

Ministère de l'Enseignement et de la Recherche Scientifique

Université des Sciences et de la Technologie
« Houari BOUMEDIENE »

38/83

2 ex

ECOLE NATIONALE POLYTECHNIQUE D'ALGER
DEPARTEMENT D'ELECTRONIQUE

الجامعة الوطنية للعلوم والتكنولوجيا
PROJET DE FIN D'ETUDES
الهندسة
INGENIORAT D'ETAT EN ELECTRONIQUE
ECOLE NATIONALE POLYTECHNIQUE
BIBLIOTHEQUE

Etude de la Fonction signe de matrice
Applications à la résolution
des équations de Riccati et de Lyapunov

Proposé et suivi par :

Mr. R. TOUMI : Dr. Ingénieur
M^{elle} M. AMINI : Dr. Ingénieur

Etudié par :

Y. LAOUAR
A. HAMZAOUI

JANVIER 1983

Ministère de l'Enseignement et de la Recherche Scientifique

38/83

Université des Sciences et de la Technologie
« Houari BOUMEDIENE »

2 ex

ECOLE NATIONALE POLYTECHNIQUE D'ALGER
DEPARTEMENT D'ELECTRONIQUE

الجامعة الوطنية للعلوم والتقنية
PROJET DE FIN D'ETUDES
الهندسة
INGENIORAT D'ETAT EN ELECTRONIQUE
ECOLE NATIONALE POLYTECHNIQUE
BIBLIOTHEQUE

Etude de la Fonction signe de matrice
Applications à la résolution
des équations de Riccati et de Lyapunov

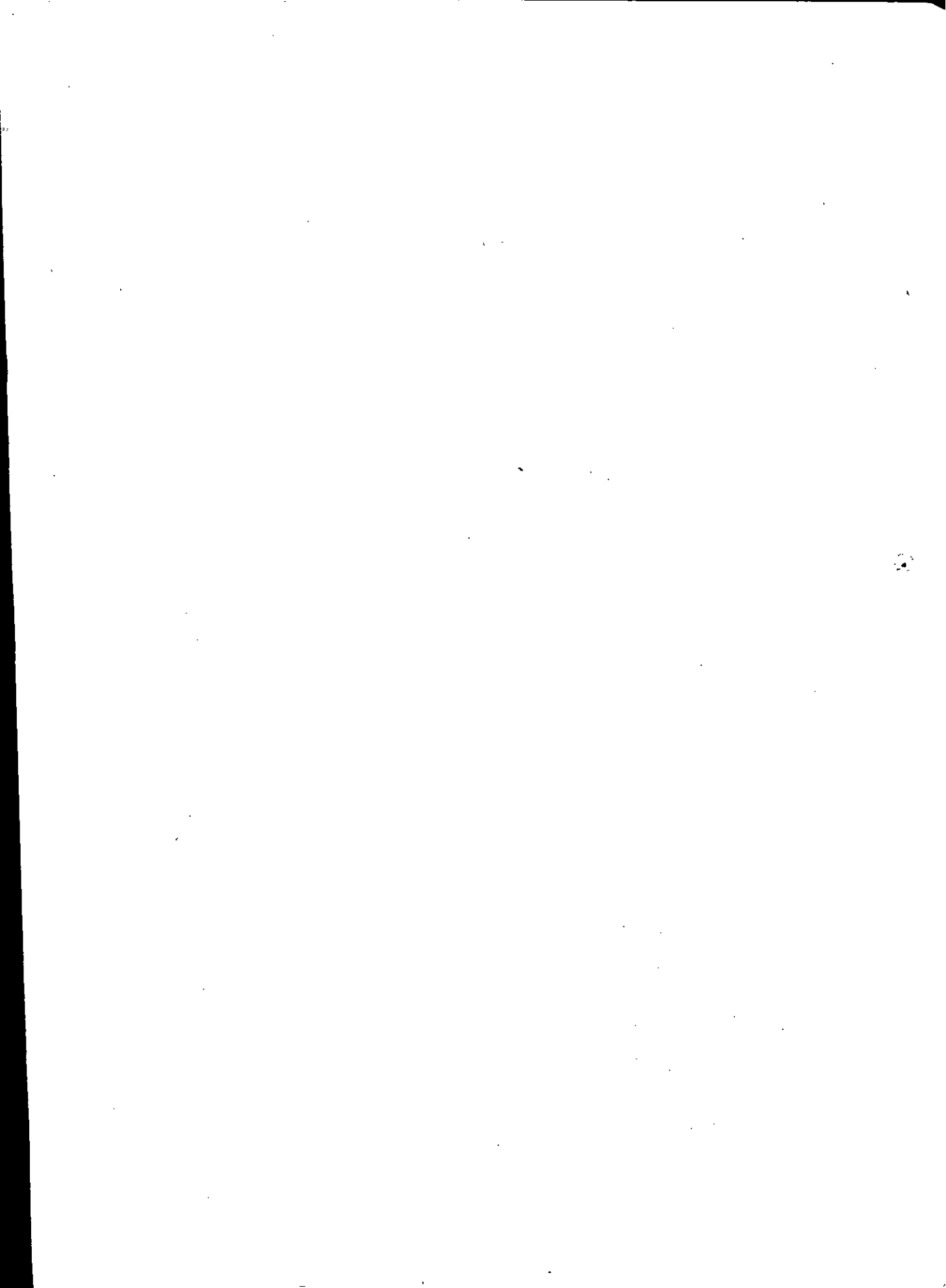
Proposé et suivi par :

Mr. R. TOUMI : Dr. Ingénieur
M^{lle} M. AMINI : Dr. Ingénieur

Etudié par :

Y. LAOUAR
A. HAMZAOU

JANVIER 1983



Ministère de l'Enseignement et de la Recherche Scientifique

Université des Sciences et de la Technologie

« Houari BOUMEDIENE »

ECOLE NATIONALE POLYTECHNIQUE D'ALGER

DEPARTEMENT D'ELECTRONIQUE

PROJET DE FIN D'ETUDES

INGENIORAT D'ETAT EN ELECTRONIQUE

Etude de la Fonction signe de matrice
Applications à la résolution
des équations de Riccati et de Lyapunov

Proposé et suivi par :

Mr. R. TOUMI : Dr. Ingénieur

M^{lle} M. AMINI : Dr. Ingénieur

Etudié par :

Y. LAOUAR

A. HAMZAOUI

JANVIER 1983

Ministère de l'Enseignement et de la Recherche Scientifique

Université des Sciences et de la Technologie
« Houari BOUMEDIENE »

ECOLE NATIONALE POLYTECHNIQUE D'ALGER
DEPARTEMENT D'ELECTRONIQUE

PROJET DE FIN D'ETUDES
INGENIORAT D'ETAT EN ELECTRONIQUE

Etude de la Fonction signe de matrice
Applications à la résolution
des équations de Riccati et de Lyapunov

Proposé et suivi par :

Mr. R. TOUMI : Dr. Ingénieur
M^{lle} M. AMINI : Dr. Ingénieur

Etudié par :

Y. LAOUAR
A. HAMZAOU

JANVIER 1983

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

«اللَّهُ نُورُ السَّمَاوَاتِ وَالْأَرْضِ مِثْلُ نُورِهِ كَمِشْكَاةٍ
فِيهَا مِصْبَاحٌ الْمِصْبَاحُ فِي زُجَاجَةٍ الزُّجَاجَةُ
كَأَنَّهُ كَوْكَبٌ ذُرِّيُّ بُرُودٍ مِنْ شَجَرَةٍ مُبَارَكَةٍ
زَيْتُونَةٍ شَرْقِيَّةٍ وَأُخْرَبِيَّةٍ يُعْادُ زَيْتُهَا يَوْمَئِذٍ
وَلَوْ لَمْ تَنْسَسْهُ نَارٌ لَوْرَعَى نُورٌ يَهْدِي اللَّهُ لِنُورِهِ
مَنْ يَشَاءُ وَيَضْرِبُ اللَّهُ الْأَمْثَالَ لِلنَّاسِ وَاللَّهُ
بِكُلِّ شَيْءٍ عَلِيمٌ»

سورة النور

كسب الله الرحمن الرحيم

«الله نور السماوات والأرض مثل نوره كمشكاة
فيها مصباح المصباح في زجاجة الزجاجة
كانها كوكب دري يوقد من شجرة مباركة
زيتونة لا شرقية ولا غربية يكاد زيتها يضيء
ولو لم يمسسه نار لنور على نور يهدي الله لنوره
من يشاء ويضرب الله الأمثال للناس والله
بكل شيء عليم»

سورة النور

E M E R C I E M E N T S

Le présent travail a été réalisé au sein de la division .V.
du C.E.N.

Pour cela, nous tenons à remercier Mr. SANSAL pour nous avoir
admis au sein de sa division.

Nous exprimons toute notre reconnaissance à Mr. TOUMI et
Mlle AMINI pour l'encadrement et l'aide considérable qui nous ont
été fourni, durant tout le semestre.

Nous tenons également à remercier tout le personnel du centre
de calcul du C.E.N. en particulier RACHID et ABDENOUR.

Le Maire et tous les employés de l'A.P.C. de M'toussa pour
la rappe.

Nos vifs remerciements aussi pour Mme. BAHLOUL H. et Monsieur
BENSLITANE. B . pour le tirage.

REMERCIEMENTS

Le présent travail a été réalisé au sein de la division .V. du C.E.N.

Pour cela, nous tenons à remercier Mr. SANSAL pour nous avoir admis au sein de sa division.

Nous exprimons toute notre reconnaissance à Mr. TOUMI et Mlle AMINI pour l'encadrement et l'aide considérable qui nous ont été fourni, durant tout le semestre.

Nous tenons également à remercier tout le personnel du centre de calcul du C.E.N. en particulier RACHID et ABDENOUR.

Le Maire et tous les employés de l'A.P.C. de M'toussa pour la drappe.

Nos vifs remerciements aussi pour Mme. BAHLOUL H. et Monsieur BENSLITANE. B . pour le tirage.

E M E R C I E M E N T S

Le présent travail a été réalisé au sein de la division .V.
du C.E.N.

Pour cela, nous tenons à remercier Mr. SANSAL pour nous avoir
admis au sein de sa division.

Nous exprimons toute notre reconnaissance à Mr. TOUMI et
Mlle AMINI pour l'encadrement et l'aide considérable qui nous ont
été fourni, durant tout le semestre.

Nous tenons également à remercier tout le personnel du centre
de calcul du C.E.N. en particulier RACHID et ABDEENOUR.

Le Maire et tous les employés de l'A.P.C. de M'toussa pour
la rappe.

Nos vifs remerciements aussi pour Mme. BAHLOUL H. et Monsieur
BENSLITANE. B . pour le tirage.

E M E R C I E M E N T S

Le présent travail a été réalisé au sein de la division .V.
du C.E.N.

Pour cela, nous tenons à remercier Mr. SANSAL pour nous avoir
admis au sein de sa division.

Nous exprimons toute notre reconnaissance à Mr. TOUMI et
Mlle AMINI pour l'encadrement et l'aide considérable qui nous ont
été fourni, durant tout le semestre.

Nous tenons également à remercier tout le personnel du centre
de calcul du C.E.N. en particulier RACHID et ABDENOUR.

Le Maire et tous les employés de l'A.P.C. de M'toussa pour
la drappe.

Nos vifs remerciements aussi pour Mme. BAHLOUL H. et Monsieur
BENSLITANE. B . pour le tirage.

IIY E II II II II II II II

XX

A MON GRAND-PERE

A MON PERE ET MA MERE QUI M'ONT TOUT DONNE

A MES FRERES ET SOEURS

A MES COUSINS ET TOUS MES AMIS

Yahia

A LA MEMOIRE DE MON PERE

A MA GRAND-MERE ET MA MERE

A MES FRERES ET SOEURS

A TOUS MES AMIS ET COUSINS

Abdelaziz

II) E II) II) II) II) II) E II)

XX

A MON GRAND-PERE

A MON PERE ET MA MERE QUI M'ONT TOUT DONNE

A MES FRERES ET SOEURS

A MES COUSINS ET TOUS MES AMIS

Yahia

A LA MEMOIRE DE MON PERE

A MA GRAND-MERE ET MA MERE

A MES FRERES ET SOEURS

A TOUS MES AMIS ET COUSINS

Abdelaziz

TABLE DES MATIERES

<u>INTRODUCTION/</u>	1
<u>CHAPITRE I/ ETUDE DE LA FONCTION SIGNE D'UNE MATRICE</u>	3
Introduction	3
1-Introduction à la fonction signe d'une matrice.	4
2-Construction d'une suite Z_k telle que $Z_\infty = \text{signe}(Z_0)$	5
2-1-Cas réel	6
2-2-Cas général: Z_0 complexe	7
3-Extension au cas matriciel	9
3-1-Algorithmme fini pour le calcul de $S = \text{signe}(A)$	11
3-2-Algorithmme de Newton accéléré.	12
3-3-Implémentation	15
3-3-1-Choix de la norme.	15
3-3-2-Critère d'arrêt.	16
<u>CHAPITRE II/ APPLICATION A LA RESOLUTION DES EQUATIONS DE RICCATI</u>	
Introduction	18
<u>Partie A/ Equations de Riccati dans les problèmes de</u>	
commande optimale avec critère quadratique.	19
1-Introduction.	19
2-Présentation du problème.	19
3-Minimisation de l'hamiltonien.	19
4-Equations canoniques réduites.	20
5-Détermination de P	21

TABLE DES MATIERES

<u>INTRODUCTION/</u>	1
<u>CHAPITRE I/ ETUDE DE LA FONCTION SIGNE D'UNE MATRICE</u>	3
Introduction	3
1-Introduction à la fonction signe d'une matrice.	4
2-Construction d'une suite Z_k telle que $Z_\infty = \text{signe}(Z_0)$..	5
2-1-Cas réel	6
2-2-Cas général: Z_0 complexe	7
3-Extension au cas matriciel	9
3-1-Algorithmme fini pour le calcul de $S = \text{signe}(A)$	11
3-2-Algorithmme de Newton accéléré.	12
3-3-Implémentation	15
3-3-1-Choix de la norme.	15
3-3-2-Critère d'arrêt.	16
<u>CHAPITRE II/ APPLICATION A LA RESOLUTION DES EQUATIONS DE RICCATI</u>	
Introduction	18
<u>Partie A/ Equations de Riccati dans les problèmes de</u>	
commande optimale avec critère quadratique.	19
1-Introduction.	19
2-Présentation du problème.	19
3-Minimisation de l'hamiltonien.	19
4-Equations canoniques réduites.	20
5-Détermination de P	21

Partie B/ Résolution de l'équation de Riccati dans

le cas continu 22

1-Introduction 22

2-Définition du problème 22

3-Résolution 23

4-Coût calcul. 25

5-Comparaison avec l'approche d'Anderson . . . 26

6-Organigramme

Partie C/ Résolution de l'équation de Riccati dans

le cas discret 28

1-Introduction 28

2-Définition du problème 28

3-Produit simple et produit "étoile" de
matrices 30

4-Construction de l'isomorphisme Ψ 34

5-Application à l'équation de Riccati discrète 36

6-Comparaisons avec d'autres méthodes. . . . 38

6-1-: Méthodes à caractère implicite 38

6-1-1-Méthode de Vaughan 38

6-1-2-Méthode de Laub ; 39

6-2-Algorithmes itératifs 40

6-2-1-Méthode du produit étoile. 40

6-2-2-Méthode de He wer 43

6-2-3-Méthode classique stabilisée 43

6-2-4-Square root d'ordre 1 44

7-Conclusions 45

8-Organigramme

Partie B/ Résolution de l'équation de Riccati dans

le cas continu 22

1-Introduction 22

2-Définition du problème 22

3-Résolution 23

4-Coût calcul 25

5-Comparaison avec l'approche d'Anderson . . . 26

6-Organigramme

Partie C/ Résolution de l'équation de Riccati dans

le cas discret 28

1-Introduction 28

2-Définition du problème 28

3-Produit simple et produit "étoile" de
matrices 30

4-Construction de l'isomorphisme Ψ 34

5-Application à l'équation de Riccati discrète 36

6-Comparaisons avec d'autres méthodes 38

6-1-: Méthodes à caractère implicite 38

6-1-1-Méthode de Vaughan 38

6-1-2-Méthode de Laub ; 39

6-2-Algorithmes itératifs 40

6-2-1-Méthode du produit étoile 40

6-2-2-Méthode de He wer 43

6-2-3-Méthode classique stabilisée 43

6-2-4-Square root d'ordre 1 44

7-Conclusions 45

8-Organigramme

CHAPITRE III/ APPLICATIONS

1-Description des programmes 48
1-1-Objet des programmes 48
1-2-Algorithmes 48
1-3-Définition des arguments 50
1-4-Structure du programme: 51
2-Applications numériques : 52
-cas continu 52
-cas discret 55
3-Applications à la stabilité des systèmes 58
3-1-Rappel sur la deuxième méthode de Lyapunov . . . 58
3-2-Etablissement de l'équation de Lyapunov . . . 60
3-3-Applications numériques. 60

CONCLUSION/63

ANNEXES/

BIBLIOGRAPHIE/

CHAPITRE III/ APPLICATIONS

1-Description des programmes 48

1-1-Objet des programmes 48

1-2-Algorithmes 48

1-3-Définition des arguments 50

1-4-Structure du programme: 51

2-Applications numériques : 52

 -cas continu 52

 -cas discret 55

3-Applications à la stabilité des systèmes 58

 3-1-Rappel sur la deuxième méthode de Lyapunov . . . 58

 3-2-Etablissement de l'équation de Lyapunov . . . 60

 3-3-Applications numériques 60

CONCLUSION/ 63

ANNEXES/

BIBLIOGRAPHIE/

INTRODUCTION

Ce travail est consacré à l'application d'un concept, la fonction signe d'une matrice, à l'élaboration d'algorithmes résolvant numériquement les équations de Riccati et de Lyapunov que l'on rencontre dans les problèmes de commande optimale et de stabilité des systèmes.

Ce problème peut paraître démodé à beaucoup de spécialistes. Il ya plus de quinze ans qu'il a été formalisé et que les premières solutions numériques ont été présentées. En fait les solutions connues et largement enseignées jusqu'à une date récente sont inefficaces dans au moins deux situations:

- lorsque la dimension du vecteur d'état est importante.
- lorsque le système "linéaire" est instable (ou voisin de l'instabilité).

L'évolution actuelle des applications de l'automatique conduit souvent à la résolution de problèmes qui peuvent notamment dépasser le cadre de la commande de systèmes intrinsèquement stables.

A ce renouvellement récent des domaines d'applications doit correspondre un renouveau des méthodes auxquelles notre travail contribue. Il présente en effet une méthode efficace de calcul direct des paramètres stationnaires optimaux d'un feedback de commande dans le cas de systèmes non nécessairement stables et de dimension quelconque.

On attire l'attention du lecteur que ces nouvelles méthodes ne sont pas le fruit du hasard, mais sont liées étroitement au problème de la réduction des calculs (temps machine) et la possibilité de commander en temps réel un processus (installations industrielles et Problème de la poursuite d'une cible) et ceci grâce à la possibilité de disposer de calculateurs puissants et rapides.

Ce concept (fonction signe) et surtout son application en automatique résulte de la convergence d'idées, d'études de plusieurs auteurs entre autres "Casti en 1977" [3] "Anderson" [4] au congrès IEEE 1977 et 1978

INTRODUCTION

Ce travail est consacré à l'application d'un concept, la fonction signe d'une matrice, à l'élaboration d'algorithmes résolvant numériquement les équations de Riccati et de Lyapunov que l'on rencontre dans les problèmes de commande optimale et de stabilité des systèmes.

Ce problème peut paraître démodé à beaucoup de spécialistes. Il ya plus de quinze ans qu'il a été formalisé et que les premières solutions numériques ont été présentées. En fait les solutions connues et largement enseignées jusqu'à une date récente sont inefficaces dans au moins deux situations:

- lorsque la dimension du vecteur d'état est importante.
- lorsque le système "linéaire" est instable (ou voisin de l'instabilité).

L'évolution actuelle des applications de l'automatique conduit souvent à la résolution de problèmes qui peuvent notamment dépasser le cadre de la commande de systèmes intrinsèquement stables.

A ce renouvellement recett des domaines d'applications doit correspondre un renouveau des méthodes auxquelles notre travail contribue. Il présente en effet une méthode efficace de calcul direct des paramètres stationnaires optimaux d'un feedback de commande dans le cas de systèmes non nécessairement stables et de dimension quelconque.

On attire l'attention du lecteur que ces nouvelles méthodes ne sont pas le fruit du hasard, mais sont liées étroitement au problème de la réduction des calculs (temps machine) et la possibilité de commander en temps réel un processus (installations industrielles et Probleme de la pour suite d'une cible) et ceci grâce à la possibilité de disposer de calculateurs puissants et rapides.

