda| N(A)}$   
 $C(A) = C(\lambda.A)$

Théorème VIII Le conditionnement général  $\gamma_{ij}$ , relatif aux deux normes  $\varphi_i$  et  $\varphi_j$ , d'une matrice A ne change pas si A est remplacée par  $Q_i A Q_j$  avec:

- $Q_i$  appartenant au groupe associé à la norme de vecteur  $\varphi_i$
- $Q_j$  appartenant au groupe associé à la norme de vecteur  $\varphi_j$

En effet si  $r_{ij} = \frac{\varphi_i(Ax)}{\varphi_j(x)} = \frac{\varphi_i(Q_i Ax)}{\varphi_j(x)}$

et puisque  $Q_j$  appartient au groupe associé à  $\varphi_j$  (cf chapitre II)

on aura  $r_{ij} = \frac{\varphi_i(Q_i A Q_j y)}{\varphi_j(Q_j y)} = \frac{\varphi_i(Q_i A Q_j y)}{\varphi_j(y)}$

donc  $r_{ij}(A, x) = r_{ij}(Q_i A Q_j, y)$

et par suite le rapport :  $\gamma_{ij}(A) = \gamma_{ij}(Q_i A Q_j)$

Corollaire Le conditionnement normalisé C(A) est invariant si A est remplacée par J A Q

Si J est soit un produit de forme  $(I - 2E_{ii}) \cdot (I - 2E_{ee}) \dots$

soit un matrice d'échange de deux lignes  $V = (I - E_{ii} - E_{jj} + E_{ij} + E_{ji})$   
 (Produit de matrices du 2° degré bien particulières)

et si Q est une matrice orthonormale ( $Q^T Q = I$ )

En effet il est clair que  $I - 2E_{ii} = J$  n'est autre que A où la i° ligne a un signe contraire. Donc quelque soit x,

$$\Phi(J.Ax) = \Phi(Ax)$$

et par suite J est du groupe associé à la norme  $\Phi(x)$  de vecteurs

De même  $V = I - E_{ii} - E_{jj} + E_{ij} + E_{ji}$  appartient à ce groupe.

Il est aussi évident que  $\rho_0(JA) = \rho_0(A)$ ,  $N(JA) = N(A)$ .

D'autre part, le rapport  $\gamma = \frac{m_A}{M_A}$  d'après le théorème précédent ne change pas si A est remplacée par JAQ.

Enfin, si  $l_i$  est la longueur d'une ligne de A

on a  $l_i^2 = A_i \cdot A_i^T$

puis  $l_i^2 = A_i \cdot Q \cdot Q^T \cdot A_i^T = (A \cdot Q)_i \cdot (A \cdot Q)_i^T$

enfin :  $l_i^2 = l_i'^2 \rightarrow \rho_0(A) = \rho_0(A \cdot Q)$

Comme

$$N^2(A) = \text{trace}(A \cdot A^T) = \text{trace}(A \cdot Q \cdot Q^T \cdot A^T) = N^2(AQ)$$

par suite le rapport  $\frac{\rho_0(A)}{N(A)} = \frac{\rho_0(AQ)}{N(AQ)}$

Application : Signalons comme conséquence de ce corollaire une

application que nous étudierons plus en détails dans le chapitre consacré à l'étude de la méthode d'orthogonalisation.

Dans la partie théorique relative à cette méthode on peut prouver le résultat suivant: Etant donnée une matrice A non singulière il existe toujours une matrice triangulaire inférieure et unitaire, T telle que  $T \cdot A = Q$ , avec Q orthogonale en ligne

soit  $Q = \begin{pmatrix} q_1 \\ q_2 \\ \vdots \\ q_n \end{pmatrix}$  les  $q_i$  étant des lignes

de cette matrice, si je pose  $d_{ii} = \|q_i\| = \sqrt{q_i \cdot q_i^T}$

et  $D = \begin{vmatrix} d_{11} & & \\ & d_{22} & \\ & & \ddots \\ & & & d_{nn} \end{vmatrix}$ ,  $Q = D \cdot Q'$

je vois que  $Q'$  est orthonormale

Donc  $T \cdot A = D \cdot Q'$

et  $A = T^{-1} \cdot D \cdot Q' = T' \cdot Q'$

Théorème IX Le conditionnement normalisé  $C(A)$  d'une matrice A (non singulière) est celui de la matrice triangulaire  $C(T)$ , qui se calcule aisément dans la méthode d'orthogonalisation.

11<sup>e</sup>) Conditionnement d'un produit de deux matrices.

Soient A et B deux matrices non singulières,

formons 
$$\frac{\Phi(A.Bx)}{\|x\|} = \frac{\Phi(Ax)}{\|x\|} \cdot \frac{\|Bx\|}{\|x\|}$$

si  $X = Bx$ .

on sait que (théorème V du Chapitre III)  $\lambda_m^{\frac{1}{2}} \leq \frac{\|Bx\|}{\|x\|} \leq \lambda_M^{\frac{1}{2}}$

si  $\lambda_1, \lambda_m$  désignent les plus grandes et plus petites valeurs propres de la matrice  $B.B^T$ .

D'où 
$$m_A \cdot \lambda_m^{\frac{1}{2}} \leq m_{AB} \leq M_{AB} \leq M_A \cdot \lambda_M^{\frac{1}{2}}$$

et 
$$\frac{m_A}{M_A} \left( \frac{\lambda_m}{\lambda_M} \right)^{\frac{1}{2}} \leq \frac{m_{AB}}{M_{AB}}, \quad \gamma_{12}(A) \cdot \gamma_{22}(B) \leq \gamma_{12}(A.B)$$

enfin

$$C_{AB} \geq \left( \frac{\lambda_m}{\lambda_M} \right)^{\frac{1}{2}} \cdot C_A \cdot \frac{\rho_0(A)}{N(A)} \cdot \frac{N(AB)}{\rho_0(AB)}$$

si l'on se rappelle que (Théorème IV du chapitre III)

$$\frac{1}{\sqrt{n}} \gamma_{22}(B) \leq \gamma_{12}(B) \leq \sqrt{n} \gamma_{22}(B)$$

on voit que :

(\*) 
$$C_{AB} \geq \frac{C_A}{\sqrt{n}} \cdot \frac{m_B}{M_B} \cdot \frac{\rho_0(A)}{N(A)} \cdot \frac{N(AB)}{\rho_0(AB)}$$

donc :

Theorème X Entre les conditionnements  $C_{AB}$  d'un produit A.B de deux matrices A et B non singulières on a la relation:

$$1 \geq C(AB) \geq \frac{C_A \cdot C_B}{\sqrt{n}} \cdot \frac{\rho_0(A) \cdot \rho_0(B)}{\rho_0(A.B)} \cdot \frac{N(A.B)}{N(A) \cdot N(B)} \quad (1)$$

En particulier, pour que le produit de deux matrices soit orthogonal en lignes il suffit que:

$$\sqrt{n} = C_A \cdot C_B \cdot \frac{\rho_0(A) \cdot \rho_0(B)}{\rho_0(AB)} \cdot \frac{N(A.B)}{N(A) \cdot N(B)}$$

Corollaire Entre les conditionnements de A et de son inverse  $A^{-1}$ , il est la relation:

(\*) 
$$1 \geq C_A \cdot C_{A^{-1}} \cdot \frac{\rho_0(A)}{N(A)} \cdot \frac{\rho_0(A^{-1})}{N(A^{-1})}$$

En effet il suffit de remarquer que  $\frac{N(I)}{\rho_0(I)} = \sqrt{n}$  et de faire  $AB=I$  dans la formule (1).

(\*) On montre aisément, sur des exemples simples, qu'il ne peut y avoir d'inégalité de sens inverse à celle-ci.

12°) Exemple de calcul d'un conditionnement normalisé, Matrice de HILBERT  $H_n$ .

Définition: On appelle matrice de HILBERT d'ordre  $n$  la matrice symétrique de type  $(n,n)$  dont le terme général est

$$h_{ij}^{(n)} = \frac{1}{i+j-1} \quad \text{pour } i=1,2,\dots,n; \quad j=1,2,\dots,n$$

c'est à dire que

$$H_n = \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \dots & \frac{1}{n} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \dots & \dots & \frac{1}{n+1} \\ \frac{1}{3} & \frac{1}{4} & \dots & \dots & \dots & \vdots \\ \frac{1}{4} & \dots & \dots & \dots & \dots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{1}{n} & \frac{1}{n+1} & \dots & \dots & \dots & \frac{1}{2n-1} \end{pmatrix}$$

Cette matrice a été très fréquemment étudiée cf. [19] [20] [21] [22] et aussi [15]. La raison en est que l'on sait, formellement, en déterminer l'inverse  $T_n$ , d'autre part l'inversion numérique de cette matrice, même pour de très faibles valeurs de  $n$  est très difficile. A ce sujet on se rapportera aux expériences de calcul de  $H_n^{-1}$ , que nous expliquerons dans le chapitre relatif à la résolution des systèmes linéaires.

On sait que si  $T_n = H_n^{-1} = (t_{ij})$

$$t_{ij} = (-1)^{i+j} \frac{(n+j-1)! (n+i-1)!}{(i+j-1) [(i-1)! (j-1)!]^2 (n-i)! (n-j)!}$$

Donc :

$$(1) \quad \left| \frac{t_{i,j+1}}{t_{ij}} \right| = \left( 1 - \frac{1}{i+j} \right) \cdot \left( 1 + \frac{n^2 - 2j^2}{j^2} \right)$$

Nous nous proposons d'évaluer  $C(H_n)$  pour des valeurs de  $n$  assez grandes.

- I°) Détermination de  $\rho_0(H_n)$ .

On a évidemment  $\rho_0^2(H_n) = \frac{1}{n^2} + \frac{1}{(n+1)^2} + \dots + \frac{1}{(2n-1)^2}$

en écrivant :  $\frac{1}{n^2} = \frac{1}{n-1} - \frac{1}{n} = \frac{1}{n^2(n-1)}$   $i=2n-1$

on voit que :  $\rho_0^2(H_n) = \frac{1}{n-1} - \frac{1}{2n-1} - \sum_{i=n}^{2n-1} \frac{1}{i^2(i-1)}$

La somme du 2° membre est inférieure à

$$\frac{n}{n^2(n-1)} = \frac{1}{n(n-1)}$$

On peut donc dire que si  $n$  augmente :  $\rho^c(H_n) \sim \frac{1}{n-1} - \frac{1}{2n-1} \sim \frac{1}{2n}$

2°) Détermination de  $N(H_n)$

$$\begin{aligned} \text{Ecrivons : } N^2(H_n) &= 1 + \left(\frac{1}{2}\right)^2 + \left(\frac{1}{3}\right)^2 + \dots + \left(\frac{1}{n}\right)^2 \\ &+ \left(\frac{1}{2}\right)^2 + \left(\frac{1}{3}\right)^2 + \dots + \left(\frac{1}{n}\right)^2 + \left(\frac{1}{n+1}\right)^2 \\ &+ \left(\frac{1}{3}\right)^2 + \dots + \left(\frac{1}{n}\right)^2 + \left(\frac{1}{n+1}\right)^2 + \left(\frac{1}{n+2}\right)^2 \\ &\dots \\ &+ \left(\frac{1}{n}\right)^2 + \left(\frac{1}{n+1}\right)^2 + \dots + \left(\frac{1}{2n-1}\right)^2 \end{aligned}$$

Donc :  $N^2(H_n) = 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n} + \varphi(n)$

Regardons :  $\varphi(n) = \frac{n-1}{(n+1)^2} + \frac{n-2}{(n+2)^2} + \dots + \frac{n-(n-1)}{(2n-1)^2}$

en écrivant sous la forme  $\varphi(n) = \frac{1}{n} \left[ \frac{1 - \frac{1}{n}}{\left(1 + \frac{1}{n}\right)^2} + \frac{1 - \frac{2}{n}}{\left(1 + \frac{2}{n}\right)^2} + \dots + \frac{1 - \frac{n-1}{n}}{\left(1 + \frac{n-1}{n}\right)^2} \right]$

On voit que si  $n \rightarrow \infty$  cela a pour limite  $\int_0^1 \frac{1-x}{(1+x)^2} dx = 1 - \log 2$

par suite, puisque  $1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n} \sim \log n + C$

on peut écrire:

$$N^2(H_n) \sim 1 - \log 2 + C + \log n = A + \log n, \quad A = 0,8841$$

3°) Détermination de  $M_{H_n} = \max_E \|A^T E\|$ , ou ici puisque  $H_n$  est symétrique  $= \max_E \|AE\|$

$$\begin{aligned} M_{H_n}^2 &= \max_E \left[ \left( \varepsilon_1 + \frac{\varepsilon_2}{2} + \frac{\varepsilon_3}{3} + \dots + \frac{\varepsilon_n}{n} \right)^2 + \left( \frac{\varepsilon_1}{2} + \frac{\varepsilon_2}{3} + \dots + \frac{\varepsilon_n}{n+1} \right)^2 + \dots + \left( \frac{\varepsilon_1}{n} + \frac{\varepsilon_2}{n+1} + \dots + \frac{\varepsilon_n}{2n-1} \right)^2 \right] \\ &= \left[ \left( 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n} \right)^2 + \left( \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n+1} \right)^2 + \dots + \left( \frac{1}{n} + \frac{1}{n+1} + \dots + \frac{1}{2n-1} \right)^2 \right] \end{aligned}$$

En remarquant que  $(z \neq 1)$

$$\frac{1}{z} + \frac{1}{z+1} + \dots + \frac{1}{z+n-1} \sim \log(z+n-1) - \log(z-1) = \log\left(1 + \frac{n}{z-1}\right)$$

On peut écrire:

$$\left(\frac{1}{2} + \dots + \frac{1}{k+1}\right)^2 + \dots + \left(\frac{1}{n} + \dots + \frac{1}{2n-1}\right)^2 \sim \left(\log\left(1 + \frac{n}{1}\right)\right)^2 + \left(\log\left(1 + \frac{n}{2}\right)\right)^2 + \dots + \left(\log\left(1 + \frac{n}{n-1}\right)\right)^2$$

Il suffit de remarquer que l'intégrale  $\int_0^1 \left[ \log\left(1 + \frac{1}{x}\right) \right]^2 dx$

a un sens, Soit  $k$  sa valeur pour pouvoir écrire

$$M_{H_n}^2 \simeq n \cdot k + (C + \log n)^2$$

4°) Détermination de  $m_{H_n} = 1/M_{\max} \|A^{-1}e_i\|$

De la formule (1) on peut déduire avec TODD [15]

$$\max_i \|T_n e_i\| \# n^{\frac{1}{2}} \max_{i,j} (|t_{ij}|) \sim A_1 n^{-\frac{1}{2}} e^{nB}, \quad B = 3,525$$

D'après ce qui précède on voit que une valeur asymptotique de  $C(H_n)$

$$\text{est } C(H_n) \simeq \alpha \cdot n^{\frac{1}{2}} (\log n)^{\frac{1}{2}} e^{-nB}$$

Cela montre la rapidité avec laquelle  $C(H_n) \rightarrow 0$  et explique le mauvais conditionnement de  $H_n$

### 13°) Cas d'une matrice d'interpolation d'opérateur différentiel

Nous allons retrouver le même genre de difficulté dans le cas d'un autre problème numérique. Soit à résoudre numériquement le problème:

$$y'' = F(x, y)$$

Résoudre l'équation différentielle  $y'' = F(x, y)$ ,  $0 \leq x \leq 1$

avec  $y(0) = y_0$  et  $y(1) = y_n$ . (problème de conditions aux limites)

Il est classique de savoir que si l'on remplace l'équation proposée par le système linéaire

$$y_{i-1} - 2y_i + y_{i+1} = h^2 F(x_i, y_i) \quad (i = 1, \dots, n)$$

$$x_i = \frac{i}{n}, \quad h = \frac{1}{n}$$

celui-ci peut s'écrire

$$(1) \quad A_n y = b$$

Avec

$$A_n = \begin{vmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & -1 & 2 & -1 & \\ & & & \ddots & \\ & & & & -1 & 2 & -1 \\ & & & & & -1 & 2 \end{vmatrix}, \quad Y = \begin{vmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{vmatrix}, \quad b = \begin{vmatrix} y_0 - h^2 F_1 \\ -h^2 F_2 \\ \vdots \\ y_{n+1} - h^2 F_n \end{vmatrix}, \quad F_u = F(x_u, y_u)$$

Il n'est pas difficile de déterminer  $A_n^{-1}$  (par exple en triangularisant)

On trouve ainsi :

$$A_n^{-1} = \frac{1}{n+1} \begin{vmatrix} 1 \cdot n & 1 \cdot (n-1) & 1 \cdot (n-2) & \dots & 1 \cdot 1 \\ 1 \cdot (n-1) & 2 \cdot (n-1) & 2 \cdot (n-2) & \dots & 2 \cdot 1 \\ 1 \cdot (n-2) & 2 \cdot (n-2) & 3 \cdot (n-2) & \dots & 3 \cdot 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 \cdot 1 & 2 \cdot 1 & 3 \cdot 1 & \dots & n \cdot 1 \end{vmatrix}$$

De cette expression on tire.

$$1^0) \quad N^2(A_n) = 4 + 1 + (n-2)(1+4+1) + 4 + 1 = 2(3n-1)$$

$$\rho^2(A) = 4 + 1 = 5$$

$$2^0) \quad M_A = \max_E \|A^T E\| = \max_E \|AE\| \quad \text{si} \quad E = \begin{vmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{vmatrix}$$

$$\|AE\|^2 = (2\varepsilon_1 - \varepsilon_2)^2 + (-\varepsilon_1 + 2\varepsilon_2 - \varepsilon_3)^2 + (\varepsilon_2 + 2\varepsilon_3 - \varepsilon_4)^2 + \dots + (-\varepsilon_{n-1} + 2\varepsilon_n)^2$$

Il est évident que le maximum est obtenu pour une suite de signes  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ , qui alternent régulièrement.

et alors,

$$M_A^2 = 3^2 + (n-2)4^2 + 3^2 = 2(8n-7)$$

Il nous reste à étudier  $m_A = \frac{1}{\max \|A^{-1}e_i\|}$

On peut écrire:

$$A_n^{-1} e_i = \frac{1}{n+1} \begin{pmatrix} 1 \cdot (n-i+1) \\ 2 \cdot (n-i+1) \\ \vdots \\ i \cdot (n-i+1) \\ i \cdot (n-i) \\ \vdots \\ i \cdot 1 \end{pmatrix} \begin{matrix} \leftarrow i^0 \\ \\ \\ \leftarrow (i+1)^0 \\ \\ \end{matrix}$$

Donc :

$$\begin{aligned} \|A_n^{-1} e_i\|^2 &= \frac{1}{(n+1)^2} \left[ (n-i+1)^2 (1+2^2+\dots+i^2) + i^2 (1+2^2+\dots+(n-i)^2) \right] \\ &= \frac{1}{(n+1)^2} \left[ (n-i+1)^2 \frac{(i+1)i(2i+1)}{6} + i^2 \frac{(n-i)(n-i+1)(2n-2i+1)}{6} \right] \\ &= \frac{i(n-i+1)}{6(n+1)^2} \left[ 2i(n^2+n+1) + (n+1)(i-2i^2) \right] \end{aligned}$$

Le 2° facteur est maximum pour  $i = \frac{n}{2} + \frac{1}{2(n+1)} \sim \frac{n}{2}$

Le 1° facteur est maximum pour  $i = \frac{n+1}{2} \sim \frac{n}{2}$

par suite  $\text{Max}_i \|A_n^{-1} e_i\|^2 \sim \frac{n^2}{4 \cdot 6(n+1)^2} \times \frac{n^3}{2} \sim \frac{n^3}{48}$

$$\text{Donc Pour } C(A_n) : C(A_n) \sim \frac{1}{\sqrt{2(8n-7)}} \cdot \sqrt{\frac{48}{n^3}} \cdot \sqrt{\frac{2}{5}(3n-1)}$$

$$C(A_n) \sim \frac{1,8}{\sqrt{n^3}}$$

On peut enconclure que l'inversion numérique du système  $A_n y = b$  est d'autant plus difficile que n est grand .Or cette inversion est engénéral utilisée si les  $F_i$  varient peu,oubien dans certaines méthodes itératives .On voit que les opérateurs aux différences fournissent des systèmes linéaires de plus en plus mal conditionnés et cela au fur et à mesure que l'on augmente l'ordre, c'est à dire si l'on désire diminuer l'erreur de méthode . Ce que l'on veut gagner d'un côté est perdu de l'autre. Cette décroissance est cependant beaucoup plus lente que dans des cas analogues à celui du paragraphe précédent.

## CHAPITRE V

### Les erreurs dans la résolution de systèmes linéaires par élimination

#### 1°) Les principes de ces calculs d'erreurs.

La résolution de systèmes du premier degré est de toute première importance en calcul numérique. La plupart des procédés numériques ne consistent-ils pas à se ramener à la solution de systèmes linéaires ? Un tel problème ne comporte aucune erreur de méthode, seules des erreurs dues aux opérations élémentaires (sommes, différences, produits et quotients) peuvent intervenir et la répercussion, l'accumulation de telles erreurs sont d'une très grande importance.

Lorsque les moyens de calculs étaient limités aux machines de bureaux, le contrôle, à chaque opération, par le calculateur, des erreurs d'arrondis se faisait très facilement. Celui-ci se rendait compte à chaque instant de la valeur de ses résultats intermédiaires, pouvait par exemple, décider en cours de calcul, de changer le nombre de chiffres significatifs conservés par lui, pour remédier à une certaine "dérive" de son calcul global. Il faut ajouter qu'à ce moment des résolutions de systèmes linéaires d'ordre 20 étaient considérés comme des curiosités et bien rares sont les calculateurs en ayant résolus.

A l'heure actuelle, les calculateurs électroniques ont permis, de par leur extrême rapidité, d'envisager et d'exécuter des résolutions de systèmes d'ordre extrêmement élevé (10.000 par exemple). Ces grands systèmes sont résolus à l'aide de programmes généraux pour ces calculateurs qui ne donnent pas de renseignements sur la valeur des résultats fournis par la machine. On peut penser combien, des résultats dépendant d'un aussi grand nombre d'opérations élémentaires peuvent être sensibles aux erreurs faites dans chacune de ces opérations.

Les calculs d'erreurs qui vont suivre sont basés, d'une part sur des hypothèses de travail relatives aux erreurs dans les opérations élémentaires, d'autre part, sur une méthode de détermination de la matrice d'erreur.

- 1°) Il y a deux façons d'utiliser une machine à calculer à programme relativement aux opérations élémentaires. L'une est dite "calcul en opérations à virgule fixe", l'autre est dite "calcul en opérations en virgule flottante". Nous ne pouvons pas entrer dans les détails mais l'on peut justifier les règles suivantes:

Règle I Si un calcul se fait en virgule fixe, tout se passe comme si il n'y avait pas d'erreur dans une somme algébrique, les erreurs dans les produits sont constantes =  $e$ , les erreurs dans les quotients sont aussi constantes =  $e'$

Si  $a$  et  $b$  sont deux nombres

$\oplus, \ominus, \odot, \oslash$  désignent les opérations-machine correspondant à  $+, -, \cdot, :$

On écrira:

$$a \oplus b = a + b$$

$$a \ominus b = a - b$$

$$a \odot b = a \cdot b + e$$

$$a \oslash b = a : b + e'$$

Pour  $e$  et  $e'$  on prendra

$$e = k \cdot \beta^{-n}$$

$$e' = k' \cdot \beta^{-n'}$$

si  $\beta$  est la base de numération en laquelle la machine calcule  $n'$  et  $n$  désignant les nombres de chiffres conservés après la virgule dans ce calcul pour les divisions et pour les multiplications. Ce sont des constantes tout au long du calcul, c'est la raison de la dénomination "en fixe" de ce genre de calcul.

Règle II Si un calcul se fait en virgule flottante, tout se passe comme si dans toute opération il y avait une erreur relative constante.

Si  $a$  et  $b$  sont deux nombres "flottants" (c'est à dire sur lesquels la machine peut effectuer ces opérations flottantes).

$\boxplus, \boxminus, \boxdot, \boxdiv$  désignent les opérations-machine correspondant à  $+, -, \cdot, :$

On écrira:

$$a \oplus b = (1 + \varepsilon)(a + b)$$

$$a \ominus b = (1 + \varepsilon)(a - b)$$

$$a \otimes b = (1 + \varepsilon')(a \cdot b)$$

$$a \oslash b = (1 + \varepsilon')(a : b)$$

Pour  $\varepsilon$  et  $\varepsilon'$  on prendra

$$\varepsilon = k_1 \cdot \beta^{-N}$$

$$\varepsilon' = k'_1 \cdot \beta^{-N}$$

si  $\beta$  est la base d'numération en laquelle la machine calcule et  $N$  est le nombre de chiffres de mantisse (ou capacité de la représentation flottante) utilisés.

Nous répétons que les règles I et II sont des hypothèses de travail bien suffisantes d'ailleurs pour dire quelque chose de valable sur des programmes généraux. Il est d'ailleurs inextricable de vouloir utiliser des règles plus fines : les efforts de VON NEUMANN et GOLDSTINE [16] le prouvent suffisamment. Nous reparlerons plus loin d'un procédé (Le principe de simulation d'opérations) capable de doser la sensibilité d'un calcul particulier aux erreurs d'arrondis.

Voici maintenant la méthode de détermination de la matrice d'erreur :

Soit à résoudre le système linéaire :

$$A x = y$$

par exemple par la règle d'élimination de GAUSS.

On a déjà dit un mot du procédé au chapitre I, ou l'on a mis en évidence le rôle important joué par des matrices du 2° degré.

Si  $M$  est une matrice  $(n, n)$  n'ayant que des termes nuls au dessous de sa diagonale dans les colonnes  $1, 2, \dots, k$ ,

Il existe une matrice du 2° degré  $f(M)$ , telle que

$$M_1 = f(M) \cdot M$$

n'ait que des termes nuls au dessous de sa diagonale dans les colonnes  $1, 2, \dots, k, k+1$ .

On peut remarquer que si le calcul effectif donne  $\bar{M}_1$  au lieu de  $M_1$  on peut déterminer  $\delta M$  telle que <sup>l'on ait</sup> exactement :

$$\bar{M}_1 = f(M + \delta M) \cdot (M + \delta M)$$

De proche en proche on montre que :



$$\text{et } y^{(k+1)} \begin{cases} y_i^{(k+1)} = y_i^{(k)} & , i \leq k \\ y_i^{(k+1)} = y_i^{(k)} - \bar{p}_i^{(k)} y_k^{(k)} & ; i = k+1, \dots, n \end{cases}$$

2°) Calcul réel . Dans le calcul réel, l'existence d'erreur dans les opérations élémentaires conduit à une autre suite de matrices:

$$\bar{A}^{(1)} = A, \bar{A}^{(2)}, \bar{A}^{(3)}, \dots, \bar{A}^{(n)} = \bar{C}$$

et une autre suite de 2° membres:

$$\bar{y}^{(1)} = y, \bar{y}^{(2)}, \bar{y}^{(3)}, \dots, \bar{y}^{(n)} = \bar{y}$$

$\bar{A}^{(1)}$  est identique à  $A^{(1)}$ ,  $\bar{y}^{(1)}$  à  $y^{(1)}$ , mais  $\bar{C}$  et  $\bar{y}$  diffèrent de  $C$  et  $y$  . Cela provient de deux raisons :

- a) les matrices  $J_k$  sont mal déterminées et remplacées par  $\bar{J}_k$

(car il y a des erreurs dans la détermination des nombres

$\bar{p}_i^{(k)}$ ) . Les termes  $\bar{p}_i^{(k)}$  de  $\bar{J}_k$  sont remplacés par  $\bar{p}_i^{(k)} = p_i^{(k)} + \varepsilon_i^k$ , ( $i = k+1, \dots, n$ )  $\varepsilon_i^k$  désignant l'erreur de division  $a_{ik}^{(k)} : a_{kk}^{(k)}$  .

- b) Le produit  $J_k \cdot A^{(k)}$  est remplacé par  $\bar{J}_k \circ \bar{A}^{(k)}$ , qui lui-même n'est réalisé qu'avec des erreurs de calcul on désignera par  $-\eta_{ij}^{(k)}$  l'erreur dans les produits  $\bar{p}_i^{(k)} \cdot \bar{a}_{kj}^{(k)}$  ( $i, j = k+1, \dots, n$ ) .

Donc les formules donnant les termes de  $\bar{A}^{(k+1)}$  sont telles que

$$\bar{a}_{ij}^{(k+1)} = \bar{a}_{ij}^{(k)} - \bar{p}_i^{(k)} \cdot \bar{a}_{kj}^{(k)} + \eta_{ij}^{(k)}$$

De même pour les 2° membres, si l'erreur dans  $\bar{p}_i^{(k)} \cdot \bar{y}_k^{(k)}$  est  $-p_i^{(k)}$

$$\bar{y}_i^{(k+1)} = \bar{y}_i^{(k)} - \bar{p}_i^{(k)} \bar{y}_k^{(k)} + p_i^{(k)}$$

Conclusions : Les formules qui permettent d'obtenir la suite des  $\bar{A}^{(k)}$ ,

et des  $\bar{y}^{(k)}$  sont :

$$\bar{A}^{(1)} = A \quad ; \quad \bar{y}^{(1)} = y$$

$$G.F. \left[ \begin{cases} \bar{p}_i^{(k)} = (\bar{a}_{ik}^{(k)} : \bar{a}_{kk}^{(k)}) + \varepsilon_i^k & (i = k+1, \dots, n) \\ \bar{a}_{ij}^{(k+1)} = \bar{a}_{ij}^{(k)} & , i \leq k, j < k \\ \bar{a}_{ik}^{(k+1)} = 0 & , i > k \quad (j = k) \\ \bar{a}_{ij}^{(k+1)} = \bar{a}_{ij}^{(k)} + \eta_{ij}^{(k)} - \bar{p}_i^{(k)} \cdot \bar{a}_{kj}^{(k)} & (i, j = k+1, \dots, n) \\ \bar{y}_i^{(k+1)} = \bar{y}_i^{(k)} & , i \leq k \\ \bar{y}_i^{(k+1)} = \bar{y}_i^{(k)} + p_i^{(k)} - \bar{p}_i^{(k)} \cdot \bar{y}_k^{(k)} & (i = k+1, \dots, n) \end{cases} \right.$$





Par suite, si les erreurs de division sont de l'ordre de  $e'$  et celles de multiplication de l'ordre de  $e$ , cette matrice est de la forme:

$$\delta A \# e' \cdot \begin{vmatrix} 0 & 0 & \dots & 0 \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 1 & 1 & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & 0 \end{vmatrix} + e \cdot \begin{vmatrix} \pi_1 & 0 & \dots & 0 \\ 0 & \pi_2 & & \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & \pi_n \end{vmatrix} + e \cdot \begin{vmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 1 & \dots & 1 \\ 0 & 1 & 2 & \dots & 2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 1 & 2 & 3 & \dots & n-1 \end{vmatrix}$$

b) Pour  $\delta y$ , la valeur exacte est:

$$\delta y = \begin{vmatrix} 0 \\ p_2^1 \\ p_3^1 + p_3^2 \\ \vdots \\ \sum_{k=1}^{j-1} p_j^{(k)} \leftarrow j^o \\ \vdots \\ p_n^1 + p_n^2 + \dots + p_n^{n-1} \end{vmatrix}$$

et la valeur, en ordre de grandeur:

$$\delta y \# e \cdot \begin{vmatrix} 0 \\ 1 \\ 2 \\ \vdots \\ n-1 \end{vmatrix} *$$

5°) Retour-arrière. Soit alors à résoudre le système triangulaire:

$$\bar{C} x = \bar{y}$$

avec:

$$\bar{C} = \begin{vmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ & a_{22} & a_{23} & \dots \\ & & a_{33} & \dots \\ & & & \ddots \\ & 0 & & & a_{nn} \end{vmatrix}, \quad \bar{y} = \begin{vmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{vmatrix}$$

Les formules sont:

$$\xi_n = \frac{y_n}{a_{nn}}$$

$$\xi_k = \frac{1}{a_{kk}} \left[ y_k - (a_{k,k+1} \xi_{k+1} + a_{k,k+2} \xi_{k+2} + \dots + a_{k,n} \xi_n) \right]; \quad (k=n-1, \dots, 1)$$

Plus précisément, pour calculer  $\xi_k$ , ( $k=n, n-1, \dots, 1$ )

on "boucle" les différentes valeurs de  $S_k^h$  par:

$$S_k^{h+1} = 0, \quad S_k^i = S_k^{i+1} + a_{k,i} \cdot \xi_i \quad (i=k, \dots, k+1)$$

(pour,  $k=n, \dots, 1$ )

\* En supposant un nombre de chiffres différent dans la manipulation des seconds membres.



3°) Calculs en fixe. Inversion d'une matrice par la méthode de GAUSS.

Nous pouvons avec les formules qui précèdent, obtenir, dès à présent, un résultat relatif à l'inversion numérique des matrices par la méthode de GAUSS.

On sait que l'inverse d'une matrice se détermine numériquement en résolvant les n systèmes linéaires

$$(1) : Ax = e_1$$

$$(2) : Ax = e_2$$

⋮

$$(n) : Ax = e_n$$

les 2° membres étant les vecteurs unités :  $e_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow i^{\circ}$

Or ces systèmes ont la même matrice de premier membre A.

D'autre part il est clair que :

Pour le système (1), les  $\mu_i^{(1)}$  de la formule donnant  $\delta y$  sont nuls

Pour le système (2), les  $\mu_i^{(1)}, \mu_i^{(2)}$  sont nuls, etc...

Pour le système (n), les  $\mu_i^{(1)}, \mu_i^{(2)}, \dots, \mu_i^{(n-1)}$  sont nuls.

Si bien que :

La solution approchée du système (1), est la solution exacte d'un système de la forme :

$$(1)' : (A + \delta A) \cdot x = e_1 + e' \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix} + e \begin{pmatrix} n-1 \\ n-2 \\ \vdots \\ n-2 \end{pmatrix}$$

La solution approchée du système (2) est la solution exacte d'un système de la forme :

$$(2)' : (A + \delta A) \cdot x = e_2 + e' \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix} + e \begin{pmatrix} n-1 \\ n-2 \\ n-3 \\ \vdots \\ n-3 \end{pmatrix}$$

etc.....

La solution approchée du système (n) est la solution exacte d'un système de la forme :

$$(n)' : (A + \delta A) \cdot x = e_n + e' \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix} + e \begin{pmatrix} n-1 \\ n-2 \\ n-3 \\ \vdots \\ 2 \\ 1 \\ 0 \end{pmatrix}$$

on en conclut le théorème:

Théorème IV

Inverser une matrice A par la méthode de GAUSS avec des erreurs de calcul c'est résoudre exactement :  $(A + \delta I)X = I + \delta I$ , au lieu de  $A \cdot X = I$

Avec

$$\delta A \# e' \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ \pi_1 & 0 & 0 & \dots & 0 \\ \pi_1 & \pi_2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \pi_1 & \pi_2 & & & \pi_{k-1} & 0 \end{pmatrix} + e \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 1 & \dots & 1 \\ 0 & 1 & 2 & \dots & 2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 1 & 2 & \dots & k-1 \end{pmatrix}$$

et

$$\delta I \# e' \begin{pmatrix} \pi_1 & \pi_1 & \dots & \pi_2 \\ \pi_2 & \pi_2 & \dots & \pi_2 \\ \pi_3 & \pi_3 & \dots & \pi_3 \\ \vdots & \vdots & \ddots & \vdots \\ \pi_k & \pi_k & \dots & \pi_k \end{pmatrix} + e \begin{pmatrix} k-1 & k-1 & \dots & k-1 \\ k-2 & k-2 & \dots & k-2 \\ k-2 & k-3 & \dots & k-3 \\ \vdots & \vdots & \ddots & \vdots \\ k-2 & k-3 & \dots & 1 & 0 & 1 \\ & & & & & 0 \end{pmatrix}$$

- Nous allons procéder à une étude tout-à-fait analogue dans le cas de calculs conduits en point décimal flottant.

4°) Calculs en virgule flottante. Triangularisation et résolution d'un système.

- 1°) Les notations utilisées seront les plus voisines possibles de celles du § 2°, dans le calcul réel, l'existence d'erreurs dans les opérations élémentaires conduit à une suite de matrices,

$$\bar{A}^{(1)} = A, \bar{A}^{(2)}, \dots, \bar{A}^{(n)} = \bar{E}$$

et une suite de seconds membres:

$$\bar{y}^{(1)} = y, \bar{y}^{(2)}, \dots, \bar{y}^{(n)} = \bar{y}$$

$\bar{A}^{(n)}$  est identique à A,  $\bar{y}^{(n)}$  à y, mais  $\bar{E}$  et  $\bar{y}$  diffèrent de E et y. Cela provient de deux raisons:

- a) Les matrices  $J_k$  sont mal déterminées et remplacées par  $\bar{J}_k$  (car il y a des erreurs dans la détermination des nombres  $f_i^{(k)}$ ). Les termes  $f_i^{(k)}$  de  $J_k$  sont remplacés par  $\bar{f}_i^{(k)} = f_i^{(k)} + \varepsilon_i^{(k)} \cdot f_i^{(k)}$   $\varepsilon_i^{(k)}$  désignant l'erreur relative de division.

d'où: 
$$\bar{f}_i^{(k)} = (\bar{a}_{ik}^{(k)} + \varepsilon_i^{(k)} \bar{a}_{ik}^{(k)}) : \bar{a}_{kk}^{(k)}$$

b) Le produit  $J_k \cdot A^{(k)}$  est remplacé par  $\bar{J}_k \cdot \bar{A}^{(k)}$ , qui lui-même n'est réalisé qu'avec des erreurs de calcul, on désignera dans ce cas par  $\eta_{ij}^{(k)}$  l'erreur relative faite dans le calcul de  $\bar{a}_{ij}^{(k+1)}$  donc on peut écrire:

$$\bar{a}_{ij}^{(k+1)} = \bar{a}_{ij}^{(k)} - \bar{f}_i^{(k)} \cdot \bar{a}_{kj}^{(k)} + \eta_{ij}^{(k)} \cdot (\bar{a}_{ij}^{(k+1)})$$

De même, pour les seconds membres, on désignera par  $\mu_i^{(k)}$  l'erreur relative dans le calcul de  $\bar{y}_i^{(k+1)}$ . Si bien que :

$$\bar{y}_i^{(k+1)} = \bar{y}_i^{(k)} - \bar{p}_i^{(k)} \bar{y}_k^{(k)} + \mu_i^{(k)} \bar{y}_i^{(k)}$$

Conclusions : Les formules qui permettent, dans ce cas, d'obtenir la suite des  $\bar{A}^{(k)}$ , et des  $\bar{y}^{(k)}$  sont :

$$\bar{A}^{(1)} = A, \quad \bar{y}^{(1)} = y$$

G.FE

$$\left\{ \begin{array}{l} \bar{p}_i^{(k)} = (\bar{a}_{ik}^{(k)} + \varepsilon_i^{(k)} \bar{a}_{ik}^{(k)}) : a_{ik}^{(k)}, \quad (i = k+1, \dots, n) \\ \bar{a}_{ij}^{(k+1)} = a_{ij}^{(k)}, \quad i \leq k, \quad j < k \\ \bar{a}_{ik}^{(k+1)} = 0, \quad i > k, \quad (j = k) \\ \bar{a}_{ij}^{(k+1)} = (\bar{a}_{ij}^{(k)} - \bar{p}_i^{(k)} \bar{a}_{ij}^{(k)}) (1 + \eta_{ij}^{(k)}) \# \bar{a}_{ij}^{(k)} + \eta_{ij}^{(k)} (\bar{a}_{ij}^{(k+1)}) - \bar{p}_i^{(k)} \bar{a}_{kj}^{(k)}, \quad (i, j = k+1, \dots, n) \\ \bar{y}_i^{(k+1)} = \bar{y}_i^{(k)}, \quad i \leq k \\ \bar{y}_i^{(k+1)} = (\bar{y}_i^{(k)} - \bar{p}_i^{(k)} \bar{y}_k^{(k)}) (1 + \mu_i^{(k)}) \# \bar{y}_i^{(k)} + \mu_i^{(k)} \bar{y}_i^{(k+1)} - \bar{p}_i^{(k)} \bar{y}_k^{(k)}, \quad (i = k+1, \dots, n) \end{array} \right.$$

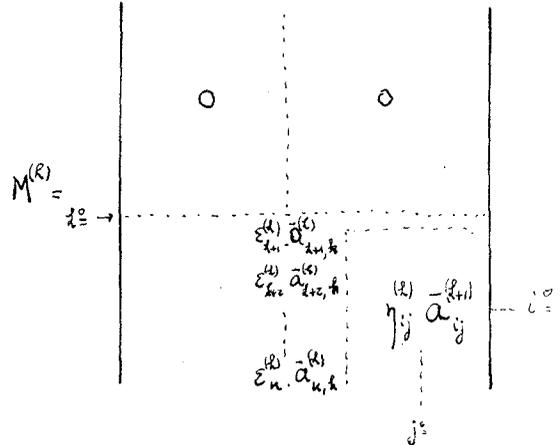
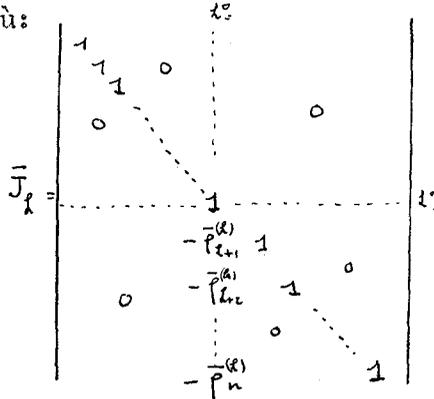
Comme nous l'avons fait nous écrivons ces formules :

$$\bar{A}^{(k+1)} = \bar{J}_k \cdot \bar{A}^{(k)}$$

avec :

$$\bar{A}^{(k)} = \bar{A}^{(k)} + M^{(k)}$$

où :



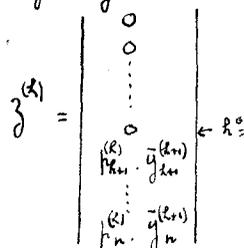
de même :

$$\bar{y}^{(k+1)} = \bar{J}_k \cdot \bar{y}^{(k)}$$

avec :

$$\bar{y}^{(k)} = \bar{y}^{(k)} + z^{(k)}$$

où :



2°) Détermination de  $\bar{C}$ .

De même que nous l'avons déjà fait, on tire

$$\bar{A}^{(k+1)} = \bar{J}_k (\bar{A}^{(k)} + M^{(k)}) = \bar{J}_k [\bar{J}_{k-1} (\bar{A}^{(k-1)} + M^{(k-1)}) + M^{(k)}]$$

Par un raisonnement identique on prouve que cela peut être écrit:

$$\bar{A}^{(k+1)} = \bar{J}_k \cdot \bar{J}_{k-1} \cdot [\bar{A}^{(k-1)} + M^{(k-1)} + M^{(k)}]$$

et par suite :

$$\bar{A}^{(n)} = \bar{C} = \bar{J}_{n-1} \cdot \bar{J}_{n-2} \cdot \dots \cdot \bar{J}_1 \cdot (A + \sum_{i=1}^{i=n-1} M^{(i)})$$

Donc,

théorème V Dans la triangularisation de la matrice  $A$ , si l'on commet dans les calculs élémentaires les erreurs relatives  $\varepsilon_i^{(k)}$  et  $\eta_{ij}^{(k)}$ , on peut considérer la matrice triangulaire obtenue  $\bar{C}$  comme résultant de la triangularisation sans erreur de la matrice  $A + \delta A$  avec  $\delta A = \sum_{i=1}^{i=n-1} M^{(i)}$  et les  $M^{(i)}$ , les valeurs indiquées plus haut.

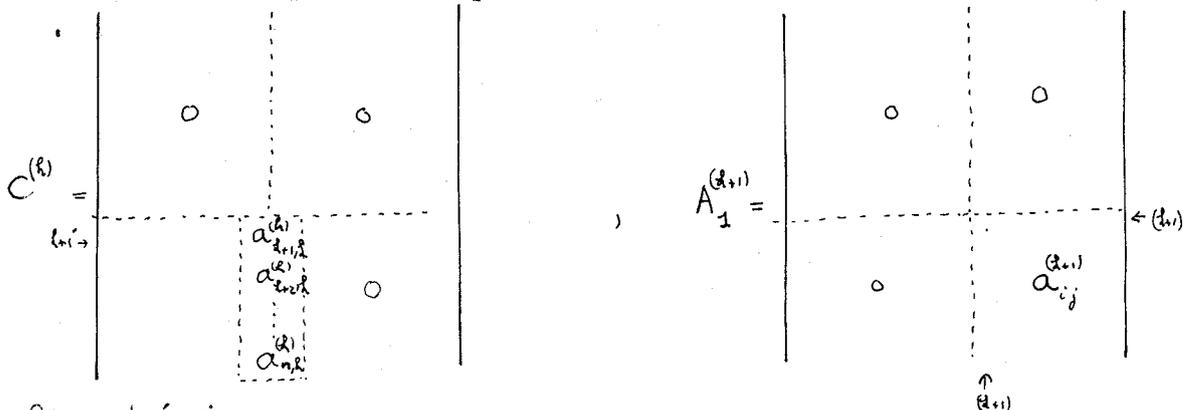
De même on voit que:  $\bar{y}^{(k+1)} = \bar{J}_k \cdot \bar{J}_{k-1} \cdot (\bar{y}^{(k)} + \delta y^{(k)} + \delta y^{(k)})$

et  $\bar{y}^{(n)} = \bar{y} = \bar{J}_{n-1} \cdot \bar{J}_{n-2} \cdot \dots \cdot \bar{J}_1 \cdot (y + \delta y)$ ,  $\delta y = \sum_{i=1}^{i=n-1} \delta y^{(i)}$

3°) Structures de la matrice  $\delta A$  et du vecteur  $\delta y$ .

Devant la complexité des formes obtenues nous ferons ici une hypothèse simplificatrice pour obtenir en ordre de grandeur les valeurs de  $\delta A$  et  $\delta y$ . Nous admettons que les erreurs relatives  $\varepsilon_i^{(k)}$ ,  $\eta_{ij}^{(k)}$  et  $\rho_i^{(k)}$  sont constantes (pour  $\varepsilon_i^{(k)}$  cela résulte de la règle II du § I) et de l'ordre de  $\varepsilon = k \cdot \beta^{-N}$  ( $N$  est la capacité de la représentation flottante utilisée,  $\beta$  la base de numération). On fait cette hypothèse dans ce §.

On voit alors que si l'on pose:



On peut écrire:

$$M^{(k)} = \varepsilon (C^{(k)} + A_1^{(k+1)})$$

$$\text{D'où } \delta A = \sum_{i=1}^{i=n-1} M^{(i)} = \varepsilon \sum_{i=1}^{i=n-1} C^{(i)} + \varepsilon \sum_{i=1}^{i=n-1} A_1^{(i)} = \varepsilon \cdot T + \varepsilon \cdot R$$

Donc :

La matrice d'erreur  $\delta A$ , à ces approximations est la somme de  $\varepsilon \cdot T + \varepsilon \cdot R$

$T$  : matrice triangulaire à diagonale nulle des "colonnes éliminées "

$R$  : matrice somme des matrices (rendues à l'ordre  $n$ ) des systèmes intermédiaires.

Pour  $\delta y$  on posera de même: 
$$\bar{y}_1^{(i+1)} = \begin{pmatrix} 0 \\ \vdots \\ \frac{1}{d_{i+1}} (R_{i+1}) \\ \vdots \\ y_n^{(i+1)} \end{pmatrix}$$

et 
$$\delta y = \varepsilon \cdot \sum_{i=1}^{i=n-1} \bar{y}_1^{(i)}$$

4°) Retour-arrière.

Si l'on doit alors résoudre le système:  $\bar{C} x = \bar{y}$ ,

(mêmes notations que précédemment)

Le calcul conduit comme on le voit immédiatement à prendre :

$$\bar{s}_r^i = \bar{s}_r^{i+1} + \alpha_{r,i} \bar{s}_i - p_{r,i} \bar{s}_r^i$$

si  $p_{r,i}$  est l'erreur relative dans cette opération.

Et par suite

$$\bar{s}_r = \frac{1}{d_{rr}} [\gamma_r - \bar{s}_r^{r+1}] + d_{r,r} \bar{s}_r$$

$d_{r,r}$  étant l'erreur relative dans la division.

d'où:

$$\bar{s}_r = \frac{1}{d_{rr}} \left[ \gamma_r + d_{rr} \cdot d_{r,r} \bar{s}_r + p_{r,r+1} \bar{s}_r^{r+1} + \dots + p_{r,n} \bar{s}_r^n - (d_{r,r+1} \bar{s}_{r+1} + \dots + \alpha_{r,n} \bar{s}_n) \right]$$

Théorème VI Dans la résolution d'un système triangulaire  $\bar{C} x = \bar{y}$ , si l'on commet les erreurs relatives  $p_{r,i}, d_{r,r}$ , on peut considérer la solution obtenue comme la solution exacte du système triangulaire  $\bar{C} x = \bar{y} + \delta y$  de même matrice de  $1^o$  membre, avec

$$\delta y = \begin{pmatrix} \alpha_{11} d_{11} \bar{s}_1 \\ \vdots \\ \alpha_{ii} d_{ii} \bar{s}_i \\ \vdots \\ \alpha_{nn} d_{nn} \bar{s}_n \end{pmatrix} + \begin{pmatrix} p_{12} \bar{s}_1^2 + p_{13} \bar{s}_1^3 + \dots + p_{1n} \bar{s}_1^n \\ \vdots \\ \vdots \\ \vdots \end{pmatrix}$$

5°) Système complet.

Il résulte de tout cela que l'on peut énoncer, comme plus haut :

Théorème VII : La solution approchée du système  $Ax = y$  et obtenue par un calcul en virgule flottante où les erreurs sont, pour la triangularisation,  $\varepsilon_i^{(k)}, \eta_{ij}^{(k)}, \rho_i^{(k)}$ , et  $d_k, \rho_{ij}^{(k)}$  pour le retour-arrière, et par la méthode de GAUSS, peut être considérée comme la solution exacte du système :

$$(A + \delta A) x = y + \delta y$$

$$\delta A = \varepsilon \cdot T + \varepsilon \cdot R$$

T matrice des colonnes éliminées,

R somme des matrices des systèmes intermédiaires

$$\delta y = \varepsilon \cdot \sum \bar{y}_i^{(k)}$$

somme des seconds membres sous-diagonaux

Obtention de bornes pour l'erreur

Soit le système  $Ax = y$ , résolu par la méthode de GAUSS,

Le calcul conduit à un vecteur solution  $x^*$ .

Celui-ci satisfait par suite du théorème III à l'égalité :

$$(A + \Delta A) \cdot x^* = y + \Delta y$$

où l'on a posé:  $\delta A = \Delta A$

$$\delta y + \delta_i y = \Delta y$$

D'où l'on tire :

$$(A + \Delta A) \cdot x^* = Ax + \Delta y$$

puis ,

$$A \cdot (x^* - x) = \Delta y - \Delta A x^*$$

D'où la formule fondamentale :

$$(F) \quad \boxed{x^* - x = A^{-1}(\Delta y - \Delta A x^*)}$$

D'autre part, soient  $\varphi_i$  et  $\varphi_j$  deux normes de vecteurs

on a :  $\varphi_i(A(x^* - x)) = \varphi_i(\Delta y - \Delta A x^*)$

Donc 
$$\frac{\varphi_i(A(x^* - x))}{\varphi_j(x^* - x)} = \frac{\varphi_i(\Delta y - \Delta A x^*)}{\varphi_j(x^* - x)}$$

, si  $m_{ij}(A), M_{ij}(A)$  ont les significations du chapitre III.

$$m_{ij}(A) \leq \frac{\varphi_i(\Delta y - \Delta A x^*)}{\varphi_j(x^* - x)} \leq M_{ij}(A)$$

Il en résulte : 
$$\frac{\varphi_i(\Delta y - \Delta A x^*)}{M_{ij}(A)} \leq \varphi_j(x^* - x) \leq \frac{\varphi_i(\Delta y - \Delta A x^*)}{m_{ij}(A)}$$

soit : 
$$\frac{1}{S_{ij}(A)} \cdot \varphi_i(\Delta y - \Delta A x^*) \leq \varphi_j(x^* - x) \leq S_{ji}(A^{-1}) \cdot \varphi_i(\Delta y - \Delta A x^*)$$

Et si l'on veut se contenter d'une majoration supérieure de l'erreur :

$$\varphi_j(x^* - x) \leq S_{ji}(A^{-1}) \cdot [\varphi_i(\Delta y) + \varphi_i(\Delta A x^*)]$$

ou en introduisant le conditionnement général  $\gamma_{ij}(A)$  :

$$\frac{\varphi_i(\Delta y - \Delta A x^*)}{S_{ij}(A)} \leq \varphi_j(x^* - x) \leq \frac{1}{\gamma_{ij}(A)} \cdot \frac{\varphi_i(\Delta y - \Delta A x^*)}{S_{ij}(A)}$$

D'où le théorème :

Théorème VIII

Si  $\varphi_i$  et  $\varphi_j$  sont deux normes de vecteurs,  $S_{ij}$  la norme de matrice associée à ces deux normes, dans cet ordre, et  $\gamma_{ij}$  le conditionnement général relatif à celles-ci, on a pour la norme  $\varphi_j$  de l'erreur commise dans la résolution d'un système linéaire :

$$(1) \quad \boxed{\frac{\varphi_i(\Delta y - \Delta A x^*)}{S_{ij}(A)} \leq \varphi_j(x^* - x) \leq \frac{1}{\gamma_{ij}(A)} \cdot \frac{\varphi_i(\Delta y - \Delta A x^*)}{S_{ij}(A)}}$$

Voici ce que donne cette formule dans les cas usuels :

1°) si  $i=j=2$  , on a 
$$\frac{\|\Delta y - \Delta A x^*\|}{\sqrt{\lambda_2}} \leq \|x^* - x\|$$

puis : 
$$(2) \quad \|x^* - x\| \leq \frac{\|\Delta y - \Delta A x^*\|}{\sqrt{\lambda_2}}$$

2°) Si  $i=1$  et  $j=2$  on trouve :

$$(2') \quad \frac{1}{M_A} \Phi(\Delta y - \Delta A x^*) \leq \|x^* - x\| \leq \frac{1}{m_A} \Phi(\Delta y - \Delta A x^*)$$

En introduisant le conditionnement normalisé  $C(A)$ ,

$$(3) \quad \frac{1}{M_A} \Phi(\Delta y - \Delta A x^*) \leq \|x^* - x\| \leq \frac{1}{C(A)} \frac{N(A)}{p_0(A)} \cdot \frac{1}{M_A} \Phi(\Delta y - \Delta A x^*)$$

Tout ce que nous venons de voir est valable quel que soit la façon dont les calculs sont conduits, je vais alors supposer qu'ils sont tout d'abord exécutés en fixe.

Utilisons tout d'abord, la formule (2)

On peut écrire :

$$\|x^* - x\| \leq \frac{1}{\sqrt{\lambda_1}} \|\Delta y\| + \frac{1}{\sqrt{\lambda_1}} \|\Delta A \cdot A^{-1} \cdot y\|$$

$$x^* \neq A^{-1} \cdot y$$

Il est rappelé que dans ces relations  $\lambda_1$  désigne la plus petite valeur propre de la matrice  $A^T.A$ .

Or il est immédiat de voir que :

$$\frac{\|\Delta A . A^{-1} y\|}{\|y\|} = \frac{\|\Delta A . X\|}{\|X\|} \cdot \frac{\|X\|}{\|AX\|}, \quad y = AX$$

puis,

$$\|\Delta A . A^{-1} y\| \leq N(\Delta A) \cdot \frac{1}{\sqrt{\lambda_1}} \|y\|$$

donc :

$$\|x^* - x\| \leq \frac{1}{\sqrt{\lambda_1}} \|\Delta y\| + \frac{1}{\lambda_1} N(\Delta A) \cdot \|y\|$$

Mais ici, d'après les formules du théorème III de ce chapitre,

$$N(\Delta A) = N(\Delta A) \leq e' \cdot N \begin{pmatrix} 0 & 0 & \dots & 0 \\ \pi_1 & 0 & & \vdots \\ \pi_1 & \pi_2 & & \vdots \\ \vdots & \vdots & & \vdots \\ \pi_1 & \pi_2 & & 0 \end{pmatrix} + e \cdot N \begin{pmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & 1 & \dots & 1 \\ \vdots & 1 & 2 & \dots & 2 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 1 & 2 & \dots & h-1 \end{pmatrix}$$

Donc

$$N(\Delta A) \leq e' \cdot R_1 + e \cdot R_2$$

$$R_1^2 = (n-1)\pi_1^2 + (n-2)\pi_2^2 + \dots + \pi_{n-1}^2$$

$$R_2^2 = (2n-3) + (2n-5) \cdot 2^2 + \dots + (n-1)^2 = \sum_{i=1}^{i=n-1} [2n - (2i+1)] i^2 = (2n-1) \sum_{i=1}^{i=n-1} i^2 - 2 \sum_{i=1}^{i=n-1} i^3$$

$$= \frac{1}{6} n(n-1)(n^2-n+1) \approx \frac{n^4}{6}$$

Puis;

$$\|\Delta y\| \leq e' \cdot R_1' + e \cdot (n-1) \sqrt{n}$$

$$R_1'^2 = \pi_1^2 + \pi_2^2 + \dots + \pi_n^2$$

Enfin :

$$(4) \quad \|x^* - x\| \leq \frac{\|y\|}{\lambda_1} R_1 \cdot e' + \frac{1}{\sqrt{\lambda_1}} R_1' \cdot e' + \frac{\|y\|}{\lambda_1} R_2 \cdot e + \frac{(n-1)\sqrt{n}}{\sqrt{\lambda_1}} \cdot e_1$$

①
②
③
④

où les termes (1) et (2) sont relatifs aux divisions dans la triangularisation, et dans le retour-arrière

(3) et (4) sont relatifs aux multiplications dans la triangularisation et le retour-arrière

Intéressons-nous tout d'abord aux erreurs dues aux produits. La chose la plus remarquable est de voir la croissance en ordre de  $n^2$  due à  $R_2$ , autrement dit si  $n$  est grand, afin de pouvoir obtenir

une erreur sur la solution due aux multiplications dans la triangularisation inférieure à  $\beta^{-m}$  donné, il suffit de prendre  $e$  de sorte que :

$$\frac{1}{\sqrt{6}} \cdot e \cdot \frac{\|y\|}{\lambda_1} h^2 \leq \beta^{-m}$$

Or, on adit que pour des calculs en fixe,  $e = k \cdot \beta^{-r}$

si l'on fait  $k = \frac{1}{2}$

on peut énoncer :

Prème IX : Pour pouvoir résoudre un système linéaire en virgule fixe et pour un ordre  $n$  assez élevé afin que ce soit le 3ème terme de la formule (4) qui prédomine, et pour assurer une erreur sur la solution inférieure à  $\beta^{-m}$ , il suffit de conduire les calculs avec  $p$  chiffres après la virgule ;  $p$  satisfaisant à :

$$r \geq m - \log_p \left( \frac{4,9}{h^2} \frac{\lambda_1}{\|y\|} \right)$$

(2.√6 # 4,9)

Comme suite à cela, étudions la partie relative aux erreurs de division. C'est à dire le 1<sup>o</sup> terme du second membre de la formule (4)

Deux cas sont à examiner :

I<sup>o</sup>)  $\sqrt{\lambda_1}$  est plus grand que  $\|y\|$ , ou bien l'ordre  $n$  est peu élevé de sorte que  $R'_1$  et  $R_1$  soient du même ordre de grandeur.

Dans ces cas les erreurs systématiques dues aux divisions ont une borne supérieure dont le facteur prépondérant est

$$R'_1 = \left( \sum_{i=1}^{i=n} \pi_i^2 \right)^{\frac{1}{2}} \neq R_1 \quad \text{pour } e_1 \text{ et } e'_1 \text{ de même ordre.}$$

Comment affaiblir le plus possible ce facteur ?

Si l'on dispose d'un système  $Ax = y$  à résoudre, un certain nombre d'opérations peuvent être faites, sans aucune erreur, qui substituent au système à résoudre, un autre dont la solution est celle du système proposé à l'ordre près des composantes: ce sont des échanges de lignes ou de colonnes dans le système donné ou dans les systèmes intermédiaires.

Par exemple, l'échange de deux lignes dans le système proposé donne, après élimination de la première inconnue un système ( d'ordre  $n-1$  ) qui diffère de celui obtenu à partir du système initial, etc... Cela provient de la dissymétrie de la méthode de GAUSS. Il a toujours été important pour le calculateur d'avoir une règle de conduite, pour organiser son calcul afin d'obtenir la meilleure solution.

Nous pouvons remarquer qu'il y a  $n^2$  façon de réorganiser, en ce sens, le système donné, puis  $(n-1)^2$  façons de réorganiser le système d'ordre  $(n-1)$  après élimination de la première inconnue, etc .... .  
 Donc on peut dire que il y a en tout:  $(n!)^2$  façons de résoudre numériquement un système donné par la méthode de GAUBS et, en ce sens.

Dans chaque résolution, les "pivots" changent et l'on peut se demander comment les choisir de sorte que la solution obtenue soit la meilleure. La formule (4) permet <sup>d'essayer</sup> de répondre à cette question: Il suffit de s'assurer que

$$R_1^2 = \pi_1^2 + \pi_2^2 + \dots + \pi_n^2$$

est minimum.

Or d'après la proposition II du chapitre I, il est clair que le produit des valeurs absolues des pivots  $|\pi_1| \cdot |\pi_2| \dots |\pi_n|$  est égal à la valeur absolue du déterminant  $|\Delta|$  de la matrice A du premier membre du système: donc ce produit est un nombre qui est toujours le même quelle que soit le procédé de résolution adopté.

Mais le minimum d'une somme de nombres positifs dont le produit est constant :

$$\pi_1^2 + \pi_2^2 + \dots + \pi_n^2$$

avec:  $\pi_1^2 \cdot \pi_2^2 \cdot \dots \cdot \pi_n^2 = \Delta^2$

a lieu si les nombres sont égaux, donc on peut dire que le facteur  $R_1$  est le plus petit possible si les pivots sont les plus voisins entre eux. D'où la règle:

Règle I Dans le cas où  $\sqrt{\lambda_1}$  est plus grand que  $\|y\|$ , ou bien que l'ordre  $n$  du système est assez faible pour que l'on puisse considérer  $R_1$  et  $R_2$  du même ordre de grandeur, le choix le plus avantageux des pivots est celui qui les assure tous le plus près de l'égalité entre-eux et forcément alors peu différents de  $\sqrt[2]{|\Delta|}$ . [23] *la val. absolue*

2°)  $\sqrt{\lambda_1}$  est plus petit que  $\|y\|$  et l'ordre est assez élevé pour que le facteur qui prédomine soit  $R_1$ .