Ce concept (fonction signe) et surtout son application en automatique résulte de la convergence d'idées, d'études de plusieurs auteurs entre autres "Casti en 1977" [3] "Anderson" [4] au congrès IEEE 1977 et 1978

et les travaux de "Beavers et Denman" en 1976; ainsi que ceux de Baraud.

Par ce travail nous nous proposons dans un premier chapitre d'étudier le concept de fonction signe de matrice. Le deuxième chapitre est relatif à l'application de ce concept à la résolution des équations de Riccati et de Lyapunov.

Quand au dernier chapitre, il fait l'objet de diverses applications et de tests numériques avec interprétations.

et les travaux de "Beavers et Denman" en 1976; ainsi que ceux de Baraud.

Par ce travail nous nous proposons dans un premier chapitre d'étudier le concept de fonction signe de matrice. Le deuxième chapitre est relatif à l'application de ce concept à la résolution des équations de Riccati et de Lyapunov.

Quand au dernier chapitre, il fait l'objet de diverses applications et de tests numériques avec interprétations.

CHAPITRE I/

ETUDE DE LA FONCTION SIGNE D'UNE MATRICE

CHAPITRE I/

ETUDE DE LA FONCTION SIGNE D'UNE MATRICE

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

Nous allons introduire dans ce chapitre, le concept de fonction signe d'une matrice et nous montrerons que l'algorithme de NEWTON permettant de calculer cette fonction, présente d'abord une convergence linéaire à évolution chaotique avant d'aborder sa phase finale du second ordre. De cette analyse découlent deux nouveaux algorithmes, l'un fini pour les matrices à spectres réels, l'autre relatif au cas général constituant une méthode de NEWTON accélérée.

Quand à l'application de cette technique, elle fera l'objet des parties qui suivent.

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

Nous allons introduire dans ce chapitre, le concept de fonction signe d'une matrice et nous montrerons que l'algorithme de NEWTON permettant de calculer cette fonction, présente d'abord une convergence linéaire à évolution chaotique avant d'aborder sa phase finale du second ordre. De cette analyse découlent deux nouveaux algorithmes, l'un fini pour les matrices à spectres réels, l'autre relatif au cas général constituant une méthode de NEWTON accélérée.

Quant à l'application de cette technique, elle fera l'objet des parties qui suivent.

I-INTRODUCTION A LA FONCTION SIGNE D'UNE MATRICE

On sait par definition qu'une matrice carrée $A(n.n)$ est semblable à sa forme de Jordan:

$$A = MJM^{-1} \tag{1}$$

où M est la matrice des vecteurs propres de A et la matrice J a la structure:

$$\begin{bmatrix} J & \dots & \dots & \dots & 0 \\ & & & & \\ & & & & \\ & & & & \\ 0 & & & & J \end{bmatrix} \quad \text{avec} \quad J_j = \begin{bmatrix} \lambda_i & 1 & & 0 \\ & \lambda_i & 1 & \\ & & \lambda_i & \dots \\ & & & \dots \\ 0 & & & & \lambda_i \end{bmatrix}$$

où le nombre de blocs de Jordan J_j associés à la valeur propre λ_i est égal au nombre de vecteurs propres lineairement independants relatif à λ_i . Soit maintenant $f(\lambda)$ d'une fonction definie sur le spectre de A . On aura

$$f(A) = Mf(J)M^{-1}$$

avec

$$f(J) = \begin{bmatrix} f(J_1) & & & \\ & 0 & & \\ & & \ddots & \\ & & & f(J_k) \\ & & & & 0 \end{bmatrix}; \quad f(J_j) = \begin{bmatrix} f(\lambda_i) & f'(\lambda_i) & \dots & \frac{f(\lambda_i)^{(n_j-1)}}{(n_j-1)!} \\ & \vdots & & \\ & & f'(\lambda_i) & \\ & & & f(\lambda_i) \end{bmatrix}$$

On définit le signe d'un nombre complexe $Z = x + iy$ comme étant

$$\text{signe}(Z) = \frac{x}{|x|}$$

Appliquons maintenant la fonction signe à une matrice; on aura

$$\text{signe}(A) = M \cdot \text{signe}(J) \cdot M^{-1} = S \tag{2}$$

avec

$$\text{signe}(J) = \begin{bmatrix} \text{signe}(J_1) & & & \\ & 0 & & \\ & & \ddots & \\ & & & \text{signe}(J_k) \\ & & & & 0 \end{bmatrix}$$

I-INTRODUCTION A LA FONCTION SIGNE D'UNE MATRICE

On sait par definition qu'une matrice carrée $A(n.n)$ est semblable à sa forme de Jordan:

$$A = MJM^{-1} \quad (1)$$

où M est la matrice des vecteurs propres de A et la matrice J a la structure:

$$\begin{bmatrix} J & \dots & \dots & \dots & 0 \\ & & & & \\ & & & & \\ & & & & \\ 0 & & & & J \end{bmatrix} \quad \text{avec} \quad J_j = \begin{bmatrix} \lambda_i & 1 & & 0 \\ & \lambda_i & 1 & \\ & & \lambda_i & \dots \\ & & & \lambda_i \end{bmatrix}$$

où le nombre de blocs de Jordan J_j associés à la valeur propre λ_i est égal au nombre de vecteurs propres lineairement independants relatif à λ_i . Soit maintenant $f(\lambda)$ d'une fonction definie sur le spectre de A . On aura

$$f(A) = Mf(J)M^{-1}$$

avec

$$f(J) = \begin{bmatrix} f(J_1) & & & \\ & 0 & & \\ & & \ddots & \\ & & & f(J_k) \end{bmatrix}; \quad f(J_j) = \begin{bmatrix} f(\lambda_i) & f'(\lambda_i) & \dots & \frac{f(\lambda_i)^{(n_j-1)}}{(n_j-1)!} \\ & & & \vdots \\ & & & f'(\lambda_i) \\ & & & f(\lambda_i) \end{bmatrix}$$

On définit le signe d'un nombre complexe $Z = x + iy$ comme étant

$$\text{signe}(Z) = \frac{x}{|x|}$$

Appliquons maintenant la fonction signe à une matrice; on aura

$$\text{signe}(A) = M \cdot \text{signe}(J) \cdot M^{-1} = S \quad (2)$$

avec

$$\text{signe}(J) = \begin{bmatrix} \text{signe}(J_1) & & & \\ & 0 & & \\ & & \ddots & \\ & & & \text{signe}(J_k) \end{bmatrix}$$

et

$$\text{signe}(J_j) = \begin{bmatrix} \text{signe}(\lambda_i) & & \\ & \theta & \\ & 0 & \text{signe}(\lambda_i) \end{bmatrix} = \pm I$$

$\pm I$ si $\text{réel}(\lambda_i) > 0$ $-I$ si $\text{réel}(\lambda_i) < 0$

Resumé et propriétés

$A = MJM^{-1}$

$S = \text{signe}(A)$

$S = MDM$ $D = \text{diag}(\dots \text{signe}(\lambda_i) \dots)$

$S = I$ si $\text{réel}(\lambda_i) > 0 \forall i$

$S = -I$ si $\text{réel}(\lambda_i) < 0 \forall i$

$S = A$ si A orthogonale symétrique

$\text{Signe}(A^{-1}) = \text{Signe}(A)$

$\text{Signe}(\alpha A) = \text{Signe}(A) \quad \forall \alpha \in \mathbb{R}^+ *$

$\text{Signe}(A^T) = S^T$

$\text{Signe}(B) = VS^{-1} \quad \text{si } B = VAV^{-1}$

2-CONSTRUCTION D'UNE SUITE Z_k TELLE QUE $Z_k = \text{Signe}(Z_0)$

On a vu ultérieurement que pour connaître le signe d'une matrice, il suffisait de connaître les signes de ses valeurs propres. Donc il nous a paru utile de voir dans ce paragraphe comment trouver le signe d'un nombre. On démontrera ici que pour trouver le signe d'un nombre (complexe ou réel), on peut utiliser une suite qui à l'infini, nous donne son signe.

Une telle suite peut être introduite en cherchant à déterminer les zéros de la fonction complexe

$$f(z) = z^2 - I$$

par la méthode de Newton. Ce qui donne l'algorithme

$$Z_{k+1} = Z_k - \frac{f(Z_k)}{f'(Z_k)}$$

c'est à dire (annexe):

$$Z_{k+1} = \frac{1}{2} \left(Z_k + \frac{I}{Z_k} \right)$$

(3)

et

$$\text{signe}(J_j) = \begin{bmatrix} \text{signe}(\lambda_i) & & \\ & \theta & \\ & & \text{signe}(\lambda_i) \end{bmatrix} = \pm I$$

$\begin{cases} +I & \text{si } \text{réel}(\lambda_i) > 0 \\ -I & \text{si } \text{réel}(\lambda_i) < 0 \end{cases}$

Resumé et propriétés

$A = M J M^{-1}$

$S = \text{signe}(A)$

$S = M D M^{-1}$ $D = \text{diag}(\dots \text{signe}(\lambda_i) \dots)$

$S = I$ si $\text{réel}(\lambda_i) > 0 \forall i$

$S = -I$ si $\text{réel}(\lambda_i) < 0 \forall i$

$S = A$ si A orthogonale symétrique

$\text{Signe}(A^{-1}) = \text{Signe}(A)$

$\text{Signe}(\alpha A) = \text{Signe}(A) \quad \forall \alpha \in \mathbb{R}^+ *$

$\text{Signe}(A^T) = S^T$

$\text{Signe}(B) = V S V^{-1}$ si $B = V A V^{-1}$

2-CONSTRUCTION D'UNE SUITE Z_k TELLE QUE $Z_k = \text{Signe}(Z_0)$

On a vu ultérieurement que pour connaître le signe d'une matrice, il suffisait de connaître les signes de ses valeurs propres. Donc il nous a paru utile de voir dans ce paragraphe comment trouver le signe d'un nombre. On démontrera ici que pour trouver le signe d'un nombre (complexe ou réel), on peut utiliser une suite qui à l'infini, nous donne son signe.

Une telle suite peut être introduite en cherchant à déterminer les zéros de la fonction complexe

$$f(z) = z^2 - I$$

par la méthode de Newton. Ce qui donne l'algorithme

$$Z_{k+1} = Z_k - \frac{f(Z_k)}{f'(Z_k)}$$

c'est à dire (annexe):

$$Z_{k+1} = \frac{1}{2} \left(Z_k + \frac{I}{Z_k} \right)$$

(3)

2-I- CAS REEL

Etudions maintenant plus en détail l'évolution de Z_K . Pour cela nous allons supposer dans un premier cas que Z_0 est réel ($Z_0 = x_0$).

Notre suite devient:

$$x_{K+1} = \frac{I}{2} \left(x_K + \frac{I}{x_K} \right) \quad (4)$$

on voit clairement que:

$$\text{signe}(x_K) = \text{signe}(x_0) \quad \forall K \quad (5)$$

On peut donc supposer pour la suite de notre raisonnement que $x_0 > 0$ et par raison de symétrie de x et I/x par rapport à I , il suffit de considérer

$$x_0 > I$$

d'où l'on peut écrire

$$I < x_{K+1} < x_K \dots \dots \dots < x_0 \quad (6)$$

Si on prend $x_K = I + \epsilon \implies x_{K+1} = I + \frac{I}{2} \cdot \frac{\epsilon^2}{1 + \epsilon}$ (7)

c'est à dire

$$x_{K+1} = I + \theta \left(\frac{\epsilon^2}{2} \right)$$

resultat qui caractérise la convergence d'ordre deux de l'algorithme de Newton.

D'après (4) et (5) on peut écrire

$$\frac{x_K}{2} < x_{K+1}$$

et sachant que :

$$x_K > I \implies I/x_K < I$$

donc

$$I/2(x_K + I) > I/2(x_K + I/x_K)$$

d'où

$$I/2(x_K + I) > x_{K+1}$$

finalement on aura la relation:

$$x_K/2 < x_{K+1} < I/2(I + x_K) \quad (8)$$

En prenant $K+1=n$ dans l'inégalité (5), on aura:

$$x_0/2^n < x_n < (I - I/2^n) + x_0/2^n \quad (9)$$

2-I- CAS REEL

Etudions maintenant plus en détail l'évolution de Z_K . Pour cela nous allons supposer dans un premier cas que Z_0 est réel ($Z_0 = x_0$).

Notre suite devient:

$$x_{K+1} = \frac{I}{2} \left(x_K + \frac{I}{x_K} \right) \quad (4)$$

on voit clairement que:

$$\text{signe}(x_K) = \text{signe}(x_0) \quad \forall K \quad (5)$$

On peut donc supposer pour la suite de notre raisonnement que $x_0 > 0$ et par raison de symétrie de x et I/x par rapport à I , il suffit de considérer $x_0 > I$

d'où l'on peut écrire

$$I < x_{n+1} < x_K \dots \dots \dots < x_0 \quad (6)$$

Si on prend $x_K = I + \epsilon \implies x_{K+1} = I + \frac{I}{2} \cdot \frac{\epsilon^2}{1 + \epsilon}$ (7)

c'est à dire

$$x_{K+1} = I + \theta \left(\frac{\epsilon^2}{2} \right)$$

resultat qui caractérise la convergence d'ordre deux de l'algorithme de Newton.

D'après (4) et (5) on peut écrire

$$\frac{x_K}{2} < x_{K+1}$$

et sachant que :

$$x_n > I \implies I/x < I$$

donc

$$I/2(x_K + I) > I/2(x_K + I/x_K)$$

d'où

$$I/2(x_K + I) > x_{K+1}$$

finalement on aura la relation:

$$x_K/2 < x_{K+1} < I/2(I + x_K) \quad (8)$$

En prenant $K+1=n$ dans l'inégalité (5), on aura:

$$x_0/2^n < x_n < (I - I/2^n) + x_0/2^n \quad (9)$$

Ainsi si X_0 est de l'ordre de 2^n , il nous faut donc n itérations pour le ramener à une valeur comprise entre 1 et 2; ce qui traduit une convergence linéaire dans un premier temps, après cela la convergence d'ordre 2 intervient.

A l'aide de la relation (6), pour $K > n+1$, on peut connaître le nombre d'itérations à effectuer pour avoir une précision donnée.

En posant $X_{n+p} = 1 + \epsilon_p$ $0 < \epsilon_0 < 1$

on aura:
$$\epsilon_p = \frac{1}{2} \frac{\epsilon_{p-1}^2}{1 + \epsilon_{p-1}}$$

En prenant $\epsilon_0 = 1$, qui correspond au cas le plus défavorable, on pourra dresser le tableau suivant:

ϵ_0	ϵ_1	ϵ_2	ϵ_3	ϵ_4	ϵ_5	ϵ_6	ϵ_7
1	0,25	0,025	$3 \cdot 10^{-4}$	$4,6 \cdot 10^{-8}$	$1,1 \cdot 10^{-15}$	$5,8 \cdot 10^{-31}$	$1,7 \cdot 10^{-61}$

Ainsi le nombre total d'itérés pour obtenir une précision machine vaut

$$N = n + p \tag{10}$$

2-2-CAS GENERAL: Z_0 COMPLEXE

Sans restreindre la généralité du raisonnement, on peut supposer que

$\rho_0 > 1$ et $0 < \theta < \frac{\pi}{2}$
 Par définition on a $Z_0 = \rho_0 (\cos \theta_0 + i \sin \theta_0)$

d'où $Z_{K+1} = \rho_{K+1} (\cos \theta_{K+1} + i \sin \theta_{K+1})$

Ainsi en utilisant (3), on obtient

$$\rho_{K+1}^2 = 1/4 (\rho_{K-1}/\rho_K)^2 + \cos^2 \theta_K = 1/4 (\rho_{K+1}/\rho_K)^2 - \sin^2 \theta_K \tag{11}$$

et

$$\operatorname{tg} \theta_{K+1} = \frac{\rho_K^2 - 1}{\rho_K^2 + 1} \operatorname{tg} \theta_K \tag{12}$$

Etant donné que

Ainsi si X_0 est de l'ordre de 2^n , il nous faut donc n itérations pour le ramener à une valeur comprise entre 1 et 2; ce qui traduit une convergence linéaire dans un premier temps, après cela la convergence d'ordre 2 intervient.

A l'aide de la relation (6), pour $K > n+1$, on peut connaître le nombre d'itérations à effectuer pour avoir une précision donnée.

En posant $X_{n+p} = 1 + \epsilon_p$ $0 < \epsilon_0 < 1$

on aura:
$$\epsilon_p = \frac{1}{2} \frac{\epsilon_{p-1}^2}{1 + \epsilon_{p-1}}$$

En prenant $\epsilon_0 = 1$, qui correspond au cas le plus défavorable, on pourra dresser le tableau suivant:

ϵ_0	ϵ_1	ϵ_2	ϵ_3	ϵ_4	ϵ_5	ϵ_6	ϵ_7
1	0,25	0,025	$3 \cdot 10^{-4}$	$4,6 \cdot 10^{-8}$	$1,1 \cdot 10^{-15}$	$5,8 \cdot 10^{-31}$	$1,7 \cdot 10^{-61}$

Ainsi le nombre total d'itérés pour obtenir une précision machine vaut

$$N = n + p \tag{10}$$

2-2-CAS GENERAL: Z_0 COMPLEXE

Sans restreindre la généralité du raisonnement, on peut supposer que

$\rho_0 > 1$ et $0 < \theta < \frac{\pi}{2}$
 Par définition on a $Z_0 = \rho_0 (\cos \theta_0 + i \sin \theta_0)$

d'où $Z_{K+1} = \rho_{K+1} (\cos \theta_{K+1} + i \sin \theta_{K+1})$

Ainsi en utilisant (3), on obtient

$$\rho_{K+1}^2 = 1/4 (\rho_{K-1}/\rho_K)^2 + \cos^2 \theta_K = 1/4 (\rho_{K+1}/\rho_K)^2 - \sin^2 \theta_K \tag{11}$$

et

$$\operatorname{tg} \theta_{K+1} = \frac{\rho_K^2 - 1}{\rho_K^2 + 1} \operatorname{tg} \theta_K \tag{12}$$

Etant donné que

$$c \left\langle \left| \frac{\rho_k^2 - 1}{\rho_k^2 + 1} \right| \right\rangle < I \quad \Rightarrow \quad |tg \epsilon_{K+1}| < |tg \epsilon_K|$$

on a ainsi

$$c \left\langle |\epsilon_{K+1}| \right\rangle < |\epsilon_K| \quad (13)$$

ce qui traduit la convergence de ϵ_K vers zero.

Etudions maintenant la décroissance de ρ en supposant $\rho \gg I$ et si

$$\rho_K > 1$$

d'après (II) on peut écrire

$$\frac{1}{4} \rho_K^2 - \frac{1}{2} \left\langle \rho_{K+1}^2 \right\rangle < \frac{1}{4} \rho_K^2 + \frac{3}{4} \left\langle \rho_K^2 \right\rangle \quad (14)$$

Si on prend $K+1=n$ et en opérant comme pour le cas réel, on aboutit à

$$\frac{1}{4} \rho_0^2 - \frac{2}{3} (1-1/4^n) \left\langle \rho_n^2 \right\rangle < \frac{1}{4} \rho_0^2 + (1-1/4^n) \quad (15)$$

et si à tous les itérés on a

$$1 \left\langle \rho_{n-1} \right\rangle < \rho_{n-2} < \dots < \rho_1 < \rho_0 \quad (16)$$

et dans l'hypothèse où

$$2^n \leq \rho_0 < 2^{n+1} \quad \rho_0^2 = \alpha 4^n \quad 1 \leq \alpha < 4$$

on pourra donc écrire

$$\sqrt{\frac{Ic}{3}} \left\langle \rho_n \right\rangle < \sqrt{I7} \quad \text{et} \quad \sqrt{\frac{I}{3}} \left\langle \rho_n \right\rangle < \sqrt{5} \quad (17)$$

ce qui traduit une convergence lineaire de facteur $\frac{I}{2}$ pour ρ .

En fait le module de Z_K decroit légèrement plus vite qu'une progression geometrique de raison $\frac{I}{2}$ lorsque ϵ est grand (ϵ voisin de $\frac{\pi}{2}$) et un peu moins vite lorsqu'on se rapproche de l'axe des réels.

Maintenant voyons la décroissance monotone de ϵ_K son évolution découle de de celle du gain $f(\rho) = \frac{\rho^2 - 1}{\rho^2 + 1}$ et de $tg \epsilon$, qui à priori ne relève d'aucun processus simple tel celui mis en évidence pour ρ .

On peut cependant déjà conclure pour un Z_0 tel que $\rho_0 \gg 1$ et $tg \epsilon_0 \gg 1$ L'évolution de Z_K se décompose en trois phase :

$$c \left\langle \frac{\rho_k^2 - 1}{\rho_k^2 + 1} \right\rangle < I \quad \Rightarrow \quad |tg \theta_{K+1}| < |tg \theta_K|$$

on a ainsi

$$c \left\langle |\epsilon_{K+1}| \right\rangle < |\epsilon_K| \quad (13)$$

ce qui traduit la convergence de ϵ_K vers zero.