Alors le meilleur choix est obtenu en essayant d'amoindrir

$$Z = (n-1)\pi_1^2 + (n-2)\pi_2^2 + \dots + \pi_n^2$$

Soit  $\mathcal{F}$  l'ensemble fini des nombres pouvant <sup>(comme)</sup> comme valeur absolue de dernier pivot dans une résolution de la famille de celles envisagées.

C'est à dire l'ensemble des nombres de la forme

$$\left| \frac{\Delta}{M} \right|$$

Où M est un mineur d'ordre  $n-1$  de la matrice A.

Si un choix de dernier pivot est imposé les

$$|\pi_1|, |\pi_2|, \dots, |\pi_{k-1}|$$

peuvent varier encore de sorte que

$$|\pi_1| \cdot |\pi_2| \cdot \dots \cdot |\pi_{k-1}| = |M|$$

Parmi ceux-ci, ceux qui minimisent  $Z$  sont tels que :

$$\sqrt{k-1} |\pi_1| = \sqrt{k-2} |\pi_2| = \dots = |\pi_{k-1}| \quad (1)$$

et dans ce cas  $Z_M$  vaudrait  $(k-1) \pi_{k-1}^2$

et d'après (1) la valeur de ces pivots serait telle que

$$|M| = \frac{|\pi_{k-1}|^{k-1}}{\sqrt{(k-1)!}}$$

ou :

$$\pi_{k-1} = [(k-1)!]^{1/(k-1)} \cdot |M|^{1/(k-1)}$$

Donc

$$Z_M = (k-1) [(k-1)!]^{1/(k-1)} \cdot |M|^{2/(k-1)}$$

Il y a donc à procéder ainsi :

Règle II Dans le cas où  $\sqrt{\lambda_1}$  est plus petit que  $\|y\|$  et l'ordre assez élevé pour que le facteur qui prédomine dans la partie de la formule (4) relative aux erreurs de divisions, le choix le plus avantageux est obtenu de la façon suivante :

1° On "organise" la matrice de sorte que le mineur  $\Delta$  qui tombe "en tête" de  $A$  (c'est à dire obtenu en supprimant la dernière ligne et la dernière colonne de  $A$ ) soit le plus petit possible.

2° Le pivot n° 1 doit être pris le plus près possible de  $[(k-1)!]^{1/(k-1)} \cdot |M|^{1/(k-1)} \cdot (k-1)^{-1/2}$

Le pivot n° 2 doit être pris le plus près possible de

$$[(k-1)!]^{1/(k-1)} \cdot |M|^{1/(k-1)} \cdot (k-2)^{-1/2}$$

etc.....

Remarquons, qu'au lieu d'utiliser la formule (4) faisant intervenir une majoration de l'erreur, on peut, d'après la formule fondamentale

$$(F) : x^* - x = A^{-1} \cdot [\Delta y - \Delta A \cdot x^*]$$

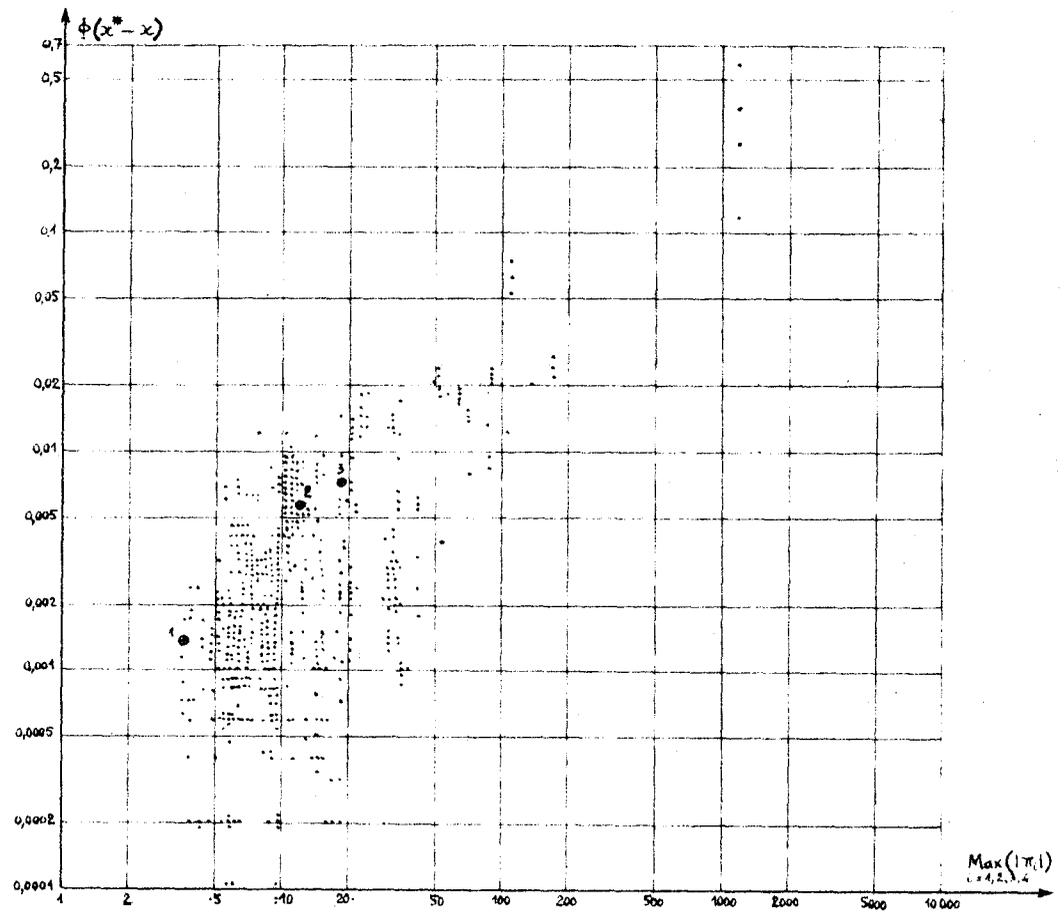
et d'après les expressions de  $\Delta A$  et  $\Delta y$  trouvées au théorème III, montrer que l'affaiblissement de  $\Delta y$  et  $\Delta A$  est réalisé si l'on évite les pivots par trop grands (ou par trop petits à cause de la constance du produit).

6°) Expériences de calcul en fixe

Afin de pouvoir juger de la validité des règles ci-dessus, nous avons fait les calculs suivants.

On a résolu le système  $Ax = y$ , avec :

A=	<table border="0" style="width: 100%; text-align: center;"> <tr> <td>0,053</td> <td>2,0135</td> <td>1,0172</td> <td>1,4521</td> </tr> <tr> <td>1,2142</td> <td>1,5676</td> <td>5,0431</td> <td>1,5834</td> </tr> <tr> <td>2,0534</td> <td>2,9415</td> <td>6,7123</td> <td>2,1421</td> </tr> <tr> <td>3,9835</td> <td>9,1121</td> <td>11,3451</td> <td>10,4520</td> </tr> </table>	0,053	2,0135	1,0172	1,4521	1,2142	1,5676	5,0431	1,5834	2,0534	2,9415	6,7123	2,1421	3,9835	9,1121	11,3451	10,4520	et y =	<table border="0" style="width: 100%; text-align: center;"> <tr> <td>12,94</td> </tr> <tr> <td>25,8123</td> </tr> <tr> <td>46,6417</td> </tr> <tr> <td>98,051</td> </tr> </table>	12,94	25,8123	46,6417	98,051
0,053	2,0135	1,0172	1,4521																				
1,2142	1,5676	5,0431	1,5834																				
2,0534	2,9415	6,7123	2,1421																				
3,9835	9,1121	11,3451	10,4520																				
12,94																							
25,8123																							
46,6417																							
98,051																							



Pour  $A \cdot A^T$  on a :  $\lambda_1 = 0,167$  ;  $\lambda_4 = 469,5$  ;  $\sqrt{\lambda_1} = 0,409$  ;  $\|y\| \# 112,7$

On est donc dans le cas où la règle I s'applique ( n =4 est faible )

Les calculs ont été exécutés en virgule fixe , à quatre décimales après celle-ci , sur la machine électronique " Gamma E.T. " de l' Université de Grenoble pour laquelle nous avons réalisé un programme capable de faire les  $(4!)^2 = 576$  résolutions de ce système . Il faut environ trois heures-machine pour terminer ce calcul. Le graphique de la page ci-contre résume les résultats obtenus. Le système a été formé de sorte que la solution exacte soit la colonne des nombres 1,2,3,4 .

Nous avons porté en abscisse le maximum de la valeur absolue des pivots relatif à une résolution, et en ordonnée, pour cette même résolution, la valeur de  $\Phi(x^* - x)$  , norme de l'erreur.

Le point marqué 1 de ce graphique est relatif à la résolution donnant les pivots les plus près possibles de  $\sqrt{|det A|}$  , c'est la résolution suivant la Règle 1 ci-dessus.

Le point marqué 2 est relatif à la règle habituelle du "pivot maximal" c'est à dire que les systèmes successifs sont toujours réorganisés de sorte que le plus grand élément , en valeur absolue, du système partiel tombe en tête.

Le point 3 est relatif à une résolution suivant la règle contraire à la précédente: placer en tête le plus petit des éléments relatifs, en valeur absolue, des systèmes intermédiaires.

Enfin le graphique montre combien les variations de l'erreur sont grandes et qu'il faut éviter à tout prix d'utiliser de trop grands pivots (et par contre coup, à cause du produit constant, de trop petits ) .

Il faut signaler que ce calcul n'a pu se faire sur la machine électronique qu'en utilisant ce que nous appelons des "simulations d'opérations "

En effet , la programmation en fixe sur un calculateur est très difficile si l'on ne connaît à l'avance la position de la virgule dans tous les résultats intermédiaires , ce qui est le cas dans un tel calcul. Par contre , la programmation en virgule flottante est très aisée . Or dans la plus part des machines l'ordre de commande de l'opération flottante est un ordre d'appel d'un sous-programme , Il suffit donc d'écrire un programme en virgule flottante, de substituer aux sous-programmes d'opérations flottantes des sous-programmes d'opérations déterminées (en fixe dans notre cas, mais pouvant , pour l'étude d'autres questions être des opérations d'autre nature: virgule fixe à capacité variable, avec ou sans arrondi , virgule flottante à capacité variable, avec ou sans arrondi) pour obtenir ce que l'on désire.

7°) Calculs en virgule flottante, Evaluation de l'erreur.

Nous chercherons ici une majoration convenable pour une erreur relative dans la solution.

Avec les mêmes notations qu'au § 5, le Théorème VIII est encore valable.

On peut écrire,

$$(1) \quad \varphi_j(x^* - x) \leq \frac{1}{s_{ij}(A)} \cdot \frac{\varphi_i(\Delta y - \Delta A x^*)}{s_{ij}(A)}$$

ou

$$(2) \quad \|x^* - x\| \leq \frac{1}{\sqrt{\lambda_1}} \|\Delta y - \Delta A x^*\|$$

ou

$$(2)' \quad \|x^* - x\| \leq \frac{1}{m_A} \Phi(\Delta y - \Delta A x^*)$$

Mais ici, il est plus difficile d'obtenir des ordres de grandeurs corrects pour  $\Delta A$  et  $\Delta y$ .

Pour  $\Delta A$ , sa structure est déterminée par ce que l'on a dit au §4,3°.

$$\Delta A = \varepsilon_1 (T + R)$$

Pour  $\delta y$ , d'après toujours ce même §, on peut remarquer (mêmes notations)

$$A_1^{(k+1)} \cdot x^* \neq \bar{y}_1^{(k+1)}$$

donc :

$$\delta y = \varepsilon \sum_{i=1}^{i=n-1} \bar{y}_i^{(k)} \neq \varepsilon R \cdot x^*$$

Pour  $\delta_1 y$  D'après §4,4°; on peut écrire

$$\delta_1 y = \varepsilon \cdot \begin{pmatrix} \alpha_{11} \bar{s}_1 \\ \alpha_{12} \bar{s}_2 \\ \vdots \\ \alpha_{1n} \bar{s}_n \end{pmatrix} + \varepsilon \cdot \begin{pmatrix} \bar{s}_2 + \bar{s}_3 + \dots + \bar{s}_n \\ \bar{s}_3 + \bar{s}_4 + \dots + \bar{s}_n \\ \vdots \\ \bar{s}_{n-1} \\ 0 \end{pmatrix}$$

Or, si :

$$\bar{C}_1 = \begin{pmatrix} a_{11} & a_{12} & 2a_{13} & 3a_{14} & \dots & (n-1)a_{1n} \\ 0 & a_{22} & a_{23} & 2a_{24} & \dots & (n-2)a_{2n} \\ 0 & \dots & a_{33} & \dots & \dots & \vdots \\ \vdots & & & & & a_{n-1,n} \\ 0 & & & & & 0 & a_{nn} \end{pmatrix}$$

Donc :

$$\delta_1 y = \varepsilon \bar{C}_1 x^*$$

Il en résulte que en ordre de grandeur on peut dire que:

$$\Delta y = \delta y + \delta_1 y = \varepsilon (R + \bar{C}_1) x^*$$

De la formule(1) on tire donc :

$$\frac{\varphi_j(x^* - x)}{\varphi_j(x^*)} \leq \frac{1}{\gamma_{ij}(A)} \cdot \varepsilon \cdot \frac{1}{S_{ij}(A)} \cdot \left( \frac{\varphi_j((R+\bar{E}_1)x^*)}{\varphi_j(x^*)} + \frac{\varphi_j((R+T)x^*)}{\varphi_j(x^*)} \right)$$

puis,

$$\frac{\varphi_j(x^* - x)}{\varphi_j(x^*)} \leq \frac{\varepsilon}{\gamma_{ij}(A)} \cdot \frac{1}{S_{ij}(A)} \left( S_{ij}(R+\bar{E}_1) + S_{ij}(R+T) \right)$$

enfin,

$$\frac{\varphi_j(x^* - x)}{\varphi_j(x^*)} \leq \frac{\varepsilon}{\gamma_{ij}(A)} \cdot \frac{2S_{ij}(R) + S_{ij}(T) + S_{ij}(\bar{E}_1)}{S_{ij}(A)}$$

qui est la majoration en erreur relative cherchée.

Dans les cas particulier qui nous intéressent,

$$\frac{\|x^* - x\|}{\|x^*\|} \leq \frac{\varepsilon}{m_A} \cdot (2M_R + M_T + M_{\bar{E}_1})$$

et

$$\frac{\|x^* - x\|}{\|x^*\|} \leq \frac{\varepsilon}{\sqrt{\lambda_1}} \left( 2\sqrt{\Lambda_R} + \sqrt{\Lambda_T} + \sqrt{\Lambda_{\bar{E}_1}} \right)$$

Si  $\Lambda_R, \Lambda_T, \Lambda_{\bar{E}_1}$  désignent les plus grandes valeurs propres des matrices  $R^T \cdot R, T^T \cdot T, \bar{E}_1^T \cdot \bar{E}_1$ .

En particulier, soit  $m_B$  le plus grand des nombres (en valeur absolue) qui peuvent être obtenus dans le calcul, c'est à dire pouvant être un élément des matrices  $A^{(k)}$ .

Puisque

$$M_R \leq n\sqrt{n} \cdot (n-1) m_B$$

et

$$M_T \leq n\sqrt{n} m_B$$

et

$$M_{\bar{E}_1} \leq n\sqrt{n} (n-1) m_B$$

on peut dire que:

$$\frac{\|x^* - x\|}{\|x^*\|} \leq \frac{\varepsilon}{m_A} \cdot m_B \cdot (2n-2 + 1 + n-1) n^{\frac{3}{2}}$$

Donc :

$$\frac{\|x^* - x\|}{\|x^*\|} \leq (3n-2) n^{\frac{3}{2}} \cdot \frac{\varepsilon}{m_A} \cdot m_B$$

Théorème X Pour pouvoir résoudre en virgule fixe un système linéaire, afin d'assurer au résultat une précision relative inférieure à  $\beta^{-p}$  il suffit de conduire les calculs avec une virgule flottante ayant une capacité de mantisse  $N$  satisfaisant à:

$$(3n-2) \cdot n^{\frac{3}{2}} \beta^{-N} \cdot \frac{m_B}{m_A} \leq \beta^{-p}$$

$$\left| N \geq p - \log_{\beta} \left( \frac{m_A}{m_B} \cdot \frac{1}{(3n-2) \cdot n^{\frac{3}{2}}} \right) \right|$$

On peut remarquer que dans cette formule intervient seulement le nombre  $n_A/m_B$  qui est presque de même nature que les conditionnements généraux que nous avons définis. On peut appeler ce nombre " conditionnement global " pour le système à résoudre

3°) Expériences de calcul en virgule flottante.

On peut se demander si les évaluations d'erreurs qui sont données dans ce qui précède ne sont pas trop peu précises en ce sens que les inégalités qu'elles font intervenir ne sont pas assez "serrées".

Il est donc indispensable de faire des expériences de calcul pour savoir ce qu'il en est.

Nous avons considéré le système suivant,  $Ax = y$ ,

avec

$$A = \begin{pmatrix} 6! & 0 & \dots & 0 \\ 0 & 7! & & \\ & & 8! & \\ & & & 9! \\ & & & & 10! \\ & & & & & \dots \\ & & & & & & 11! \end{pmatrix} \cdot \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} \\ \frac{1}{2} & \frac{1}{3} & \dots & & & \frac{1}{7} \\ \frac{1}{3} & & & & & \\ \vdots & & & & & \\ \frac{1}{6} & \dots & & & & \frac{1}{11} \end{pmatrix}, \quad y = \begin{pmatrix} 6! \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

La solution exacte de ce système n'est autre que la première colonne :

$$x = \begin{pmatrix} 36 \\ -630 \\ 3.360 \\ -7.560 \\ 7.560 \\ -2.772 \end{pmatrix}$$

de l'inverse de la matrice  $H_6$  de HILBERT;

Dans ce cas, on trouve :

$$n_A \# 2$$

$$m_B \# 3 \cdot 10^6$$

$$(3n-2) \cdot n^2 \# 2,3 \cdot 10^2$$

Si bien que la formule du théorème X s'écrit :

$$N \geq p - \log_{10} \left( \frac{2}{3} \cdot 10^{-8} \cdot \frac{1}{2,3} \right) \# p + 8 + 0,65$$

On voit d'après cela que le nombre de chiffres significatifs à la mantisse du vecteur solution doit être à peu près :

$$p \# N - 8,65$$

# H<sub>6</sub>

Pas de chiffre  
significatif

Valeur exacte	3	4	5	6	7	8	9	9
							PDF	PDF avec arrondi
36	22,9	12,38	15,456	49,3277	36,526 11	36,067 416	36,009 469 4	35,996 275 0
- 630	-205	- 60,26	- 82,547	- 998,547	- 644,152 3	- 631,850 95	- 630,253 404	- 629,901 619
3 360	477	- 43,12	- 167,93	5 808,43	3 451,857	3 372,223 6	3 361,643 28	3 359,366 98
-7 560	-252	480,8	1 289,0	-13 852,7	-7 791,998	-7 591,270 2	-7 564,149 10	-7 558,441 51
7 560	-219	-662,7	-1 935,7	14 447,7	7 810,528	7 594,100 6	7 564,479 93	7 558,293 47
-2 772	174	272,6	886,1	- 5 469,26	-2 869,055	-2 785,313 9	-2 773,735 48	-2 771,341 69
Max. erreur relative sur une compos.	?	?	?	1	# $\frac{1}{30}$	$\frac{3}{750}$ # $\frac{4}{1000}$	$\frac{4}{7560}$ # $\frac{5}{10000}$	$\frac{2}{7560}$ # $\frac{2}{10000}$

Tableau des résultats des différentes résolutions au sujet du 2 8

Ce qui revient à dire que pour résoudre le système avec quelque chance d'avoir des résultats avec un chiffre décimal exact, il faut atteindre les mantisses de capacité 8 au moins. La table ci-dessus indique des résultats moins pessimistes. On constate, sur cette table, que c'est à partir de 7 chiffres significatifs (en décimal) de mantisse que l'on a une erreur relative, sur la solution, de l'ordre de  $1/10$ . En fait, le tableau montre que l'on a  $P \approx N-6$ , au lieu de la formule précédente. Cela prouve que notre formule n'est pas trop "lâche". Il faut signaler, au point de vue de la réalisation de ces expériences, que c'est encore le principe de simulation d'opérations qui nous a servi pour l'exécution en machine de ces calculs. On peut noter le gain assez faible qu'apporte l'utilisation d'une virgule flottante à 9 chiffres de mantisse et avec arrondi sur l'utilisation sans arrondi, mais avec toujours 9 chiffres de mantisse.

9°) Evolution des conditionnements de systèmes successifs d'une élimination.

Dans ce dernier paragraphe, nous nous proposons d'étudier ce que l'on peut dire des conditionnements des matrices obtenues comme premiers membres des systèmes provenant d'une élimination.

Soit à résoudre le système  $Ax = y$

Eliminons la première inconnue  $x_1$  ; on peut considérer le système  $A' \cdot x' = y'$  où la matrice est la partie carrée

de la matrice  $A^{(2)}$ , c'est à dire que:

$$a'_{ij} = a_{ij} - \frac{a_{i1} \cdot a_{1j}}{a_{11}} \quad (i, j = 2, \dots, n)$$

et l'ordre en étant  $n-1$ ,

$$x' = \begin{pmatrix} x_2 \\ \vdots \\ x_n \end{pmatrix}, \quad y' = \begin{pmatrix} y_2^{(2)} \\ y_3^{(2)} \\ \vdots \\ y_n^{(2)} \end{pmatrix}$$

On va déterminer les nombres  $m_{A'}, M_{A'}$  pour  $A'$  d'ordre  $n-1$

Posons :

$$\Phi_A(x_1, x_2, \dots, x_n) = \sum_{i=1}^{i=n} |f_i(x_1, \dots, x_n)|$$

si

$$f_i(x_1, \dots, x_n) = \sum_{j=1}^{j=n} a_{ij} x_j$$

puis

$$\Phi_{A'}(x_2, \dots, x_n) = \sum_{i=1}^{i=n} \left| f_i - \frac{a_{i1}}{a_{11}} f_1 \right| \quad (2)$$

On a donc;

$$\Phi_{A'}(x_2, \dots, x_n) \equiv \Phi_A \left( -\frac{1}{a_{11}} (a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n), x_2, x_3, \dots, x_n \right)$$

et par suite

$$\frac{\Phi_{A'}(x_2, \dots, x_n)}{\sqrt{x_2^2 + \dots + x_n^2}} \geq \frac{\Phi_A \left( -\frac{1}{a_{11}} (\dots), x_2, x_3, \dots, x_n \right)}{\sqrt{\frac{1}{a_{11}^2} (\dots)^2 + x_2^2 + \dots + x_n^2}} \geq m_A$$

D'où :

$$m_A \leq m_{A'} \quad (I)$$

D'autre part, d'après (1)

$$\Phi_{A'}(x_2, \dots, x_n) \leq \Phi_A(x_1, \dots, x_n) + |f_1| \cdot \left[ \frac{|a_{21}| + \dots + |a_{n1}|}{|a_{11}|} - 1 \right]$$

donc:

$$\frac{\Phi_{A'}(x_2, \dots, x_n)}{\sqrt{x_2^2 + \dots + x_n^2}} \leq \frac{\Phi_A(x_1, \dots, x_n)}{\sqrt{x_1^2 + \dots + x_n^2}} + \frac{|f_1|}{\sqrt{x_1^2 + \dots + x_n^2}} \cdot \left[ \frac{|a_{21}| + \dots + |a_{n1}|}{|a_{11}|} - 1 \right]$$

et si l'on fait  $x_1 = 0$  dans cette inégalité valable quels que soient

les  $x_1, x_2, \dots, x_n$ ,

$$\frac{\Phi_{A'}(x_2, \dots, x_n)}{\sqrt{x_2^2 + \dots + x_n^2}} \leq \frac{\Phi_A(0, x_2, \dots, x_n)}{\sqrt{0 + x_2^2 + \dots + x_n^2}} + \frac{|a_{12}x_2 + \dots + a_{1n}x_n|}{\sqrt{x_2^2 + \dots + x_n^2}} \cdot \left[ \frac{|a_{21}| + \dots + |a_{n1}|}{|a_{11}|} - 1 \right]$$

I) si  $\left[ \frac{|a_{21}| + \dots + |a_{n1}|}{|a_{11}|} - 1 \right] < 0$ , il est clair qu'on peut choisir les  $x_2, \dots, x_n$

de sorte que  $|a_{12}x_2 + \dots + a_{1n}x_n| = 0$  et par conséquent:

$$\frac{\Phi_{A'}}{\|x'\|} \leq \frac{\Phi_A(0, x_2, \dots, x_n)}{\sqrt{0 + x_2^2 + \dots + x_n^2}} \leq M_A$$

d'où:

$$M_{A'} \leq M_A \quad (II)$$

2) Si  $\left[ \frac{|a_{21}| + \dots + |a_{n1}|}{|a_{11}|} - 1 \right] \gg 0$ , puisque (Inégalité de SCHWARZ)

$$\frac{|a_{12}x_2 + \dots + a_{1n}x_n|}{\sqrt{x_2^2 + \dots + x_n^2}} \leq \sqrt{a_{12}^2 + \dots + a_{1n}^2}$$

on a donc,

$$\frac{\Phi_{A'}}{\|x'\|} \leq \frac{\Phi_A(0, x_2, \dots, x_n)}{\sqrt{x_2^2 + \dots + x_n^2}} + \sqrt{a_{12}^2 + \dots + a_{1n}^2} \left[ \frac{|a_{21}| + \dots + |a_{n1}|}{|a_{11}|} - 1 \right]$$

or si dans  $\frac{\Phi_A(x)}{\|x\|}$  on fait  $x_1 = 0, x_i = a_{1i} \quad (i=2, \dots, n)$

cela donne:  $\sqrt{a_{12}^2 + \dots + a_{1n}^2} \leq M_A$

et par suite:  $\frac{\Phi_{A'}}{\|x'\|} \leq M_A + M_A \left[ \frac{|a_{21}| + \dots + |a_{n1}|}{|a_{11}|} - 1 \right]$

d'où  $M_{A'} \leq M_A \cdot \frac{|a_{21}| + \dots + |a_{n1}|}{|a_{11}|} \quad (II)'$

Mais les deux cas peuvent se résumer dans cette même formule:

(III)  $M_{A'} \leq M_A \cdot K, \quad \frac{1}{K} \frac{m_A}{M_A} \leq \frac{m_{A'}}{M_{A'}}$

où  $K = \max \left\{ 1, \frac{|a_{21}| + \dots + |a_{n1}|}{|a_{11}|} \right\}$

D'où le théorème:

théorème XII: Pour être sûr que l'élimination de  $x_1$ , conduite à un système de conditionnement général  $\frac{m_{A'}}{M_{A'}}$  supérieur à celui  $\frac{m_A}{M_A}$  du système initial, il suffit que l'on ait:

$$|a_{11}| \gg |a_{21}| + |a_{31}| + \dots + |a_{n1}|$$



CHAPITRE VI

Les erreurs dans la résolution  
des systèmes linéaires par orthogonalisation .

1° Solution d'un système linéaire par orthogonalisation en ligne :

Soit  $A \cdot x = y$  (1) : un système linéaire

$A = (a_{ij})$  ( $i, j = 1, \dots, n$ ) <sup>matrice</sup> non singulière, et à éléments réels.

On désigne toujours par  $a_k$  (un seul indice), la  $k$ -ième ligne de A,

$$a_k = A_k = (a_{k1}, a_{k2}, \dots, a_{kn})$$

$y_i$ , ( $i = 1, \dots, n$ ) sont les composantes du vecteur du 2° membre.

La méthode consiste à remplacer le système donné par

$$H \cdot x = Y \quad (2)$$

avec  $H = A^{(n)} = S_{n-1} \cdot S_{n-2} \cdot \dots \cdot S_2 \cdot A$  ;  $Y = S_{n-1} \cdot S_{n-2} \cdot \dots \cdot S_1 \cdot y$   
H orthogonale en lignes. (cf. Chapitre I)

La solution du système  $H \cdot x = Y$  est alors très facile, en effet l'orthogonalité en lignes de H se traduit, sous forme matricielle, par:

$$H \cdot H^T = D$$

où D est une matrice diagonale, dont les éléments diagonaux sont les carrés des longueurs des lignes de H .

Si l'on pose

$$x = H^T \cdot z$$

le système (1) est remplacé par:

puis :

$$\begin{cases} HH^T z = Y \\ x = H^T z \end{cases}$$

Enfin:

$$(3) \quad x = H^T \cdot D^{-1} \cdot Y$$

Pour écrire les formules de calcul nous poserons:

$$H = (h_{ij}) \quad , \quad (i, j = 1, 2, \dots, n)$$

puis,  $h_k = (h_{k1}, h_{k2}, \dots, h_{kn})$

et  $\gamma = (\gamma_i) \quad , \quad (i=1, \dots, n)$

alors on a :

$$(I) \quad \begin{cases} h_1 = a_1 \\ h_2 = a_2 + \lambda_{21} \cdot h_1 \\ h_3 = a_3 + \lambda_{31} \cdot h_1 + \lambda_{32} \cdot h_2 \\ \vdots \\ h_n = a_n + \lambda_{n1} \cdot h_1 + \dots + \lambda_{n,n-1} \cdot h_{n-1} \end{cases}$$

$$(II) \quad \begin{cases} \gamma_1 = y_1 \\ \gamma_2 = y_2 + \lambda_{21} \gamma_1 \\ \gamma_3 = y_3 + \lambda_{31} \gamma_1 + \lambda_{32} \gamma_2 \\ \vdots \\ \gamma_n = y_n + \lambda_{n1} \gamma_1 + \dots + \lambda_{n,n-1} \gamma_{n-1} \end{cases}$$

$$(III) \quad \lambda_{ij} = - \frac{a_i \cdot h_j^T}{h_j \cdot h_j^T} \quad \begin{matrix} (i > j ; & j = 1, \dots, n-1) \\ & i = 2, \dots, n \end{matrix}$$

D'après ces formules on voit que si l'on pose

$$\Lambda = \begin{vmatrix} 0 & 0 & 0 & \dots & 0 \\ \lambda_{21} & 0 & \dots & \dots & 0 \\ \lambda_{31} & \lambda_{32} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \lambda_{n1} & \lambda_{n2} & \dots & \lambda_{n,n-1} & 0 \end{vmatrix}$$

Les formules ci-dessus permettent d'écrire:

$$(I)' \quad H = A + \Lambda \cdot H$$

$$(II)' \quad \gamma = y + \Lambda \cdot \gamma$$

et si  $T = (I - \Lambda)^{-1}$ , triangulaire inférieure et unitaire,

$$(I)'' \quad H = T \cdot A$$

$$(II)'' \quad \gamma = T \cdot y$$

La matrice  $T$  n'est autre que le produit:

$$T = S_{n-1} \cdot S_{n-2} \cdot \dots \cdot S_2 \cdot S_1 \cdot A$$

Comme nous l'avons déjà fait nous allons étudier les erreurs de calcul pour cette méthode en séparant le cas où le calcul est conduit en virgule fixe de celui conduit en virgule flottante.

2.9 Etude des erreurs en calcul en virgule fixe .

Nous désignerons par  $\epsilon_{ij}$  l'erreur de calcul d'un  $\lambda_{ij}$  :

$$\bar{\lambda}_{ij} = -\frac{a_i \cdot \bar{h}_j}{\bar{h}_j \cdot \bar{h}_j} + \epsilon_{ij}$$

et par suite la détermination de H et Y conduit à  $\bar{H}, \bar{Y}$  tels

$$(I)_e \left\{ \begin{array}{l} \bar{h}_1 = a_1 \\ \bar{h}_2 = a_2 + \bar{\lambda}_{21} \circ \bar{h}_1 \\ \bar{h}_3 = a_3 + \bar{\lambda}_{31} \circ \bar{h}_1 + \bar{\lambda}_{32} \circ \bar{h}_2 \\ \vdots \\ \bar{h}_n = a_n + \bar{\lambda}_{n1} \circ \bar{h}_1 + \dots + \bar{\lambda}_{nn-1} \circ \bar{h}_{n-1} \end{array} \right. , \quad (II)_e \left\{ \begin{array}{l} \bar{y}_1 = y_1 \\ \bar{y}_2 = y_2 + \bar{\lambda}_{21} \circ \bar{y}_1 \\ \vdots \\ \bar{y}_n = y_n + \bar{\lambda}_{n1} \circ \bar{y}_1 + \dots + \bar{\lambda}_{nn-1} \circ \bar{y}_{n-1} \end{array} \right.$$

si nous supposons (calcul en fixe) que nous ne faisons pas d'erreur dans les additions algébriques .

Posons  $\bar{\lambda}_{ij} \circ \bar{h}_j = \bar{\lambda}_{ij} \cdot \bar{h}_j + \eta_{ij}$

$\bar{\lambda}_{ij} \circ \bar{y}_j = \bar{\lambda}_{ij} \cdot \bar{y}_j + \rho_{ij}$

où  $\eta_{ij}$  désigne une ligne dont les valeurs sont les erreurs de multiplication dans l'évaluation du produit:

$$\bar{\lambda}_{ij} \cdot \bar{h}_j$$

Les formules de calcul réel s'écrivent donc:

$$(I)'_e \left\{ \begin{array}{l} \bar{h}_1 = a_1 \\ \bar{h}_2 = a_2 + (\lambda_{21} + \epsilon_{21}) \cdot \bar{h}_1 + \eta_{21} \\ \bar{h}_3 = a_3 + (\lambda_{31} + \epsilon_{31}) \cdot \bar{h}_1 + (\lambda_{32} + \epsilon_{32}) \cdot \bar{h}_2 + \eta_{31} + \eta_{32} \\ \vdots \\ \bar{h}_n = a_n + (\lambda_{n1} + \epsilon_{n1}) \cdot \bar{h}_1 + (\lambda_{n2} + \epsilon_{n2}) \cdot \bar{h}_2 + \dots + (\lambda_{nn-1} + \epsilon_{nn-1}) \cdot \bar{h}_{n-1} + \eta_{n1} + \eta_{n2} + \dots + \eta_{nn-1} \end{array} \right.$$

$$(II)'_e \left\{ \begin{array}{l} \bar{y}_1 = y_1 \\ \bar{y}_2 = y_2 + (\lambda_{21} + \epsilon_{21}) \cdot \bar{y}_1 + \rho_{21} \\ \vdots \\ \bar{y}_n = y_n + (\lambda_{n1} + \epsilon_{n1}) \cdot \bar{y}_1 + (\lambda_{n2} + \epsilon_{n2}) \cdot \bar{y}_2 + \dots + (\lambda_{nn-1} + \epsilon_{nn-1}) \cdot \bar{y}_{n-1} + \rho_{n1} + \rho_{n2} + \dots + \rho_{nn-1} \end{array} \right.$$

Par suite les  $\bar{H}$  et  $\bar{Y}$  auxquels on arrive satisfont aux équations:

$$(I)''_e \quad \bar{H} = A + (\Lambda + E) \cdot \bar{H} + K$$

$$(II)''_e \quad \bar{Y} = y + (\Lambda + E) \cdot \bar{Y} + P$$

Si l'on pose :

$$E = \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ \varepsilon_{21} & 0 & 0 & \dots & 0 \\ \varepsilon_{31} & \varepsilon_{32} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \varepsilon_{n1} & \varepsilon_{n2} & \dots & \dots & \varepsilon_{nn} \end{pmatrix}, \quad K = \begin{pmatrix} 0 \\ \eta_{21} \\ \eta_{31} + \eta_{32} \\ \vdots \\ \eta_{n1} + \eta_{n2} + \dots + \eta_{nn} \end{pmatrix}, \quad P = \begin{pmatrix} 0 \\ \rho_{21} \\ \rho_{31} + \rho_{32} \\ \vdots \\ \rho_{n1} + \rho_{n2} + \dots + \rho_{nn} \end{pmatrix}$$

Pour la matrice K il faut noter que chaque symbole représente une ligne.

On peut remarquer qu'en ordre de grandeur, ces matrices peuvent être écrites

$$E \# e', \quad K \# e, \quad P \# e.$$

$$\begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ 1 & 0 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & \dots & 1 \end{pmatrix}, \quad \begin{pmatrix} 0 & 0 & 0 & 0 & \dots & 0 \\ 1 & 1 & 1 & 1 & \dots & 1 \\ 2 & 2 & 2 & \dots & \dots & 2 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ n-1 & n-1 & \dots & \dots & n-1 & n-1 \end{pmatrix}, \quad \begin{pmatrix} 0 \\ 1 \\ 2 \\ \vdots \\ n-1 \end{pmatrix}$$

Si comme plus haut,  $e'$  et  $e$  désignent les erreurs en division et en multiplication.

Donc la solution de  $Ax = y$ , rigoureusement la même que celle du système  $Hx = Y$ , est remplacée par celle de :  $\bar{H}x = \bar{Y}$

Posons encore  $\bar{H} = H + \delta H$  et  $\bar{Y} = Y + \delta Y$ .

Il en résulte que :

$$\delta H = (\Lambda + E) \cdot (H + \delta H) - \Lambda \cdot H + K$$

$$\delta Y = (\Lambda + E) \cdot (Y + \delta Y) - \Lambda \cdot Y + P$$

d'où en négligeant les termes de  $E \cdot \delta H$  et ceux de  $E \cdot \delta Y$ ,

$$\delta H = E \cdot H + \Lambda \cdot \delta H + K$$

$$\delta Y = E \cdot Y + \Lambda \cdot \delta Y + P$$

et puisque  $T = (I - \Lambda)^{-1}$ , on en déduit le théorème :

Théorème I.

Si  $T$  est la matrice triangulaire inférieure unitaire qui dans la procédé d'orthogonalisation permet d'écrire  $H = T \cdot A$  avec  $H$  orthogonale en lignes, le calcul numérique, en virgule fixe, de cette matrice conduit à une matrice  $\bar{H}$  et le système à résoudre  $Hx = Y$  équivalent à  $Ax = y$  est remplacé par  $(H + \delta H)x = Y + \delta Y$  avec,

$$\delta H = T \cdot (E \cdot H + K)$$

$$\delta Y = T \cdot (E \cdot Y + P)$$

Il est bien clair que la matrice  $\bar{H}$  n'est pas orthogonale en ligne, et cependant nous allons traiter le système comme si il était à matrice de premier membre orthogonale en lignes. Cela conduit aux opérations

1°) Déterminer les éléments diagonaux de  $(H + \delta H) \cdot (H + \delta H)^T = (H + \delta H)(H^T + (\delta H)^T)$

Je pose  $\bar{D}$  pour la matrice obtenue en ne prenant que les éléments diagonaux de ce produit, comme on peut écrire :

$$Z = H \cdot H^T + \delta H \cdot H^T + H \cdot (\delta H)^T$$

(négligeant les éléments de  $\delta H \cdot H^T$ )

il en résulte quasi l'on pose encore :  $\bar{D} = D + \delta D$

$$\delta D = \text{diag} [ \delta H \cdot H^T + H \cdot (\delta H)^T ]$$

2°) Cela fait il reste pour déterminer la solution la formule:

$$x^* = \bar{H}^T \circ (\bar{D}^{-1} \circ \bar{Y})$$

- a) Je vais poser  $\bar{D}^{-1} \circ \bar{Y} = \bar{D}^{-1} \cdot \bar{Y} + Q$ , avec

$$Q = \begin{pmatrix} q_1 \\ q_2 \\ \vdots \\ q_n \end{pmatrix}, \quad Q \neq e' = \begin{pmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

Si  $q_i$  désignent les erreurs de division dans le calcul de  $\bar{D}^{-1} \cdot \bar{Y}$

En négligeant les termes du 2° ordre

par rapport à  $\delta D$ ,

$$\begin{aligned} \bar{D}^{-1} &= (D + \delta D)^{-1} = (I + D^{-1} \delta D)^{-1} \cdot D^{-1} = (I - D^{-1} \delta D) \cdot D^{-1} \\ &= D^{-1} - D^{-1} \delta D \cdot D^{-1} \end{aligned}$$

d'où :

$$\bar{D}^{-1} \circ \bar{Y} = (D^{-1} - D^{-1} \delta D \cdot D^{-1}) \cdot (Y + \delta Y) + Q$$

enfin

$$\bar{D}^{-1} \circ \bar{Y} = D^{-1} \cdot Y + D^{-1} \delta Y - D^{-1} \delta D \cdot D^{-1} \cdot Y + Q$$

- b) Il reste à calculer le produit  $\bar{H}^T \circ (\bar{D}^{-1} \circ \bar{Y})$

Là encore on posera  $\bar{H}^T \circ (\bar{D}^{-1} \circ \bar{Y}) = \bar{H}^T \cdot (\bar{D}^{-1} \circ \bar{Y}) + R$

la matrice étant d'ordre n, le vecteur d'erreur  $R$  aura pour composantes les erreurs sur les composantes du produit et:

$$R = \begin{pmatrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{pmatrix}, \quad R \neq e = \begin{pmatrix} n \\ n \\ \vdots \\ n \end{pmatrix}$$

Il en résulte que :

$$x^* = x + \delta x = (H^T + (\delta H)^T) \cdot (D^{-1} \gamma + D^{-1} \delta \gamma - D^{-1} \delta D \cdot D^{-1} \gamma + Q) + R$$

Donc :

Théorème II. La méthode de résolution d'un système d'équations linéaires par le procédé d'orthogonalisation en lignes, conduit à une solution  $x^* = x + \delta x$  avec

$$(I) \quad \delta x = (\delta H)^T \cdot D^{-1} \gamma - H^T D^{-1} \delta D \cdot D^{-1} \gamma + H^T D^{-1} \delta \gamma + H^T Q + R$$

les matrices et les vecteurs  $\delta H, \delta D, \delta \gamma, Q, R$  ayant été précédemment définis, les calculs étant conduits en virgule fixe.

Cette formule permet d'écrire :

$$\delta x = \delta (H^T \cdot D^{-1} \gamma) + H^T Q + R$$

Or :

$$H^T D^{-1} = H^{-1}$$

donc :

$$\delta x = H^{-1} \delta \gamma - H^{-1} \delta H \cdot H^{-1} \gamma + H^T Q + R$$

Mais si l'on revient aux expressions indiquées plus haut pour  $\delta H$  et  $\delta \gamma$  cela se réduit à la formule très simple :

$$\delta x = H^{-1} [T(E\gamma + P) - T(EH + K)H^{-1}\gamma] + H^T Q + R$$

$$(II) \quad \delta x = A^{-1} \cdot [P - Kx^*] + H^T Q + R$$

$x^*$  désigne la solution obtenue.

Comme utilisation de cette dernière formule nous l'écrivons en ordre de grandeur avec les expressions données pour les vecteurs et matrices  $K, P, Q, R$ .

$$\delta x = e \cdot A^{-1} \cdot \begin{bmatrix} 0 \\ (1 - \xi_1, -\xi_2, \dots, -\xi_n) \\ 2(1 - \xi_1, -\xi_2, \dots, -\xi_n) \\ \vdots \\ (n-1)(1 - \xi_1, \dots, -\xi_n) \end{bmatrix} + e' \cdot H^T \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + ne \begin{bmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

Si  $\xi_1, \xi_2, \dots, \xi_n$ , sont les composantes de la solution:  $x^*$

Le fait le plus remarquable à signaler est que , à cet ordre d'approximation et en calcul en fixe, la matrice B ne figure plus dans l'expression de l'erreur, cela d'après la formule II .

Comme application posons-nous la question suivante:

Etant donné un système linéaire,  $Ax = y$  , résolvons-le une première fois en fixe, par orthogonalisation, cela conduit à une solution  $x_1^*$  , puis échangeons les positions de deux équations du système cela donne une autre solution etc.. Il est clair que je peux ainsi résoudre  $n!$  systèmes linéaires qui ne diffèrent entre-eux qu'à une permutation près des équations, donc qui ont parfaitement le même vecteur solution  $x$ .

Or l'expérience prouve que la valeur du vecteur solution-approchée  $x^*$  n'est pas la même et cela à cause des erreurs de calcul . Quelle est donc la solution la meilleure ?

D'après la dernière formule écrite on voit que la norme ( $\varphi_x$ , par exemple) du premier vecteur donnant  $\delta x$  ne va pas beaucoup varier pour  $n$  peu élevé, par contre pour le second  $H^T \cdot Q$ , les choses sont différentes . il est clair que l'on aura toujours le plus grand intérêt à choisir la résolution qui conduit à affaiblir le plus possible une norme de  $H^T$  .

Or

$$H \cdot H^T = D$$

Puis comme :  $H = T \cdot A$  ,  $\text{Det}(H) = \Delta = \text{Det}(A)$  ,  $\text{Det}(D) = \Delta^2$

On en déduit que :  $N^2(H^T) = \text{trace}(H \cdot H^T) = \text{trace}(D)$

Si l'on pose  $e_i^2 = h_i \cdot h_i^T$  ,  $e_i > 0$

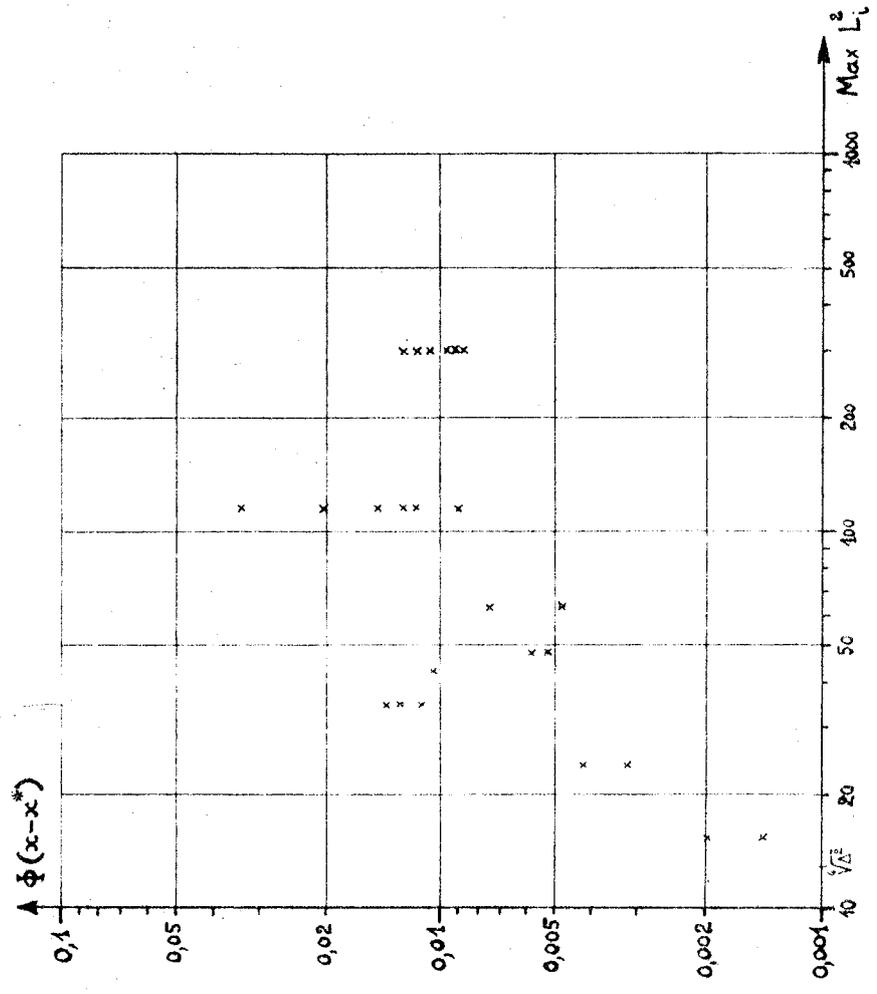
$$N^2(H^T) = N^2(H) = e_1^2 + e_2^2 + \dots + e_n^2$$

Donc puisque le produit  $e_1^2 \cdot e_2^2 \dots e_n^2 = \text{Det}(D) = \Delta^2$  est constant dans toute nos résolutions ,  $N(H^T)$  sera minimum si les longueurs des lignes de  $H$  sont toutes presque égales entre-elles.

D'où:

Règle Afin de minimiser l'erreur de calcul en virgule fixe dans la solution d'un système linéaire par la méthode d'orthogonalisation , il suffit d'organiser le système de sorte que les différentes <sup>longueurs des</sup> lignes de  $H$  soient les plus près possible de l'égalité entre-elles et par suite près de :  $\sqrt[n]{\Delta}$

Afin de vérifier cette règle nous avons fait résoudre au calculateur



Resolutions d'un même système linéaire par la méthode d'orthogonalisation en lignes.

électronique les 24 systèmes linéaires obtenus par toutes les permutations des équations à partir du système déjà utilisé, d'ordre 4 au chapitre V, § 6 qui nous a servi à vérifier la règle du choix des pivots dans une résolution par élimination. Les résultats sont indiqués sur la figure ci-contre. On a porté en abscisse sur ce graphique  $M_{\max} (e_i^2)$  pour une résolution et en ordonnée, la valeur de la norme  $\Phi(x^2-z)$  de l'erreur correspondant à cette résolution.

3°) Etude des erreurs en calcul en virgule flottante .

De même qu'au § 3, nous sommes conduits à trouver une matrice  $\bar{H}$  et un 2° membre  $\bar{Y}$  définis par les relations (mêmes notations)

$$(I)'_f \left\{ \begin{array}{l} \bar{h}_1 = a_1 \\ \bar{h}_2 = a_2 + (\lambda_{21} + \varepsilon_{21}) \cdot \bar{h}_1 + \xi_2 \cdot \bar{h}_2 \\ \vdots \\ \bar{h}_n = a_n + (\lambda_{n1} + \varepsilon_{n1}) \cdot \bar{h}_1 + (\lambda_{n2} + \varepsilon_{n2}) \cdot \bar{h}_2 + \dots + (\lambda_{nn-1} + \varepsilon_{nn-1}) \cdot \bar{h}_{n-1} + \xi_n \cdot \bar{h}_n \end{array} \right.$$

$$(II)'_f \left\{ \begin{array}{l} \bar{y}_1 = a_1 \\ \bar{y}_2 = a_2 + (\lambda_{21} + \varepsilon_{21}) \cdot \bar{y}_1 + \pi_2 \cdot \bar{y}_2 \\ \vdots \\ \bar{y}_n = a_n + (\lambda_{n1} + \varepsilon_{n1}) \cdot \bar{y}_1 + \dots + (\lambda_{nn-1} + \varepsilon_{nn-1}) \cdot \bar{y}_{n-1} + \pi_n \cdot \bar{y}_n \end{array} \right.$$

Cette fois  $\xi_i$  et  $\pi_i$  désignent les erreurs relatives faites sur les calculs des  $\bar{h}_i$  et  $\bar{y}_i$ .

On peut écrire ces formules sous forme matricielle :

$$(I)''_f \quad \bar{H} = A + (\Lambda + E) \cdot \bar{H}' + K_1 \cdot \bar{H}$$

$$(II)''_f \quad \bar{Y} = y + (\Lambda + E) \cdot \bar{Y}' + P_1 \cdot \bar{Y}$$

Où l'on a posé :

$$K_1 = \begin{vmatrix} 0 & & & & \\ & \xi_2 & & & \\ & & \xi_3 & & \\ & & & \ddots & \\ & & & & \xi_n \\ & 0 & & & \end{vmatrix}, \quad P_1 = \begin{vmatrix} 0 & & & & \\ & \pi_2 & & & \\ & & \pi_3 & & \\ & & & \ddots & \\ & & & & \pi_n \\ & 0 & & & \end{vmatrix}$$

Il est à remarquer que les formules  $(I)'_i$  ne sont qu'approchées car le fait d'écrire que l'erreur absolue dans le calcul de la ligne  $\bar{h}_i$  est proportionnelle à cette ligne, elle-même, est assez schématique. Malgré cela, nous nous contenterons de cette hypothèse simplificatrice et nous poursuivrons dans ce cas l'étude des erreurs. Il est à noter que dans l'étude faite en virgule fixe les formules  $(I)'_i$  sont exactes.

On notera aussi les valeurs en ordre de grandeur des matrices:

$$K_1 = \varepsilon \begin{vmatrix} 0 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 & & \\ & & & & & 0 \\ & & & & & & & & & 1 \end{vmatrix}, \quad P_1 = \varepsilon \begin{vmatrix} 0 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 & & \\ & & & & & 0 \\ & & & & & & & & & 1 \end{vmatrix}$$

Cela posé, si  $\bar{H} = H + \delta H$  et  $\bar{Y} = Y + \delta Y$

On trouve de la même façon qu'au § 3,  $\delta H, \delta Y$ :

D'où:

Théorème III

Si T est la matrice triangulaire inférieure et unitaire qui dans le procédé d'orthogonalisation permet d'écrire  $H = T.A$ , avec H orthogonale en lignes, le calcul numérique en virgule flottante, de cette matrice conduit à une matrice  $\bar{H}$  et le système à résoudre  $Ax = y$ , qui a la même solution que  $Hx = Y$  est remplacé par  $(H + \delta H) \cdot x = Y + \delta Y$ , avec:

$$\delta H = T.(E + K_1).H$$

$$\delta Y = T.(E + P_1).Y$$

Si l'on continue à suivre le même raisonnement qu'au § 3, nous voyons que:  $x^* = \bar{H}^T \square (\bar{D}^{-1} \square \bar{Y})$

-1°) On peut encore écrire

$$\bar{D}^{-1} \square \bar{Y} = \bar{D}^{-1} \cdot \bar{Y} + Q_1 \cdot (\bar{D}^{-1} \bar{Y})$$

où

$$Q_1 = \begin{vmatrix} q'_1 & & & \\ & q'_2 & & \\ & & \ddots & \\ & & & q'_k & & \\ & & & & & 0 \end{vmatrix} \quad \# \quad \varepsilon \cdot I$$

-2°) Puis pour calculer le produit  $\bar{H}^T \square (\bar{D}^{-1} \square \bar{Y}) = x^*$

On aura:

$$\bar{H}^T \square (\bar{D}^{-1} \square \bar{Y}) = \bar{H}^T \cdot (\bar{D}^{-1} \square \bar{Y}) + R_1 \cdot [\bar{H}^T (\bar{D}^{-1} \square \bar{Y})]$$

Avec :

$$R_1 = \begin{vmatrix} r_{11} & & & \\ & r_{22} & & \\ & & \ddots & \\ & & & r_{nn} \end{vmatrix} \# \varepsilon \cdot I, \quad R_1 \cdot [\bar{H}^T \cdot (\bar{D}^{-1} \bar{Y})] = R_1 \cdot x^*$$

Enfin, après le même calcul que celui déjà fait au paragraphe 3, on trouve :

$$\delta x = S (H^T \cdot D^{-1} Y) + H^T Q \cdot D^{-1} Y + R_1 \cdot x^*$$

D'où :

Théorème IV : La méthode de résolution d'un système d'équations linéaires par le procédé d'orthogonalisation en lignes, en virgule flottante, conduit à une solution

$$x^* = x + \delta x \quad \text{avec :}$$

$$\delta x = A^{-1} [P \cdot H - K \cdot H] x^* + H^T Q \cdot D^{-1} H x^* + R_1 x^*$$

les matrices  $P, K, Q, R_1$ , ayant été précédemment définies.

En ordre de grandeur, on peut écrire :

et, par suite :

$$\begin{aligned} \delta x &\# 2\varepsilon A^{-1} H x^* + 2\varepsilon x^* \\ \delta x &\# 2\varepsilon (T + I) x^* = 2\varepsilon T_1 x^* \quad (1) \end{aligned}$$

$$T_1 = T + I, \quad \text{triangulaire}$$

Ces formules sont, certes, assez grossières, mais leur extrême simplicité peut justifier leur utilisation.

De la formule (1), on tire :  $\frac{\Phi(\delta x)}{\|x^*\|} \leq 2\varepsilon \cdot M_{T_1}$

$M_{T_1} = \max_E [T_1^T E]$ , facile à calculer, et, par suite, si l'on veut obtenir une précision relative de l'ordre de  $\beta^{-r}$ ,  $2\varepsilon$  étant de l'ordre de  $\beta^{-k}$ , il suffit d'avoir :

$$\beta^{-k} \cdot M_{T_1} \leq \beta^{-r}$$

ou  $k \geq r + \log_p (M_{T_1})$  d'où le résultat :

Théorème V : Pour obtenir pour la solution d'un système linéaire par la méthode d'orthogonalisation, en point décimal flottant, au moins  $\mu$  chiffres significatifs exacts aux nombres-solution, il suffit de conduire les calculs en point décimal flottant de capacité  $k$  telle :

$$k \geq \mu + \log_p (M_{T_1})$$



CHAPITRE VII

Les normes générales et les procédés itératifs de résolution des systèmes linéaires .

1°) Notations

Dans ce qui suit , pour débiter, nous ne considérons que l'espace affine  $\mathbb{R}^n$  . La solution d'un système linéaire : (0)  $A \cdot x = y$  sera désignée par  $\omega$  , pour la colonne, et le point de l'espace ayant ces coordonnées sera désigné par  $\Omega$  .

Ce point  $\Omega$  est le point commun aux hyper-plans (ou plans ) dont les équations sont:

$$\begin{cases} f_1(x) = A_1 \cdot x - y_1 = 0 \\ f_2(x) = A_2 \cdot x - y_2 = 0 \\ \vdots \\ f_n(x) = A_n \cdot x - y_n = 0 \end{cases}$$

Si  $x_0$  est un "point" de  $\mathbb{R}^n$  le vecteur :

$$r(x_0) = A \cdot x_0 - y$$

est dit "vecteur résidu" en  $x_0$  .

Il est évident que  $r(x_0) = 0$  , si, et seulement si  $x_0 = \omega$  .  
Si non, peut-on obtenir facilement un vecteur (ou un point)  $x_1$  tel que:

$$\varphi(x_1) < \varphi[r(x_0)]$$

Si  $\varphi$  est une certaine norme sur  $\mathbb{R}^n$  ? C'est la question fondamentale posée dans toutes les méthodes de résolution, itératives .

Si  $\varphi(x)$  est une certaine norme je poserai:

$$\varphi(x) = F(x_1, x_2, \dots, x_n)$$

si  $x$  est la colonne :  $x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$

Soit  $x_0$  un point ou  $F(x_1, \dots, x_n) = \varphi(x)$  est différentiable on désignera par  $p(x_0)$  , le vecteur-colonne:

$$p(x) = \begin{pmatrix} \left(\frac{\partial F}{\partial x_1}\right)_0 \\ \left(\frac{\partial F}{\partial x_2}\right)_0 \\ \vdots \\ \left(\frac{\partial F}{\partial x_n}\right)_0 \end{pmatrix} = \text{grad}(\varphi(x))$$

Soit  $x_0$  un point de  $\mathbb{R}^n$ , et considérons le lieu  $\mathcal{B}_\omega^0(A)$  des points  $x$  de  $\mathbb{R}^n$  tels:

$$\varphi(Ax - y) \leq \varphi(Ax_0 - y)$$

puisque  $A\omega = y$ , il est clair que c'est aussi la lieu des points tels:  $\varphi(A(x-\omega)) \leq \varphi(A(x_0-\omega))$

si je pose  $x' = A(x-\omega)$  je vois que ces lieux ne sont autres que les transformés des boules:

$$\varphi(x') \leq \varphi(x'_0)$$

de centre l'origine du système d'axes de  $\mathbb{R}^n$ , dans la transformation définie par:

$$x' = A(x-\omega), \quad x = \omega + A^{-1}x'$$

c'est à dire le produit d'une translation qui les centre en  $\omega$  et d'une affinité définie par  $A$ .

$\Pi_\omega^0(A)$  Désignera la frontière de ce lieu:  $\mathcal{B}_\omega^0(A)$ .

2°) Rapprochement de la solution .

En calcul numérique les seules normes utilisées sont les trois normes  $\varphi_1(x)$ ,  $\varphi_2(x)$ ,  $\varphi_\infty(x)$  déjà étudiées, et le but de la méthode itérative cherchée, doit être de pouvoir obtenir  $x^*$  tel que  $\varphi_1(x^* - \omega) \leq \varepsilon$ , ou  $\varphi_2(x^* - \omega) \leq \varepsilon$ , ou  $\varphi_\infty(x^* - \omega) \leq \varepsilon$ ,  $\varepsilon$  étant donné.

Or nous disposons pour définir  $\omega$  des plans donnés :

$$f_1(x) = 0, \quad f_2(x) = 0, \quad \dots, \quad f_n(x) = 0$$

et aussi d'autant de plans que l'on veut: il suffit de remplacer une des équations ci-dessus par une combinaison linéaire convenable de celles-ci pour obtenir un système ayant la même solution. Si  $z$  est un vecteur colonne quelconque,

$$(1) \quad z^T \cdot (Ax - y) = 0$$

est aussi l'équation d'un plan passant par  $\omega$ .

Or de nos trois normes, une seule, la norme euclidienne  $\varphi_2(x)$  a cette propriété qu'il <sup>est</sup> immédiat de trouver dans un plan arbitraire passant par  $\omega$  un point  $x_1$  plus rapproché, au sens de cette norme, de  $\omega$ . Ce point n'est autre que la projection orthogonale de  $x_0$  sur ce plan.

Soit  $x_0$  donné, cherchons la projection orthogonale  $x_1$  de ce point sur le plan :

$$z^T \cdot (Ax - y) = 0$$

Ce dernier est perpendiculaire au vecteur  $A^T z$  donc  $x_1$ , projection cherchée est tel qu'il existe un scalaire  $\lambda$  avec le quel on puisse écrire:

$$x_1 = x_0 + \lambda \cdot A^T z$$

Mais comme  $x_1$  est dans le plan, c'est que  $\lambda$  satisfait à:

$$z^T \cdot (Ax_0 + \lambda AA^T z - y) = 0$$

d'où :

$$\lambda = - \frac{z^T \cdot (Ax_0 - y)}{z^T AA^T z} = - \frac{z^T \cdot z(x_0)}{\|A^T z\|^2}$$

et par suite :

$$(1) \quad x_1 = x_0 - \frac{z^T \cdot z(x_0)}{\|A^T z\|^2} \cdot A^T z$$

3°) Décomposition d'une norme .