Etudions maintenant la décroissance de ρ en supposant $\rho \gg I$ et si

$$\rho_K > 1$$

d'après (II) on peut écrire

$$\frac{1}{4} \rho_K^2 - \frac{1}{2} \langle \rho_{K+1}^2 \rangle < \frac{1}{4} \rho_K^2 + \frac{3}{4} \langle \rho_K^2 \rangle \quad (14)$$

Si on prend $K+1=n$ et en opérant comme pour le cas réel, on aboutit à

$$\frac{1}{4} \rho_0^2 - \frac{2}{3} (1-1/4^n) \langle \rho_n^2 \rangle < \frac{1}{4} \rho_0^2 + (1-1/4^n) \quad (15)$$

et si à tous les itérés on a

$$1 \langle \rho_{n-1} \rangle \langle \rho_{n-2} \rangle \dots \langle \rho_1 \rangle \langle \rho_0 \rangle \quad (16)$$

et dans l'hypothèse où

$$2^n \ll \rho_0 < 2^{n+1} \quad \rho_0^2 = \alpha 4^n \quad 1 \ll \alpha < 4$$

on pourra donc écrire

$$\sqrt{\frac{I_0}{3}} \langle \rho_n \rangle < \sqrt{I_7} \quad \text{et} \quad \sqrt{\frac{I}{3}} \langle \rho_n \rangle < \sqrt{5} \quad (17)$$

ce qui traduit une convergence lineaire de facteur $\frac{I}{2}$ pour ρ .

En fait le module de Z_K decroit légèrement plus vite qu'une progression geometrique de raison $\frac{I}{2}$ lorsque θ est grand (θ voisin de $\frac{\pi}{2}$) et un peu moins vite lorsqu'on se rapproche de l'axe des réels.

Maintenant voyons la décroissance monotone de θ_K son évolution découle de de celle du gain $f(\rho) = \rho^2 - 1 / \rho^2 + 1$ et de $tg \theta$, qui à priori ne relève d'aucun processus simple tel celui mis en évidence pour ρ .

On peut cependant déjà conclure pour un Z_0 tel que $\rho_0 \gg 1$ et $tg \theta_0 \gg 1$ L'évolution de Z_K se décompose en trois phase :

1^{re} PHASE/ décroissance de ρ (du type progression geometrique de raison $I/2$) à ϵ approximativement constant.

2^{de} PHASE/ décroissance de ϵ , ρ oscillant autour de I .

3^{de} PHASE/ convergence du second ordre simultané de ϵ vers zero et ρ vers I .

3-EXTENSION AU CAS MATRICIEL

Soit A une matrice carrée (n.n) dont la partie réelle des valeurs propres non nulle: .Par analogie au cas scalaire definissons la suite:

$$A_{K+I} + I = I/2(A_K + A_K^{-1}) \quad A_0 = A \quad (I8)$$

Supposons d'abord que A est diagonalisable, c'est à dire qu'on peut écrire:

$$A_{K+I} = M D_{K+I} M^{-1} = I/2 M (D_K + D_K^{-1}) M^{-1} \quad \text{avec} \quad A_0 = M D_0 M^{-1}$$

par identification on aura donc

$$D_{K+I} = I/2 (D_K + D_K^{-1}) = \text{diag}(\dots I/2(\lambda_i + I/\lambda_i) \dots)$$

de sorte qu'on ait

$$D_\infty = \text{diag}(\dots \text{signe}(\lambda_i) \dots)$$

En vertu de ce qu'on a vu au paragraphe precedent ,on peut écrire:

$$\text{signe}(A) = S = M D M^{-1} = \lim_{K \rightarrow \infty} A_K$$

Maintenant supposons que notre matrice ne peut être mise que sous la forme de Jordan. Soit (J) un bloc de Jordan, et définissons la suite:

$$J_{K+I} = I/2 (J_K + J_K^{-1}) \quad \text{avec} \quad J_0 = J$$

1^{re} PHASE/ décroissance de ρ (du type progression geometrique de raison $1/2$) à ϵ approximativement constant.

2^{de} PHASE/ décroissance de ϵ , ρ oscillant autour de 1.

3^{de} PHASE/ convergence du second ordre simultané de ϵ vers zero et ρ vers 1.

3-EXTENSION AU CAS MATRICIEL

Soit A une matrice carrée (n.n) dont la partie réelle des valeurs propres non nulle. Par analogie au cas scalaire définissons la suite:

$$A_{K+1} + I = 1/2(A_K + A_K^{-1}) \quad A_0 = A \quad (I8)$$

Supposons d'abord que A est diagonalisable, c'est à dire qu'on peut écrire:

$$A_{K+1} = M D_{K+1} M^{-1} = 1/2 M (D_K + D_K^{-1}) M^{-1} \quad \text{avec} \quad A_0 = M D_0 M^{-1}$$

par identification on aura donc

$$D_{K+1} = 1/2 (D_K + D_K^{-1}) = \text{diag}(\dots I/2(\lambda_i + 1/\lambda_i) \dots)$$

de sorte qu'on ait

$$D_{\infty} = \text{diag}(\dots \text{signe}(\lambda_i) \dots)$$

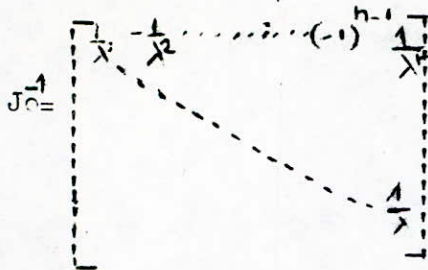
En vertu de ce qu'on a vu au paragraphe précédent, on peut écrire:

$$\text{signe}(A) = S = M D_{\infty} M^{-1} = \lim_{K \rightarrow \infty} A_K$$

Maintenant supposons que notre matrice ne peut être mise que sous la forme de Jordan. Soit (J) un bloc de Jordan, et définissons la suite:

$$J_{K+1} = 1/2 (J_K + J_K^{-1}) \quad \text{avec} \quad J_0 = J$$

d'où



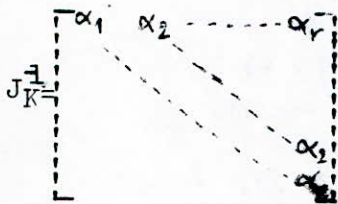
or

$$J_I = I/2(J_0 + J_0^{-1})$$

par identification on aura:

$$a_1^I = I/2(\lambda + I/\lambda) \quad a_2^I = I/2(I - I/\lambda^2) \quad a_i^I = (-I)^{i-1} \cdot I/2\lambda^2 \quad i=3, \dots, r$$

Posons



puisque $J_K J_K^{-1} = I$

$$\Rightarrow \alpha_1 = I/a_1 \quad \alpha_2 = -(I/a_1)a_2\alpha_1; \quad \alpha_r = -I/a_1 \cdot (a_2\alpha_{r-1} + \dots + a_r\alpha_1)$$

Sachant que $J_{K+I} = I/2(J_K + J_K^{-1})$, il en découle:

$$a_1^{K+1} = I/2(a_1^K + I/a_1^K)$$

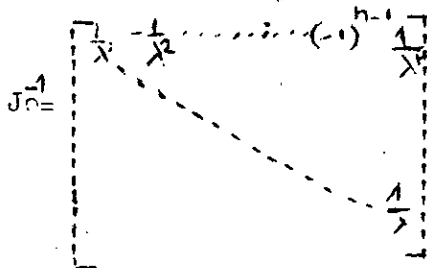
$$a_2^{K+1} = I/2 \cdot a_2^K (I - I/(a_1^K)^2)$$

$$a_r^{K+1} = I/2 \cdot a_r^K \left[I - I/(a_1^K)^2 \right] - (I/2a_1^K) (a_2^K \alpha_{r-1}^K + \dots + a_r^K \alpha_1^K)$$

Cette structure montre clairement que lorsque a_1^K tend vers $\pm I$, alors a_i^K , $i=2, \dots, r$ tendent vers zero, ce qui donne:

$$J_{00} = \begin{bmatrix} \text{signe}(\lambda) & & & \\ & 0 & & \\ & & 0 & \\ & & & \text{signe}(\lambda) \end{bmatrix} = \text{signe}(J)$$

d'où



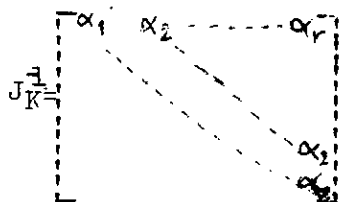
or

$$J_I = I/2(J_0 + J_{2I}^{-1})$$

par identification on aura:

$$a_1^I = I/2(\lambda + I/\lambda) \quad a_2^I = I/2(I - I/\lambda^2) \quad a_i^I = (-I)^{i-1} \cdot I/2\lambda^2 \quad i=3, \dots, r$$

Posons



puisque $J_K J_K^{-1} = I$

$$\Rightarrow \alpha_1 = I/a_1 \quad \alpha_2 = -(I/a_1)a_2\alpha_1; \quad \alpha_r = -I/a_1 \cdot (a_2\alpha_{r-1} + \dots + a_r\alpha_1)$$

Sachant que $J_{K+I} = I/2(J_K + J_K^{-1})$, il en découle:

$$a_1^{K+1} = I/2(a_1^K + I/a_1^K)$$

$$a_2^{K+1} = I/2 \cdot a_2^K (I - I/(a_1^K)^2)$$

$$a_r^{K+1} = I/2 \cdot a_r^K \left[I - I/(a_1^K)^2 \right] - (I/2a_1^K) (a_2^K \alpha_{r-1}^K + \dots + a_r^K \alpha_1^K)$$

Cette structure montre clairement que lorsque a_1^K tend vers $\pm I$, alors a_i^K , $i=2, \dots, r$ tendent vers zero, ce qui donne:

$$J_{00} = \begin{bmatrix} \text{signe}(\lambda) & & & \\ & 0 & & \\ & & 0 & \\ & & & \text{signe}(\lambda) \end{bmatrix} = \text{signe}(J)$$

3-I-ALGORITHME FINI POUR LE CALCUL DE $S = \text{signe}(A)$

Nous allons considérer dans ce paragraphe un cas particulier à savoir les matrices régulières à spectre réel ; on aura donc la suite :

$$\lambda_{K+1} = \frac{I}{2} \left(\lambda_K + \frac{I}{\lambda_K} \right) \quad \text{avec } \lambda_0 = \max(|\lambda_{\max}|, |\frac{1}{\lambda_{\min}}|)$$

et le nombre d'itérés pour calculer S sera donné comme pour le cas réel c'est à dire :

$$N = n + r$$

Soit

$$s = |\lambda_{\max}| \quad t = |\lambda_{\min}| \quad \text{et } \alpha = \frac{I}{\sqrt{st}}$$

définissons la transformation

$$A^* = \alpha A$$

de façon à avoir

$$s \left(\frac{s}{t} \right)^{1/2} ; \quad t \left(\frac{t}{s} \right)^{1/2} \implies s = \frac{I}{t^*} \quad (I9)$$

Appliquons alors une itération de (I8) pour $K=0$, il vient :

$$A_1^* = \frac{I}{2} (A_0^* + (A_0^*)^{-1})$$

d'après (I9), A_0^* a deux valeurs propres de module inverse, alors A_1^* aura deux valeurs propres de même module maximum :

$$\lambda_{\max}(A_1^*) = \frac{I}{2} (s^* + t^*) = \frac{I}{2} \alpha (s + t)$$

Plus généralement nous avons :

$$s_K = |\lambda_{\max}(A_K^*)|, \quad t_K = |\lambda_{\min}(A_K^*)|, \quad \alpha_K = \frac{I}{\sqrt{s_K t_K}}$$

et (20)

$$A_{K+1}^* = \frac{I}{2} \left[\alpha_K A_K^* + \frac{I}{\alpha_K} (A_K^*)^{-1} \right] \quad A_0 = A_0 = A$$

En supposant que A_K possède au moins $(K+1)$ valeurs propres de même module maximum. L'axiome de récurrence permet d'affirmer, qu'en faisant une itération (c-à-d A_{K+1}), celle-ci conduit à au moins $(K+2)$ valeurs propres de même module maximum :

3-I-ALGORITHME FINI POUR LE CALCUL DE $S = \text{signe}(A)$

Nous allons considérer dans ce paragraphe un cas particulier à savoir les matrices régulières à spectre réel ; on aura donc la suite :

$$\lambda_{K+1} = \frac{I}{2} \left(\lambda_K + \frac{I}{\lambda_K} \right) \quad \text{avec } \lambda_0 = \max(|\lambda_{\max}|, |\frac{1}{\lambda_{\min}}|)$$

et le nombre d'itérés pour calculer S sera donné comme pour le cas réel c'est à dire :

$$N = n + r$$

Soit

$$s = |\lambda_{\max}| \quad t = |\lambda_{\min}| \quad \text{et } \alpha = \frac{I}{\sqrt{st}}$$

définissons la transformation

$$A^* = \alpha A$$

de façon à avoir

$$s \left(\frac{s}{t} \right)^{1/2} ; \quad t \left(\frac{t}{s} \right)^{1/2} \implies s = \frac{I}{t^*} \quad (19)$$

Appliquons alors une itération de (18) pour $K=0$, il vient :

$$A_1^* = \frac{I}{2} (A_0^* + (A_0^*)^{-1})$$

d'après (19), A_0^* a deux valeurs propres de module inverse, alors A_1^* aura deux valeurs propres de même module maximum :

$$\lambda_{\max}(A_1^*) = \frac{I}{2} (s^* + t^*) = \frac{I}{2} \alpha (s + t)$$

Plus généralement nous avons :

$$s_K = |\lambda_{\max}(A_K^*)|, \quad t_K = |\lambda_{\min}(A_K^*)|, \quad \alpha_K = \frac{I}{\sqrt{s_K t_K}} \quad (20)$$

et

$$A_{K+1}^* = \frac{I}{2} \left[\alpha_K A_K^* + \frac{I}{\alpha_K} (A_K^*)^{-1} \right] \quad A_0 = A_0^* = A$$

En supposant que A_K possède au moins $(K+1)$ valeurs propres de même module maximum. L'axiome de récurrence permet d'affirmer, qu'en faisant une itération (c-à-d A_{K+1}^*), celle-ci conduit à au moins $(K+2)$ valeurs propres de même module maximum :

$$|\lambda_{\max}(A_{K+I}^*)| = \frac{1}{2} \frac{|\lambda_{\max}(A_K^*)| + |\lambda_{\min}(A_K^*)|}{\sqrt{|\lambda_{\max}(A_K^*) \cdot \lambda_{\min}(A_K^*)|}}$$

on peut donc affirmer qu'il existe $1 \leq n-1$ tel que A_1^* possède n valeurs propres de même module.

On peut alors écrire:

$$S = \text{signe}(A) = \frac{A_1^*}{|\lambda(A_1^*)|} \quad (21)$$

Ainsi il est possible d'énoncer le théorème suivant:

Soit A (n,n) régulière à spectre réel, $S = \text{signe}(A)$ est obtenue par l'algorithme fini:

$$\begin{aligned} A_{K+I}^* &= \frac{1}{2} \left[\alpha_K A_K^* \frac{1}{\alpha_K} (A_K^*)^{-1} \right] & A_0 &= A \\ \alpha_K &= \frac{1}{\sqrt{|\lambda_{\max}(A_K^*) \cdot \lambda_{\min}(A_K^*)|}} \\ S &= \frac{1}{|\lambda(A_1^*)|} \cdot A_1^* \end{aligned} \quad (22)$$

3-2-ALGORITHME DE NEWTON ACCELERE

Introduction

La méthode qu'on a vue précédemment n'est malheureusement utilisable que si A est diagonalisable ; à ceci il faut ajouter le fait le plus fondamentale que l'algorithme ne s'étend pas au cas des matrices à spectre complexe.

Il est néanmoins possible d'élaborer à partir du théorème vu précédemment une procédure itérative répondant à ces deux objectifs et qui constituera ce qu'on appellera une version accélérée de la méthode de Newton.

Algorithme de Newton accéléré

Considérons d'abord le problème des valeurs propres complexes.

Soient Z_1 et Z_2 deux de ces valeurs. Définissons alors la transformation

$$|\lambda_{\max}(A_{K+I}^*)| = \frac{1}{2} \frac{|\lambda_{\max}(A_K^*)| + |\lambda_{\min}(A_K^*)|}{\sqrt{|\lambda_{\max}(A_K^*) \cdot \lambda_{\min}(A_K^*)|}}$$

on peut donc affirmer qu'il existe $1 \leq n-1$ tel que A_1^* possède n valeurs propres de même module.

On peut alors écrire:

$$S = \text{signe}(A) = \frac{1}{|\lambda(A_1^*)|} \quad (21)$$

Ainsi il est possible d'énoncer le théorème suivant:

Soit A (n,n) régulière à spectre réel, $S = \text{signe}(A)$ est obtenue par l'algorithme fini:

$$A_{K+I}^* = \frac{1}{2} \left[\alpha_K A_K^* \frac{1}{\alpha_K} (A_K^*)^{-1} - I \right] \quad A_0 = A$$

$$\alpha_K = \frac{1}{\sqrt{|\lambda_{\max}(A_K^*) \cdot \lambda_{\min}(A_K^*)|}}$$

$$S = \frac{1}{|\lambda(A_1^*)|} \cdot A_1^* \quad (22)$$

3-2-ALGORITHME DE NEWTON ACCELERE

Introduction

La méthode qu'on a vue précédemment n'est malheureusement utilisable que si A est diagonalisable ; à ceci il faut ajouter le fait le plus fondamentale que l'algorithme ne s'étend pas au cas des matrices à spectre complexe.

Il est néanmoins possible d'élaborer à partir du théorème vu précédemment une procédure itérative répondant à ces deux objectifs et qui constituera ce qu'on appellera une version accélérée de la méthode de Newton.

Algorithme de Newton accéléré

Considérons d'abord le problème des valeurs propres complexes.

Soient Z_1 et Z_2 deux de ces valeurs. Définissons alors la transformation

$$Z_i^* = \alpha Z \quad i=1,2$$

puis

$$Z_i' = \frac{1}{2} \left(Z_i + \frac{1}{Z_i} \right) \quad \text{et} \quad \tilde{Z}_i = \frac{1}{2} \left(Z_i^* + \frac{1}{Z_i^*} \right)$$

on a donc

$$\begin{aligned} \rho_i'^2 &= |Z_i'|^2 = \frac{1}{4} \left(\rho_i + \frac{1}{\rho_i} \right)^2 - \sin^2 \epsilon_i \\ \tilde{\rho}_i^2 &= |\tilde{Z}_i|^2 = \frac{1}{4} \left(\alpha \rho_i + \frac{1}{\alpha \rho_i} \right)^2 - \sin^2 \epsilon_i \\ \operatorname{tg} \tilde{\epsilon}_i &= \frac{\alpha^2 \rho_i^2 - 1}{\alpha^2 \rho_i^2} \operatorname{tg} \epsilon_i \\ \operatorname{tg} \tilde{\epsilon}_i &= \frac{\rho_i^2 - 1}{\rho_i^2} \operatorname{tg} \epsilon_i \end{aligned}$$

On voit clairement (en général) qu'il n'existe pas de α fonction de ρ_i, ϵ_i tel que $\tilde{\epsilon}_1 = \tilde{\epsilon}_2$ et $\tilde{\rho}_1 = \tilde{\rho}_2$ soient vérifiés simultanément.

Une première approximation consiste par extension du cas réel à faire abstraction des arguments ϵ_1 et ϵ_2 et ne raisonner que sur les modules ρ_1 et ρ_2 .

Le mieux qu'on puisse faire est d'égaliser la contribution de ρ_1 et ρ_2 dans les modules $\tilde{\rho}_{1,2}$ d'où l'on prend:

$$\alpha = \frac{1}{\sqrt{\rho_1 \rho_2}} \quad (23)$$

d'où l'on a

$$\begin{aligned} |\tilde{\rho}_1 - \tilde{\rho}_2| &< |\rho_1' - \rho_2'| \\ 0 &< |\tilde{\rho}_1^2 - \tilde{\rho}_2^2| < 1 \end{aligned} \quad (24)$$

et une égalisation du facteur de réduction de $\operatorname{tg} \epsilon$ qui découle des relations

$$\begin{aligned} \rho_1^* \rho_2^* &= \alpha \rho_1 \cdot \alpha \rho_2 = 1 \\ \Psi(\rho) &= \frac{\rho^2 - 1}{\rho^2 + 1} = -\Psi\left(\frac{1}{\rho}\right) \\ \left. \begin{aligned} \operatorname{tg} \tilde{\epsilon}_1 &= \Psi(\rho_1^*) \operatorname{tg} \epsilon_1 \\ \operatorname{tg} \tilde{\epsilon}_2 &= -\Psi(\rho_1^*) \operatorname{tg} \epsilon_2 \end{aligned} \right\} \quad (25) \end{aligned}$$

$$Z_i^* = \sqrt{Z} \quad i=1,2$$

puis

$$Z_i' = \frac{1}{2} \left(Z_i + \frac{1}{Z_i} \right) \quad \text{et} \quad \tilde{Z}_i = \frac{1}{2} \left(Z_i^* + \frac{1}{Z_i^*} \right)$$

on a donc

$$\rho_i'^2 = |Z_i|^2 = \frac{1}{4} \left(\rho_i + \frac{1}{\rho_i} \right)^2 - \sin^2 \epsilon_i$$

$$\tilde{\rho}_i^2 = |\tilde{Z}_i|^2 = \frac{1}{4} \left(\alpha \rho_i + \frac{1}{\alpha \rho_i} \right)^2 - \sin^2 \tilde{\epsilon}_i$$

$$\operatorname{tg} \tilde{\epsilon}_i = \frac{\alpha^2 \rho_{i-1}^2 \operatorname{tg} \epsilon_i}{\alpha^2 \rho_{i+1}^2}$$

$$\operatorname{tg} \epsilon_i = \frac{\rho_{i-1}^2}{\rho_{i+1}^2} \cdot \operatorname{tg} \epsilon_i$$

On voit clairement (en général) qu'il n'existe pas de α fonction de ρ_i, ϵ_i tel que $\tilde{\epsilon}_1 = \tilde{\epsilon}_2$ et $\tilde{\rho}_1 = \tilde{\rho}_2$ soient vérifiées simultanément.