Définition : Nous dirons qu'une norme  $\varphi$  a été décomposée, si l'on a pu trouver une application  $z : x \longrightarrow z(x)$  de  $\mathbb{R}^n$  dans lui-même, telle que l'on ait quel que soit  $x$  :

$$\varphi(x) = z^T(x) \cdot x$$

Il est clair que:

- 1°) Pour  $\Phi(x) = \varphi$ , on peut écrire :  $\Phi(x) = z^T(x) \cdot x$

avec :

$$z(x) = m(x) = \begin{pmatrix} \text{sig}(\xi_1) \\ \text{sig}(\xi_2) \\ \vdots \\ \text{sig}(\xi_n) \end{pmatrix}$$

où la fonction  $\text{sig}(\xi)$  est définie par:

$$\begin{aligned} \text{sig}(\xi) &= +1 & , \xi \geq 0 \\ \text{sig}(\xi) &= -1 & , \xi < 0 \end{aligned}$$

- 2°) Pour  $\mathcal{M}(x) = \varphi_\infty$  on peut écrire:  $\mathcal{M}(x) = z^T(x) \cdot x$

avec :

$$z(x) = \begin{pmatrix} 0 \\ \vdots \\ \text{sig}(\xi_j) \\ \vdots \\ 0 \end{pmatrix} = p(x) \quad , \quad |\xi_j| = \max_{j=1, \dots, n} |\xi_j|$$

3°) Pour  $\|x\| = \varphi_2$ , on peut évidemment écrire :  $\|x\| = \frac{x^T}{\|x\|} \cdot x$

avec :  $z(x) = \frac{x}{\|x\|}$

Plus généralement si B est une matrice symétrique définie positive on sait (24) que  $\sqrt{x^T B x} = \varphi(x)$  est une norme. alors si l'on remarque que :

$$(Bx)^T \cdot x = x^T B x$$

on voit que pour cette norme  $\varphi(x) = z^T(x) \cdot x$  avec  $z(x) = \frac{Bx}{\sqrt{x^T B x}}$

4°) Enfin pour une norme de HOLDER  $\varphi_r(x) = \left( \sum_{i=1}^{i=n} |\xi_i|^r \right)^{\frac{1}{r}}$ ,  $r \geq 1$

on peut écrire  $\varphi_r(x) = z^T(x) \cdot x$

avec

$$z(x) = \frac{1}{\left( \sum_{i=1}^{i=n} |\xi_i|^r \right)^{\frac{r-1}{r}}} \cdot \begin{pmatrix} \alpha_j |\xi_1| \cdot |\xi_1|^{r-1} \\ \alpha_j |\xi_2| \cdot |\xi_2|^{r-1} \\ \vdots \\ \alpha_j |\xi_n| \cdot |\xi_n|^{r-1} \end{pmatrix}$$

Conclusions: Toutes les normes de l'analyse numérique sont très facilement décomposables comme on vient de le voir.

Une remarque importante pour la suite est celle-ci :

Dans tous les cas précédents on a, pour tout x

$$\|z(x)\| \leq K$$

avec  $K = n$  pour le premier cas, dans les autres  $(n, \| \cdot \|)$ ;  $K = 1$

4°) Procédé itératif associé à une décomposition d'une norme.

Dès que l'application  $z$ , en général très simple, a été mise en évidence, nous avons un procédé qui permet de faire correspondre à un  $x_0$  un plan passant par  $\omega$  : celui qui correspond à la valeur  $z(x_0)$  de  $z$ .

Partant de  $x_0$  arbitraire, on prend donc :

$$x_1 = x_0 - \frac{\varphi(x_0)}{\|A^T \cdot z(x_0)\|^2} \cdot A^T \cdot z(x_0)$$

cela, d'après la formule (1) du § 2, et le choix de  $z(x)$ .

Je pose pour simplifier l'écriture :

$$V(x) = A^T \cdot z(x)$$

donc, alors :

$$x_1 = x_0 - \varphi(z(x_0)) \cdot \frac{V(x_0)}{\|V(x_0)\|^2}$$

L'itération consiste à prendre

$$x_n = x_{n-1} - \varphi(z(x_{n-1})) \cdot \frac{V(x_{n-1})}{\|V(x_{n-1})\|^2}$$

Étudions la convergence du procédé.

Je remarque que l'on peut écrire, à cause de l'orthogonalité :

$$\|\omega - x_n\|^2 = \|\omega - x_{n-1}\|^2 - \|x_n - x_{n-1}\|^2 = \|\omega - x_{n-1}\|^2 \cdot [1 - L(x_{n-1})]$$

$$\text{avec } L(x_{n-1}) = \frac{\varphi^2(z(x_{n-1}))}{\|x_{n-1} - \omega\|^2} \cdot \frac{1}{\|V(x_{n-1})\|^2}$$

$$\text{Or } \varphi(z(x_{n-1})) = \varphi(A(x_{n-1} - \omega))$$

$$\text{ou } L(x_{n-1}) = \frac{\varphi^2(A(x_{n-1} - \omega))}{\|x_{n-1} - \omega\|^2} \cdot \frac{1}{\|V(x_{n-1})\|^2}$$

et si  $m_A$  désigne  $\min_{x \neq 0} \frac{\varphi(Ax)}{\|x\|}$ , on peut dire que

$$L(x_{n-1}) \geq \frac{m_A^2}{\|V(x_{n-1})\|^2}$$

D'autre part, (25)

$$\|V(x_{n-1})\|^2 = \|A^T \cdot z(x_{n-1})\|^2 \leq N^2(A^T) \cdot \|z(x_{n-1})\|^2 \leq N^2(A) \cdot K^2$$

en désignant toujours par  $N(A) = \left(\sum_{i,j} a_{ij}^2\right)^{\frac{1}{2}}$  qui est une norme sur l'ensemble des matrices carrées d'ordre  $n$ , et puisque sur cet ensemble toutes les normes de matrices sont équivalentes, il existe, si  $M_A$  désigne la norme  $\max_{x \neq 0} \left(\frac{\varphi(Ax)}{\|x\|}\right)$  (cf chap II) une constante  $R$ , indépendante de  $A$ , telle que

$$N(A) \leq R \cdot M_A$$

Donc,

$$\|V(x_{n-1})\|^2 \leq K^2 \cdot R^2 \cdot M_A^2 = R^2 \cdot M_A^2$$

puis,

$$L(x_{n-1}) \geq \frac{1}{R^2} \cdot \left(\frac{m_A}{M_A}\right)^2$$

enfin,

$$\|\omega - x_n\| \leq \|\omega - x_{n-1}\| \cdot \sqrt{1 - \frac{1}{R^2} \cdot \left(\frac{m_A}{M_A}\right)^2}$$

et de là :

$$\|\omega - x_p\| \leq \|\omega - x_0\| \cdot \left(\sqrt{1 - \frac{1}{R^2} \cdot \left(\frac{m_A}{M_A}\right)^2}\right)^p$$

Et en introduisant le conditionnement général  $\gamma(A)$  par rapport aux deux normes  $\varphi$  et  $\|x\|$ , cela peut s'écrire :

$$\|w - x_p\| \leq \|w - x_0\| \cdot \left( \sqrt{1 - \frac{\gamma^2}{\kappa^2}} \right)^p$$

d'où le théorème :

**Théorème I** Etant donnée une norme  $\varphi$  sur  $R^n$ , le procédé itératif associé à une décomposition :  $\varphi(x) = z^T(x) \cdot x$  de cette norme, avec  $\|z^T(x)\|$  borné quelque soit  $x$ , est toujours convergent. Cette convergence est assurée par l'inégalité :

$$\|w - x_p\| \leq \|w - x_0\| \cdot \left( \sqrt{1 - \frac{\gamma^2}{\kappa^2}} \right)^p$$

qui montre le rôle essentiel joué dans le procédé par le conditionnement général  $\gamma$  de  $A$  par rapport aux deux normes  $\varphi$  et  $\|x\|$ . [26]

De ce théorème on peut déduire autant de procédés itératifs que l'on veut, nous allons en étudier quelques-uns dans ce qui va suivre.

5°) Procédé associé à la norme  $\varphi_\infty = \mathcal{M}(x)$ .

On sait que dans ce cas, on peut écrire  $\mathcal{M}(x) = z^T(x) \cdot x$

avec :

$$z(x) = \begin{pmatrix} 0 \\ \vdots \\ \text{sig}(z_i) \\ \vdots \\ 0 \end{pmatrix}, \quad |z_i| = \text{Max}_{j=1, \dots, n} |z_j|$$

Le plan sur lequel, pour trouver  $x_n$ , on doit projeter orthogonalement  $x_{n-1}$ , n'est autre que

$$z^T(Ax - y) = 0,$$

c'est à dire  $f_i(x) = 0$ ,  $i$  désignant le rang de la composante de plus grande valeur absolue de  $z(x_{n-1})$  :

On voit donc que, au cours de l'itération seuls les plans donnés sont utilisés, le procédé est donc celui des projections sur les plans donnés, mais au lieu de projeter sur ces plans dans un ordre imposé comme dans la méthode de KACZMARZ [27] classique (cf aussi BODEWIG [28]) à chaque étape on choisit, un peu comme dans la méthode de relaxation de SOUTHWELL [29], le plan qui correspond au plus fort résidu en module.

La formule de calcul est :

$$x_n = x_{n-1} - \mathcal{M}(z(x_{n-1})) \cdot \text{sig}(z_i(x_{n-1})) \cdot \frac{A_i^T}{\|A_i\|^2} = x_{n-1} - z_i(x_{n-1}) \cdot \frac{A_i^T}{\|A_i\|^2}$$

et ici :

$$L(x_{n-1}) = \frac{\mathcal{M}^2(z(x_{n-1}))}{\|x_{n-1} - \omega\|^2} \cdot \frac{1}{\|A_i\|^2}$$

donc si  $\mu^0 = \max_i \|A_i\|$

en utilisant comme au chapitre IV,  $m'_A$  pour désigner :  $\frac{\max_{\forall \neq 0} (M(Ax))}{\|x\|}$   
 la formule déterminant le rapprochement devient :

$$\| \omega - x_p \| \leq \| \omega - x_0 \| \cdot \left( \sqrt{1 - \left( \frac{m'_A}{\mu^0} \right)^2} \right)^p$$

Enfin si le système de plans donné est orthogonal le procédé doit théoriquement, s'arrêter au bout de  $n$  tours au plus. Il est à peine besoin de remarquer que dans une utilisation pratique de ce procédé il serait recommandé de transformer le système donné :  $Ax = y$  en un système tel que les lignes de la matrice  $A$  soient toutes de même longueur  $1$ , avec cela les formules de calcul sont :

$$x_{p-1} \rightarrow x_p \begin{cases} \sum_i^{(p)} = \sum_i^{(p-1)} - \rho_i^{(p-1)} \cdot a_{i1} \\ \sum_i^{(p)} = \sum_i^{(p-1)} - \rho_i^{(p-1)} \cdot a_{i2} \\ \vdots \\ \sum_i^{(p)} = \sum_i^{(p-1)} - \rho_i^{(p-1)} \cdot a_{in} \end{cases}$$

$$|\rho_i^{(p-1)}| = \max_{j=1, \dots, n} |\rho_j^{(p-1)}|$$

$$r(x_p) = r(x_{p-1}) - \rho_i^{(p-1)} \cdot A \cdot A^T$$

6°) Procédé associé à la norme  $\varphi_A(x) = \Phi(x)$  .

On a vu que dans ce cas, on peut écrire :  $\Phi(x) = z^T(x) \cdot x$

avec :

$$z(x) = \begin{pmatrix} m_j(\xi_1) \\ m_j(\xi_2) \\ \vdots \\ m_j(\xi_n) \end{pmatrix} = m(x)$$

Les plans sur lesquels pour trouver  $x_p$ , on doit projeter orthogonalement  $x_{p-1}$ , sont les plans passant par  $\omega$  et perpendiculaires aux vecteurs :

$$A^T \cdot E$$

si  $E$  parcourt l'ensemble des vecteurs de composantes de valeur absolue  $1$  .

La formule de calcul est : 
$$x_p = x_{p-1} - \frac{\Phi(x_{p-1})}{\|A^T m(x_{p-1})\|^2} \cdot A^T m(x_{p-1})$$

étudions encore ici la convergence de ce procédé,

$$L(x_{p-1}) = \frac{\Phi^2(r(x_{p-1}))}{\|x_{p-1} - \omega\|^2} \cdot \frac{1}{\|A^T m(x_{p-1})\|^2}$$

et par suite

$$L(x_{p-1}) \geq \frac{m_A^2}{\|A^T m(x_{p-1})\|^2}$$

Or on a vu au chapitre IV que  $M_A = \max_E \|AE\|$ , donc

$$L(x_{p-1}) \geq \frac{m_A^2}{M_A^2} = \gamma^2, \quad \|x_p - \omega\| \leq \|x_{p-1} - \omega\| \cdot \left(\sqrt{1 - \gamma^2}\right)^k$$

Nota au § précédent on aurait pu aussi bien remarquer que  $\mu = M'_A$ ,

donc tomber sur une formule très analogue:

$$\|x_p - \omega\| \leq \|x_{p-1} - \omega\| \cdot \left(\sqrt{1 - \gamma'^2}\right)^k$$

ou  $\gamma'$  est le conditionnement défini par les deux normes  $\mathcal{M}_0$  et  $\|\cdot\|$ .

Or il est facile de voir que, puisque l'on sait que les conditionnements sont à peu près équivalents, en gros les deux procédés sont de même convergence.

Les formules de calcul sont les suivantes :

$$x_p = x_{p-1} - \Phi(r(x_{p-1})) \cdot \frac{A^T E_{p-1}}{\|A^T E_{p-1}\|^2}$$

$$E_{p-1} = m(x_{p-1})$$

$$r(x_p) = r(x_{p-1}) - \Phi(r(x_{p-1})) \cdot \frac{AA^T E_{p-1}}{\|A^T E_{p-1}\|^2}$$

Cette méthode a été effectivement utilisée sur la machine électronique GAMMA A.E.T de l'université de Grenoble. Ont été étudiés des systèmes de 4, 6, 8 et 12 équations. La convergence est très faible dans le cas d'un système ayant comme matrice de premier membre la matrice  $H_0$  de HILBERT. Cela n'a rien de surprenant d'après ce que nous avons dit au chapitre IV, sur les conditionnements de telles matrices. Par contre, sur un système d'ordre 8, où la méthode d'itération classique de GAUSS-SEIDEL ne donne pas de résultats, et citée par FREEMANN [30], à partir du vecteur initial nul, la méthode conduit au bout de 50 tours d'itération à 3 chiffres significatifs exacts. Dans le cas d'un système d'ordre 12, cinq chiffres ont été obtenus après 250 tours.

7°) Le procédé associé à la norme  $\varphi_1(x) = \|x\|$ .

Nous pouvons bien entendu, écrire  $\|x\| = \frac{x^T}{\|x\|} \cdot x$   
avec :  $z(x) = \frac{x}{\|x\|}$

Le procédé consiste à projeter orthogonalement  $x_{p-1}$  sur le plan :

$$z^T(x_{p-1}) \cdot (Ax - y) = (Ax_{p-1} - y)^T \cdot (Ax - y) = 0$$

Nous obtenons ainsi la formule :

$$x_p = x_{p-1} - \frac{\|z(x_{p-1})\|^2}{\|A^T z(x_{p-1})\|^2} \cdot A^T z(x_{p-1})$$

et par suite,

$$L(x_{p-1}) = \frac{\|z(x_{p-1})\|^2}{\|x_{p-1} - \omega\|^2} \cdot \frac{\|z(x_{p-1})\|^2}{\|A^T z(x_{p-1})\|^2}$$

Donc,

$$L(x_{p-1}) \geq \frac{\lambda_1}{\lambda_n}$$

Si  $\lambda_1$  et  $\lambda_n$  sont les plus petite et plus grande valeur propre de  $A \cdot A^T$  ( ou de  $A^T \cdot A$  ).

Par suite la formule donnant une évaluation de la convergence est dans ce cas encore :

$$\|x_p - \omega\| \leq \|x_0 - \omega\| \cdot \left( \sqrt{1 - \gamma^2} \right)^p$$

On peut donc affirmer:

Théorème II Dans les trois procédés associés aux trois normes  $\varphi_1, \varphi_2, \varphi_\infty$ , on peut dire que :

$$\|x_p - \omega\| \leq \|x_0 - \omega\| \cdot \left( \sqrt{1 - \gamma^2} \right)^p$$

si  $\gamma$  est le conditionnement général associé d'une des trois normes et de la norme euclidienne.

Dans ce cas les formules de calcul sont les suivantes:

$$x_p = x_{p-1} - \frac{\|z(x_{p-1})\|^2}{\|A^T z(x_{p-1})\|^2} \cdot A^T z(x_{p-1})$$

$$z(x_p) = z(x_{p-1}) - \frac{\|z(x_{p-1})\|^2}{\|A^T z(x_{p-1})\|^2} \cdot A \cdot A^T z(x_{p-1})$$

Nous allons étudier maintenant les relations entre les procédés que nous venons d'exposer et les méthodes de gradient.

### 3°) Méthode du gradient .

Avec les notations du § I , soit  $x_0$  un point de  $R^n$  , par ce point il passe une seule des frontières  $\Pi_\omega^\circ(A)$  de la région  $B_\omega^\circ(A)$  , cette surface admet pour équation ;

$$\varphi(x) = \varphi(x_0)$$

Supposons que  $x_0$  soit tel qu'en ce point il existe un plan tangent à  $\Pi_\omega^\circ(A)$  cela revient à dire que  $\varphi(x)$  soit différentiable en  $x_0$ .

$$\text{Or } \frac{\partial \varphi(x)}{\partial x_i} = \sum_{j=1}^{j=n} \frac{\partial \varphi}{\partial X_j} \frac{\partial X_j}{\partial x_i} = (A^T)_i \cdot p(x_0) ; (X_j = f_j(x))$$

Donc le plan tangent est le plan perpendiculaire au vecteur :

$$r_1(x_0) = A^T \cdot p(x_0) = A^T \cdot p(x_0)$$

Le plan  $z^T(Ax - y) = 0$  passant par  $\omega$  sera parallèle au plan tangent si :

$$r_1(x_0) = A^T \cdot z, \quad z = p(x_0)$$

Supposons qu'il en soit ainsi et cherchons l'intersection de ce plan et de la normale en  $x_0$  à  $\Pi_\omega^\circ(A)$ .

Ce point d'intersection  $x_1$  est défini par :

$$x_1 = x_0 + \lambda \cdot r_1$$

(  $\lambda$  scalaire ) et par le fait d'être dans ce plan ce qui permet d'obtenir :

$$z^T(x_0 + \lambda A r_1) = 0$$

d'où

$$\lambda = - \frac{z^T \cdot x_0}{r_1^T \cdot A^T z}$$

et par suite :

$$x_1 = x_0 - \frac{z^T \cdot x_0}{r_1^T A^T z} \cdot r_1$$

Mais si l'on remarque :

$$r_1 = A^T z = A^T \cdot p(x_0)$$

cela peut s'écrire :

$$x_1 = x_0 - \frac{z^T \cdot x_0}{\|r_1\|^2} \cdot r_1 = x_0 - \frac{p^T(x_0) \cdot x_0}{\|p_1\|^2} \cdot p_1$$

Prenons nos différentes normes :

1°) Pour  $\varphi_0(x) = \mathcal{M}_0(x)$  , on a :  $r_1(x) = p(x)$

$$x_1 = x_0 - \frac{\mathcal{M}_0(x_0) \cdot A^T p(x_0)}{\|A^T p(x_0)\|^2}$$

ce qui redonne bien la formule du §5

2°) Pour  $\varphi_1(x) = \Phi(x)$  , on a :  $r_1(x) = m(x)$

et la formule se réduit bien à :  $x_p = x_{p-1} - \frac{\Phi(x_{p-1})}{\|A^T m(x_{p-1})\|^2} \cdot A^T m(x_{p-1})$

qui est celle du § 6  
3°) pour  $\varphi(x) = \|x\|$ , on a  $r(x) = \frac{x}{\|x\|}$

et par suite encore :

$$x_p = x_{p-1} - \frac{\|r(x_{p-1})\|^2}{\|A^T r(x_{p-1})\|^2} \cdot A^T r(x_{p-1})$$

On en conclut que pour les trois cas qui viennent d'être cités la méthode de la décomposition de la norme coïncide avec celle du gradient .

Toutes les fois qu'une norme est dérivable dans une certaine région contenant  $x_0$ , on peut écrire pour tout  $\lambda > 0$ ,

$$\varphi(\lambda x_0) = \lambda \varphi(x_0)$$

la fonction  $F(\lambda) = \varphi(\lambda x_0)$  est dérivable pour  $\lambda = 1$ , donc :

$$\varphi(x_0) = \left( \frac{\partial \varphi}{\partial \xi_1} \right)_0 \cdot \xi_1^0 + \dots + \left( \frac{\partial \varphi}{\partial \xi_n} \right)_0 \cdot \xi_n^0 = \text{grad } \varphi \cdot x$$

cela prouve comment on peut localement découvrir une décomposition d'une norme .

### 9°) Généralisation de la méthode

Nous allons voir comment on peut généraliser la méthode que l'on vient d'étudier au cas où l'on ne fait pas forcément des projections orthogonales de  $x_{p-1}$  sur les plans passant par la solution et qui sont associés à ces points.

Je me restreins au cas de l'utilisation de la norme  $\Phi(x)$ .  
Considérons les  $2^n$  régions déterminées dans  $R^n$  par les plans

$$A_j \cdot x - y_j \equiv f_j(x) = 0$$

du système donné. Soit  $R_j$ , ( $j=1, 2, \dots, 2^n$ ), une telle région, c'est un lieu de points de  $R^n$  pour les quels les signes des  $f_j(x)$  restent les mêmes .

Soient  $\alpha_j$ , ( $j=1, 2, \dots, 2^n$ ) un ensemble de  $2^n$  directions de  $R^n$  définies par des vecteurs  $\alpha_j$ , unitaires,  $\|\alpha_j\|=1$  et associés bi-univoquement aux régions  $R_j$  : à  $\alpha_j$  correspondra  $R_j$  de même indice .

D'après ce que l'on a vu plus haut, si  $x_0$  est un point de la région  $R_j$ , le plan :

$$\pi_j : m^T(x_0) \cdot (Ax - y) = 0$$

est le plan passant par  $\omega$  et parallèle à la "face" du polyèdre

$\prod_{\omega}^{\circ}(A)$ , soit  $\pi_j'$ , qui passe par  $x_0$ .

Au point  $x_0$ , faisons correspondre le point  $x_1$ , où la parallèle à  $\alpha_j$  menée par  $x_0$  coupe  $\Pi_j$ . (On choisit  $\alpha_j$  non parallèle à  $\Pi_j$ )

on a :  $x_1 = x_0 + t\alpha_j$ .

et  $t$  est tel que :  $m^T(x_1) \cdot (A(t\alpha_j + x_0) - y) = 0$

ou bien :  $t \cdot m^T \cdot A \alpha_j + m^T(x_0) \cdot (Ax_0 - y) = 0$

et si l'on remarque que :  $m^T(x_0) \cdot (Ax_0 - y) = \bar{\Phi}_1(x_0) = \bar{\Phi}(x_0)$

cela donne :  $t = -\bar{\Phi}_1(x_0) / m^T A \alpha_j$  :

$$x_1 = x_0 - \frac{\bar{\Phi}_1(x_0)}{m^T A \alpha_j} \cdot \alpha_j$$

formule très simple .

La question importante est de savoir si le procédé est convergent et comment il converge dans ce cas .

Il est clair que :

$$\| \omega - x_1 \|^2 = \| x_1 - x_0 \|^2 + \| x_0 - \omega \|^2 - 2(x_1 - x_0)^T \cdot (x_0 - \omega)$$

donc :

$$\| \omega - x_1 \|^2 = \| \omega - x_0 \|^2 - \| x_1 - x_0 \|^2 + 2(x_1 - x_0)^T \cdot (x_0 - \omega)$$

et ensuite successivement :

$$\begin{aligned} \| \omega - x_1 \|^2 &= \| \omega - x_0 \|^2 - \frac{\bar{\Phi}_1^2(x_0)}{(m^T A \alpha_j)^2} - 2 \frac{\bar{\Phi}_1(x_0)}{(m^T A \alpha_j)} \alpha_j^T \cdot (x_0 - \omega) \\ &= \| \omega - x_0 \|^2 - \frac{\bar{\Phi}_1^2(x_0)}{(m^T A \alpha_j)^2} - 2 \frac{\bar{\Phi}_1(x_0)}{(m^T A \alpha_j)} \cdot \alpha_j^T \left[ x_0 - \omega - \frac{\bar{\Phi}_1(x_0)}{m^T A \alpha_j} \alpha_j \right] \\ &= \| \omega - x_0 \|^2 \left[ 1 - \frac{\bar{\Phi}_1(x_0)}{(m^T A \alpha_j)^2 \| \omega - x_0 \|^2} \left[ 2 m^T A \alpha_j \alpha_j^T (x_0 - \omega) - \bar{\Phi}_1(x_0) \right] \right] \\ &= \| \omega - x_0 \|^2 \left[ 1 - L(x_0) \right] \end{aligned}$$

avec

$$L(x_0) = \frac{\bar{\Phi}_1(x_0)}{(m^T A \alpha_j)^2 \| \omega - x_0 \|^2} \left[ 2 (m^T A \alpha_j) \cdot \alpha_j^T (x_0 - \omega) - \bar{\Phi}_1(x_0) \right]$$

et puisque :  $\bar{\Phi}_1(x_0) = m^T A (x_0 - \omega)$

$$L(x_0) = \frac{\bar{\Phi}_1(x_0)}{(m^T A \alpha_j)^2 \| x_0 - \omega \|^2} \left[ 2 (m^T A \alpha_j) \cdot \alpha_j - A^T m \right]^T \cdot (x_0 - \omega)$$

Remarquons que si l'on fait  $\alpha_j = \frac{A^T m}{\|A^T m\|}$  dans cette formule on trouve bien  $L(x_j) = \frac{\Phi_j(x_0)}{\|A^T m\|^2 \|x_0 - \omega\|^2}$  qui correspond au procédé déjà vu.

Posons

$$k_j = 2 \cdot (m^T A \alpha_j) \cdot \alpha_j - A^T m$$

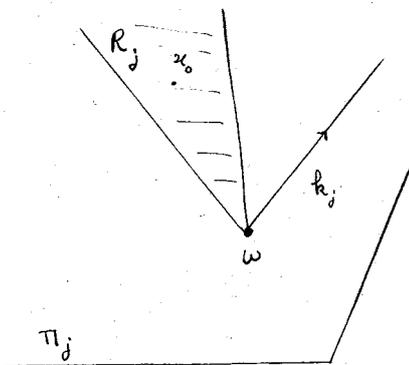
Alors

$$L(x_0) = \frac{\Phi_j(x_0)}{(m^T A \alpha_j)^2 \|x_0 - \omega\|^2} k_j^T \cdot (x_0 - \omega)$$

La question, pour nous permettre d'affirmer la convergence, peut se poser ainsi :

Peut-on choisir les  $\alpha_j$  de sorte que  $L(x)$  soit toujours  $> 0$  et le plus grand possible ? ( il est toujours  $\leq 1$  )

1°) Il est clair que l'on peut toujours faire un choix des  $k_j$  de sorte que  $k_j^T (x_0 - \omega) > 0$ . En effet si l'on remarque que  $R_j$  est un demi-cône convexe tout entier dans l'un des deux demi-espaces déterminés par le plan  $\pi_j$ . Il suffit de choisir un  $k_j$  ( convenablement orienté ) de sorte que en balayant  $R_j$   $x_0 - \omega$ , ne devienne jamais perpendiculaire à celui-ci. Il suffit de choisir dans la portion du cône supplémentaire de  $R_j$  qui contient le vecteur  $\frac{A^T m}{\|A^T m\|}$  que nous savons convenir.



2°) les  $k_j$  étant choisis selon la règle que nous venons de voir il nous faut déterminer les  $\alpha_j$  cela se fait en résolvant les équations:

$$2 (m^T A \alpha_j) \cdot \alpha_j = k_j + A^T m$$

on tire  $2 (m^T A \alpha_j)^2 = m^T A k_j + \|A^T m\|^2$

D'où  $S_j = m^T A \alpha_j = \sqrt{\frac{\|A^T m\|^2 + m^T A k_j}{2}}$

ayant  $S_j$  on en déduit:

$$\alpha_j = \frac{1}{2 S_j} [k_j + A^T m]$$

et pour satisfaire à la condition  $\|\alpha_j\| = 1$

il suffit de prendre  $k_j$  de sorte que :

$$1 = \|\alpha_j\|^2 = \frac{\|k_j\|^2 + \|A^T m\|^2 + 2 k_j^T A^T m}{2 \|A^T m\|^2 + 2 k_j^T A^T m}$$

$$\|k_j\| = \|A^T m\|$$

Cela fait :

$$L(x_0) = \frac{\Phi_1(x_0)}{\|x_0 - \omega\|} \cdot \frac{k_j^T(x_0 - \omega)}{\|x_0 - \omega\|} \cdot \frac{2}{m^T A k_j + \|A^T m\|^2}$$

et puisque  $k_j^T \cdot (x_0 - \omega) \geq \|k_j\| \cdot \|x_0 - \omega\| \cdot \cos p_j$

$p_j$  étant le maximum de l'angle aigu que fait  $k_j$  avec  $x_0 - \omega$  si  $x_0$  parcourt la région  $R_j$ .

donc 
$$L(x) \geq \frac{\Phi_1(x_0)}{\|x_0 - \omega\|} \|k_j^T\| \cos p_j \cdot \frac{2}{m^T A k_j + \|A^T m\|^2} \geq \frac{\Phi_1(x_0)}{\|x_0 - \omega\|} \cos p_j$$

soit  $m = \min_j (\cos p_j)$

il en résulte :

$$L(x_0) \geq$$

et puisque 
$$\frac{\Phi_1(x)}{\|x - \omega\|} \geq m_A$$

on a :

$$L(x) \geq m_A \cdot m$$

enfin :

$$\|x_1 - \omega\| \leq \|x_0 - \omega\| \cdot \sqrt{1 - m_A \cdot m}$$

Si l'on a déterminé les  $\alpha_j$  par les conditions précédentes on peut donc énoncer :

Théorème III

Le procédé de résolution relatif à la norme  $\Phi$  converge encore si les projections orthogonales sont remplacées par des projections parallèles à des directions  $\alpha_j$  convenablement choisies.

Ce qui fait l'intérêt de ce théorème est le corollaire suivant:

Corollaire

En calcul numérique le procédé relatif à la décomposition de la norme  $\Phi$  est pratiquement insensible aux erreurs de détermination des vecteurs  $A^T m$ .

En effet si le calcul fournit, au lieu de  $A^T m$ , un vecteur  $\overline{A^T m}$  on peut dire, vu que les erreurs sont tout de même assez faibles, que l'on a remplacé les directions par des directions  $\alpha_j$  voisines, Donc que le procédé sera celui-ci-dessus, avec des  $\alpha_j = \overline{A^T m}$  convenables; Seule la rapidité de la convergence peut en être affectée.

CHAPITRE VIII

Calcul du polynôme caractéristique,  
des valeurs et vecteurs propres.

1°) Réduction à une forme de FROBENIUS.

Nous utiliserons dans ce qui va suivre des matrices antisymétriques  $K$  :

$$K = X \cdot Y^T$$

avec :

$$X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}, \quad Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

nous poserons :  $\lambda = Y^T X = X^T Y = x_1 y_1 + x_2 y_2 + \dots + x_n y_n$

le polynôme minimal pour  $K$  est alors :  $m(u) = u^2 - \lambda u$

Si l'on se donne donc  $M = aI + bK$  son inverse :  $M^{-1} = a'I + b'K$  avec  $a'$  et  $b'$ , déterminés par les formules du Chap. I, la transmuée d'une matrice  $A = (a_{ij})$  ( $i, j = 1, 2, \dots$ )

par  $M$  sera si  $p = -\frac{a}{b} \neq 0$  et  $p \neq \lambda$  :

$$A' = A + \frac{\lambda - p}{p^2 - \lambda p} K \cdot A + \frac{p}{p^2 - \lambda p} A \cdot K - \frac{1}{p^2 - \lambda p} K \cdot A \cdot K$$

ou :

$$(II) \quad A' = A - \frac{1}{p} K \cdot A + \frac{1}{p - \lambda} A \cdot K - \frac{1}{p(p - \lambda)} K \cdot A \cdot K$$

Dans le cas que nous traitons :

$$AK = (AX) \cdot Y^T, \quad KA = X \cdot (Y^T A)$$

$$KAK = XY^T AX Y^T = (Y^T AX) \cdot XY^T = (Y^T AX) \cdot K$$

Nous écrirons cette formule:

a) En explicitant le terme général  $a'_{ij}$  de  $A'$

$$a'_{ij} = a_{ij} - \frac{1}{p} x_i (a_{1j} y_1 + a_{2j} y_2 + \dots + a_{nj} y_n) + \frac{1}{p-\lambda} y_j (a_{i1} x_1 + \dots + a_{in} x_n) - \frac{1}{p(p-\lambda)} (Y^T A X) x_i y_j$$

$$a'_{ij} = a_{ij} - \frac{1}{p} x_i (Y^T A_j) + \frac{1}{p-\lambda} y_j (A_i \cdot X) - \frac{1}{p(p-\lambda)} (Y^T A X) x_i y_j$$

b) En explicitant une colonne (de rang  $j$ ) de  $A'$

$$A'_{.j} = A_{.j} - \frac{1}{p} (Y^T A_{.j}) \cdot X + \frac{1}{p-\lambda} y_j \cdot (A X) - \frac{1}{p(p-\lambda)} (Y^T A X) \cdot y_j X$$

D'autre part, on sait ([31],[32]) que si l'équation caractéristique

d'une matrice admet des racines simples,

il est possible, par des opérations rationnelles dans le corps de ses coefficients, de trouver une matrice  $\mathcal{F}$  semblable à celle-ci et de la forme " de FROBENIUS " :

$$\mathcal{F} = \begin{pmatrix} 0 & 0 & \dots & \dots & p_1 \\ 1 & 0 & & & p_2 \\ 0 & 1 & 0 & & p_3 \\ \vdots & & \ddots & \ddots & \vdots \\ 0 & & & 1 & 0 \\ \vdots & & & & \ddots \\ 0 & \dots & \dots & \dots & 1 & p_n \end{pmatrix}$$

c'est à dire qu'il existe une matrice  $\mathcal{M}$  dont les éléments s'obtiennent par des opérations rationnelles à partir de ceux de  $A$ , non singulière et telle que :

$$\mathcal{F} = \mathcal{M} \cdot A \cdot \mathcal{M}^{-1}$$

Nous allons prouver qu'il est possible de trouver une matrice telle que  $\mathcal{M}$  par un produit de matrices du 2° degré très simples et qui permettent de programmer effectivement un tel calcul .

1) Première transmutation

Faisons  $y_i = 0$  et laissons les  $x_i$  et  $y_2, y_3, \dots, y_n$  indéterminés pour l'instant .

On aura :

$$A'_{.1} = A_{.1} - \frac{1}{p} (Y^T A_{.1}) \cdot X$$

Remarquons d'abord qu'il est impossible de choisir les  $x_i$  et  $y_i$  qui restent de sorte que les composantes de  $A'_{.1}$  soient toutes nulles, car dans ce cas on aurait :

$$X = \frac{\rho}{(\gamma^T A_{.1})} \cdot A_{.1}$$

d'où :

$$\lambda = \gamma^T X = \frac{\rho}{(\gamma^T A_{.1})} \cdot (\gamma^T A_{.1}) = \rho$$

ce qui est incompatible avec les utilisations des formules .

Par contre, et dans le but de commencer la mise en forme de FROBENIUS, peut-on faire en sorte que :  $A'_{.1} = e_2$  ?

(  $e_i$  désigne toujours le  $i^{\text{e}}$  vecteur de la base fondamentale de  $K^n$  ) .

Pour réaliser cela il suffit :

1°) de prendre  $\rho = \gamma^T A_{.1}$  .

2°) de prendre  $x_1 = a_{11}$ ,  $x_3 = a_{31}$ ,  $x_4 = a_{41}$ , ...,  $x_n = a_{n1}$

alors

$$X = A_{.1} - a_{21} e_2 + x_2 e_2$$

puis

$$A'_{.1} = A_{.1} - X = (a_{21} - x_2) \cdot e_2 = e_2$$

si l'on fait  $x_2 = a_{21} - 1$  ;  $X = A_{.1} - e_2$

Voyons si ce choix donne  $\rho \neq 0$  et  $\rho \neq \lambda$  :

Il est clair que <sup>l'arbitraire</sup> laissé sur les  $y_2, y_3, \dots, y_n$  permet (si  $A_{.1} \neq 0$  ce qui est supposé,  $A$  n'étant pas à première colonne nulle) d'obtenir  $\rho = \gamma^T A_{.1} \neq 0$

Ensuite,

$$\lambda = \gamma^T X = \gamma^T (A_{.1} - e_2)$$

Donc  $\rho = \lambda + y_2$

il faut prendre  $y_2 \neq 0$ ,  $y_2 = 1$  par exemple .

Cela nous conduit à prendre :

$$\gamma = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = e_2$$

d'où  $\rho = a_{21}$  (si  $a_{21} \neq 0$ ) et  $\lambda = a_{21} - 1$

Alors  $A'_{.1} = e_2$  .

Pour les autres colonnes de  $A'$  on aura :

$$A'_{.2} = A_{.2} - \frac{a_{22}}{a_{21}} (A_{.1} - e_2) + A (A_{.1} - e_2) - \frac{1}{a_{21}} (e_2^T A (A_{.1} - e_2)) \cdot (A_{.1} - e_2)$$

$$A'_{.j} = A_{.j} - \frac{a_{2j}}{a_{21}} (A_{.1} - e_2) \quad (j = 3, \dots, n)$$

la formule qui donne  $A'_{.2}$  se simplifie en remarquant que :  $A \cdot e_2 = A_{.2}$

et que  $e_2^T \cdot A \cdot (A_{.1} - e_2) = e_2^T \cdot A \cdot A_{.1} - e_2^T \cdot A_{.2} = (A \cdot A_{.1})_2 - a_{22}$

en désignant par  $(A \cdot A_{.1})_i$  la  $i^{\text{e}}$  composante de  $(A \cdot A_{.1})$

Alors ,  $A'_{.2} = A_{.2} - \frac{1}{a_{21}} [a_{22} + (A \cdot A_{.1})_2 - a_{22}] (A_{.1} - e_2) + A \cdot A_{.1} - A_{.2}$

cela nous fournit les formules très simples :

$$A' : \begin{cases} A'_{.1} = e_2 \\ A'_{.2} = A \cdot A_{.1} - \frac{(A \cdot A_{.1})_2 \cdot (A_{.1} - e_2)}{a_{21}} \\ A'_{.j} = A_{.j} - \frac{(A_{.j})_2 \cdot (A_{.1} - e_2)}{a_{21}} \end{cases}$$

2) Transmutation générale

Supposons être parvenu à transmuier  $A$  en  $A^{(k)}$ , de la forme ci-dessous :

$$A^{(k)} = \begin{pmatrix} 0 & 0 & \dots & 0 & x & x & \dots & x \\ 1 & 0 & \dots & 0 & x & \dots & \dots & x \\ 0 & 1 & & & & & & \\ \vdots & & \ddots & & & & & \\ 0 & & & 0 & x & \dots & \dots & x \\ \vdots & & & & & & & \\ 0 & & & 1 & x & \dots & \dots & x \\ \vdots & & & & & & & \\ 0 & & & 0 & x & & & x \end{pmatrix}, \quad k \leq n-2$$

C'est à dire que les  $k$  premières colonnes de  $A^{(k)}$  sont à termes nuls sauf les termes sur la première parallèle au dessous de la diagonale qui sont des 1.

Nous allons montrer qu'il existe un choix des  $p, x, y$  tels que la transmutation qu'ils définissent assure à la transmuée de  $A^{(k)}$  la forme  $A^{(k+1)}$ .

La formule devant donner les colonnes de  $A^{(k+1)}$  est encore :

$$A_{.j}^{(k+1)} = A_{.j}^{(k)} - \frac{1}{p} \cdot (Y^T A_{.j}^{(k)}) \cdot X + \frac{1}{p-x} y_j A^{(k)} X - \frac{1}{p(p-x)} (Y^T A^{(k)} X) \cdot y_j X$$

Faisons ici  $y_1 = y_2 = \dots = y_{k+1} = 0$ . Il est évident que  $Y^T \cdot A_{.j}^{(k)} = 0$  pour  $j=1, 2, \dots, k$ , puisque les  $A_{.j}^{(k)}$  ( $j=1, 2, \dots, k$ ) n'ont que des 0 dans les lignes  $k+2, k+3, \dots, n$ .

Donc :  $A_{.j}^{(k+1)} = A_{.j}^{(k)}$  ( $j=1, 2, \dots, k$ )

Les  $k$  premières colonnes de  $A^{(k+1)}$  sont celles de  $A^{(k)}$ .

Prenons ensuite :

$$Y = \begin{pmatrix} 0 \\ \vdots \\ 1 \leftarrow k+2 \\ \vdots \\ 0 \end{pmatrix} = e_{k+2}$$

Alors :  $A_{\cdot, k+1}^{(k+1)} = A_{\cdot, k+1}^{(k)} - \frac{1}{p} (e_{k+2}^T \cdot A_{\cdot, k+1}^{(k)}) \cdot X$

et pour faire en sorte que  $A_{\cdot, k+1}^{(k+1)} = e_{k+2}$

il suffit de prendre :  $p = e_{k+2}^T \cdot A_{\cdot, k+1}^{(k)} = a_{k+2, k+1}^{(k)} \quad (\neq 0)$

et  $X = A_{\cdot, k+1}^{(k)} - e_{k+2}$

alors

$$\lambda = Y^T X = a_{k+2, k+1}^{(k)} - 1 = p - 1 \neq p$$

si bien que les formules donnant  $A_{\cdot, j}^{(k+1)}$ , ( $j = k+1, \dots, n$ ) sont :

$$A_{\cdot, k+1}^{(k+1)} = e_{k+2}$$

$$A_{\cdot, k+2}^{(k+1)} = A_{\cdot, k+2}^{(k)} - \frac{1}{a_{k+2, k+1}^{(k)}} (e_{k+2}^T \cdot A_{\cdot, k+2}^{(k)}) (A_{\cdot, k+1}^{(k)} - e_{k+2}) + A_{\cdot, k+1}^{(k)} (A_{\cdot, k+1}^{(k)} - e_{k+2}) - \frac{1}{a_{k+2, k+1}^{(k)}} \left[ e_{k+2}^T \cdot A_{\cdot, k+1}^{(k)} (A_{\cdot, k+1}^{(k)} - e_{k+2}) \right] \cdot (A_{\cdot, k+1}^{(k)} - e_{k+2})$$

$$A_{\cdot, j}^{(k+1)} = A_{\cdot, j}^{(k)} - \frac{a_{k+2, j}^{(k)}}{a_{k+2, k+1}^{(k)}} (A_{\cdot, k+1}^{(k)} - e_{k+2}) \quad ; \quad (j = k+3, \dots, n)$$

si l'on simplifie la 2<sup>o</sup> formule comme on l'a déjà fait,

$$A_{\cdot, k+2}^{(k+1)} = A_{\cdot, k+2}^{(k)} - \frac{1}{a_{k+2, k+1}^{(k)}} \left[ a_{k+2, k+2}^{(k)} + (A_{\cdot, k+1}^{(k)} \cdot A_{\cdot, k+1}^{(k)})_{k+2} - a_{k+2, k+2}^{(k)} \right] (A_{\cdot, k+1}^{(k)} - e_{k+2}) + A_{\cdot, k+1}^{(k)} \cdot A_{\cdot, k+1}^{(k)} - A_{\cdot, k+1}^{(k)}$$

Il en résulte les formules de calcul :

$$A^{(0)} = A$$

$$A^{(k)} \rightarrow A^{(k+1)} \quad \left\{ \begin{array}{l} A_{\cdot, j}^{(k+1)} = e_{j+1} \quad (j = 1, 2, \dots, k+1) \\ A_{\cdot, k+2}^{(k+1)} = A_{\cdot, k+2}^{(k)} - \frac{(A_{\cdot, k+1}^{(k)} \cdot A_{\cdot, k+1}^{(k)})_{k+2}}{a_{k+2, k+1}^{(k)}} [A_{\cdot, k+1}^{(k)} - e_{k+2}] \\ A_{\cdot, j}^{(k+1)} = A_{\cdot, j}^{(k)} - \frac{(A_{\cdot, j}^{(k)})_{k+2}}{a_{k+2, k+1}^{(k)}} \cdot [A_{\cdot, k+1}^{(k)} - e_{k+2}] \quad (j = k+3, \dots, n) \end{array} \right.$$

Et l'on voit que  $A^{(k-1)}$  aura bien la forme de FROBENIUS  $(e_2, e_3, \dots, e_n, \pi) = \mathcal{F}$

avec  $\pi = A_{\cdot, n}^{(k-1)}$ .

Donc le polynôme caractéristique de  $A$  n'est autre que :

$$F(\lambda) = (-1)^k \left[ \lambda^n - \lambda^{n-1} (a_{n, n}^{(k-1)}) - \lambda^{n-2} (a_{n-1, n}^{(k-1)}) \dots - (a_{1, n}^{(k-1)}) \right]$$



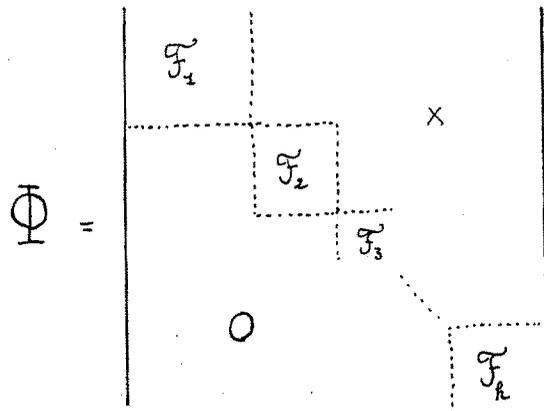
Puis  $A^{(k)} = V_{k+2, l} \cdot A_1^{(k)}$  : matrice identique à  $A_1^{(k)}$  sauf échange de la  $k+2^{\text{e}}$  ligne et de la  $l^{\text{e}}$  ligne .

Donc, dans  $A^{(k)}$  le terme  $a_{k+2, k+1}^{(k)} = a_{l, k+1} \neq 0$ , elle est bien de la même forme que  $A^{(k)}$ , enfin elle est semblable à  $A^{(k)}$ ,  $A^{(k)} = V_{k+2, l} \cdot A_1^{(k)} \cdot V_{k+2, l}^{-1}$ . Si l'on tombe sur un cas b) il suffit de remplacer  $A^{(k)}$  par  $A^{(k)}$  et de poursuivre l'algorithme .

Ce qui précède prouve le théorème [33]

Théorème II

Etant donnée une matrice carrée d'ordre  $n$ , à éléments dans un corps  $\mathcal{R}$ ,  $A$ , il est possible de la transmuter en une matrice de la forme de FROBENIUS générale, c'est à dire de la forme :



Les  $\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_h$  étant de forme de Frobenius (Simple):

$$\mathcal{F}_i = \begin{pmatrix} 0 & 0 & \dots & \pi_i \\ 1 & 0 & \dots & \pi_i \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} \quad ; \text{ ou bien } : \mathcal{F}_i = 0$$

Ces transmutations sont toutes de caractère rationnel en le corps  $\mathcal{R}$  et l'on peut écrire :

$$\Phi = \mathcal{M} \cdot A \cdot \mathcal{M}^{-1}$$

La matrice  $\mathcal{M}$  étant un produit de matrices du 2° degré de la forme soit  $-\rho I + \chi_i \cdot e_i^T$ , soit  $V_{ij}$ .

De la démonstration qui précède on déduit la possibilité de toujours calculer le polynôme caractéristique de  $A$  par cette méthode .

Il en résulte aussi un autre résultat intéressant :

Théorème III Si dans le corps  $\mathcal{R}$  auquel appartiennent les éléments  $a_{ij}$  d'une matrice  $A$  le polynôme caractéristique de celle-ci est irréductible, Il existe toujours une suite de matrices du 2° degré dont le produit transmue  $A$  en une forme simple de FROBENIUS .

C'est une simple conséquence de la démonstration précédente si l'on remarque que dans l'hypothèse du théorème on ne peut pas tomber dans la

circonstance du cas particulier a)

La méthode que nous venons de décrire est de programmation facile, surtout sur une machine possédant un programme d'assemblage matriciel.

Dans le passage de  $A^{(k)}$  vers  $A^{(k+1)}$  il nous faut :

- 1) Calculer  $A^{(k)} \cdot A^{(k)}$ , ce qui vu la forme de  $A^{(k)}$  exige  $(n-k) \cdot n$  produits
- 2) Le calcul des termes des nouvelles colonnes ( $k+2$  à  $n$ ) demande  $(n-k-1) \cdot n$  produits.
- 3) Il y a enfin à exécuter  $n-k-1$  divisions par le même facteur:  $a_{k+2, k+2}$

Le "coût" total est donc de  $n-k-1$  divisions et de  $n(2(n-k)-1)$  multiplications dans le passage  $A^{(k)} \rightarrow A^{(k+1)}$ .

La transmutation complète doit exiger :

$$\sum_{k=0}^{h-2} (n-k-1) = 1+2+\dots+(n-1) = \frac{n(n-1)}{2} = N_D \quad \text{divisions}$$

$$n \sum_{k=0}^{h-1} [2(n-k)-1] = n \cdot (h^2-1) = N_M \quad \text{multiplications.}$$

Ce nombre est de l'ordre de  $n^3$ .

### 2°) Comparaison avec la méthode de DANILEWSKI.

Le procédé que l'on vient de voir ressemble beaucoup au procédé de DANILEWSKI [34], [35], qui ne coûte que de l'ordre de  $n^3$  multiplications aussi. Il est facile de voir que ces procédés résultent tous les deux de transmutations successives de la matrice donnée par des matrices du 2° degré.

En effet notre procédé consiste à prendre au premier pas :

$$A^{(1)} = \left( I - \frac{1}{a_{21}} K \right) \cdot A \cdot \left( I + K \right)$$

avec

$$K = \begin{pmatrix} 0 & a_{11} & 0 & \dots & 0 \\ 0 & a_{21} & 0 & \dots & 0 \\ 0 & a_{31} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n1} & 0 & \dots & 0 \end{pmatrix}$$

comme il résulte aisément de ce que l'on a établi.

Le procédé de DANILEWSKI, utilise

$$A^{(1)} = \left( I - \frac{1}{a_{21}} K' \right) \cdot A \cdot \left( I + \frac{1}{a_{21}} K' \right)$$

avec  $K'^2 = 0$

et

$$K' = \begin{vmatrix} 0 & a_{11} & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & a_{21} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n1} & 0 & \dots & 0 \end{vmatrix}$$

on peut remarquer que :

$$\left(1 - \frac{1}{a_{21}} K\right) = \begin{vmatrix} 1 & \dots & 0 \\ \frac{1}{a_{21}} & 1 & \dots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & 1 \end{vmatrix} \cdot \left(1 - \frac{1}{a_{21}} K'\right)$$

On voit que la différence des deux méthodes est très peu sensible ; elles ne sont, en fait, que des cas particuliers de la méthode de transmutation générale par des matrices du 2° degré. D'autres choix de matrices K sont possibles, et cela conduira, pour la recherche d'une forme de FROBENIUS équivalente à une matrice donnée, à des algorithmes voisins de ces deux que l'on vient de présenter.

3°) Transmutation en une forme Triple-diagonale .

Nous allons voir qu'il est possible de transmuter une matrice donnée en une matrice de forme triple-diagonale:

$$B_D = \begin{vmatrix} \alpha_1 & \beta'_1 & & & \\ \beta_1 & \alpha_2 & \beta'_2 & & 0 \\ & \beta_2 & \alpha_3 & \beta'_3 & \\ & & & & \ddots \\ 0 & & & & & \beta'_{n-1} \\ & & & & & \beta_{n-1} & \alpha_n \end{vmatrix}$$

n'ayant que des termes nuls sauf dans la diagonale principale et dans les deux parallèles à celle-ci immédiatement au dessus et au dessous .

Nous utiliserons des matrices

$$M = aI + b.K \quad , \quad M^{-1} = a'I + b'K$$

avec K , de la forme  $K = X.Y^T$ .

Nous allons nous astreindre à ne prendre ici que des X et Y tels :

$$X = \begin{pmatrix} 0 \\ \vdots \\ x_e \\ \vdots \\ x_m \\ \vdots \\ 0 \end{pmatrix} \quad , \quad Y = \begin{pmatrix} 0 \\ \vdots \\ y_e \\ \vdots \\ y_m \\ \vdots \\ 0 \end{pmatrix}$$

où l et m sont deux indices différents  $l < m$  , pour fixer les idées .

Alors  $K = (k_{ij})$  a que 4 termes non nuls :

$$k_{ee} = x_e \cdot y_e \quad , \quad k_{me} = x_m \cdot y_e$$

$$k_{el} = x_e \cdot y_l \quad , \quad k_{ml} = x_m \cdot y_l$$

Nous avons ici :  $\lambda = x_e y_e + x_m y_m$

Soit  $A = (a_{ij})$  une matrice donnée, si on la transmue par  $M$  en  $A' = M.A.M^{-1}$  on a, vu que :

$$a'_{ij} = a_{ij} - \frac{1}{p} x_i (a_{ij} y_e + \dots + a_{mj} y_m) + \frac{1}{p-\lambda} y_j (a_{ie} x_e + \dots + a_{im} x_m) - \frac{(Y^T A X)}{p(p-\lambda)} x_i y_j$$

Il en résulte que  $A'$  est identique à  $A$ , sauf pour les lignes et colonnes de rang  $l$  ou  $m$ .

Si l'on pose

$$R = a_{ee} x_e y_e + a_{em} x_m y_e + a_{me} x_e y_m + a_{mm} x_m y_m$$

la formule ci-dessus donne :

$$a'_{ij} = a_{ij} - \frac{1}{p} x_i [a_{ej} y_e + a_{mj} y_m] + \frac{1}{p-\lambda} y_j [a_{ie} x_e + a_{im} x_m] + \frac{R x_i y_j}{p(p-\lambda)}, (i, j = l, m)$$

Soit  $r$  un indice tel que  $r < l < m$  par exemple,

$$a'_{mr} = a_{mr} - \frac{1}{p} x_m [a_{er} y_e + a_{mr} y_m]$$

et son "symétrique " :

$$a'_{rm} = a_{rm} + \frac{1}{p-\lambda} y_m [a_{re} x_e + a_{rm} x_m]$$

On dispose des paramètres  $p, x_e, x_m, y_e, y_m$ , peut-on déterminer ceux-ci afin que: 1° ) Ces deux termes soient nuls

2° ) Que les conditions de possibilité  $p \neq 0$  et  $p \neq \lambda$  soient assurées ?

$$(\lambda = x_e y_e + x_m y_m).$$

Il suffit donc de choisir  $p, x_e, x_m, y_e, y_m$ , de sorte que :

$$0 = p \cdot a_{mr} - x_m [a_{er} y_e + a_{mr} y_m]$$

$$0 = (p-\lambda) a_{rm} + y_m [a_{re} x_e + a_{rm} x_m]$$

de là on tire :

$$a_{mr} [p - x_m y_m] - x_m y_e a_{er} = 0$$

$$a_{rm} [p - x_e y_e] + x_e y_m a_{re} = 0$$

Par suite :

$$x_m [y_e a_{er} + y_m a_{mr}] = p \cdot a_{mr}$$

$$x_e [y_e a_{rm} - y_m a_{re}] = p \cdot a_{rm}$$

Il est facile de voir, d'après les dernières formules écrites et sous la condition que l'on n'est pas simultanément :

$$a_{ee} = a_{re} = a_{mr} = a_{rm} = 0$$

( Dans ce dernier cas l'on aurait déjà  $a_{mr} = a_{rm} = 0$  et, par suite, il est inutile de se poser la question de trouver une matrice semblable à A où ces termes sont nuls : c'est A elle-même ) que l'on peut toujours trouver  $y_e$  et  $y_m$  et prendre  $p$  de sorte que :

$$y_e a_{er} + y_m a_{mr} \neq 0, \quad y_e a_{rm} - y_m a_{re} \neq 0$$

Avec un tel choix pour  $y_e$  et  $y_m$  et  $p$ ,

$$x_m = p \cdot a_{mr} / (y_e a_{er} + y_m a_{mr})$$

$$x_e = p \cdot a_{rm} / (y_e a_{rm} - y_m a_{re})$$

Dans une programmation générale, on peut faire  $y_e = y_m = p = 1$  et cela fournit, dans le but d'éviter des calculs :

$$x_m = a_{mr} / (a_{er} + a_{mr}) = \xi_m$$

$$x_e = a_{rm} / (a_{rm} - a_{re}) = \xi_e$$

Si  $a_{er} + a_{mr} = 0$  ou  $a_{rm} - a_{re} = 0$ , il est commode de choisir encore  $p=1$  et de prendre pour  $y_e$  et  $y_m$  des nombres tels que  $y_e = \varepsilon$  et  $y_m = \varepsilon'$  avec  $|\varepsilon| = |\varepsilon'| = 1$ , il est clair que l'on arrivera toujours à en déterminer de cette sorte et convenables.

Dans ces cas les formules de calcul sont les suivantes :

$$A \xrightarrow{\textcircled{H}} A' \left\{ \begin{array}{l} a'_{ij} = a_{ij} \quad (i \neq e, m; j \neq e, m) \\ a'_{ij} = a_{ij} - x_i [a_{ej} y_e + a_{mj} y_m] + \frac{1}{1 - x_m y_m - x_e y_e} \left[ y_j (a_{ie} x_e + a_{im} x_m) - R x_i y_j \right] \\ R = a_{ee} x_e y_e + a_{em} x_m y_e + a_{me} x_e y_m + a_{mm} x_m y_m \end{array} \right.$$



De cela il résulte le théorème suivant:

Théorème IV

Etant donnée une matrice A quelconque il existe toujours une suite de matrices du 2° degré dont le produit  $\mathcal{M}$  est tel qu'il transmue A en une forme  $\mathcal{L}_D$  triple diagonale;

$$\mathcal{L}_D = \mathcal{M} \cdot A \cdot \mathcal{M}^{-1}$$

Pour achever la détermination du polynôme caractéristique de A' une fois  $\mathcal{L}_D$  trouvée, celui-ci est

$$F(\lambda) = \text{Dét} \begin{vmatrix} \alpha_1 - \lambda & \beta'_1 & & & \\ \beta_1 & \alpha_2 - \lambda & & & \\ & & \ddots & & \\ & & & \alpha_k - \lambda & \beta'_{k-1} \\ & & & \beta_{k-1} & \alpha_k - \lambda \end{vmatrix}$$

Pour en faire le calcul, on pose

$$f_i(\lambda) = \text{Dét} \begin{vmatrix} \alpha_1 - \lambda & \beta'_1 & & & \\ \beta_1 & \alpha_2 - \lambda & & & \\ & & \ddots & & \\ & & & \alpha_i - \lambda & \beta'_{i-1} \\ & & & \beta_{i-1} & \alpha_i - \lambda \end{vmatrix}$$

Si bien que  $f_n(\lambda) = F(\lambda)$ ,  $f_1(\lambda) = \alpha_1 - \lambda$  et on peut aussi poser  $f_0(\lambda) = 1$ .

Alors en développant par rapport à la dernière ligne on trouve :

$$f_i(\lambda) = (\alpha_i - \lambda) \cdot f_{i-1}(\lambda) - \beta_{i-1} \beta'_{i-1} f_{i-2}(\lambda)$$

Cette méthode est à rapprocher avec celle de LANCZOS ou de GIVENS, qui pour chercher une matrice équivalente à une matrice donnée et de forme triple diagonale utilisent des techniques bien différentes. GIVENS [38] [39] Transmue A par des "rotations", c'est à dire comme dans la méthode de JACOBI  $A' = \Omega^T A \Omega$  est telle que  $\Omega^{(*)}$  soit orthonormale, la méthode de GIVENS est exclusivement, sous sa forme primitive, réservée au cas de matrices symétriques et dans ce cas particulièrement commode puisque la suite des polynômes  $f_i(\lambda)$  ci-dessus est une suite de STURM du polynôme caractéristique.

La méthode de LANCZOS est très différente, elle se propose de construire une base, dans laquelle la transformation, définie dans la base fondamentale de  $R^n$  par A, soit définie par une matrice de forme triple-diagonale. On trouvera un exposé de cette méthode classique en [40] [41]. . .

(\*) Remarquons que les matrices des "rotations"  $\Omega$  de GIVENS sont de la forme:  $I + X \cdot Y^T$ .

4°) Etude des erreurs dans la méthode de réduction à une forme de FROBENIUS .

Nous allons revenir à la méthode exposée au §4 et étudier rapidement, en utilisant un mode de raisonnement analogue à celui fait pour l'étude des erreurs dans la résolution de systèmes linéaires, les erreurs de calcul pouvant arriver dans l'utilisation de ce procédé .

Pour simplifier l'écriture des formules de transmutation de  $A^{(k)}$  en  $A^{(k+1)}$  posons :  $A^{(k)} = A$ ,  $A^{(k+1)} = A'$ ,  $\pi = a_{k+2, k+1}^{(k)}$ .

$$X = A_{k+1}^{(k)} - e_{k+2}$$

on a vu que alors :

$$A' = \left( I - \frac{1}{\pi} X \cdot e_{k+2}^T \right) \cdot A \cdot \left( I + X \cdot e_{k+2}^T \right)$$

D'autre part, pour la détermination de  $A'$  il ne faut faire des calculs que pour les éléments  $a_{i,j}^{(k+1)}$  de  $A^{(k+1)}$  situés dans les colonnes de rang  $k+2, k+3, \dots, n$ .

On peut donc poser, si  $\bar{A}'$  est la matrice réellement obtenue à la place de  $A'$  :

$$\bar{A}' = A' + \delta A' = A' + \begin{pmatrix} 0 & \dots & 0 & \varepsilon_{1, k+2} & \dots & \varepsilon_{1, n} \\ 0 & \dots & 0 & \varepsilon_{2, k+2} & \dots & \varepsilon_{2, n} \\ \vdots & & \vdots & \vdots & & \vdots \\ 0 & \dots & 0 & \varepsilon_{n, k+2} & \dots & \varepsilon_{n, n} \end{pmatrix}$$

les  $\varepsilon_{ij}$  étant les erreurs de calcul faites dans la détermination des termes  $a_{ij}^{(k+1)}$  de  $A^{(k+1)}$ .

Or il est facile de déterminer  $\delta A$  de sorte que :

$$\bar{A}' = \left( I - \frac{1}{\pi} X \cdot e_{k+2}^T \right) \cdot (A + \delta A) \cdot \left( I + X \cdot e_{k+2}^T \right)$$

c'est à dire :

$$\delta A' = \left( I - \frac{1}{\pi} X \cdot e_{k+2}^T \right) \cdot \delta A \cdot \left( I + X \cdot e_{k+2}^T \right)$$

et puisque  $I - \frac{1}{\pi} X \cdot e_{k+2}^T$  est l'inverse de  $I + X \cdot e_{k+2}^T$ ,

$$\delta A = \left( I + X \cdot e_{k+2}^T \right) \delta A' \left( I - \frac{1}{\pi} X \cdot e_{k+2}^T \right)$$

Or si  $\delta A'$  a la forme indiquée ( ses  $k+1$  premières colonnes nulles),

Il est très facile de voir que

$$\delta A = \delta A' + X e_{k+2}^T \delta A' - \frac{1}{\Pi} \delta A' X \cdot e_{k+2}^T - \frac{1}{\Pi} X e_{k+2}^T \delta A' X e_{k+2}^T$$

admet aussi ses  $k+1$  premières colonnes nulles .. Cette remarque entraîne l'importante constatation : Si l'on veut soumettre  $A + \delta A$  à la transmutation lui donnant la forme  $(A + \delta A)^{(k+1)}$ , on fera, puisque cette transmutation ne dépend que de la  $k+1$ ° colonne de la matrice, qui est la même pour  $A$  et pour  $A + \delta A$ , la même transmutation, c'est à dire que :

$\bar{A}'$ , obtenue par une transmutation de  $A$ , avec des erreurs de calcul peut être considérée comme obtenue par transmutation de  $A + \delta A$ , sans erreurs de calcul,  $\delta A$  ayant la forme obtenue plus haut.