Une première approximation consiste par extension du cas réel à faire abstraction des arguments ϵ_1 et ϵ_2 et ne raisonner que sur les modules ρ_1 et ρ_2 .

Le mieux qu'on puisse faire est d'égaliser la contribution de ρ_1 et ρ_2 dans les modules $\tilde{\rho}_{1,2}$ d'où l'on prend:

$$\alpha = \frac{1}{\sqrt{\rho_1 \rho_2}} \quad (23)$$

d'où l'on a

$$|\tilde{\rho}_1 - \tilde{\rho}_2| < |\rho_1' - \rho_2'|$$

$$0 < |\tilde{\rho}_1^2 - \tilde{\rho}_2^2| < 1 \quad (24)$$

et une égalisation du facteur de réduction de $\operatorname{tg} \epsilon$ qui découle des relations

$$\rho_1^* \rho_2^* = \alpha \rho_1 \cdot \alpha \rho_2 = 1$$

$$\Psi(\rho) = \frac{\rho^2 - 1}{\rho^2 + 1} = -\Psi\left(\frac{1}{\rho}\right)$$

$$\operatorname{tg} \tilde{\epsilon}_1 = \Psi(\rho_1^*) \operatorname{tg} \epsilon_1$$

$$\operatorname{tg} \tilde{\epsilon}_2 = -\Psi(\rho_1^*) \operatorname{tg} \epsilon_2$$

(25)

CONCLUSION

On peut dire que la transformation α n'annule plus l'écart $\tilde{Z}_1 - \tilde{Z}_2$, mais contribue à réduire l'écart en norme de \tilde{Z}_1 et \tilde{Z}_2 par rapport à une itération de Newton (Z'_1 et Z'_2) et égalise en valeur absolue le facteur de réduction de la tangente des arguments ϵ_1 et ϵ_2 .

Donc si on fait jouer à $\lambda_{\max}(A_K^*)$ et $\lambda_{\min}(A_K^*)$ le rôle de Z_1 et Z_2 dans l'algorithme précédent. En raison des propriétés d'invariance de la fonction signe on a :

$$\forall K \text{ signe}(A_K^*) = \text{signe}(A_K) = \text{signe}(A)$$

c'est à dire

$$\lim_{K \rightarrow \infty} A_K^* = \lim_{K \rightarrow \infty} A_K = S$$

Ceci étant, on peut dire que l'algorithme accéléré se comporte dans son stade ultime exactement comme l'algorithme de Newton, l'accélération de la convergence au sens de la réduction du nombre d'itéré pour atteindre S avec une précision donnée, intervient donc essentiellement dans l'évolution de A_K^* correspondant aux phases 1 et 2 du cas scalaire.

Pour éviter le calcul des valeurs propres on est amené à faire l'approximation suivante :

$$\alpha_K = \sqrt{\frac{\|A_K^{*-1}\|}{\|A_K^*\|}} \quad (27)$$

ceci est légitime car :

$$|\lambda_{\max}(A)| \leq \|A\| \quad \text{et} \quad |\lambda_{\min}(A)| \geq 1/\|A^{-1}\|$$

L'algorithme de Newton accéléré est donc en définitive, défini par la procédure suivante.

$$\begin{aligned} A_0^* &= A \\ A_{K+1}^* &= \frac{1}{2} \left[\alpha_K A_K^* + \frac{1}{\alpha_K} A_K^{*-1} \right] \\ \alpha_K &= \sqrt{\frac{\|A_K^{-1}\|}{\|A_K\|}} \end{aligned} \quad (28)$$

CONCLUSION

On peut dire que la transformation α n'annule plus l'écart $\tilde{Z}_1 - \tilde{Z}_2$, mais contribue à réduire l'écart en norme de \tilde{Z}_1 et \tilde{Z}_2 par rapport à une itération de Newton (Z'_1 et Z'_2) et égalise en valeur absolue le facteur de réduction de la tangente des arguments ϵ_1 et ϵ_2 .

Donc si on fait jouer à $\lambda_{\max}(A_K^*)$ et $\lambda_{\min}(A_K^*)$ le rôle de Z_1 et Z_2 dans l'algorithme précédent. En raison des propriétés d'invariance de la fonction signe on a :

$$\forall K \text{ signe}(A_K^*) = \text{signe}(A_K) = \text{signe}(A)$$

c'est à dire

$$\lim_{K \rightarrow \infty} A_K^* = \lim_{K \rightarrow \infty} A_K = S$$

Ceci étant, on peut dire que l'algorithme accéléré se comporte dans son stade ultime exactement comme l'algorithme de Newton, l'accélération de la convergence au sens de la réduction du nombre d'itéré pour atteindre S avec une précision donnée, intervient donc essentiellement dans l'évolution de A_K^* correspondant aux phases 1 et 2 du cas scalaire.

Pour éviter le calcul des valeurs propres on est amené à faire l'approximation suivante :

$$\alpha_K = \sqrt{\frac{\|A_K^{*-1}\|}{\|A_K^*\|}} \quad (27)$$

ceci est légitime car :

$$|\lambda_{\max}(A)| \leq \|A\| \quad \text{et} \quad |\lambda_{\min}(A)| \geq 1/\|A^{-1}\|$$

L'algorithme de Newton accéléré est donc en définitive, défini par la procédure suivante.

$$\begin{aligned} A_0^* &= A \\ A_{K+1}^* &= \frac{1}{2} \left[\alpha_K A_K^* + \frac{1}{\alpha_K} A_K^{*-1} \right] \\ \alpha_K &= \sqrt{\frac{\|A_K^{*-1}\|}{\|A_K^*\|}} \end{aligned} \quad (28)$$

3-3-IMPLEMENTATION

On a vu précédemment que l'algorithme de Newton accéléré, était défini par

$$\begin{aligned} \Delta c^* &= A \\ \Delta_{K+1}^* &= \frac{1}{2} \left[\alpha_K \Delta_K^* + \frac{1}{\alpha_K} \Delta_{K-1}^* \right] \\ \alpha_K &= \sqrt{\frac{\|\Delta_K^*\|}{\|\Delta_{K-1}^*\|}} \end{aligned}$$

Nous remarquons que le calcul de α_K passe par celui des normes de matrices; ainsi s'impose un choix judicieux de ces dernières.

3-3-1-Choix de la norme

On a par définition

$$\|A\|_p = \max_{\|x\|_p=1} \|\Delta x\|_p$$

ce qui donne

$$\begin{aligned} p=1 & \quad \|A\|_1 = \max_j \left(\sum_{i=1}^n |a_{ij}| \right) \\ p=2 & \quad \|A\|_2 = \sqrt{\max(\lambda^T \lambda)} \\ p=\infty & \quad \|A\|_\infty = \max_i \left(\sum_{j=1}^n |a_{ij}| \right) \\ p=F & \quad \|A\|_F = \sqrt{\sum_{i,j} a_{ij}^2} \end{aligned}$$

on en déduit que:

$$\begin{aligned} \|A_1\| &= \|A^T\|_\infty \\ \|A_2\| &\leq \|A\|_F \leq \begin{cases} \sqrt{n} \cdot \|A\|_1 \\ \sqrt{n} \cdot \|A\|_\infty \end{cases} \\ \frac{1}{\sqrt{n}} \cdot \|A\|_p &\leq \|A\|_F \leq \sqrt{n} \cdot \|A\|_p \end{aligned}$$

notre choix doit ainsi satisfaire deux critères:

- le coût calcul
- la recherche du plus petit majorant de $|\lambda_{\max}|$

On remarque que le plus petit majorant de $|\lambda_{\max}|$ est $\|A\|_2$, tandis que le critère coût calcul est satisfait par $\|A\|_1$ ou $\|A\|_\infty$.

Un compromis sera ainsi donné par :

$$\|A\| = \min(\|A\|_1, \|A\|_\infty) \quad (29)$$

3-3-IMPLEMENTATION

On a vu précédemment que l'algorithme de Newton accéléré, était défini par

$$\begin{aligned} \Delta_K^* &= A \\ \Delta_{K+1}^* &= \frac{1}{2} \left[\alpha_K \Delta_K^* + \frac{1}{\alpha_K} \Delta_{K-1}^* \right] \\ \alpha_K &= \sqrt{\frac{\|\Delta_K^*\|}{\|\Delta_{K-1}^*\|}} \end{aligned}$$

Nous remarquons que le calcul de α_K passe par celui des normes de matrices; ainsi s'impose un choix judicieux de ces dernières.

3-3-1-Choix de la norme

On a par définition

$$\|A\|_p = \max_{\|x\|_p=1} \|Ax\|_p$$

ce qui donne

$$\begin{aligned} p=1 & \quad \|A\|_1 = \max_j \left(\sum_{i=1}^n |a_{ij}| \right) \\ p=2 & \quad \|A\|_2 = \sqrt{\max(A^T A)} \\ p=\infty & \quad \|A\|_\infty = \max_i \left(\sum_{j=1}^n |a_{ij}| \right) \\ p=F & \quad \|A\|_F = \sqrt{\sum_{i,j} a_{ij}^2} \end{aligned}$$

on en déduit que:

$$\begin{aligned} \|A\|_1 &= \|A^T\|_\infty \\ \|A\|_2 &\leq \|A\|_F \leq \begin{cases} \sqrt{n} \cdot \|A\|_1 \\ \sqrt{n} \cdot \|A\|_\infty \end{cases} \\ \frac{1}{\sqrt{n}} \cdot \|A\|_p &\leq \|A\|_F \leq \sqrt{n} \cdot \|A\|_p \end{aligned}$$

notre choix doit ainsi satisfaire deux critères:

- le coût calcul
- la recherche du plus petit majorant de $|\lambda_{\max}|$

On remarque que le plus petit majorant de $|\lambda_{\max}|$ est $\|A\|_2$, tandis que le critère coût calcul est satisfait par $\|A\|_1$ ou $\|A\|_\infty$.

Un compromis sera ainsi donné par :

$$\|A\| = \min(\|A\|_1, \|A\|_\infty) \quad (29)$$

3-3-2-Critère d'arrêt

Il convient de définir un critère d'arrêt tel que

$$\lim_{K \rightarrow \infty} \Lambda_K^* = S \quad \left\{ \begin{array}{l} \alpha_K \rightarrow 1 \\ \|\Lambda_{K+1}^* - \Lambda_K^*\| \rightarrow 0 \\ \frac{1}{n} \cdot \text{trace}(\Lambda_K^*)^2 \rightarrow 1 \end{array} \right.$$

On en déduit les conditions suivantes classées par coût croissant dont chacune peut constituer un test de convergence :

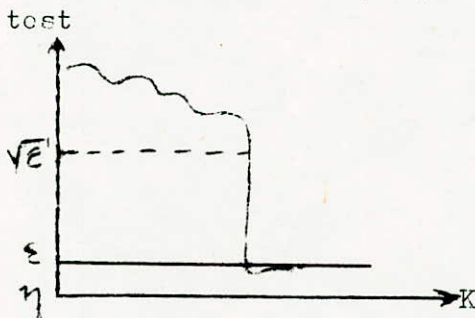
$$\begin{array}{l} 1-|\alpha_K-1| \leq \varepsilon \\ 2- \|\Lambda_{K+1}^* - \Lambda_K^*\| \leq \varepsilon \|\Lambda_K^*\| \\ 3- \left| \frac{1}{n} \cdot \text{trace}(\Lambda_K^*)^2 - 1 \right| < \varepsilon \end{array} \quad \varepsilon < \eta$$

η étant la précision machine

Remarque: $\alpha_0=1$, pour une matrice orthogonale ; donc le 1^{er} test ne s'applique pas à la première itération car $\|\Lambda_1^* - \Lambda_0^*\| \neq 0$.

Le 1^{er} et 3^{er} tests sont équivalents. En raison du coût calcul du 3^{er} test (n^3 multiplications) en comparaison avec le 1^{er} (n^2), on a ainsi éliminé le 3^{er}. On aura le test suivant

$$\text{"test"} = \max(|\alpha_K - 1|, \frac{\|\Lambda_{K+1}^* - \Lambda_K^*\|}{\|\Lambda_K^*\|}) \quad \text{fin} \quad (30)$$



La variable test évolue d'une façon chaotique avant de décroître jusqu'à une valeur de l'ordre de grandeur de η . La difficulté réside ainsi dans la fixation d'un seuil ε . Or la convergence finale a une nature quadratique; donc la variable test décroît très rapidement à partir d'un certain moment pour arriver sur une asymptote basse ε dépendant de η et de Λ .

3-3-2-Critère d'arrêt

Il convient de définir un critère d'arrêt tel que

$$\lim_{K \rightarrow \infty} \Lambda_K^* = S \quad \left\{ \begin{array}{l} \alpha_K \rightarrow 1 \\ \|\Lambda_{K+1}^* - \Lambda_K^*\| \rightarrow 0 \\ \frac{1}{n} \cdot \text{trace}(\Lambda_K^*)^2 \rightarrow 1 \end{array} \right.$$

On en déduit les conditions suivantes classées par coût croissant dont chacune peut constituer un test de convergence :

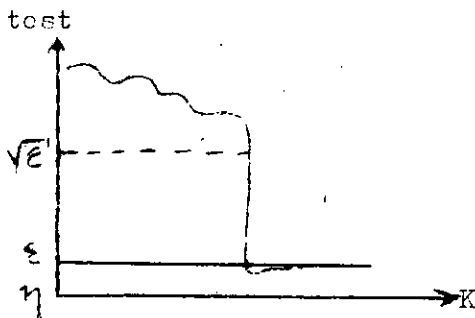
$$\begin{array}{l} 1-|\alpha_K-1| \leq \varepsilon \\ 2- \|\Lambda_{K+1}^* - \Lambda_K^*\| \leq \varepsilon \|\Lambda_K^*\| \\ 3- \left| \frac{1}{n} \cdot \text{trace}(\Lambda_K^*)^2 - 1 \right| < \varepsilon \end{array} \quad \varepsilon < \eta$$

η étant la précision machine

Remarque: $\alpha_0=1$, pour une matrice orthogonale ; donc le 1^{er} test ne s'applique pas à la première itération car $\|\Lambda_1^* - \Lambda_0^*\| \neq 0$.

Le 1^{er} et 3^{es} tests sont équivalents. En raison du coût calcul du 3^{es} test (n^3 multiplications) en comparaison avec le 1^{er} (n^2), on a ainsi éliminé le 3^{es}. On aura le test suivant

$$\text{"test"} = \max(|\alpha_K - 1|, \frac{\|\Lambda_{K+1}^* - \Lambda_K^*\|}{\|\Lambda_K^*\|}) \quad \text{fin} \quad (30)$$



La variable test évolue d'une façon chaotique avant de décroître jusqu'à une valeur de l'ordre de grandeur de η . La difficulté réside ainsi dans la fixation d'un seuil ε . Or la convergence finale a une nature quadratique; donc la variable test décroît très rapidement à partir d'un certain moment pour arriver sur une asymptote basse ε dépendant de η et de Λ .

Il faut donc s'assurer d'abord qu'on ait atteint la zone de convergence quadratique. Pour cela on introduit un seuil suffisamment grand vis à vis de η pour qu'il soit atteint, par exemple $\sqrt{\eta}$.

La détection de l'asymptote se fait sur la base de deux valeurs consécutives du "test":

$$\text{test}_{K+1} \simeq \text{test}_K^2 \quad \text{et} \quad \text{test}_{K+1} \simeq \text{test}_K$$

d'où l'on aura le critère d'arrêt suivant:

$$\text{Si } \left[\text{test}_{K+1} < \sqrt{\eta} \text{ et } \text{test}_{K+1} > \frac{1}{2} \text{test}_K \right] \text{ fin} \quad (31)$$

Enfin pour prévenir tout bouclage intempestif de l'algorithme, le nombre d'itérations doit être limité sur la base de la valeur maximum possible de N:

$$n(\varphi) = \left\lfloor \log_2 \max(\|A\|, \|A^{-1}\|) \right\rfloor$$

et

$$q(\eta) < 1 + \left\lfloor \log_2 \left(\frac{c,64}{\eta} \right) \right\rfloor \quad \text{avec} \quad \left(\text{tr} \varphi < \frac{2}{\pi \cdot \eta} \left(\frac{c,64}{\eta} \right) \right)$$

En conséquence on a implémenté le test supplémentaire:

$$N = \left\lfloor \log_2 \max(\|A\|, \|A^{-1}\|) \right\rfloor + q_m(\eta) + p(\eta) \quad (32)$$

où $q_m(\eta)$ et $p(\eta)$ sont des constantes machines introduites une fois pour toutes.

Il faut donc s'assurer d'abord qu'on ait atteint la zone de convergence quadratique. Pour cela on introduit un seuil suffisamment grand vis à vis de η pour qu'il soit atteint, par exemple $\sqrt{\eta}$.

La détection de l'asymptote se fait sur la base de deux valeurs consécutives du "test":

$$\text{test}_{K+1} \simeq \text{test}_K^2 \quad \text{et} \quad \text{test}_{K+1} \simeq \text{test}_K$$

d'où l'on aura le critère d'arrêt suivant:

$$\text{Si } \left[\text{test}_{K+1} < \sqrt{\eta} \text{ et } \text{test}_{K+1} > \frac{1}{2} \text{test}_K \right] \text{ fin} \quad (31)$$

Enfin pour prévenir tout bouclage intempestif de l'algorithme, le nombre d'itérations doit être limité sur la base de la valeur maximum possible de N:

$$n(\rho) = \left\lfloor \log_2 \max(\|A\|, \|A^{-1}\|) \right\rfloor$$

et

$$g(\rho) < 1 + \left\lfloor \log_2 \left(\frac{c,64}{\eta} \right) \right\rfloor \quad \text{avec} \quad \left(\text{trco} < \frac{2}{\pi \cdot \eta} \left(\frac{c,64}{\eta} \right) \right)$$

En conséquence on a implémenté le test supplémentaire:

$$N = \left\lfloor \log_2 \max(\|A\|, \|A^{-1}\|) \right\rfloor + g_m(\eta) + p(\eta) \quad (32)$$

où $g_m(\eta)$ et $p(\eta)$ sont des constantes machines introduites une fois pour toutes.

CHAPITRE II/

APPLICATION A LA RESOLUTION DES EQUATIONS DES EQUATIONS DE RICCATI

Partie A/

Equations de Riccati dans les problèmes de commande optimale avec critère quadratique.

Partie B/

Résolution de l'équation de Riccati dans le cas continu

Partie C/

Résolution de l'équation de Riccati dans le cas discret

CHAPITRE II/

APPLICATION A LA RESOLUTION DES EQUATIONS DES EQUATIONS DE RICCATI

Partie A/

Equations de Riccati dans les problèmes de commande optimale avec critère quadratique.

Partie B/

Résolution de l'équation de Riccati dans le cas continu

Partie C/

Résolution de l'équation de Riccati dans le cas discret

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

Après avoir étudié le concept de fonction de matrice, nous procédons dans ce chapitre à son application pour la résolution des équations de Riccati apparaissant dans les problèmes d'optimisation déterministe de la commande.

Pour cela, il nous a paru utile de diviser ce chapitre en trois parties:

- la première traitera de la commande optimale des systèmes linéaires avec critère quadratique; ceci pour montrer l'origine des équations de Riccati.
- la deuxième aura pour but de donner un algorithme (utilisant la fonction signe) résolvant l'équation de Riccati dans le cas continu et donner l'organigramme.
- quand à la dernière partie, elle fera l'objet d'un algorithme (utilisant la fonction signe) résolvant l'équation de Riccati dans le cas discret et présentation de l'organigramme.

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

Après avoir étudié le concept de fonction de matrice, nous procédons dans ce chapitre à son application pour la résolution des équations de Riccati apparaissant dans les problèmes d'optimisation déterministe de la commande.