Revenons aux notations générales, Partons de la matrice  $A$ , au lieu de la suite de matrices semblables :

$$A^0 = A, A^{(1)}, A^{(2)}, \dots, A^{(k-1)}$$

telles que :

$$A^{(k+1)} = M_k \cdot A^{(k)} \cdot M_k^{-1}$$

on obtient par le calcul effectif la suite

$$\bar{A}^0 = A, \bar{A}^{(1)}, \bar{A}^{(2)}, \dots, \bar{A}^{(k-1)}$$

telles que :

$$\bar{A}^{(k+1)} = M_k \cdot (\bar{A}^{(k)} + \delta A^{(k)}) \cdot M_k^{-1}$$

les  $\delta A^{(k)}$  étant des matrices ayant leur  $k+1$  premières colonnes nulles et se calculant comme on l'a vu plus haut à partir des erreurs de calcul sur la détermination de  $A'$ .

Or : 
$$\bar{A}^{(k)} = M_{k-1} \cdot (\bar{A}^{(k-1)} + \delta A^{(k-1)}) \cdot M_{k-1}^{-1}$$

et :

$$\bar{A}^{(k+1)} = M_k \cdot [M_{k-1} \cdot (\bar{A}^{(k-1)} + \delta A^{(k-1)}) \cdot M_{k-1}^{-1} + \delta A^{(k)}] \cdot M_k^{-1}$$

Considérons le produit :

$$M_k \cdot [M_{k-1} \cdot (\bar{A}^{(k-1)} + \delta A^{(k-1)}) \cdot M_{k-1}^{-1} + \delta A^{(k)}] \cdot M_k^{-1}$$

On peut l'écrire,

$$M_k \cdot M_{k-1} \cdot [\bar{A}^{(k-1)} + \delta A^{(k-1)} + M_{k-1}^{-1} \delta A^{(k)} M_{k-1}] \cdot M_{k-1}^{-1} \cdot M_k^{-1}$$

et ensuite,

$$\bar{A}^{(k+1)} = M_k \dots M_0 \cdot [A^0 + \delta A^0 + M_0^{-1} \delta A^1 M_0 + M_0^{-1} M_1^{-1} \delta A^2 M_1 M_0 + \dots] \cdot M_0^{-1} \cdot M_1^{-1} \dots M_k^{-1}$$

enfin,

$$\bar{A}^{(k+1)} = M_0 \left[ A + \sum_{i=0}^{k-1} M_0^{-1} \dots M_{i+1}^{-1} \delta A^{(i)} M_{i+1} \dots M_0 \right] M_0^{-1}$$

on peut donc énoncer :

Théorème V Si l'on transforme une matrice A donnée, carrée, d'ordre n, en une matrice de forme de FROBENIUS, selon la méthode ci-dessus exposée, on obtient une forme qui peut être considérée comme rigoureusement semblable à la matrice  $A + \delta A$ , avec :

$$\delta A = \sum_{i=0}^{i=k-2} M_0^{-1} \dots M_{i-1}^{-1} \delta A^{(i)} M_{i-1} \dots M_0$$

Il nous faut dire un mot de la façon dont on peut déduire de la connaissance de  $\delta A$  une estimation des erreurs faites sur les valeurs propres et vecteurs propres.

5°) Variations des éléments propres correspondant à la variation  $\delta A$  de A.

Soit A une matrice carrée, d'ordre n, ayant ses n valeurs propres  $\lambda_1, \lambda_2, \dots, \lambda_n$ , distinctes.

Soient  $u_i$  et  $v_i$  les vecteurs propres à droite et à gauche pour A et correspondant à la valeur propre  $\lambda_i$  :

$$(1) \quad A u_i = \lambda_i u_i$$

$$(2) \quad v_i^T A = \lambda_i v_i^T \quad (A^T v_i = \lambda_i v_i)$$

Je les supposerai "normalisés" de sorte que :

$$(3) \quad u_i^T u_i = 1$$

$$(4) \quad v_i^T u_i = 1$$

On désignera par T la matrice d'ordre n ayant pour colonnes les vecteurs colonnes  $u_1, u_2, \dots, u_n$ . les lignes de  $T^{-1}$  sont les  $v_i^T$ .

Soit donc  $A + \delta A$  une matrice légèrement variée de A,  $\delta A$  ayant par exemple, la valeur déterminée au § précédent.

$\lambda_i + \delta \lambda_i$  désignera la valeur propre de  $A + \delta A$  voisine de  $\lambda_i$  de A et  $u_i + \delta u_i$  et  $v_i + \delta v_i$  les vecteurs propres à droite et à gauche correspondant à celle-ci.

On doit avoir :

$$\begin{cases} (1)' & (A + \delta A - (\lambda_i I + \delta \lambda_i I)) \cdot (u_i + \delta u_i) = 0 \\ (2)' & (\sigma_i^T + (\delta \sigma_i)^T) \cdot (A + \delta A - (\lambda_i I + \delta \lambda_i I)) = 0 \end{cases}$$

écrivons :

$$(1)'' \quad (\delta A - \delta \lambda_i I) \cdot u_i + (A - \lambda_i I) \cdot \delta u_i + \underline{\delta A \cdot \delta u_i - \delta \lambda_i \cdot \delta u_i} = 0$$

$$(2)'' \quad \sigma_i^T (\delta A - \delta \lambda_i I) + (\delta \sigma_i)^T (A - \lambda_i I) + \underline{(\delta \sigma_i)^T \delta A - \delta \lambda_i \cdot (\delta \sigma_i)^T} = 0$$

Si l'on néglige les parties soulignées comme étant d'ordre supérieur :

$$(1)''' \quad (\delta A - \delta \lambda_i I) \cdot u_i + (A - \lambda_i I) \delta u_i = 0$$

$$(2)''' \quad \sigma_i^T (\delta A - \delta \lambda_i I) + (\delta \sigma_i)^T (A - \lambda_i I) = 0$$

Multiplions (1)''' par  $\sigma_k^T$  (à gauche)

$$\sigma_k^T \delta A u_i - (\delta \lambda_i) \sigma_k^T u_i + \sigma_k^T A \delta u_i - \lambda_i \sigma_k^T \delta u_i = 0$$

$$\sigma_k^T \delta A u_i - (\delta \lambda_i) \sigma_k^T u_i + (\lambda_k - \lambda_i) \cdot \sigma_k^T \delta u_i = 0$$

de même en multipliant (2)''' par  $u_k$  (à droite)

$$\sigma_i^T \delta A u_k - (\delta \lambda_i) \sigma_i^T u_k + (\lambda_k - \lambda_i) (\delta \sigma_i)^T u_k = 0$$

Or :

$$\sigma_k^T u_i = \delta_{ki}$$

( $\delta_{ki}$  : symbole de KRONECKER), d'où :

$$(I) \begin{cases} k=i & : \sigma_i^T \delta A u_i = \delta \lambda_i \\ k \neq i & : \sigma_k^T \delta u_i = \frac{1}{\lambda_i - \lambda_k} \sigma_k^T \delta A u_i \end{cases}$$

$$(II) \begin{cases} k=i & : \sigma_i^T \delta A u_i = \delta \lambda_i \\ k \neq i & : (\delta \sigma_k)^T u_i = \frac{1}{\lambda_k - \lambda_i} \sigma_k^T \delta A u_i \end{cases}$$

Ces relations obtenues, soit à

déterminer  $\delta\lambda_i, \delta u_i, \delta\sigma_i$ , en supposant connus  $\delta A, u_i, \sigma_i$ .

1°) Les  $\delta\lambda_i$  sont déterminés par :

$$\delta\lambda_i = \sigma_i^T \cdot \delta A \cdot u_i$$

2°) Les  $\delta u_i$ , pour  $i$  fixé, sont tels que :

$$\sigma_k^T \cdot \delta u_i = \frac{1}{\lambda_i - \lambda_k} \cdot \sigma_k^T \delta A u_i$$

cela donne pour  $k=1, 2, \dots, i-1, i+1, \dots, n$  seulement  $n-1$  équations avec  $n$  inconnues : les composantes de  $\delta u_i$ , soient :  $\delta u_{i1}, \delta u_{i2}, \dots, \delta u_{in}$ .

mais si l'on pense qu'en permanence :  $\sigma_k^T \cdot u_i = \delta_{ki}$  c'est à dire :

$$(\sigma_k^T + (\delta\sigma_k)^T) \cdot (u_i + \delta u_i) = \delta_{ki} ;$$

cela fait :

$$(\delta\sigma_k)^T \cdot u_i + \sigma_k^T (\delta u_i) = 0$$

en négligeant toujours les termes du 2° ordre  $(\delta\sigma_k)^T \delta u_i$ .

Donc le système fournissant les  $\delta u_i$  et  $\delta\sigma_i$

( $2 \cdot n^2$  inconnues) s'écrit :

$$(3) \begin{cases} \sigma_k^T \cdot \delta u_i = b_{ki} & , k \neq i \\ (\delta\sigma_k)^T \cdot u_i = -b_{ki} & , k \neq i \end{cases}$$

$$(4) \quad (\delta\sigma_k)^T \cdot u_i + \sigma_k^T \cdot \delta u_i = 0$$

où l'on a posé :

$$b_{ki} = \frac{1}{\lambda_i - \lambda_k} \cdot \sigma_k^T \delta A u_i \quad (i \neq k)$$

La remarque importante à faire est celle-ci :

les relations (4) (pour  $i \neq k$ ) sont toujours satisfaites si (3)

le sont

, seules les relations (4) pour  $i=k$  sont indépendantes de ces relations (3). Notre système à  $2 \cdot n^2$  inconnues est donc :

$$(3) \begin{cases} \sigma_k^T \cdot \delta u_i = b_{ki} & , k \neq i \\ (\delta\sigma_k)^T \cdot u_i = -b_{ki} & , k \neq i \end{cases}$$

$$(4)' \quad (\delta\sigma_k)^T \cdot u_i + \sigma_k^T \cdot \delta u_i = 0$$

Ce qui nous fait  $2n^2 - n$  équations seulement. Il nous manque donc, en ce point de notre étude,  $n$  conditions pour résoudre le problème.

Si A est symétrique, le cas est classique et a été traité par JACOBI [42]

Voyons le cas général :

Ecrivons que les  $\delta u_i$  sont tels que  $u_i^T \cdot u_i = 1$  se maintienne.

Cela donne les  $n$  conditions désirées :

$$(5) \quad u_i^T \cdot (\delta u_i) = 0$$

Donc, le système devant fournir les  $\delta u_i$  et  $\delta v_i$  sera :

$$(S) \quad (3) \quad \begin{cases} \sigma_l^T \cdot (\delta u_i) = \beta_{li} & (i \neq l) \\ u_l^T \cdot (\delta v_i) = -\beta_{li} & (i = 1, \dots, n) \end{cases} ; (4) \quad u_i^T \cdot (\delta v_i) + v_i^T \cdot (\delta u_i) = 0 ; (5) \quad u_i^T \cdot (\delta u_i) = 0$$

Nous allons indiquer comment on peut résoudre ce système .

Commençons par déterminer les  $\delta u_i$ .

Ils satisfont à :

$$\begin{cases} \sigma_l^T \cdot (\delta u_i) = \beta_{li} \\ u_i^T \cdot (\delta u_i) = 0 \end{cases}$$

Or on sait que

$$T^{-1} = \begin{pmatrix} \sigma_1^T \\ \sigma_2^T \\ \vdots \\ \sigma_n^T \end{pmatrix}$$

donc :

$$E_{ii} \cdot T^{-1} = \begin{pmatrix} 0 \\ \vdots \\ \sigma_i^T \\ \vdots \\ 0 \end{pmatrix}$$

ensuite

$$T = (u_1, \dots, u_n) \quad , \quad T \cdot E_{ii} = (0, 0, \dots, u_i, 0 \dots 0) \quad , \quad E_{ii} T^T = \begin{pmatrix} 0 \\ \vdots \\ u_i^T \\ \vdots \\ 0 \end{pmatrix}$$

Il en résulte que le système donnant les composantes de  $\delta u_i$  sera :

$$(T^{-1} - E_{ii} T^{-1} + E_{ii} T^T) \cdot \delta u_i = \beta_i \quad , \quad [I - E_{ii} (I - T^T T)] \cdot T^{-1} \delta u_i = \beta_i$$

Or il est très facile de voir que la matrice  $E_{ii} (I - T^T T)$  est de carré nul :

$$(T^T T)_{ii} = 1$$

donc

$$[I - E_{ii} (I - T^T T)]^{-1} = [I + E_{ii} (I - T^T T)]$$

et par suite

$$\delta u_i = T \cdot \left[ I + E_{ii} (I - T^T T) \right] \cdot \beta_i$$

Enfin le système donnant les  $\delta u_i$  peut s'écrire :

$$T^T \cdot \delta \sigma_i = \beta'_i$$

si

$$\beta'_i = - \begin{pmatrix} \beta_i \\ \vdots \\ \sigma_i^T \delta u_i \leftarrow \beta_i \\ \vdots \end{pmatrix}$$

$$\delta \sigma_i = (T^{-1})^T \cdot \beta'_i$$

Ce qui permet de calculer les éléments cherchés.

Enfin signalons une majoration des  $|\delta \lambda_i|$  qui résulte de

$$|\delta \lambda_i| = | \sigma_i^T \delta A u_i |$$

Si  $\mathcal{E}(T^{-1})$  désigne plus grand des éléments de  $T^{-1}$ , (on sait que c'est une norme de la matrice  $T^{-1}$ ), et  $M_A = \text{Max} \frac{\Phi(Ax)}{\|x\|}$  on peut écrire

$$|\delta \lambda_i| \leq \mathcal{E}(T^{-1}) \cdot M_{\delta A}$$

ce qui fixe une borne supérieure pour la variation des valeurs propres.

6°) Cas particulier d'une matrice de forme de FROBENIUS

Soit la matrice

$$A = \begin{pmatrix} 0 & 0 & \dots & \dots & \mu_n \\ 1 & 0 & & & \mu_{n-1} \\ 0 & 1 & 0 & & \mu_{n-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & 1 & \mu_1 \end{pmatrix}$$

Je suppose que son équation caractéristique :

$$F(\lambda) = (-1)^n \cdot \left[ \lambda^n - \mu_1 \lambda^{n-1} - \mu_2 \lambda^{n-2} - \dots - \mu_n \right] = 0$$

n'a que des racines simples  $\lambda_1, \lambda_2, \dots, \lambda_n$ .

Il est facile de voir que si l'on pose :

$$\begin{aligned} f_0(\lambda) &= \mu_n \\ f_1(\lambda) &= \mu_n + \lambda \mu_{n-1} \\ \vdots \\ f_i(\lambda) &= \mu_n + \lambda \mu_{n-1} + \dots + \mu_{n-i} \cdot \lambda^i \end{aligned}$$

c'est à dire que les  $f_i(\lambda)$  peuvent être calculés (Schéma de HORNER) 131  
 par :  $f_i(\lambda) = f_{i-1}(\lambda) + \lambda^i \cdot r_{n-i}$

le vecteur propre  $u_i$  de cette matrice correspondant à  $\lambda_i$  est de composantes

$$u_i = \begin{pmatrix} \xi_u \cdot f_0(\lambda_i) / \lambda_i \\ \xi_u \cdot f_1(\lambda_i) / \lambda_i^2 \\ \vdots \\ \xi_u \cdot f_{n-1}(\lambda_i) / \lambda_i^n \\ \xi_u \cdot 1 \end{pmatrix}, (\xi_u \neq 0, \text{ qq.})$$

Le vecteur propre  $v_i$ , à gauche est alors :

$$v_i = \begin{pmatrix} \eta_i \\ \eta_i \cdot \lambda_i \\ \vdots \\ \eta_i \cdot \lambda_i^{n-1} \\ \eta_i \cdot \lambda_i^n \end{pmatrix}, (\eta_i \neq 0, \text{ qq.})$$

supposons que  $A$  varie de sorte que

$$\delta A = \begin{pmatrix} \circ & \circ & \dots & \circ \delta p_n \\ \circ & \circ & \dots & \circ \delta p_{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ \circ & \dots & \circ & \delta p_1 \end{pmatrix}$$

ce qui revient à dire que les coefficients de l'équation caractéristique varient des  $\delta p_i$ . Il résulte des formules ci-dessus que :

$$\delta \lambda_i = \xi_u \cdot \eta_i (\delta p_n + \lambda_i \delta p_{n-1} + \dots + \lambda_i^{n-1} \delta p_1)$$

Or si l'on écrit que  $v_i^T \cdot u_i = 1$ , on trouve que  $\xi_u, \eta_i$  doivent être tels :

$$\frac{\xi_u \cdot \eta_i}{\lambda_i} \left[ (n-1) p_n + (n-2) p_{n-1} \lambda_i + \dots + p_1 \lambda_i^{n-2} + \lambda_i^n \right] = 1$$

ce qui donne :

$$\delta \lambda_i = \frac{\lambda_i \left[ \delta p_n + \lambda_i \delta p_{n-1} + \dots + \lambda_i^{n-1} \delta p_1 \right]}{n p_n + (n-1) p_{n-1} \lambda_i + \dots + \lambda_i^{n-1} p_1}$$

ou encore : 
$$\frac{\delta \lambda_i}{\lambda_i} = \frac{\delta p_n + \lambda_i \delta p_{n-1} + \dots + \lambda_i^{n-1} \delta p_1}{n p_n + (n-1) p_{n-1} \lambda_i + \dots + \lambda_i^{n-1} p_1}$$

et si l'on pose :

$$T^n f\left(\frac{x}{T}\right) = \Phi(x, T)$$

on a  $n \Phi(x, T) = x \Phi'_x + T \Phi_T$ , donc cela redonne la relation classique :

$$\lambda_i \cdot f'(\lambda_i) = n p_n + (n-1) p_{n-1} \lambda_i + \dots + \lambda_i^{n-1} p_1$$

$$\delta \lambda_i = \frac{\delta p_n + \lambda_i \delta p_{n-1} + \dots + \lambda_i^{n-1} \delta p_1}{f'(\lambda_i)}$$



BIBLIOGRAPHIE

On trouve une bibliographie très complète sur l'analyse numérique linéaire dans l'ouvrage :

"Simultaneous linear equations and the détermination of eigenvalues " édité par J.PAIGE et O.TAUSSKY . N.B.S. App.Math.Series .29 (1953)

Cet ouvrage comporte aussi une classification des différents ouvrages cités pour les différentes méthodes utilisées . Il ne cite ,malheureusement que des ouvrages publiés avant 1951 .

Pour une bibliographie générale très récente et assez complète on pourra se rapporter à :

" Introduction to matrix analysis " par R.BELLMAN . Mac.GRAW-HILL .(1960).

Les ouvrages cités dans la liste qui suit sont ceux pour les quels un renvoi figure dans le texte .

- [1] SPERNER et SCHEIRER "Matrix calculus and Modern Algebra " Chelsea .Pub. Comp . New-York .( 1953) Chap . .
  - [2] N.GASTINEL "Utilisation de matrices vérifiant une équation du 2° degré pour la transmutation de matrices " C.R.A.S 7 mars 1960 .
  - [3] E.BODEWIG " Matrix calculus " North-Holland Pub.Comp. Amsterdam (1956) P ages 10 et suiv .
  - [4] N.BOURBAKI " Topologie Générale " Chap IX ,§3 , Déf 5 , p 42 .
  - [5] N.BOURBAKI " Espaces Vectoriels Topologiques " Chap I,§3, Corol.I, ThI p 36 .
  - [6] S.BANACH " Théorie des opérations linéaires " Chelsea Pub.Comp? N.Y (1955) ch IV .
  - [7] N.BOURBAKI "Topologie générale " Chap IX,§3, Prop 7,p 43 .
  - [8] A.M.OSTROWSKI "Uber Normen von Matrizen " Math .Zeitsch. Bd 69 S p 2 (1955)
- T.E.EASTERFIELD "Matix norms and Vector Measures " Duke Math J. Vol 24 pp 663-671 (1955) .

- [9] E. BODEWIG Loc .Cit. p II7
- [10] N. GASTINEL "Conditionnement d'un système d'équations linéaires " C.R.A.S.  
T 248,p2 707-2 709, II mai 1959 .
- [11] G.H.HARDY,J.E.LITTLEWOOD,G.POLYA "Inequalities" Cambridge Univ.Press (1934
- [12] E. BODEWIG Loc.Cit .p IIS
- [13] A.M. TURING "Rounding-off errors in matrix processes " Quart,J.mech.appl.  
math .I,1948,p287-308 .
- [14] A.T.LONSETH "The propagation of error in linear problems" Trans .Amér.  
Math .Soc. 62 p 193-212.
- [15] J.TODD " The condition of finite segments of the HILBERT matrix"  
N.B.S. Appl.Math Series.62.39 ,pII2 (1954).
- [16] J.von NEUMANN,H.H.GOLDSTINE "Numerical inverting of matrices of high order  
Bul.Amer.Math.Soc.Vol.53,No II,pp IO2I-IO99.(1947)  
Formule (7 -5') p IO93.
- [17] J.TODD Loc.Cit. p III formule (5-4)
- [18] A.S.HOUSEHOLDER "Introduction to numerical analysis " Mc Graw-Hill (1956)  
p 44 , formule (2-08-22) .
- [19] A.R.COLLAR "On the reciprocation of certain matrices" Proc.Roy.Soc .  
Edinburgh,59,p195-208 (1939).
- [20] A.R. COLLAR "On the reciprocal of a segment of a generalised HILBERT  
matrix " Proc.Cambridge Phil? Soc. 47,pII-17 (1951)
- [21] E.H.LINFOOT , W.M.SHEPHERD "On a set of linear equations"Quart.J.Math.  
10,84 (1939)
- [22] E.LUKACS ,I.R.SAVAGE "Tables of inverses of finite segments of the  
Hilbert matrix" N.B.S. Appl.Math.Series .39 .pIO5 (1954)
- [23] N.GASTINEL "Sur le choix des pivots dans l'élimination de GAUSS pour  
la résolution d'un système linéaire " C.R.A.S. T250,p275-277  
II janv.1960 .
- [24] R.BELLMAN "Matrix analysis " Mc GRAW-HILL .N.Y.(1960) Chap III.
- [25] E. BODEWIG Loc.Cit. p38 .

- [26] N.GASTINEL "Procédé itératif pour la résolution numérique d'un système d'équations linéaires " C.R.A.S. T 246 ,p 2571 -2574 .  
5 mai 1958.
- [27] KACZMARZ Bull .Inter.Académie Polonaise des Sciences A 1937,355-7
- [28] E.BODEWIG Loc.Cit. p 163.
- [29] J.ALLEN "Relaxations methods "
- [30] FREEMANN Phil. Mag.34,1943,p 409-416.
- [31] SPERNER et SCHEIRER Loc.Cit. Chap IV
- [32] J.M.WEDDERBURN "Lectures on matrices " Am.Math.Soc.Colloq.Publ. Vol 17  
(1934) .
- [33] N.GASTINEL Loc.Cit (2)
- [34] R.A.BUCKINGHAM "Numerical Methods" Pitman ed .London .(1957)  
p 375 .
- [35] H.WAYLAND "Expansion of determinantal equations into polynomial form"  
Quart.Appl.Math.2 (1945) pp 277-306.
- [36] N.GASTINEL Loc.Cit.(2).
- [38] P.A.WHITE "The computation of eigenvalues and eigenvectors of a matrix"  
Jour.Soc.Indus.Appl.Math. Vol 6,No 4 ,(Dec 1958) p 418 .
- [39] W.GIVENS "Computation of plane unitary rotations transforming a general  
matrix to triangular form".  
J.Soc.Indus.Appl.Math. Vol 6,No1,(Mars 1958 ) P26-50 .
- W;GIVENS "Numérical computation of the characteristic values of a réal  
symétric matrix " Oak Ridge Nat.Lab.Rep. No 1574 (1954) .
- [40] P.A.WHITE Loc.Cit. p 425.
- [41] C.LANCZOS "An iteration method for the solution of the eigenvalue problem  
of linear differential and integral operators " J.Res.Nat.B.S.  
45, (1950) pp 255-282 .
- [42] C.G.JACOBI ,J;de Crelle Vol 30,(1846) pp51-95 ,Voir aussi : E.BODEWIG  
Loc.Cit. p 283.
- [43] BOURBAKI . Esp . Vect . Topo. Chap II,§ 5 ,Prop 3 ,PP 94 et suiv.



TABLE DES MATIERES

<u>Introduction</u> , .....	I-III
<u>Chapitre I; Matrices du 2° degré</u> , .....	1
1°) Anneaux de matrices du 2° degré, .....	1
2°) Détermination des matrices du 2° degré, de type (n,n) sur le corps C, .....	2
3°) Utilisation de matrices $X.Y^T$ antiscales, ....	5
<u>Chapitre II : Normes de vecteurs et de matrices</u> , .....	14
1°) Normes de vecteurs, .....	14
2°) Normes de matrices, .....	16
3°) Normes de matrices associées à des normes de vecteurs, .....	19
4°) Groupes de matrices associés à des normes de vecteurs ou de matrices, .....	23
5°) Normes de matrices complètement invariantes, ..	25
6°) Normes de matrices usuelles, .....	26
<u>Chapitre III : Les conditionnements numériques</u> , .....	27
1°) Notion de conditionnement, .....	27
2°) Examen des différentes définitions de condi- tionnement, .....	28
3°) Définition générale des conditionnements numé- riques, .....	29
4°) Cas d'applications linéaires, .....	30
5°) Comparaison des différents conditionnements généraux d'une matrice, .....	32

<u>Chapitre IV : Le conditionnement normalisé <math>C(A)</math>, .....</u>	35
1°) Notations, .....	35
2°) Etude géométrique du rapport pour A normée en lignes, .....	35
3°) Etude analytique, .....	36
4°) Etude du rapport pour A à termes réels non singulière, .....	39
5°) Le conditionnement normalisé $C(A)$ , .....	40
6°) Détermination de et de $C(A)$ , .....	41
7°) Le conditionnement normalisé complémentaire, ?.	43
8°) Relation entre ces deux conditionnements, .....	45
9°) Relation de $C(A)$ et les nombres de VON NEUMANN et de TURING, .....	45
10°) Opérateurs laissant $C(A)$ invariant, .....	49
11°) Conditionnement d'un produit de deux matrices, .....	51
12°) Exemple de calcul d'un conditionnement norma- lisé, matrice de HILBERT : $H_n$ , .....	52
13°) Cas d'une matrice d'interpolation d'opérateur différentiel, .....	54

Chapitre V : Les erreurs dans la résolution des systèmes linéaires

<u>par élimination, .....</u>	57
1°) Les principes de ces calculs d'erreurs, .....	57
2°) Calculs en fixe. Triangularisation et résolu- tion d'un système, .....	60
3°) Calculs en fixe. Inversion d'une matrice par la méthode de GAUSS, .....	66
4°) Calculs en virgule flottante. Triangularisation et résolution d'un système, .....	67
5°) Obtention de bornes pour l'erreur, .....	71
6°) Expériences de calculs en fixe, .....	77
7°) Calculs en virgule flottante. Evaluation de l'erreur, .....	79
8°) Expériences de calcul en virgule flottante, ...	81

9°) Evolution des conditionnements de systèmes successifs d'une élimination, .....	83
<u>Chapitre VI - Les erreurs dans la résolution des systèmes linéaires par</u>	
<u>orthogonalisation, .....</u>	86
1°) Solution d'un système linéaire par orthogonalisation.	86
2°) Etude des erreurs en virgule fixe, .....	88
3°) Etude des erreurs en virgule flottante, .....	94
<u>Chapitre VII - Les normes générales et les procédés itératifs de résolu-</u>	
<u>tion des systèmes linéaires, .....</u>	97
1°) Notations, .....	97
2°) Rapprochement de la solution, .....	98
3°) Décomposition d'une norme, .....	99
4°) Procédé itératif associé à une décomposition d'une norme, .....	100
5°) Procédé associé à la norme $\varphi_0$ , .....	102
6°) Procédé associé à la norme $\varphi_1$ , .....	103
7°) Procédé associé à la norme $\varphi_2$ , .....	105
8°) Méthode du gradient, .....	106
9°) Généralisation de la méthode, .....	107
<u>Chapitre VIII - Calcul du polynôme caractéristique, des valeurs et vec-</u>	
<u>teurs propres, .....</u>	111
1°) Réduction à une forme de FROBENIUS, .....	111
2°) Comparaison avec la méthode de DANILEWSKI, .....	116
3°) Transmutation en une forme triple diagonale, .....	119
4°) Etude des erreurs dans la méthode de réduction à une forme de FROBENIUS, .....	124
5°) Variation des éléments propres correspondant à la variation $\delta A$ de $A$ , .....	126
6°) Cas particulier d'une matrice de forme de FROBENIUS	130
<u>Bibliographie, .....</u>	132-134



II ° THESE

LE THEOREME DE STONE -WEIERSTRASS



LE THEOREME DE STONE - WEIERSTRASS

---:---

I. - INTRODUCTION -

En 1885, WEIERSTRASS [1] a énoncé le théorème célèbre :

"Etant donné une fonction  $f$ , de la variable réelle, et à valeurs réelles, dont le domaine de définition est un segment  $[a \quad b]$ , et qui est continue sur ce segment, il existe toujours un polynôme qui l'approche d'aussi près que l'on veut uniformément, c'est-à-dire  $\varepsilon$  étant donné, il existe un polynôme  $P(x)$  tel que :

$$|f(x) - P(x)| \leq \varepsilon \quad \text{pour tout } x \in [a \quad b] \text{ " .}$$

La démonstration originale de WEIERSTRASS provient du fait, apparemment fort éloigné de la question, que la fonction

$$\psi(x, k) = \frac{1}{k \cdot \sqrt{\pi}} \int_{-\infty}^{+\infty} f(u) \cdot e^{-\left(\frac{u-x}{k}\right)^2} \cdot du$$

tend uniformément vers  $f(x)$  pour  $x \in [a \quad b]$ , si  $k \rightarrow 0$ .

Ayant prouvé que  $\psi(z, k)$  est une fonction entière de  $z$  (dans le plan complexe), on en déduit pour  $\psi(x, k)$  ( $x$  réel) une série uniformément convergente dans  $[a \quad b]$ , il s'ensuit qu'il existe un polynôme tel

$$|\psi(x, k) - f(x)| \leq \frac{\varepsilon}{2}$$

et puisque  $k$  peut être choisi

assez petit afin que  $|f(x) - \psi(x, k)| \leq \frac{\varepsilon}{2}$  pour tout  $x$ , la proposition en résulte.

Les très nombreuses méthodes de démonstration du théorème qui furent présentées par la suite peuvent se diviser en deux catégories (d'après E. BOREL [2]) les unes : de caractère "transcendant" (c'est-à-dire faisant appel à la théorie des fonctions analytiques et du genre de celle de WEIERSTRASS lui-même), les autres : de nature élémentaire.