Pour cela, il nous a paru utile de diviser ce chapitre en trois parties:

- la première traitera de la commande optimale des systèmes linéaires avec critère quadratique; ceci pour montrer l'origine des équations de Riccati.
- la deuxième aura pour but de donner un algorithme (utilisant la fonction signe) résolvant l'équation de Riccati dans le cas continu et donner l'organigramme.
- quand à la dernière partie, elle fera l'objet d'un algorithme (utilisant la fonction signe) résolvant l'équation de Riccati dans le cas discret et présentation de l'organigramme.

EQUATIONS DE RICCATI DANS LES PROBLEMES DE COMMANDE

OPTIMALE AVEC CRITERE QUADRATIQUE

1-INTRODUCTION

La minimisation d'un critère quadratique constitue l'un des moyens de parvenir à la détermination d'une structure de commande par retour d'état pour les systèmes multidimensionnels.

En effet, un critère quadratique permet d'exprimer d'une manière convenable les qualités globales recherchées par la commande tant en assurant le meilleur compromis entre certaines performances, représentées par des termes de pondération faisant intervenir les sorties ou les variables d'état, et une économie d'énergie.

Un autre avantage, non moins négligeable de la méthode quadratique est de conduire à des développements mathématiques nombreux et puissants.

2-PRESENTATION DU PROBLEME

Soient le système linéaire invariant et gouvernable

$$\dot{X}(t) = A.X(t) + B.U(t)$$

$$Y(t) = C.X(t)$$

(33)

X: état , U: commande , Y: sortie

et le coût (en supposant l'instant final infini)

$$J = \frac{1}{2} \int_0^{\infty} \left[\langle X(t) , Q.X(t) \rangle + \langle U(t) , R.U(t) \rangle \right] dt \quad (34)$$

3-MINIMISATION DE L'HAMILTONNIEN

L'hamiltonien associé à (33) et (34) est

$$H = \frac{1}{2} \langle X(t) , Q.X(t) \rangle + \frac{1}{2} \langle U(t) , R.U(t) \rangle \\ + \langle p(t) , A.X(t) + B.U(t) \rangle$$

EQUATIONS DE RICCATI DANS LES PROBLEMES DE COMMANDE

OPTIMALE AVEC CRITERE QUADRATIQUE

1-INTRODUCTION

La minimisation d'un critère quadratique constitue l'un des moyens de parvenir à la détermination d'une structure de commande par retour d'état pour les systèmes multidimensionnels.

En effet, un critère quadratique permet d'exprimer d'une manière convenable les qualités globales recherchées par la commande tant en assurant le meilleur compromis entre certaines performances, représentées par des termes de pondération faisant intervenir les sorties ou les variables d'état, et une économie d'énergie.

Un autre avantage, non moins négligeable de la méthode quadratique est de conduire à des développements mathématiques nombreux et puissants.

2-PRESENTATION DU PROBLEME

Soient le système linéaire invariant et gouvernable

$$\dot{X}(t) = A.X(t) + B.U(t)$$

$$Y(t) = C.X(t) \tag{33}$$

X: état , U: commande , Y: sortie

et le coût (en supposant l'instant final infini)

$$J = \frac{1}{2} \int_0^{\infty} \left[\langle X(t) , Q.X(t) \rangle + \langle U(t) , R.U(t) \rangle \right] dt \tag{34}$$

3-MINIMISATION DE L'HAMILTONNIEN

L'hamiltonien associé à (33) et (34) est

$$H = \frac{1}{2} \langle X(t) , Q.X(t) \rangle + \frac{1}{2} \langle U(t) , R.U(t) \rangle \\ + \langle p(t) , A.X(t) + B.U(t) \rangle$$

$$H = \frac{1}{2} \langle X(t), Q \cdot X(t) \rangle + \frac{1}{2} \langle U(t), R \cdot U(t) \rangle + \langle p(t), A \cdot X(t) \rangle + \langle p(t), B \cdot U(t) \rangle$$

où p est le Vecteur adjoint solution de l'équation :

$$\dot{p}(t) = -\nabla_x H$$

d'où

$$\dot{p}(t) = -Q \cdot X(t) - A^T p(t)$$

Le long de la trajectoire optimale, nous devons avoir

$$\nabla_u H = 0 \implies R \cdot U + B^T \cdot p(t) = 0$$

d'où

$$U(t) = -R^{-1} B^T p(t) \tag{35}$$

R^{-1} existe du moment que R est définie positive.

Remarque

$\nabla_u H = 0$ implique seulement l'existence d'un extrémum. Si de plus on a

$\frac{\partial^2 H}{\partial u^2}$ défini positif, alors cet extrémum s'identifie à un minimum.

Or

$$\frac{\partial^2 H}{\partial u^2} = R \text{ qui est définie positive}$$

alors $U(t)$ minimalise bien l'hamiltonien

4-EQUATIONS CANONIQUES REDUITES

En remplaçant $U(t)$ par sa valeur dans (33), on aura

$$\dot{X}(t) = A \cdot X(t) - B R^{-1} B^T p(t)$$

$$\dot{p}(t) = -Q \cdot X(t) - A^T \cdot p(t)$$

qui sont les équations canoniques réduites.

Soit
$$V(t) = B R^{-1} B^T$$

on aura donc

$$H = \frac{1}{2} \langle X(t), Q \cdot X(t) \rangle + \frac{1}{2} \langle U(t), R \cdot U(t) \rangle \\ + \langle p(t), A \cdot X(t) \rangle + \langle p(t), B \cdot U(t) \rangle$$

où p est le Vecteur adjoint solution de l'équation :

$$\dot{p}(t) = -\nabla_x H$$

d'où

$$\dot{p}(t) = -Q \cdot X(t) - A^T p(t)$$

Le long de la trajectoire optimale, nous devons avoir

$$\nabla_u H = 0 \quad \Longrightarrow \quad R \cdot U + B^T \cdot p(t) = 0$$

d'où

$$U(t) = -R^{-1} B^T p(t) \quad (35)$$

R^{-1} existe du moment que R est définie positive.

Remarque

$\nabla_u H = 0$ implique seulement l'existence d'un extrémum. Si de plus on a

$\frac{\partial^2 H}{\partial u^2}$ défini positif, alors cet extrémum s'identifie à un minimum.

Or

$$\frac{\partial^2 H}{\partial u^2} = R \quad \text{qui est définie positive}$$

alors $U(t)$ minimalise bien l'hamiltonien

4-EQUATIONS CANONIQUES REDUITES

En remplaçant $U(t)$ par sa valeur dans (33), on aura

$$\dot{X}(t) = A \cdot X(t) - B R^{-1} B^T p(t)$$

$$\dot{p}(t) = -Q \cdot X(t) - A^T \cdot p(t)$$

qui sont les équations canoniques réduites.

Soit $V(t) = B R^{-1} B^T$

on aura donc

$$\begin{bmatrix} \dot{X}(t) \\ \dot{p}(t) \end{bmatrix} = \begin{bmatrix} A & -V \\ -Q & -A^T \end{bmatrix} \begin{bmatrix} X(t) \\ p(t) \end{bmatrix} \quad (36)$$

On a donc un système de $2n$ équations différentielles homogènes. Il admet une solution unique lorsque sont fixées les $2n$ conditions aux limites.

Or n conditions sont fournies par $X(t_1)$, les autres sont données par $p(t_2) = 0$ ($t_2 = \infty$).

5-détermination de P

On voit aisément que $p(t)$ et $X(t)$ sont reliées par la relation suivante:

$$p(t) = P.X(t)$$

d'où

$$\frac{d}{dt} p(t) = \frac{d}{dt} (P.X(t)) = P.\dot{X}(t) \quad (37)$$

D'après (36), on peut écrire

$$\dot{X}(t) = \begin{bmatrix} A - VP \end{bmatrix} X(t)$$

$$\dot{p}(t) = \begin{bmatrix} -Q - A^T P \end{bmatrix} X(t)$$

en remplaçant $\dot{X}(t)$ par sa valeur dans (37), on aboutit à:

$$\dot{p}(t) = P \begin{bmatrix} A - VP \end{bmatrix} X(t)$$

$$\begin{bmatrix} -Q - A^T P \end{bmatrix} X(t) = P \begin{bmatrix} A - VP \end{bmatrix} X(t)$$

d'où finalement

$$PA + A^T P - PVP + Q = 0 \quad (38)$$

c'est l'équation de Riccati matricielle.

Ainsi pour avoir la commande optimale d'un système linéaire avec critère quadratique, il suffit de calculer la solution P de l'équation de Riccati pour avoir

$$U(t) = -R^{-1} B^T P.X(t) \quad (39)$$

$$\begin{bmatrix} \dot{X}(t) \\ \dot{p}(t) \end{bmatrix} = \begin{bmatrix} A & -V \\ -Q & -A^T \end{bmatrix} \begin{bmatrix} X(t) \\ p(t) \end{bmatrix} \quad (36)$$

On a donc un système de $2n$ équations différentielles homogènes. Il admet une solution unique lorsque sont fixées les $2n$ conditions aux limites.

Or n conditions sont fournies par $X(t_1)$, les autres sont données par $p(t_2) = 0$ ($t_2 = \infty$).

5-détermination de P

On voit aisément que $p(t)$ et $X(t)$ sont reliées par la relation suivante:

$$p(t) = P.X(t)$$

d'où

$$\frac{d}{dt} p(t) = \frac{d}{dt} (P.X(t)) = P.\dot{X}(t) \quad (37)$$

D'après (36), on peut écrire

$$\dot{X}(t) = \begin{bmatrix} A - VP \end{bmatrix} X(t)$$

$$\dot{p}(t) = \begin{bmatrix} -Q - A^T P \end{bmatrix} X(t)$$

en remplaçant $\dot{X}(t)$ par sa valeur dans (37), on aboutit à:

$$\dot{p}(t) = P \begin{bmatrix} A - VP \end{bmatrix} X(t)$$

$$\begin{bmatrix} -Q - A^T P \end{bmatrix} X(t) = P \begin{bmatrix} A - VP \end{bmatrix} X(t)$$

d'où finalement

$$PA + A^T P - PVP + Q = 0 \quad (38)$$

c'est l'équation de Riccati matricielle.

Ainsi pour avoir la commande optimale d'un système linéaire avec critère quadratique, il suffit de calculer la solution P de l'équation de Riccati pour avoir

$$U(t) = -R^{-1} B^T P.X(t) \quad (39)$$

RESOLUTION DE L'EQUATION DE RICCATI (CAS CONTINU)

A L'AIDE DE LA FONCTION SIGNE DE MATRICE

1-INTRODUCTION

Après avoir introduit le concept de la fonction signe de matrice et le problème de commande optimale avec critère quadratique, nous allons maintenant aborder l'application de ce concept à la résolution de l'équation de Riccati introduite en termes de commande optimale des systèmes continus.

2-DEFINITION DU PROBLEME

Soit le système invariant à instant final infini:

$$\dot{X}(t) = A.X(t) + B.U(t)$$

et le critère

$$J = \frac{1}{2} \int_0^{\infty} (U^T R U + X^T Q X) dt. \quad (40)$$

La commande par retour d'état linéaire minimisant J est

$$\hat{U} = -R^{-1} B^T P X(t)$$

avec P solution de l'équation de Riccati

$$PA + A^T P - P B R^{-1} B^T P + Q = 0 \quad (41)$$

satisfaisant les hypothèses suivantes

$$\left. \begin{array}{l} R^T = R \quad R \gg 0 \quad Q^T = Q \quad Q \gg 0 \\ (A, B) \text{ stabilisable} \\ (C, A) \text{ détectable} \quad C: C^T C = Q \end{array} \right\} \quad (42)$$

On sait alors que P^T est aussi solution de (41) et que

$$P^T = P \quad P \gg 0 \quad P \text{ unique} \quad (43)$$

$$\text{réel } \lambda_i[\tilde{A}] < 0 \quad \text{avec} \quad \tilde{A} = A - B R^{-1} B^T P \quad (44)$$

L'équation de Riccati peut alors s'écrire

$$P \tilde{A} + \tilde{A}^T P + Q = 0 \quad (45)$$

RÉSOLUTION DE L'ÉQUATION DE RICCATI (CAS CONTINU)

A L'AIDE DE LA FONCTION SIGNE DE MATRICE

1-INTRODUCTION

Après avoir introduit le concept de la fonction signe de matrice et le problème de commande optimale avec critère quadratique, nous allons maintenant aborder l'application de ce concept à la résolution de l'équation de Riccati introduite en termes de commande optimale des systèmes continus.

2-DEFINITION DU PROBLÈME

Soit le système invariant à instant final infini:

$$\dot{X}(t) = A.X(t) + B.U(t)$$

et le critère

$$J = \frac{1}{2} \int_0^{\infty} (U^T R U + X^T Q X) dt. \quad (40)$$

La commande par retour d'état linéaire minimisant J est

$$\hat{U} = -R^{-1} B^T P X(t)$$

avec P solution de l'équation de Riccati

$$PA + A^T P - P B R^{-1} B^T P + Q = 0 \quad (41)$$

satisfaisant les hypothèses suivantes

$$\left. \begin{array}{l} R^T = R \quad R \gg 0 \quad Q^T = Q \quad Q \gg 0 \\ (A, B) \text{ stabilisable} \\ (C, A) \text{ détectable} \quad C: C^T C = Q \end{array} \right\} \quad (42)$$

On sait alors que P^T est aussi solution de (41) et que

$$P^T = P \quad P \gg 0 \quad P \text{ unique} \quad (43)$$

$$\text{réel } \lambda_i[\tilde{A}] < 0 \quad \text{avec} \quad \tilde{A} = A - B R^{-1} B^T P \quad (44)$$

L'équation de Riccati peut alors s'écrire

$$P \tilde{A} + \tilde{A}^T P + Q = 0 \quad (45)$$

Cela exprime donc que le système bouclé est asymptotiquement stable. Si on remplace l'hypothèse, (C, A) détectable, par (C, A) reconstructible (observable), on aura $P > 0$.

3-RESOLUTION

On sait que la résolution de (40) par la méthode des variations, nous conduit à introduire un état adjoint λ tel que

$$\begin{bmatrix} \dot{X} \\ \dot{\lambda} \end{bmatrix} = \begin{bmatrix} A & -BR^{-1}B^T \\ -Q & -A^T \end{bmatrix} \begin{bmatrix} X \\ \lambda \end{bmatrix}$$

et la relation

$$\lambda = P.X$$

conduit alors à l'équation de Riccati.

Soit donc la matrice

$$H = \begin{bmatrix} A & -BR^{-1}B^T \\ -Q & -A^T \end{bmatrix} \quad \begin{matrix} \uparrow \\ 2n \\ \downarrow \end{matrix} \quad (46)$$

et la transformation

$$Z = \begin{bmatrix} C & I \\ -I & C \end{bmatrix} \quad \begin{matrix} \uparrow \\ 2n \\ \downarrow \end{matrix} \quad (47)$$

La matrice $H(2n, 2n)$ est hamiltonienne si et seulement si

$$H = ZH^T Z \quad (48)$$

c'est à dire, si on partitionne H comme suit

$$H = \begin{bmatrix} H_1 & H_{12} \\ H_{21} & H_2 \end{bmatrix}$$

on aura

$$H_1 = -H_2^T, \quad H_{12} = H_{12}^T, \quad H_{21} = H_{21}^T$$

Il est ainsi aisé de vérifier que H est une matrice hamiltonienne; de plus on a d'après (47):

Cela exprime donc que le système bouclé est asymptotiquement stable. Si on remplace l'hypothèse, (C, A) détectable, par (C, A) reconstructible (observable), on aura $P > 0$.

3-RESOLUTION

On sait que la résolution de (40) par la méthode des variations, nous conduit à introduire un état adjoint λ tel que

$$\begin{bmatrix} \dot{X} \\ \dot{\lambda} \end{bmatrix} = \begin{bmatrix} A & -BR^{-1}B^T \\ -Q & -A^T \end{bmatrix} \begin{bmatrix} X \\ \lambda \end{bmatrix}$$

et la relation

$$\lambda = P.X$$

conduit alors à l'équation de Riccati.

Soit donc la matrice

$$H = \begin{bmatrix} A & -BR^{-1}B^T \\ -Q & -A^T \end{bmatrix} \quad \begin{matrix} \uparrow \\ 2n \\ \downarrow \end{matrix} \quad (46)$$

et la transformation

$$Z = \begin{bmatrix} C & I \\ -I & C \end{bmatrix} \quad \begin{matrix} \uparrow \\ 2n \\ \downarrow \end{matrix} \quad (47)$$

La matrice $H(2n, 2n)$ est hamiltonienne si et seulement si

$$H = ZH^T Z \quad (48)$$

c'est à dire, si on partitionne H comme suit

$$H = \begin{bmatrix} H_1 & H_{12} \\ H_{21} & H_2 \end{bmatrix}$$

on aura

$$H_1 = -H_2^T, \quad H_{12} = H_{12}^T, \quad H_{21} = H_{21}^T$$

Il est ainsi aisé de vérifier que H est une matrice hamiltonienne; de plus on a d'après (47):

$$Z^{-1} = Z^T = -Z$$

Remarques

- les valeurs propres d'une matrice hamiltonienne apparaissent par paires opposées.
- si $(v_1, v_2)^T$ est un vecteur propre droit, $(-v_2, v_1)^T$ est un vecteur propre gauche.

Maintenant on est en mesure d'énoncer le théorème suivant:

THEOREME: étant donné le système invariant à temps final infini défini par les équations (40), et les hypothèses (42), alors la matrice

$$H = \begin{bmatrix} A & -BR^{-1}B^T \\ -Q & -A^T \end{bmatrix}$$

est telle que

1- réel $\lambda[H] \neq 0$

2- $\lambda i[\tilde{A}] = \lambda i[A - BR^{-1}B^T] = \lambda i[H] \quad i: \text{réel } \lambda i < 0$

3-

$$H = \begin{bmatrix} W_1 & W_{12} \\ W_{21} & W_2 \end{bmatrix} \begin{bmatrix} J & \\ & O \end{bmatrix} \begin{bmatrix} V_1 & V_{12} \\ V_{21} & V_2 \end{bmatrix}$$

4- $P = W_2 W_{12}^{-1} = -V_{12}^{-1} V_1$

5- $\tilde{A} = W_{12} (-J) W_2^{-1}$

où la troisième relation exprime la décomposition de H sous forme de Jordan, et la cinquième relation celle qui en découle pour \tilde{A} (système bouclé).

Remarque: la démonstration se trouve en annexe III

Ainsi d'après ce théorème on a une solution explicite pour

$$P = W_2 W_{12}^{-1} = -V_{12}^{-1} V_1$$

Mais si H n'est pas diagonalisable le problème est pratiquement insoluble. Pour contourner cette difficulté, on va aborder le problème comme suit:

Soit

$$H = U \begin{bmatrix} \tilde{A} & O \\ O & -\tilde{A} \end{bmatrix} U^{-1}$$

$$Z^{-1} = Z^T = -Z$$

Remarques

- les valeurs propres d'une matrice hamiltonienne apparaissent par paires opposées.
- si $(v_1, v_2)^T$ est un vecteur propre droit, $(-v_2, v_1)^T$ est un vecteur propre gauche.

Maintenant on est en mesure d'énoncer le théorème suivant:

THEOREME: étant donné le système invariant à temps final infini défini par les équations(40), et les hypothèses(42), alors la matrice

$$H \equiv \begin{bmatrix} A & -BR^{-1}B^T \\ -Q & -A^T \end{bmatrix}$$

est telle que

1- réel $\lambda[H] \neq 0$

2- $\lambda i[\tilde{A}] = \lambda i[A - BR^{-1}B^T] = \lambda i[H] \quad i: \text{réel } \lambda i < 0$

3-

$$H = \begin{bmatrix} W_1 & W_{12} \\ W_{21} & W_2 \end{bmatrix} \begin{bmatrix} J & 0 \\ 0 & -J \end{bmatrix} \begin{bmatrix} V_1 & V_{12} \\ V_{21} & V_2 \end{bmatrix}$$

4- $P = W_2 W_{12}^{-1} = -V_{12}^{-1} V_1$

5- $\tilde{A} = W_{12} (-J) W_2^{-1}$

où la troisième relation exprime la décomposition de H sous forme de Jordan, et la cinquième relation celle qui en découle pour \tilde{A} (système bouclé).

Remarque: la démonstration se trouve en annexe III

Ainsi d'après ce théorème on a une solution explicite pour

$$P = W_2 W_{12}^{-1} = -V_{12}^{-1} V_1$$

Mais si H n'est pas diagonalisable, le problème est pratiquement insoluble. Pour contourner cette difficulté, on va aborder le problème comme suit:

Soit

$$H = U \begin{bmatrix} \tilde{A} & 0 \\ 0 & -\tilde{A} \end{bmatrix} U^{-1}$$

d'où

$$S = \text{signe}(H) = U \begin{bmatrix} \text{signe}(\tilde{\lambda}) & 0 \\ 0 & \text{signe}(-\tilde{\lambda}) \end{bmatrix} U^{-1} \quad (49)$$

\tilde{A} étant asymptotiquement stable, donc
réel $\lambda_i[\tilde{A}] < 0$

d'où

$$\text{signe}(\tilde{\lambda}) = -1 \quad S = U \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} U^{-1} \quad (50)$$

Introduisons alors la matrice

$$F = \frac{1}{2}(1 + S) \quad (51)$$

c'est à dire

$$F = U \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} U^{-1}$$

en remplaçant U par sa valeur, il vient

$$F = \begin{bmatrix} VP & -V \\ -(1-PV)P & 1-PV \end{bmatrix} = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix}$$

d'où finalement par identification, on aboutit à

$$P = -F_{12}^{-1} \cdot F_{21} \quad (52)$$

On obtient ainsi une expression explicite pour P comparable à celle du théorème avec l'avantage d'éliminer la factorisation de Jordan.

Toute fois ce résultat n'est valable que si V est régulière.

D'après l'équation de Lyapunov (annexe), V est a priori symétrique non négative. La matrice A étant stable, V est donnée par

$$V = \int_0^{\infty} e^{\tilde{A}t} B^T B^{-1} B^T e^{-\tilde{A}t} dt$$

Or V doit être régulière, donc définie positive. Ceci n'est vrai que si la paire (\tilde{A}, B) est commandable, c'est à dire la paire (A, B) commandable

d'où

$$S = \text{signe}(H) = U \begin{bmatrix} \text{signe}(\tilde{A}) & 0 \\ 0 & \text{signe}(-\tilde{A}) \end{bmatrix} U^{-1} \quad (49)$$

\tilde{A} étant asymptotiquement stable, donc
réel $\lambda_i[\tilde{A}] < 0$

d'où

$$\text{signe}(\tilde{A}) = -1 \quad S = U \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} U^{-1} \quad (50)$$

Introduisons alors la matrice

$$F = \frac{1}{2}(1 + S) \quad (51)$$

c'est à dire

$$F = U \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} U^{-1}$$

en remplaçant U par sa valeur, il vient

$$F = \begin{bmatrix} VP & -V \\ -(1-PV)P & 1-PV \end{bmatrix} = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix}$$

d'où finalement par identification, on aboutit à

$$P = -F_{12}^{-1} \cdot F_{11} \quad (52)$$

On obtient ainsi une expression explicite pour P comparable à celle du théorème avec l'avantage d'éliminer la factorisation de Jordan.

Toutefois ce résultat n'est valable que si V est régulière.

D'après l'équation de Lyapunov (annexe), V est a priori symétrique non négative. La matrice \tilde{A} étant stable, V est donnée par

$$V = \int_0^{\infty} e^{\tilde{A}t} B B^{-1} B^T e^{\tilde{A}t} dt$$

Or V doit être régulière, donc définie positive. Ceci n'est vrai que si la paire (\tilde{A}, B) est commandable, c'est à dire la paire (A, B) commandable

En résumé on aura l'algorithme de calcul suivant

Calculer $S = \text{signe}(H)$

$$\text{Poser } F = \frac{1}{2}(I + S) = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix} \quad (53)$$

Calculer $P = -F_{12}^{-1} \cdot F_{11}$

Il ne nous reste plus maintenant qu'à évaluer le coût calcul de cet algorithme et le comparer à l'approche d'Anderson, qui jusqu'ici s'avère être la plus efficace.

4-COUT CALCUL

Le coût calcul de cet algorithme s'évalue comme suit:

- n^3 multiplications pour le calcul de l'inverse d'une matrice (n'étant sa dimension).

-K: nombre d'itérations nécessaires pour calculer $S = \text{signe de } (H)$

Donc le coût calcul est de l'ordre de $K(2n)^3 = 8Kn^3$. D'autre part le calcul de P nécessite approximativement $2n^3$ multiplications.

D'où le coût total de l'algorithme:

$$C_m = (8K + 2)n^3 \quad (54)$$

Compte tenu de la nature du problème de la commande que l'on traite, les pôles du système bouclé (valeurs propres de H) sont situés dans le plan complexe au voisinage de l'axe des réels; on peut donc admettre que

$$|\arg \theta| < t \quad t \approx 1$$

θ argument des pôles complexes.

Ce qui permettra un bon conditionnement de la matrice H vis à vis du calcul de sa matrice signe.

En résumé on aura l'algorithme de calcul suivant

Calculer $S = \text{signe}(H)$

$$\text{Poser } F = \frac{1}{2}(I + S) = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix} \quad (53)$$

Calculer $P = -F_{12}^{-1} \cdot F_{11}$

Il ne nous reste plus maintenant qu'à évaluer le coût calcul de cet algorithme et le comparer à l'approche d'Anderson, qui jusqu'ici s'avère être la plus efficace.

4-COUT CALCUL

Le coût calcul de cet algorithme s'évalue comme suit:

- n^3 multiplications pour le calcul de l'inverse d'une matrice (n'étant sa dimension).

-K: nombre d'iterations nécessaires pour calculer $S = \text{signe de } (H)$

Donc le coût calcul est de l'ordre de $K(2n)^3 = 8Kn^3$. D'autre part le calcul de P nécessite approximativement $2n^3$ multiplications.

D'où le coût total de l'algorithme:

$$C_m = (8K + 2)n^3 \quad (54)$$

Compte tenu de la nature du problème de la commande que l'on traite, les pôles du système bouclé (valeurs propres de H) sont situés dans le plan complexe au voisinage de l'axe des réels; on peut donc admettre que

$$|\text{tg } \theta| < t \quad t \approx 1$$

θ argument des pôles complexes.

Ce qui permettra un bon conditionnement de la matrice H vis à vis du calcul de sa matrice signe.

En fait dans la plupart des situations, le nombre d'itérations K varie de 5 à 20 en fonction du spectre du système bouclé indépendamment de la dimension n du problème et du spectre du système en boucle ouverte.

5-COMPARAISON AVEC L'APPROCHE D'ANDERSON:

Le principe de l'algorithme d'ANDERSON consiste à poser:

$$H_K = \begin{bmatrix} A_K & B_K \\ -C_K & -A_K^T \end{bmatrix} \quad (55)$$

d'où

$$A_{K+1} = \frac{1}{2}(A_K + A_K^{-1})$$

$$B_{K+1} = \frac{1}{2}(B_K + B_K^{-1})$$

$$C_{K+1} = \frac{1}{2}(C_K + C_K^{-1})$$

Calculons A_K^{-1} , B_K^{-1} et C_K^{-1}

Soit

$$M = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \quad M^{-1} = \begin{bmatrix} A & B^{-1} & E & F \\ C & D & G & H \end{bmatrix}$$

donc

$$E = (A - BD^{-1}C)^{-1}$$

$$G = -D^{-1}C(A - BD^{-1}C)^{-1}$$

$$F = -A^{-1}B(D - CA^{-1}B)^{-1}$$

$$H = (D - CA^{-1}B)^{-1}$$

d'où

$$A_{K+1} = \frac{1}{2} \left[A_K + \left[A_K + B_K (A_K^T)^{-1} C_K \right]^{-1} \right]$$

$$B_{K+1} = \frac{1}{2} \left[B_K + \left[A_K + B_K (A_K^T)^{-1} C_K \right]^{-1} B_K (A_K^T)^{-1} \right]$$

En fait dans la plupart des situations, le nombre d'itérations K varie de 5 à 20 en fonction du spectre du système bouclé indépendamment de la dimension n du problème et du spectre du système en boucle ouverte.

5-COMPARAISON AVEC L'APPROCHE D'ANDERSON:

Le principe de l'algorithme d'ANDERSON consiste à poser:

$$H_K = \begin{bmatrix} A_K & B_K \\ -C_K & -A_K^T \end{bmatrix} \quad (55)$$

d'où

$$A_{K+1} = \frac{1}{2}(A_K + A_K^{-1})$$

$$B_{K+1} = \frac{1}{2}(B_K + B_K^{-1})$$

$$C_{K+1} = \frac{1}{2}(C_K + C_K^{-1})$$

Calculons A_K^{-1} , B_K^{-1} et C_K^{-1}

Soit

$$M = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \quad M^{-1} = \begin{bmatrix} A & B^{-1} & E & F \\ C & D & G & H \end{bmatrix}$$

donc

$$E = (A - B D^{-1} C)^{-1}$$

$$G = -D^{-1} C (A - B D^{-1} C)^{-1}$$

$$F = -A^{-1} B (D - C A^{-1} B)^{-1}$$

$$H = (D - C A^{-1} B)^{-1}$$

d'où

$$A_{K+1} = \frac{1}{2} \left[A_K + \left[A_K + B_K (A_K^T)^{-1} C_K \right]^{-1} \right]$$

$$B_{K+1} = \frac{1}{2} \left[B_K + \left[A_K + B_K (A_K^T)^{-1} C_K \right]^{-1} B_K (A_K^T)^{-1} \right]$$

$$C_{K+1} = \frac{1}{2} \left[C_K + (A_K^T)^{-1} C_K \left[A_K + B_K (A_K^T)^{-1} C_K \right]^{-1} \right]$$

Comparativement à cette approche d'ANDERSON [4][5], notre algorithme ne tient pas compte de la structure de l'hamiltonien pour calculer sa matrice signe.

L'avantage de cette approche (d'ANDERSON) pourrait être le gain en place mémoire et en coût calcul.

Cependant cette économie de mémoires est inférieure à $2n^2$ ($6n^2 + n$ au lieu de $8n^2$) et celle en multiplication de n^3 ($7n^3$ au lieu de $8n^3$).

Mais vu que l'implémentation est beaucoup moins simple, cette économie est relativement faible. De plus cette solution a recours à 2 inversions de matrices $(A_K^T)^{-1}$ et $\left[A_K + B_K (A_K^T)^{-1} C_K \right]^{-1}$ dont il faut connaître le conditionnement au cours des itérations, en particulier la matrice du système en boucle ouverte.

$$C_{K+1} = \frac{1}{2} \left[C_K + (A_K^T)^{-1} C_K \left[A_K + B_K (A_K^T)^{-1} C_K \right]^{-1} \right]$$

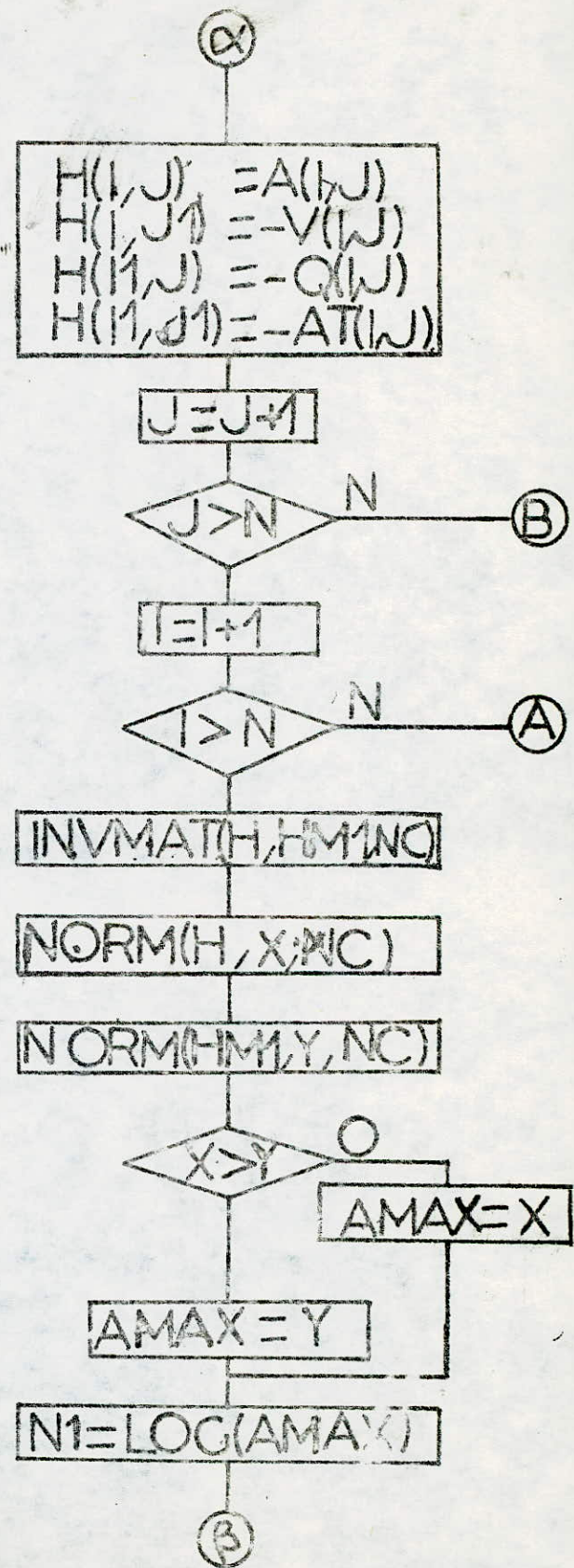
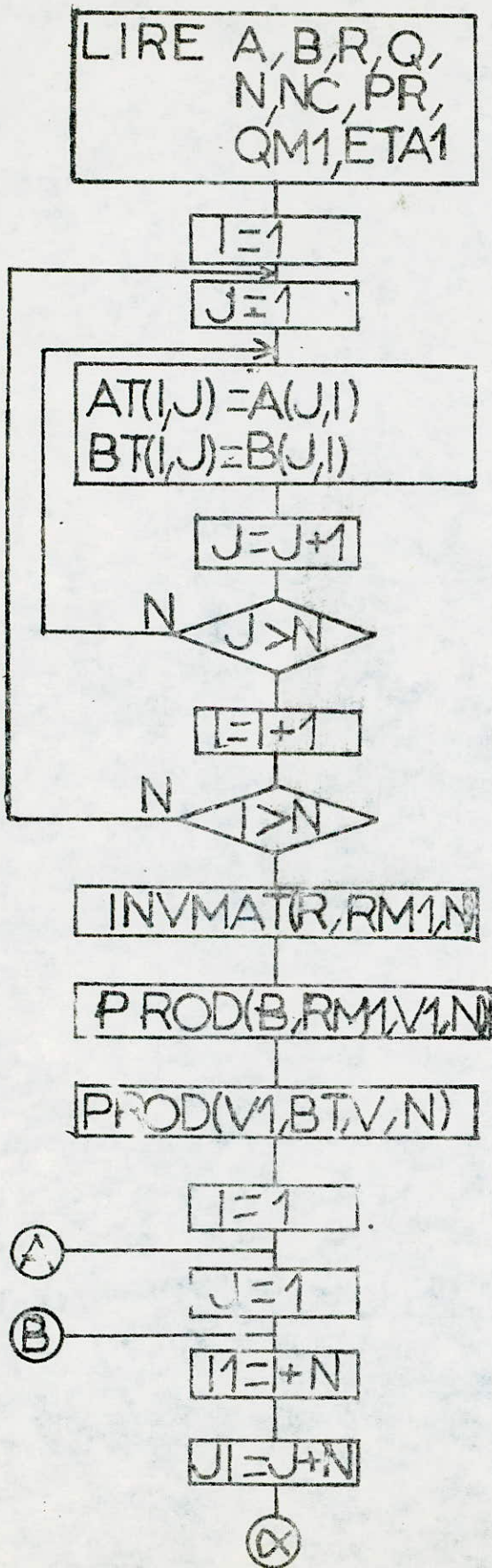
Comparativement à cette approche d'ANDERSON [4][5], notre algorithme ne tient pas compte de la structure de l'hamiltonien pour calculer sa matrice signe.

L'avantage de cette approche (d'ANDERSON) pourrait être le gain en place mémoire et en coût calcul.

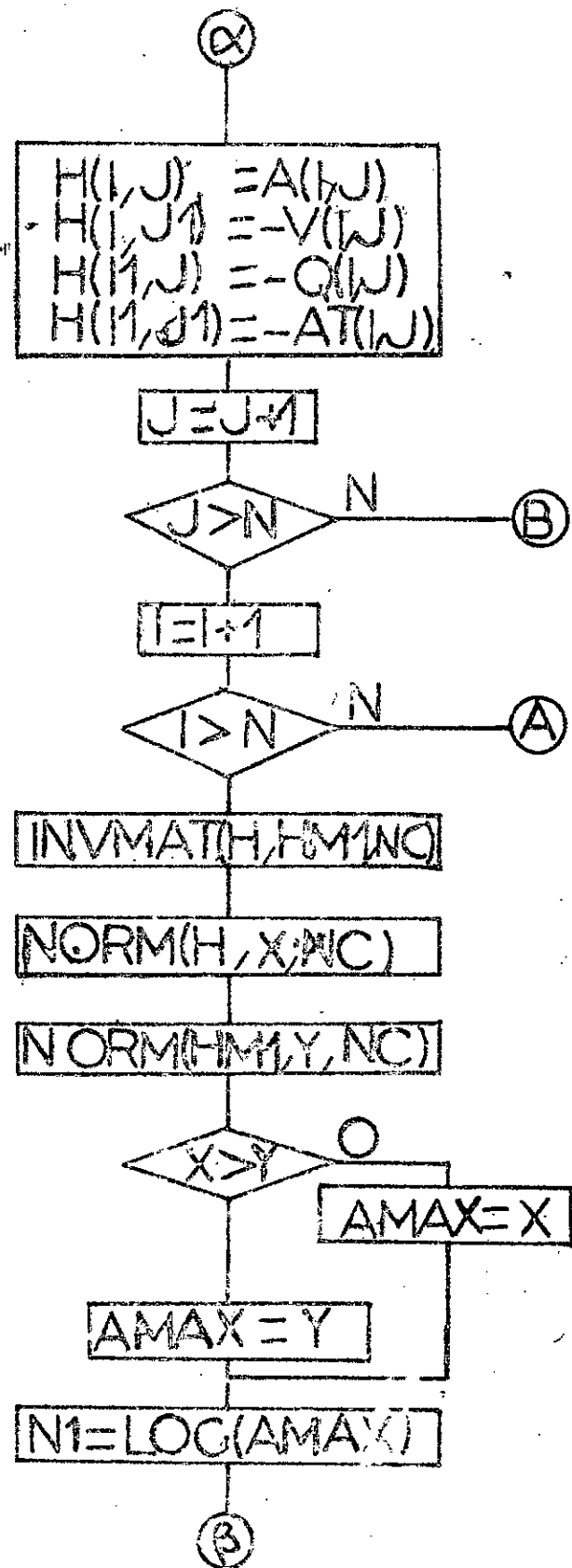
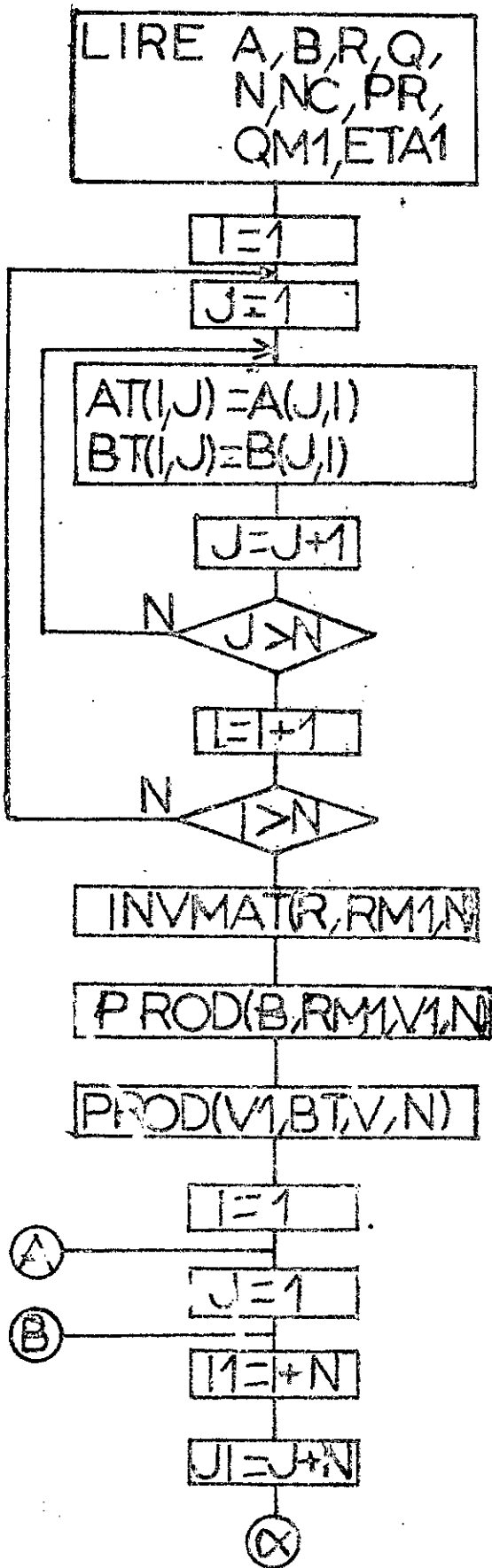
Cependant cette économie de mémoires est inférieure à $2n^2$ ($6n^2 + n$ au lieu de $8n^2$) et celle en multiplication de n^3 ($7n^3$ au lieu de $8n^3$).