Les démonstrations de PICARD [3], LERCH[4], VOLTERRA[5], font appel à la théorie des développements en série de FOURIER. Celle de VOLTERRA, par exemple, procède de la façon suivante : il est facile de montrer qu'étant donné  $f$  sur  $[a \quad b]$ , il existe une subdivision finie

$D = (a = x_0, \dots, x_1, \dots, x_n = b)$  de  $[a \quad b]$  telle que si l'on considère la ligne polygonale  $L_p$  dont les côtés joignent les points  $M_i (x_i, f(x_i))$  et  $M_{i+1} (x_{i+1}, f(x_{i+1}))$ , cette ligne brisée  $L_p$  puisse être considérée comme le graphe d'une fonction  $g$ , continue sur  $[a \quad b]$ , qui approche uniformément  $f$  et de sorte que

$$|f(x) - g(x)| \leq \varepsilon \quad \text{pour tout } x \in [a \quad b], \text{ quelque soit } \varepsilon.$$

La question revient donc à l'approche de la fonction  $g(x)$ .

VOLTERRA prolonge la définition de  $g(x)$  afin que dans  $[a \quad c]$ ,  $c > b$ ,  $g(a) = g(c)$ , et considère ce prolongement comme la partie du graphe d'une fonction  $g_1(x)$  périodique et de période  $(c-a)$ . Or  $g_1(x)$  est continue et n'admet qu'un nombre fini de maxima et minima dans une période ; un théorème classique de la théorie des séries de FOURIER montre que l'on peut écrire :

$$g_1(x) = a_0 + \sum_{n=1}^{\infty} \left( a_n \cos \frac{2n\pi}{T} x + b_n \sin \frac{2n\pi}{T} x \right), \quad T = c - a$$

et la série trigonométrique converge uniformément dans  $[a \quad c]$ .

Il existe,  $\varepsilon$  étant donné, un rang  $N$  tel que dès que  $p > N$  et pour tout  $x \in [a \quad c]$

$$\left| g_1(x) - \left\{ a_0 + \sum_{k=1}^p \left( a_k \cos \frac{2k\pi}{T} x + b_k \sin \frac{2k\pi}{T} x \right) \right\} \right| \leq \frac{\varepsilon}{2}$$

Cela prouve déjà l'approximation uniforme de  $f(x)$  par des polynômes trigonométriques. Pour passer aux polynômes ordinaires, il suffit de remarquer que  $\cos \frac{2\pi}{T} kx$ ,  $\sin \frac{2\pi}{T} kx$  sont égales à des séries entières qui convergent uniformément dans  $[a \quad c]$  ; en ne prenant qu'un nombre fini de termes de ces séries,

On en conclut : il existe un polynôme  $P(x)$  tel que

$$\left| \left\{ a_0 + \sum_{k=1}^p \left( a_k \cos \frac{2k\pi}{T} kx + b_k \sin \frac{2k\pi}{T} kx \right) \right\} - P(x) \right| \leq \frac{\varepsilon}{2} \quad \text{quel que soit } x \in [a \quad b]$$

De là résulte la proposition.

Dans ces démonstrations de caractère transcendant, il faut signaler la démonstration d'HILBERT[6], qui varie peu de la démonstration de WEIERSTRASS,

mais a le grand mérite de donner explicitement un polynôme d'approximation.

En effet, HILBERT prouve que, si  $0 < \alpha < a < b < \beta < 1$

On peut toujours par un changement de variable supposer que  $[a \quad b]$  tombe dans  $[0 \quad 1]$  strictement), les polynômes en  $x$  :

$$P_{2n}(x) = \frac{\int_{\alpha}^{\beta} f(u) [1 - (u-x)^2]^n}{\int_{-1}^{+1} (1-u^2)^n}$$

pour  $n \rightarrow \infty$ , convergent uniformément sur  $[a \quad b]$  vers  $f(x)$ .

Les démonstrations classées dans la catégorie "élémentaire" par E. BOREL, sont celles de RUGE [7], LEBESGUE [8], MITTAG-LEFFLER [9], BERNSTEIN [10]. Comme dans la démonstration de VOLTERRA, le problème est ramené à l'approximation uniforme sur  $[a \quad b]$ , de la fonction  $g$ , dont le graphe est une ligne brisée inscrite dans le graphe de  $f$  en  $[a \quad b]$ .

Posons  $y_i = f(x_i)$ .

Soit 
$$\varphi_i(x) = y_{i-1} + \frac{x-x_{i-1}}{x_i-x_{i-1}} (y_i - y_{i-1}) \quad ;$$

Soit (1) 
$$\varphi(x) = \varphi_i(x) + \sum_{i=1}^{i=h-1} [\varphi_{i+1}(x) - \varphi_i(x)] \alpha(x-x_i) \quad ;$$

$\alpha(x)$  étant la fonction échelon :

$$\alpha(x) = \begin{cases} 1 & \text{pour } x \geq 0 \\ 0 & \text{pour } x < 0 \end{cases}$$

D'après (1), il suffirait de trouver une approximation uniforme de l'échelon entre  $-(b-a)$  et  $+(b-a)$  par des polynômes. Or, cela est impossible, (l'oscillation de  $\alpha$  autour de 0 est égale à 1), mais il est possible de trouver une fonction  $\beta$  telle que :

$$|\alpha(x) - \beta(x)| \leq \varepsilon'$$

pour tout  $x \in [-(b-a) - \eta] \cup [+\eta \quad (b-a)]$

quel que soit  $\varepsilon', \eta$  donnés à l'avance.

Si 
$$\psi(x) = \varphi_i(x) + \sum_{i=1}^{i=h-1} [\varphi_{i+1}(x) - \varphi_i(x)] \beta(x-x_i) \quad , \text{ alors}$$

$$|\varphi(x) - \psi(x)| \leq 2h M \varepsilon' \quad (\text{si } |\varphi(x)| \leq M \text{ sur } [a \quad b])$$

pour  $x$  en dehors des intervalles  $[x_i - \eta \quad x_i + \eta]$

Mais la fonction continue  $\varphi_{i+1}(x) - \varphi_i(x)$  est nulle en  $x_i$ , on peut choisir  $\eta$  assez petit de sorte que dans  $[x_i - \eta \quad x_i + \eta]$ ,  $|\varphi_{i+1}(x) - \varphi_i(x)| \leq \varepsilon'$

et dans cet intervalle  $|\varphi(x) - \psi(x)| \leq \varepsilon' M_1 + 2h M \varepsilon'$

si  $|\alpha(x) - \beta(x)| \leq M$ , pour tout  $x \in [-(b-a) \quad +(b-a)]$ , on peut donc dire qu'on peut approcher  $\varphi(x)$  par  $\psi(x)$  uniformément, donc  $f(x)$  par  $\psi(x)$ .

RUNGE prend pour  $\beta(x)$  la fraction rationnelle  $\frac{1}{1 + (1-x)^{2k}}$ . Après avoir ramené  $(b-a)$  à être  $< 1$  (par un changement de variables), il prouve que  $\beta(x)$  répond à la question pour  $k$  assez grand et montre ainsi que l'on peut trouver une approximation de  $f$  sur  $[a \quad b]$  par une fraction rationnelle. Il indique une méthode [7] pour passer d'une approximation par des fractions rationnelles à une approximation par des polynômes.

MITTAG-LEFFLER propose de prendre  $\beta(x) = 1 - 2^{-(1-x)^k}$  (l'intervalle  $[-(b-a) \quad +(b-a)]$  toujours ramené à être dans  $[-1 \quad +1]$ ) pour  $k$  assez grand, c'est une fonction entière, qui à son tour sera approchée uniformément par un polynôme.

Enfin, LEBESGUE remarque que  $\alpha(x)$  est telle que

$$x \cdot [2\alpha(x) - 1] = \begin{cases} x & \text{pour } x \geq 0 \\ -x & \text{pour } x < 0 \end{cases} = +|x| = + (x^2)^{\frac{1}{2}}$$

si l'on remarque  $|x| = (1+z)^{\frac{1}{2}} = 1 + \frac{1}{2}z - \frac{1.3}{2.4}z^2 + \dots$

si  $z = x^2 - 1$ , la série converge pour  $|z| < 1$ , et même pour  $z = \pm 1$ , elle converge absolument et uniformément pour  $-1 \leq z \leq 1$ , donc il existe un polynôme,  $\varepsilon'$  étant donné tel  $||x| - P(z)| \leq \varepsilon'$  pour tout  $x \in [-1 \quad +1]$

Il est clair que pour 0 cela donne  $|P(0)| \leq \varepsilon'$  donc :  $||x| - (P(x) - P(0))| \leq 2\varepsilon'$

On peut, si  $\varepsilon = 2\varepsilon'$ , trouver un polynôme sans terme constante tel que

$$||x| - Q(x)| \leq \varepsilon$$

Dès lors, puisque  $|x[2\alpha(x) - 1] - Q(x)| \leq \varepsilon$

si  $xQ'(x) = Q(x)$ ,  $|2\alpha(x) - 1 - Q'(x)| \leq \frac{\varepsilon}{\eta}$  si  $\eta \leq |x|$

On est bien arrivé à trouver, si  $\eta$  est donné, un polynôme :  $\frac{1}{2}(1 + Q'(x)) = \beta_1(x)$ , tel :

$$|\alpha(x) - \beta_1(x)| \leq \frac{\varepsilon}{2\eta} \quad \text{qui assure l'approximation uniforme de } \alpha(x)$$

dans  $[-1 \quad -\eta] \cup [\eta \quad +1]$ .

LEBESGUE ne présentait pas sa démonstration tout à fait de cette façon, mais elle s'y ramène. Nous verrons le parti qu'en a tiré STONE pour la démonstration de

son théorème.

Enfin, les généralisations diverses du théorème au cas, soit des sommes trigonométriques, soit des fonctions à plusieurs variables, furent obtenues par des procédés analogues par LEBESGUE [3], MONTEL [10], DE LA VALLEE POUSSIN [12], BERNSTEIN [10].

II. - GENERALITES SUR LA DEMONSTRATION DE STONE . - [13], [14]

Il y a trois remarques qui furent le départ de la démonstration de STONE et qui assurent la pleine généralité du théorème par lui prouvé :

- 1) la première est que le problème est de trouver l'approximation uniforme pour les fonctions continues (à valeurs réelles) dont le domaine de définition est un ensemble compact S .
- 2) la deuxième idée est que l'ensemble des fonctions continues devant assurer cette approximation forment une algèbre de fonctions continues.
- 3) Etant donnés deux points  $x_1, x_2$  distincts de S, quelconques, il est évident qu'il y a : 1°) des fonctions continues  $f$  sur S telles que  $f(x_1) \neq f(x_2)$  .  
 2°) si l'approximation de  $f$  est réalisée, pour tout  $\epsilon$ , une fonction  $g \in A$  telle que  $|f-g| \leq \epsilon$  pour tout  $x \in S$ ,  
 donc  $|f(x_1) - g(x_1)| \leq \epsilon$  et  $|f(x_2) - g(x_2)| \leq \epsilon$   
 et l'hypothèse  $f(x_1) \neq f(x_2)$  implique que  $g(x_1) \neq g(x_2)$   
 pour un choix de  $\epsilon$  .

Donc, les fonctions de A devront vérifier cette condition nécessaire : pour tout couple  $x_1, x_2$  de points distincts de S, il existe au moins  $g \in A$  telle que  $g(x_1) \neq g(x_2)$ , on dira que A sépare les points de S

Le théorème de STONE montre que ces conditions sont généralement suffisantes : en voici l'énoncé :

Théorème : Si S est un espace compact, A une algèbre de fonctions continues à valeurs réelles qui sépare les points de S, c'est-à-dire si  $x_1, x_2$  sont distincts ( $x_1 \neq x_2$ ) et quelconques, il existe toujours une  $f \in A$  telle que  $f(x_1) \neq f(x_2)$ . Alors :

La fermeture uniforme  $\bar{A}$  de A est :

- ou bien l'ensemble de toutes les fonctions continues sur S
- ou bien l'ensemble des fonctions continues qui s'annulent en un point déterminé de S.

Le plan de la démonstration est le suivant :

1°) on prouve que l'ensemble  $\bar{A}$  est nécessairement réticulé  
 c'est-à-dire si  $f, g \in \bar{A}$ , les fonctions  $f \vee g = \text{Max}(f, g)$ ,  
 $f \wedge g = \text{Min}(f, g) \in \bar{A}$ , par un lemme (L)  
 (Il est clair que si  $A$  est une algèbre de fonctions continues en  $S$ ,  
 $\bar{A}$  est aussi une algèbre de fonctions continues).

2°) On prouve ensuite que si un ensemble  $A_1$  de fonctions continues en  $S$   
 est réticulé, et qu'une fonction  $f$  continue est telle qu'elle  
 puisse être approchée en tout couple  $x_1, x_2$  uniformément par une  
 fonction de  $A_1$  (quels que soient  $x_1, x_2, \epsilon$ ) il existe  $g \in A_1$   
 telle  

$$|f(x_1) - g(x_1)| < \epsilon, \quad |f(x_2) - g(x_2)| < \epsilon$$
  
 alors  $f$  peut être approchée uniformément dans tout  $S$ , c'est-à-dire  
 $f \in \bar{A}_1$ .

III. - DEMONSTRATION? -

a) Un ensemble de fonctions  $f$  est une algèbre sur  $R$  :

- 1) si c'est un espace vectoriel :
  - $f \pm g \in B, \lambda f \in B$  quels que soient  $f, g \in B$  et  $\lambda \in R$
- 2) si un produit  $f.g$  est défini distributif par rapport aux opérations sur l'espace vectoriel

On dira que cette algèbre est uniformément fermée si toute fonction  
 adhérente à  $B$  y appartient (au sens de la topologie de la convergence uniforme)

Lemme I : toute algèbre uniformément fermée de fonctions bornées sur un ensemble  $S$   
 (pas forcément compact), est clos pour les opérations de treillis Max. et  
Min., c'est-à-dire est réticulée.

On sait que

$$\begin{cases} \text{Max}(f, g) = \frac{1}{2}(f + g + |f - g|) \\ \text{Min}(f, g) = \frac{1}{2}(f + g - |f - g|) \end{cases}$$

je veux prouver, selon les hypothèses, que si  $f, g \in A$  :  $\text{Max}(f, g)$  et  $\text{Min}(f, g) \in A$

D'après ce qui est écrit plus haut, il suffit de montrer (A est un espace vectoriel sur R) que si une fonction  $f \in A$ , la fonction  $|f| \in A$  aussi.

Par hypothèse, toutes les fonctions de A sont bornées - je prends -  $\|f\| = \max_{x \in S} |f(x)| < 1$

$$\max_{x \in S} |f(x)| < 1 \quad (\text{sinon je remplace } f \text{ par } \lambda f$$

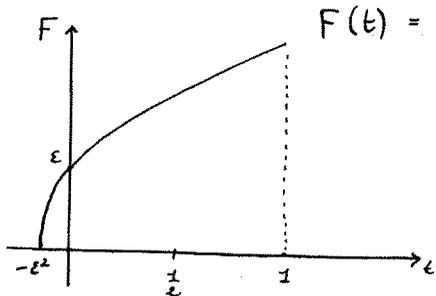
où  $\lambda$  est une constante  $\in R$ )

-Il suffit donc de voir si,  $f \in A$  et  $\|f\| < 1$  implique  $|f| \in A$ .

Puisque A est uniformément fermée, il suffit de montrer que  $|f|$  peut être uniformément approchée par des fonctions de A (hypothèses), c'est-à-dire que, quel que soit  $\varepsilon'$ , il existe  $g \in A$  telle  $|g - |f|| < \varepsilon'$ , pour tout  $x \in S$ .

-l'idée est d'étudier le même problème dans un cas particulier.

Considérons la fonction  $F: t \rightarrow F(t) = +\sqrt{t+\varepsilon^2} = (t+\varepsilon^2)^{\frac{1}{2}}$ , ( $t \in [-\varepsilon^2, +\infty)$ )  
et formons le développement de TAYLOR de  $F(t)$  autour de  $\frac{1}{2}$



$$F(t) = F\left(\frac{1}{2}\right) + \left(t - \frac{1}{2}\right)F'\left(\frac{1}{2}\right) + \frac{1}{2!} \left(t - \frac{1}{2}\right)^2 F''\left(\frac{1}{2}\right) + \dots$$

il est clair que ce développement est convergent et uniformément convergent pour t entre  $0 \leq t \leq 1$   
Si je pose  $t = x^2$ , je peux donc, étant donné  $\varepsilon$  trouver un polynôme  $P(x^2)$  tel que

$$|P(x^2) - (x^2 + \varepsilon^2)^{\frac{1}{2}}| < \varepsilon, \quad \forall x \in [-1, +1]$$

mais alors  $|P(0) - \varepsilon| < \varepsilon$ ,  $|P(0)| < 2\varepsilon$  si l'on pose  $Q(x) = P(x) - P(0)$

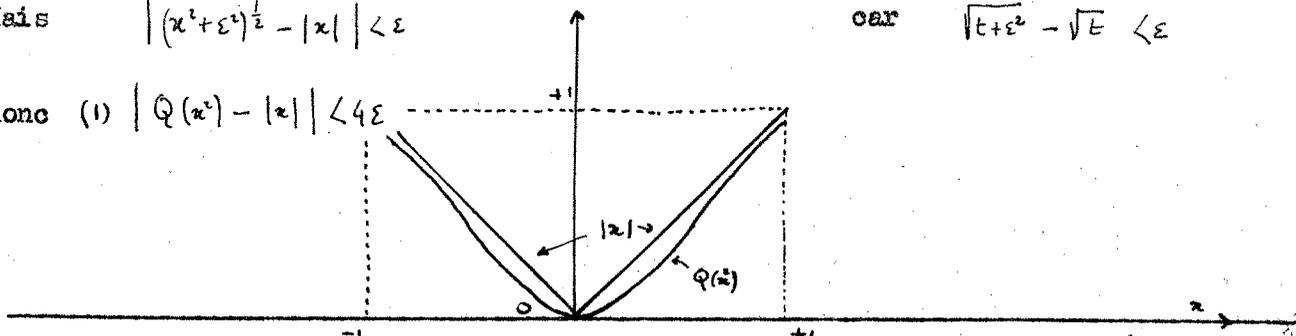
(c'est-à-dire le polynôme  $P(x)$  sauf son terme constant)

$$|Q(x^2) - (x^2 + \varepsilon^2)^{\frac{1}{2}}| \leq |P(x^2) - (x^2 + \varepsilon^2)^{\frac{1}{2}}| + |P(0)| < 3\varepsilon$$

Mais  $|(x^2 + \varepsilon^2)^{\frac{1}{2}} - |x|| < \varepsilon$

$$\text{car } \sqrt{t+\varepsilon^2} - \sqrt{t} < \varepsilon$$

donc (1)  $|Q(x^2) - |x|| < 4\varepsilon$



Si  $f$  est donné et puisque  $Q$  n'apas de terme constant  $\alpha f \in A$ ,  $\alpha f^2 \in A$ ...  
... $Q(f^2) \in A$  et puisque  $\|f\| \leq 1$ , donc (1) implique:

$$|Q(f^2) - |f|| < \varepsilon' \quad \text{pour } \forall x \in S, \quad \varepsilon' = 4\varepsilon$$

Q.E.D.

b) Prouvons le lemme II . S est ici un compact et A n'a même pas besoin d'être une algèbre.

Mais supposons  $A \subset \mathcal{C}(S, \mathbb{R})$  réticulé . ( $\mathcal{C}(S, \mathbb{R})$  : ensemble des applications continues de S dans  $\mathbb{R}$ )

Lemme II : Si une fonction  $f \in \mathcal{C}(S, \mathbb{R})$  a la propriété de pouvoir être approchée d'aussi près que l'on veut par une fonction de A, en deux points et quels que soient ces deux points p, q de S, alors  $f \in \overline{A}$ .

C'est à dire il existe une fonction de A qui l'approche uniformément en tout point.

Par hypothèse, quels que soient  $\varepsilon$ , et p, q,  $\exists f_{p,q} \in A$  telle

$$(1) \quad |f(p) - f_{p,q}(p)| < \varepsilon \quad \text{et} \quad |f(q) - f_{p,q}(q)| < \varepsilon \quad (1')$$

Soit

$$U_{p,q} = \{ x; f_{p,q}(x) < f(x) + \varepsilon \} = \{ x; f_{p,q}(x) - f(x) < \varepsilon \}$$

puisque  $f_{p,q} - f$  est continue et que l'ensemble des réels  $< \varepsilon$  est un ouvert  $O_\varepsilon$  de  $\mathbb{R}$ ,  $U_{p,q}$  est l'image réciproque de  $O_\varepsilon$  ; il est ouvert dans S.

Il contient p, q de part (1) et (1').

De même :

"  $V_{p,q} = \{ x; f_{p,q}(x) > f(x) - \varepsilon \}$  est un ouvert de S qui contient p et q.

- Supposons q fixé, et faisons varier dans S, p. Les  $U_{p,q}$  recouvrent tout S, il existe un recouvrement (propriété de compacité de S) fini de S par des  $U_{p_i,q}$  et  $p_1, \dots, p_k$  les points qui permettent cela pour q fixé.

Si  $x \in S$ , il existe donc  $U_{p_i,q}$  tel  $x \in U_{p_i,q}$  c'est-à-dire tel  $f_{p_i,q}(x) < f(x) + \varepsilon$

soit  $f_q(x) = \text{Min} [f_{p_1,q}, f_{p_2,q}, \dots, f_{p_k,q}]$ , on a  $f_q \in A$  (hypothèse que A est réticulé)

et aussi  $f_q(x) < f(x) + \varepsilon$  quel que soit  $x \in S$ .

De plus, soit  $W_q = \bigcap_{i=1, \dots, k} V_{p_i,q}$  c'est un ouvert (int. finie d'ouverts)

contient q (chaque  $V_{p_i,q}$  contient q). Si  $x \in W_q$  on a évidemment :

$f_{p_i,q}(x) > f(x) - \varepsilon$  quel que soit  $p_i$  ( $i=1, \dots, k$ ), donc :

$$f_q(x) > f(x) - \varepsilon$$

En résumé, on vient de trouver, pour tout q et  $\varepsilon$ , un  $W_q$  et une fonction  $f_q \in A$  tels :

- 1) quel que soit  $x \in S$  :  $f_q(x) < f(x) + \varepsilon$
- 2) pour  $x \in W_q$  :  $f_q(x) > f(x) - \varepsilon$
- 3)  $W_q$  est ouvert et contient q

Les  $W_q$ , pour  $q$  variant dans  $S$ , recouvrent  $S$ , donc il existe (propriété de compacité de  $S$ ) un recouvrement fini de  $S$  par de tels  $W_1; W_{q_1}, W_{q_2}, \dots, W_{q_m}$

Prenons ;

$$f_\varepsilon = \text{Max} [f_{q_1}, f_{q_2}, \dots, f_{q_m}] \quad , \quad f_\varepsilon \in A$$

(hypothèse que  $A$  est réticulé)

Quel que soit  $x$ : 1)  $f_\varepsilon(x) < f(x) + \varepsilon$

2) pour  $x$  donné, il existe un indice  $j$  tel  $x \in W_{q_j}$  et tel

$$f_{q_j}(x) > f(x) - \varepsilon \quad \text{donc} \quad f_\varepsilon(x) > f(x) - \varepsilon$$

Par suite, on vient de trouver  $f_\varepsilon \in A$  tel que  
ou

$$f(x) - \varepsilon < f_\varepsilon(x) < f(x) + \varepsilon$$

$$|f(x) - f_\varepsilon(x)| < \varepsilon \quad , \quad \forall x \in S$$

Q.E.D.

C) La démonstration des deux lemmes obtenus, passons à la démonstration du théorème de STONE lui-même.

I°) Supposons que pour tout point  $x \in S$  il existe une  $f \in A$  telle que  $f(x) \neq 0$  (c'est le cas si les fonctions constantes sur  $S$ ,  $\in A$ ).

Alors si  $x_1 \neq x_2$ , il existe un  $f$  tel  $0 \neq f(x_1) \neq f(x_2) \neq 0$

En effet, il existe  $f_1 \in A$  telle que :  $f_1(x_1) \neq f_1(x_2)$ ,

(hypothèse). De plus, il existe  $f_2 \in A$  telle :  $f_2(x_2) \neq 0$

- Supposons, parexemple, que  $f_1(x_1) \neq 0$  on aura :  $0 \neq f_1(x_1) \neq f_1(x_2)$

si  $f_1(x_2) = 0$ , considérons la fonction  $\alpha \cdot \frac{f_1(x)}{f_1(x_1)} + f_2(x) = \varphi(x)$

elle est de  $A$ ,

( $\alpha$  nombre réel.)

Il est clair que je peux choisir  $\alpha$  de sorte que : 1)  $\varphi(x_1) \neq \varphi(x_2)$

2)  $f_2(x_2) + \alpha \neq 0$

- Conclusion : Dans le cas où, pour tout  $x$ , il existe une  $f \in A$  telle que  $f(x) \neq 0$  et si  $A$  sépare  $S$ , il existe pour tout couple  $x_1 \neq x_2$ , une  $f \in A$  telle

$$0 \neq f(x_1) \neq f(x_2) \neq 0$$

Posons :  $f(x_1) = \lambda_1$ ,  $f(x_2) = \lambda_2$  ;  $0 \neq \lambda_1 \neq \lambda_2 \neq 0$

Cette remarque faite,

a) soit  $\bar{A}$  la fermeture uniforme de l'algèbre  $A$ , c'est une algèbre, et puisque ces fonctions sont à domaine de définition sur un compact  $S$ , chacune est bornée. Le lemme II prouve qu'elle est réticulée.

b) D'autre part, soit  $\varphi$  une fonction continue quelconque, et posons :

$$\varphi(x_1) = a \quad \text{et} \quad \varphi(x_2) = b, \text{ étant 2 points quelconques.}$$

Soit  $f$  la fonction déterminée plus haut. Il est facile de déterminer un polynôme en  $X$  tel  $\alpha X^2 + \beta X$  tel  $\begin{cases} \alpha X_1^2 + \beta X_1 = a \\ \alpha X_2^2 + \beta X_2 = b \end{cases}$

(On prend :  $\alpha = \frac{bX_1 - aX_2}{X_1X_2(X_2 - X_1)}, \beta = \frac{bX_1^2 - aX_2^2}{X_1X_2(X_2 - X_1)}$ )

alors :  $g = \alpha f^2 + \beta f \in A$  et pour  $x_1, x_2$  prend les mêmes valeurs que  $\varphi$ .

Il est donc clair que  $\varphi$  fonction continue quelconque est "approchée" en chaque couple aussi près que l'on veut par des fonctions de  $A$  : le lemme I  $\Rightarrow \varphi \in \bar{A}$  qui est donc identique à  $\mathcal{C}(S, \mathbb{R})$

II°) Supposons que pour un point particulier  $x_0$ , toutes les  $f \in A$  s'annulent. Je dois prouver que si  $g$  est une fonction qui est continue et s'annule pour  $x = x_0$ ,  $g \in \bar{A}$

Supposons augmenter l'ensemble  $A$  des fonctions constantes sur  $S$ . et considérons l'algèbre  $A_1$  engendrée par  $A$  et ces fonctions. L'ensemble  $A_1$  ainsi obtenu est une algèbre qui sépare les points de  $S$ , et pour tout  $x$ , il existe bien une  $f \in A$  telle que  $f(x) \neq 0$ . Donc, pour  $A_1$ , la première partie s'applique. Toute fonction de  $A_1$  est de la forme  $f + c$ , ( $f \in A$ ) et  $g$  étant donné ainsi que  $\varepsilon$ , il existe une  $f + c$  telle que sur tout  $S$ ,  $|f + c - g| < \frac{\varepsilon}{2}$  donc en  $x_0$  :  $|c| < \frac{\varepsilon}{2}$

si bien  $|f(x) - g(x)| < \varepsilon, \forall x \in S$  Q.E.D.

IV. - APPLICATIONS. -

Nous citerons pour terminer les applications du théorème de STONE, et, pour commencer, le théorème général de WEIERSTRASS qui en résulte directement.

I. Théorème de WEIERSTRASS - (approximation par des polynômes).

Si  $X$  est un ensemble fermé borné de  $\mathbb{R}^n$ , toute fonction  $f$  continue sur  $X$  peut être approchée uniformément sur  $X$  par des polynômes en  $(x_1, \dots, x_n)$ . Si  $X$  contient l'origine  $(0, \dots, 0)$  la fonction  $f$  peut être approchée uniformément par des polynômes s'annulant à l'origine si, et seulement si,  $f$  s'annule elle-même à l'origine.

## II. Théorème d'approximation par des fonctions trigonométriques -

Soit  $f$  une fonction à valeurs réelles de la variable  $\theta$  réelle et périodique, de période  $2\pi$  et continue.

$f$  peut être uniformément approchée par des polynômes trigonométriques de la forme :

$$f(\theta) = \frac{a_0}{2} + \sum_{k=1}^{k=N} (a_k \cos k\theta + b_k \sin k\theta)$$

Le raisonnement est classique (déjà fait par HILBERT) : sur le cercle unité ( $x_1^2 + x_2^2 = 1$ ), si l'on pose :

$$F(x_1, x_2) = r \cdot f(\theta), \quad r = +\sqrt{x_1^2 + x_2^2}$$

$F(x_1, x_2)$  est une fonction continue sur ce cercle qui est fermé, borné, donc compact.

Il existe un polynôme  $P(x_1, x_2)$  qui approche  $F(x_1, x_2)$  sur le cercle uniformément.

Donc :

$$P(\cos \theta, \sin \theta) = \sum_{i,j \leq N} c_{ij} \cos^i \theta \sin^j \theta, \quad \text{approche } f(\theta) \text{ uniformément.}$$

Il suffit de transformer les produits  $\cos^i \theta, \sin^j \theta$  en  $\sin$  faisant intervenir que des expressions de la forme  $\cos p\theta, \sin q\theta$  pour obtenir le théorème.

Citons enfin des théorèmes d'approximations très analogues, obtenus en "rendant compact"  $\mathbb{R}$  ou une partie de  $\mathbb{R}$ .

## III. Théorème d'approximation par des fonctions de LAGUERRE -

Toute fonction continue, à valeurs réelles, définie sur  $0 \leq x < +\infty$

et qui est telle que  $\lim_{x \rightarrow \infty} f(x) = 0$  peut être approchée

uniformément par des fonctions de la forme :  $e^{-\alpha x} p(x)$

( $p(x)$  (polynôme en  $x$ )) (fonctions de LAGUERRE).

## IV. - Théorème d'approximation par des fonctions d'HERMITE -

Toute fonction continue, à valeurs réelles, définie sur  $-\infty < x < +\infty$

et telle que  $\lim_{x \rightarrow \pm \infty} f(x) = 0$

peut être approchée uniformément par des fonctions de la forme :

$e^{-x^2} p(x)$  ( $P(x)$  : polynôme en  $x$ ) : fonctions

d'HERMITE.

Enfin, indiquons, pour terminer, que STONE [15] a étendu son théorème dans des conditions un peu plus larges : fonctions à valeurs dans  $\mathbb{C}$ , et définies sur des espaces localement compacts.

Il donne, comme exemples d'applications, en topologie générale :

- une démonstration du théorème de LEBESGUE-URYSONN sur l'extension de la définition d'une fonction continue sur une partie fermée d'un espace topologique à tout cet espace.
  - une démonstration du théorème de DIEUDONNE indiquant que si  $X$  est le produit d'une famille  $X_i$  ( $i \in I$ ), d'espaces compacts, toute fonction continue à valeurs réelles sur  $X$  peut être approchée uniformément par des sommes finies de produits finis de fonctions d'une variable sur  $X$ .
-