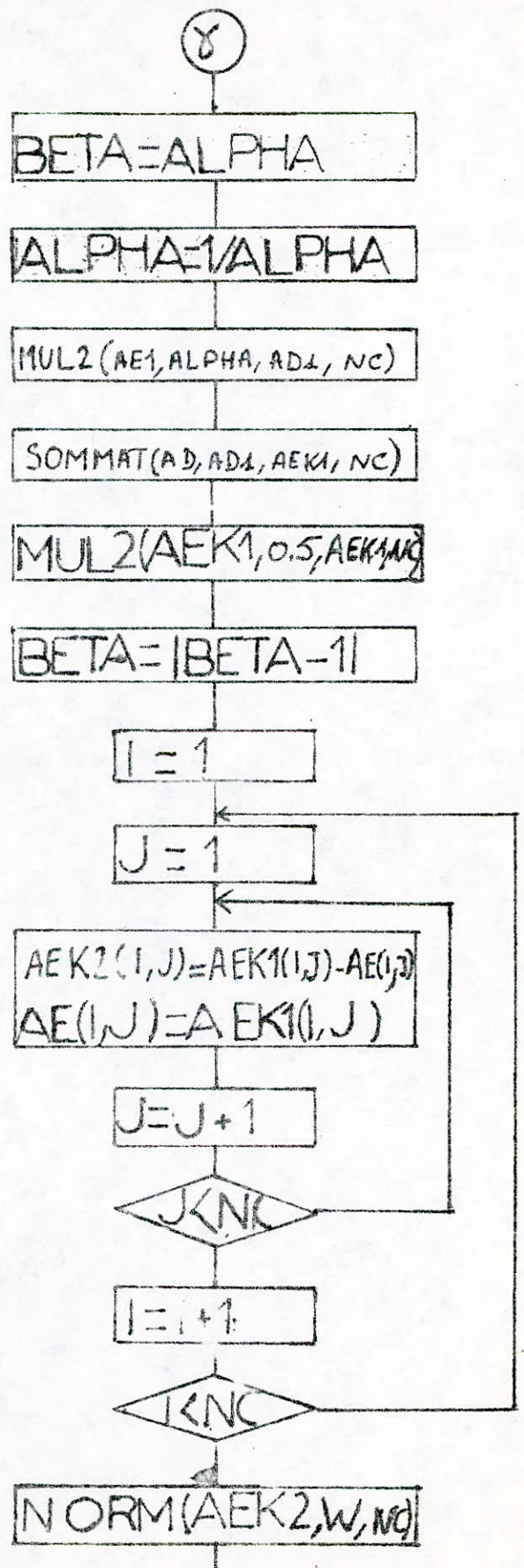
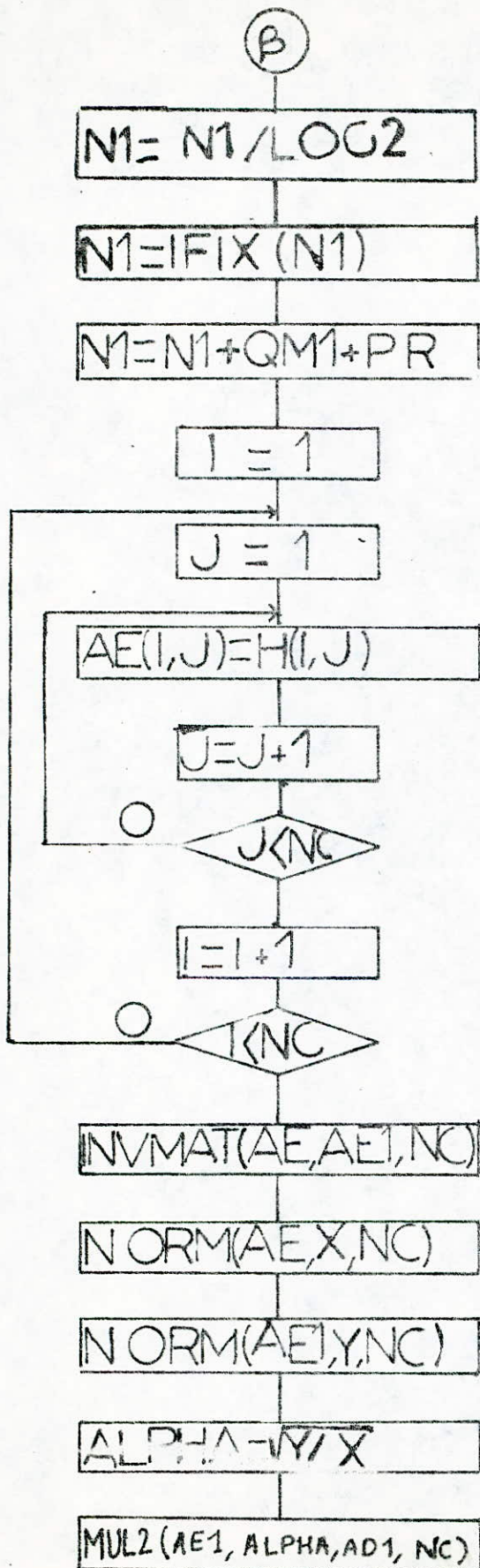
Mais vu que l'implimentation est beaucoup moins simple, cette économie est relativement faible. De plus cette solution a recours à 2 inversions de matrices $(A_K^T)^{-1}$ et $\left[A_K + B_K (A_K^T)^{-1} C_K \right]^{-1}$ dont il faut connaître le conditionnement au cours des itérations, en particulier la matrice du système en boucle ouverte.

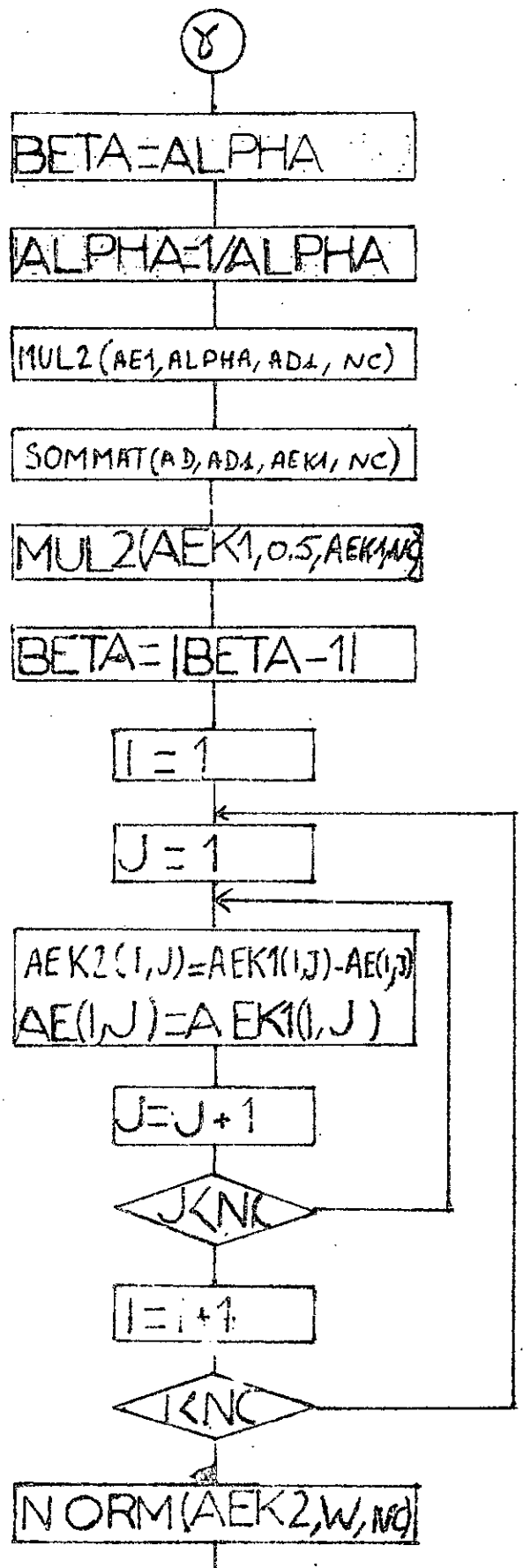
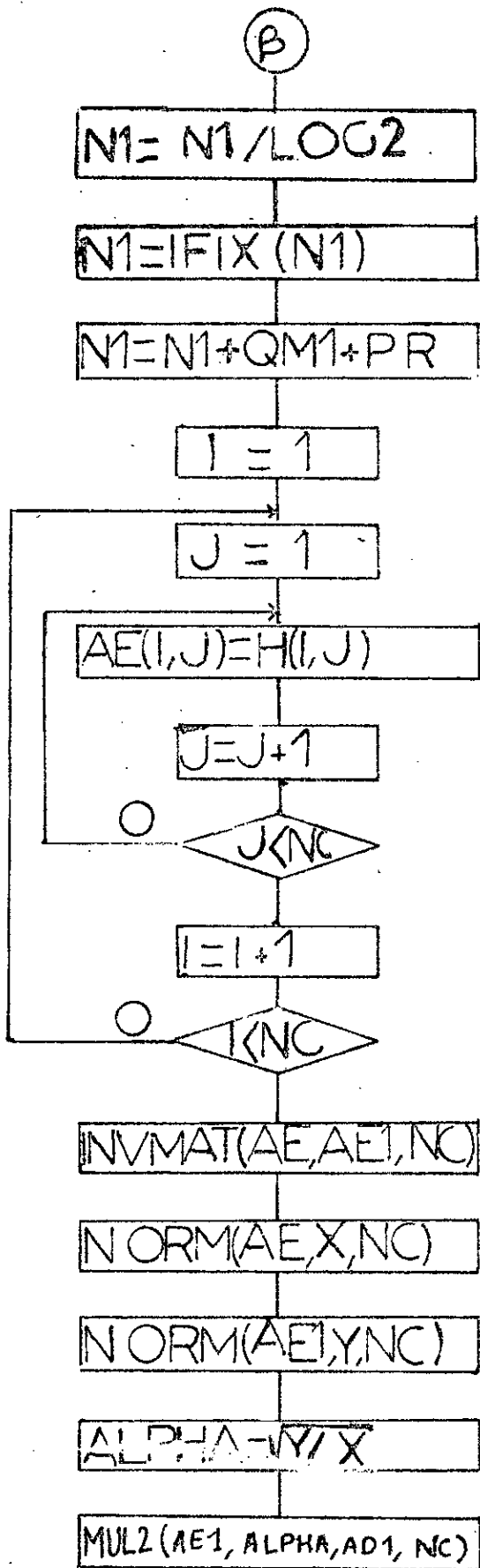
ORGANIGRAMME PRINCIPAL

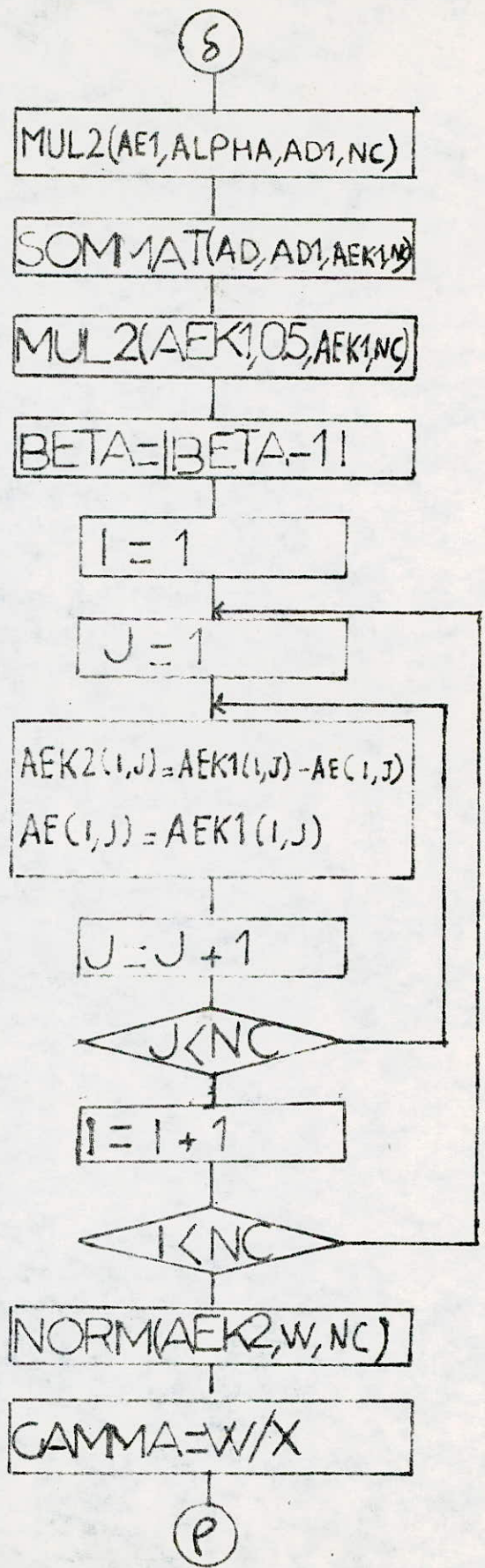
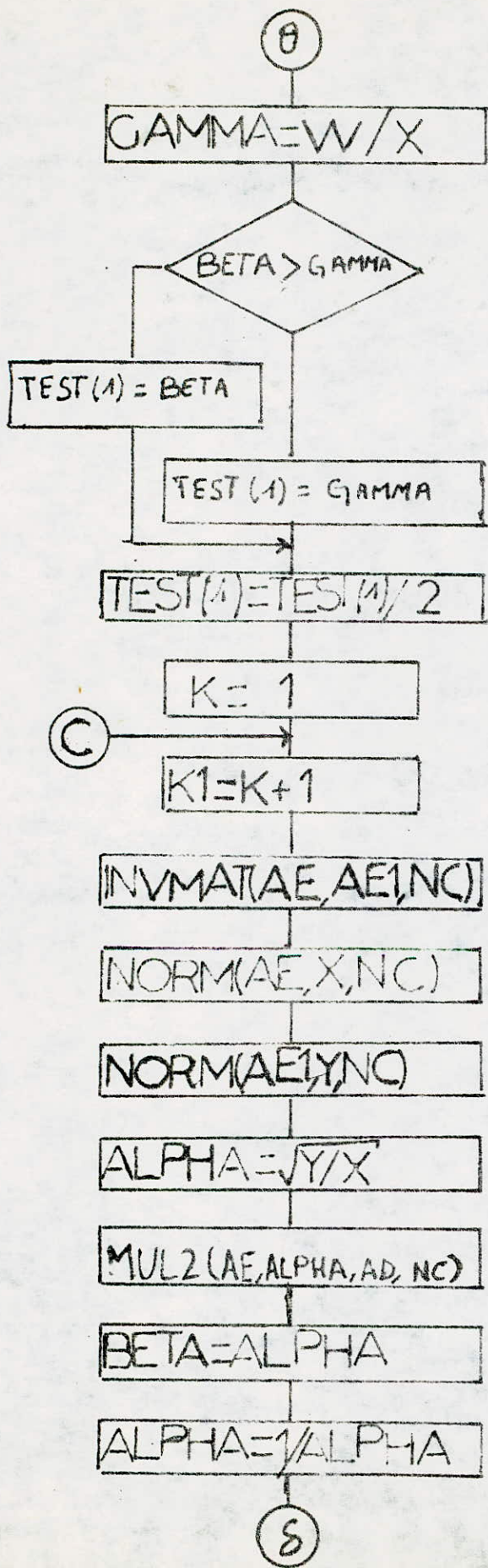


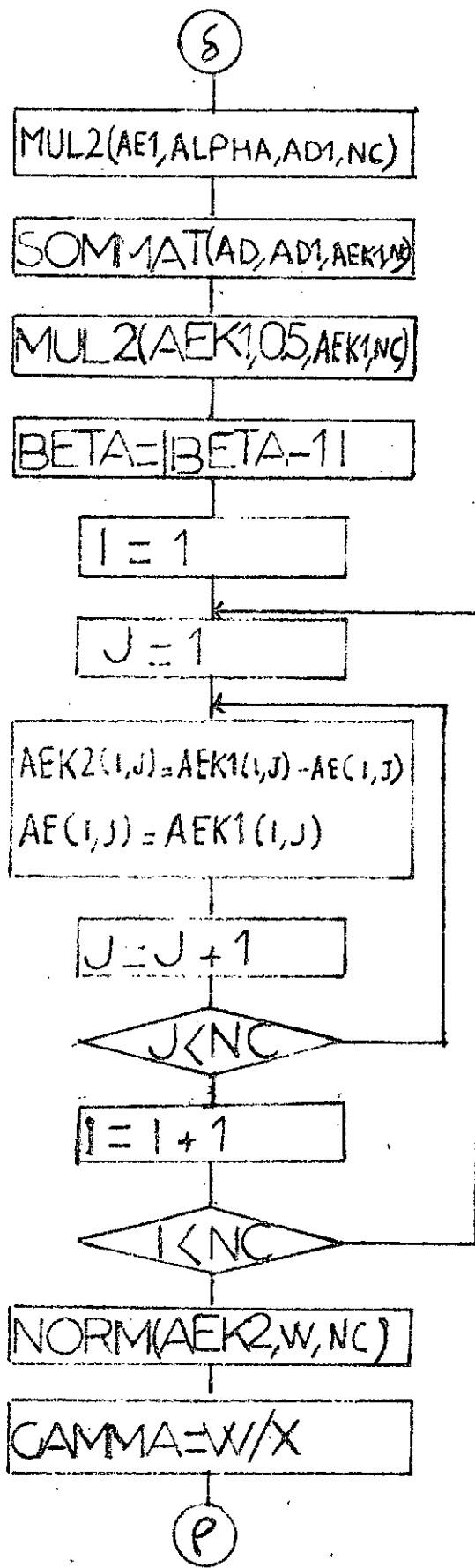
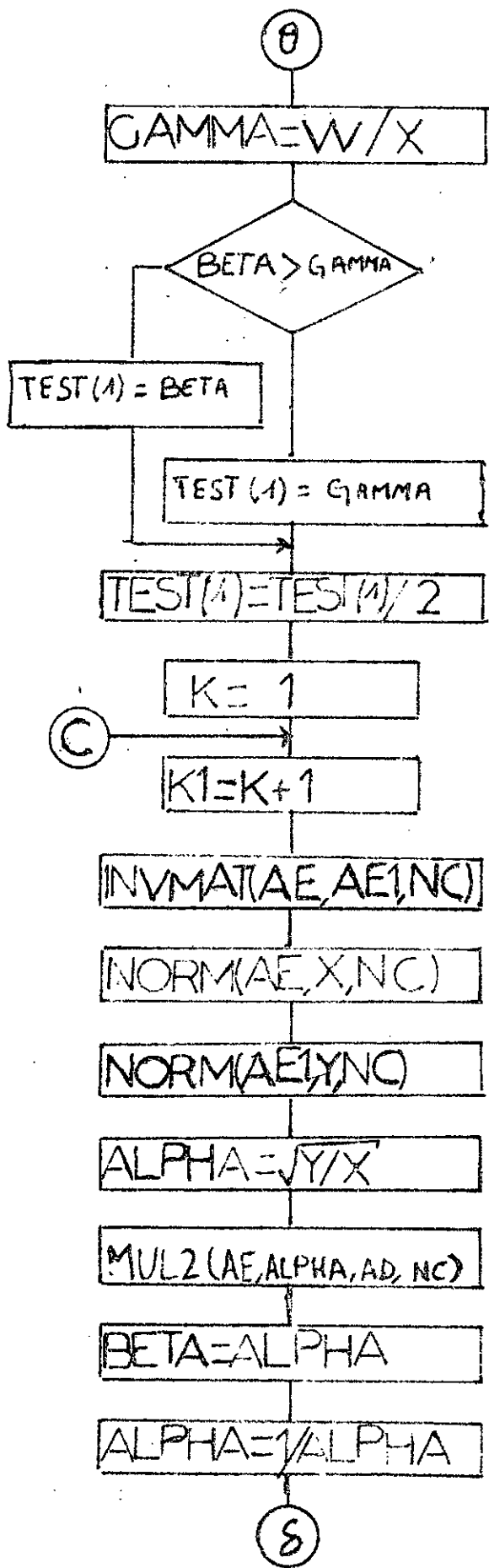
ORGANIGRAMME PRINCIPAL

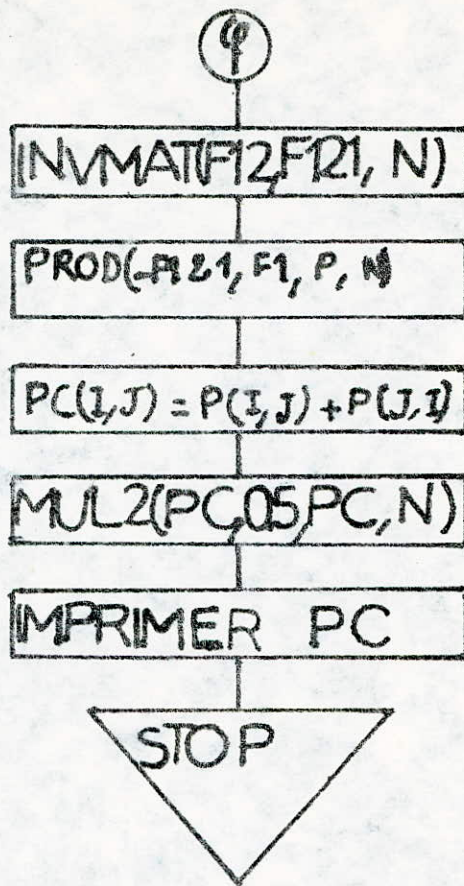
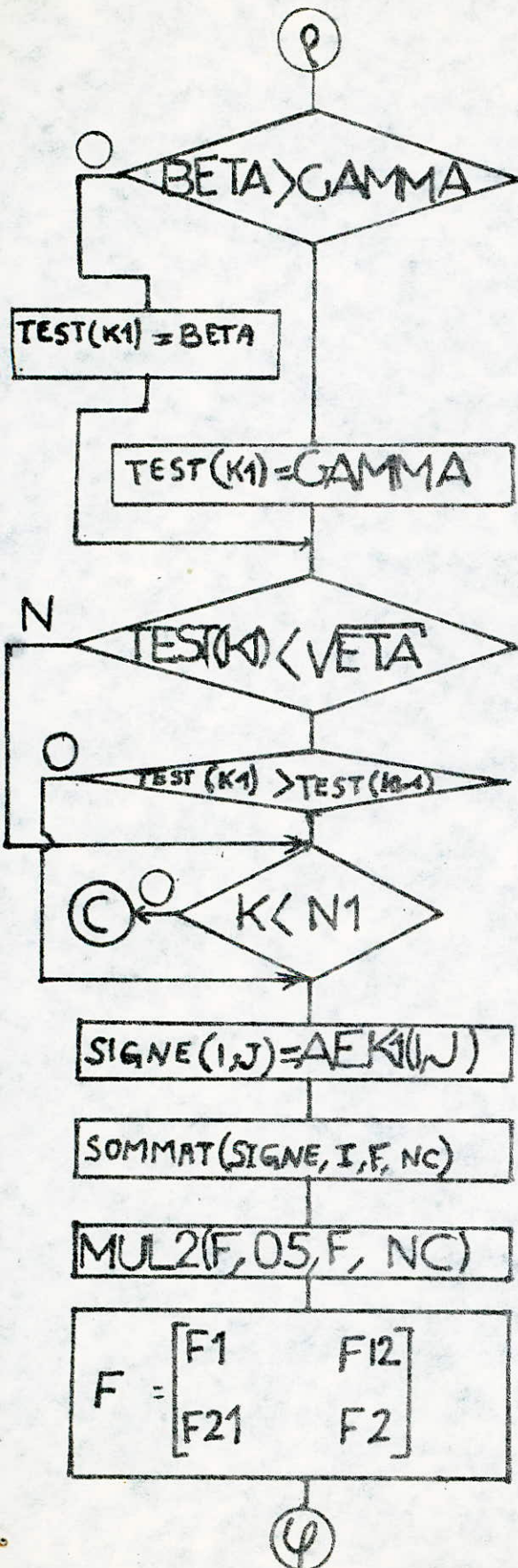


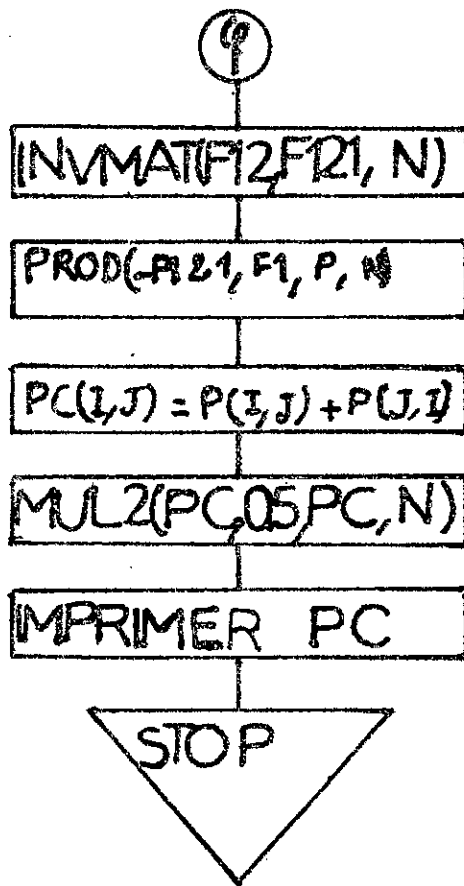
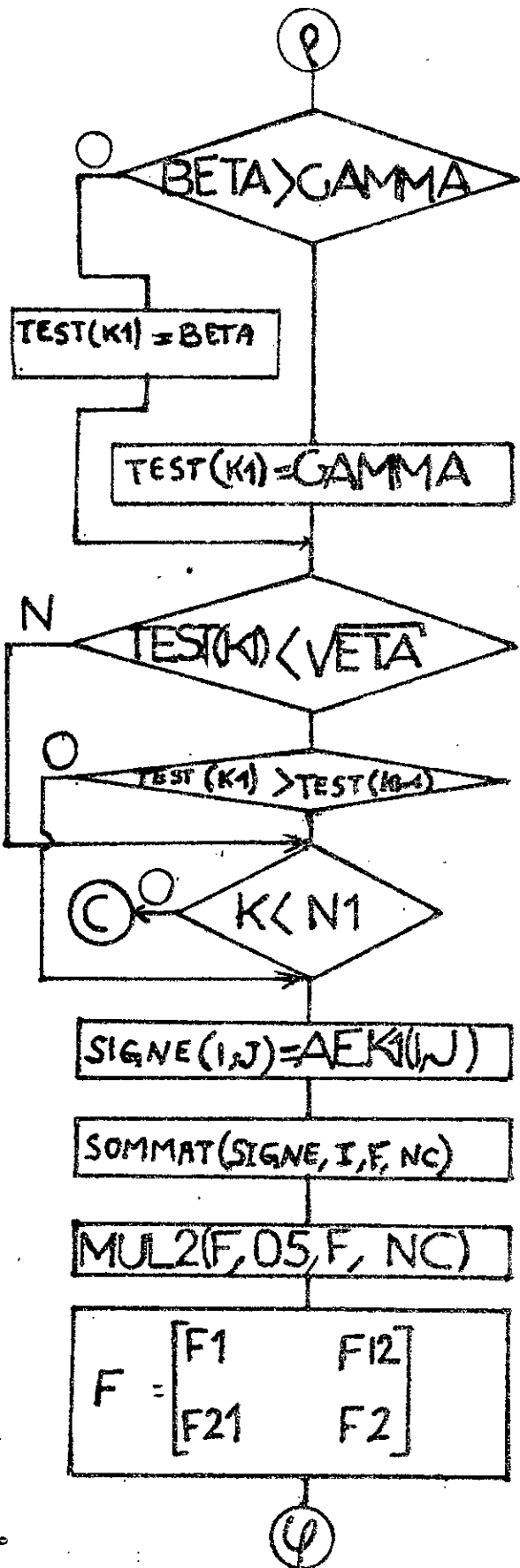




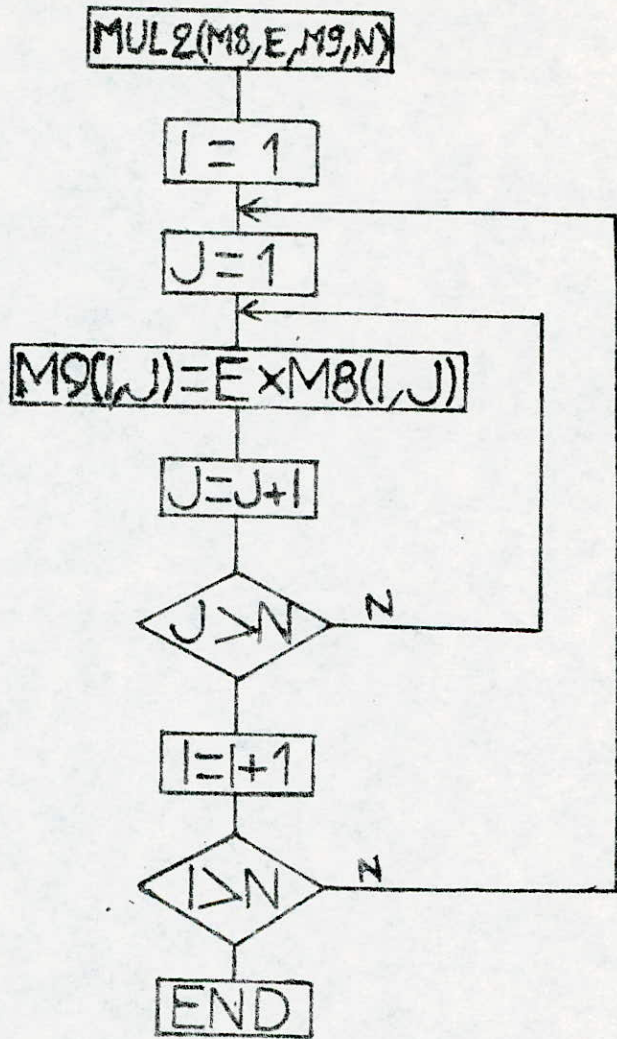
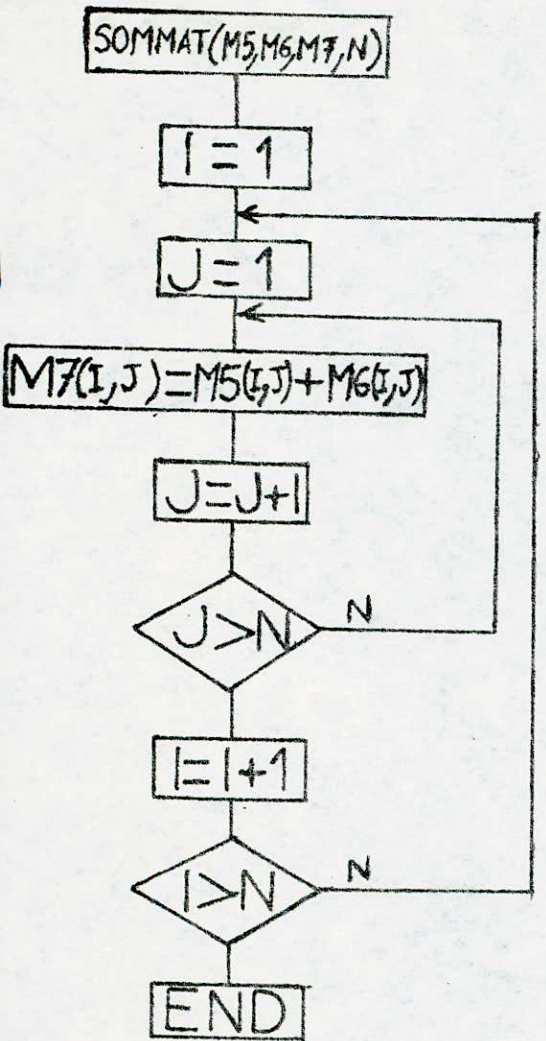




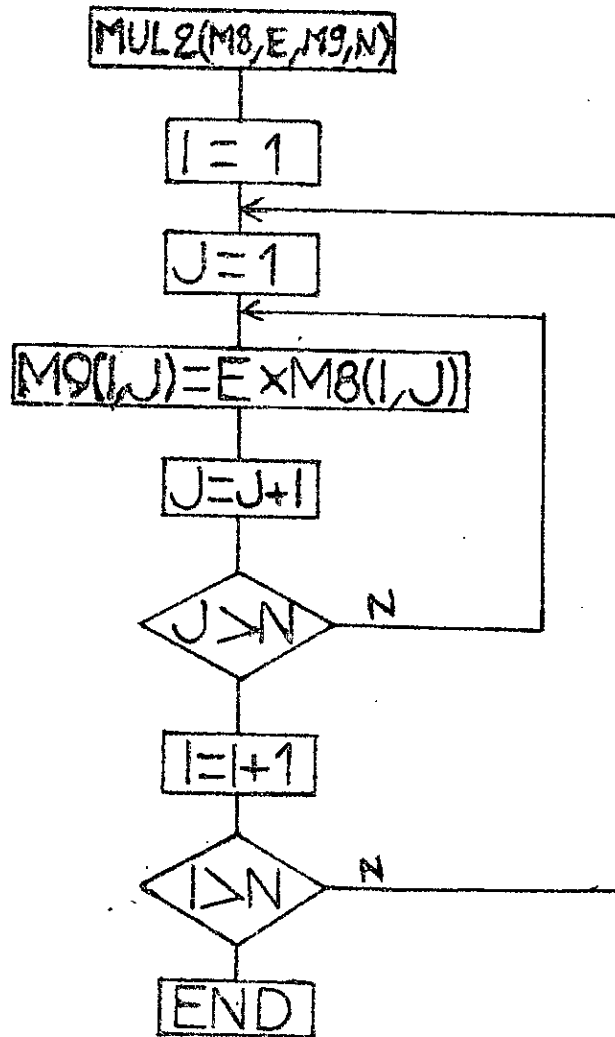
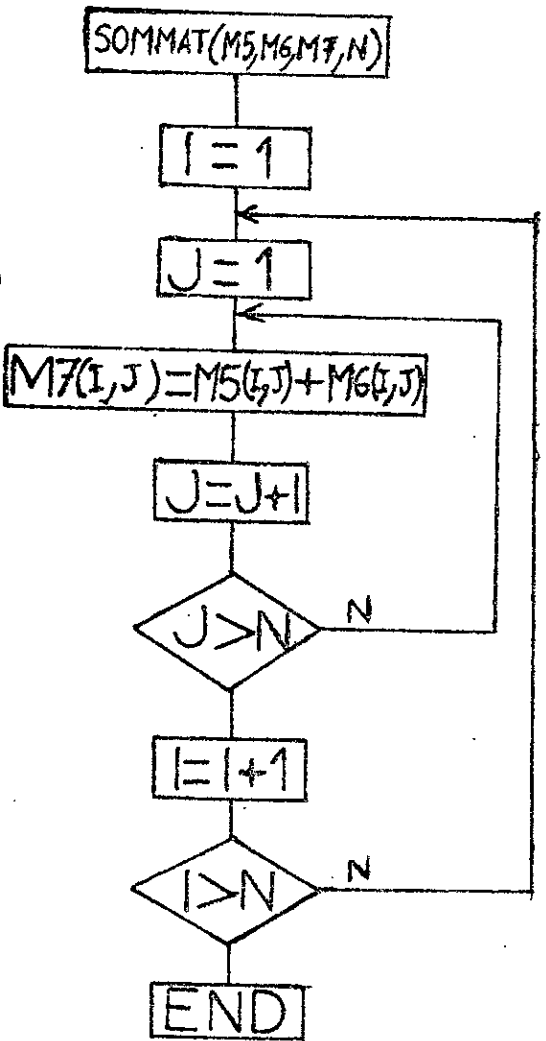




SOUS PROGRAMMES



SOUS PROGRAMMES



PROD(X,Y,Z,N)

I = 1

K = 1

Z(I,K) = 0

J = 1

$Z(I,K) = Z(I,K) + X(I,J) \times Y(J,K)$

J = J + 1

J > N

K = K + 1

K > N

I = I + 1

I > N

END

NORM(A,N,Z,N)

Z1 = 0

Z2 = 0

Z3 = 0

J = 1

TAB(J) = 0

J = J + 1

J < N

I = 1

TABL(I) = C

I = I + 1

I < N

J = 1

I = 1

^

α

β

PROD(X,Y,Z,N)

I = 1

K = 1

Z(I,K) = 0

J = 1

$Z(I,K) = Z(I,K) + X(I,J) \times Y(J,K)$

J = J + 1

J > N

K = K + 1

K > N

I = I + 1

I > N

END

NORM(A,N,Z,N)

Z1 = 0
Z2 = 0
Z3 = 0

J = 1

TAB(J) = 0

J = J + 1

J < N

I = 1

TABL(I) = C

I = I + 1

I < N

J = 1

I = 1

^

α

β

(A)

TABL(J) = TABL(J) + AN(I, J)

I = I + 1

J < N

J = J + 1

J < N

I = 1

Z1 > TABL(I)

Z1 = TABL(J)

J = J + 1

J < N

I = 1

J = 1

(B)

(B)

TABL(I) = TABL(I) + AN(I, J)

J = J + 1

J < N

I = I + 1

I < N

I = 1

Z2 > TABL(I)

Z2 = TABL(I)

I = I + 1

I < N

Z1 < Z2

Z = Z1

Z = Z2

END

(A)

TABL(J) = TABL(J) + A(I, J)

I = I + 1

J < N

I = J + 1

J < N

I = 1

Z1 > TABL(I)

Z1 = TABL(J)

J = J + 1

J < N

I = 1

J = 1

(B)

(B)

TABL(I) = TABL(I) + A(I, J)

J = J + 1

J < N

I = I + 1

I < N

I = 1

Z2 > TABL(I)

Z2 = TABL(I)

I = I + 1

I < N

Z1 < Z2

Z = Z1

Z = Z2

END

RESOLUTION DE L'EQUATION DE RICCATI

DANS LE CAS DISCRET

I-INTRODUCTION

On considère ici la résolution numérique de l'équation algébrique de Riccati,

$$A^T P A - P - A^T P B (B^T P B + R)^{-1} B^T P A + Q = 0$$

introduite en termes de commande des systèmes discrets; on développe, pour cela, le concept de produit étoile et les propriétés qui s'y rattachent.

L'application de ces résultats conduit à une méthode déduite de la solution apportée au problème continu ainsi qu'à un algorithme itératif du type "Square root" ($P = L.L^T$) présentant une convergence d'ordre 2. La première méthode qui est en fait une nouvelle voie constituant un traitement unifié des équations de Riccati à l'aide de la fonction signe de matrice, est alors comparée à cet algorithme ainsi qu'à d'autres méthodes, les unes considérées déjà comme classiques, les autres un peu plus récentes.

2-DEFINITION DU PROBLEME

On se propose de résoudre l'équation algébrique de Riccati

$$P = A^T P A - A^T P B (B^T P B + R)^{-1} B^T P A + Q \quad (56)$$

introduite ici dans le cadre de la commande optimale des systèmes discrets linéaires invariants.

Soit donc le système

$$X_{t+1} = A.X_t + B.U_t \quad (57)$$

et le critère quadratique à minimiser $J = \frac{1}{2} \sum_{t=0}^{\infty} (U_t^T R U_t + X_t^T Q X_t)$ (58)

moyennant les hypothèses suivantes:

$$(A, B) \text{ stabilisable} \quad (59a)$$

$$(C, A) \text{ détectable} \quad C^m C = Q \succcurlyeq 0 \quad (59b)$$

RESOLUTION DE L'EQUATION DE RICCATI

DANS LE CAS DISCRET

I-INTRODUCTION

On considère ici la résolution numérique de l'équation algébrique de Riccati,

$$A^T P A - P - A^T P B (B^T P B + R)^{-1} B^T P A + Q = 0$$

introduite en termes de commande des systèmes discrets; on développe, pour cela, le concept de produit étoile et les propriétés qui s'y rattachent.

L'application de ces résultats conduit à une méthode déduite de la solution apportée au problème continu ainsi qu'à un algorithme itératif du type "Square root" ($P = L L^T$) présentant une convergence d'ordre 2. La première méthode qui est en fait une nouvelle voie constituant un traitement unifié des équations de Riccati à l'aide de la fonction signe de matrice, est alors comparée à cet algorithme ainsi qu'à d'autres méthodes, les unes considérées déjà comme classiques, les autres un peu plus récentes.

2-DEFINITION DU PROBLEME

On se propose de résoudre l'équation algébrique de Riccati

$$P = A^T P A - A^T P B (B^T P B + R)^{-1} B^T P A + Q \quad (56)$$

introduite ici dans le cadre de la commande optimale des systèmes discrets linéaires invariants.

Soit donc le système

$$X_{t+1} = A X_t + B U_t \quad (57)$$

et le critère quadratique à minimiser $J = \frac{1}{2} \sum_{t=0}^{\infty} (U_t^T R U_t + X_t^T Q X_t)$ (58)

moyennant les hypothèses suivantes:

$$(A, B) \text{ stabilisable} \quad (59a)$$

$$(C, A) \text{ détectable } C^m C = Q \succeq C \quad (59b)$$

$$B^T Q B + R = (B^T Q B + R)^T \quad (59c)$$

On sait $\begin{bmatrix} I \\ \end{bmatrix}$ que

$$P + P^T \geq 0 \quad P \text{ unique} \quad (60)$$

$$|\lambda_i(\tilde{A})| < 1 \quad \tilde{A} = A - B (B^T P B + R)^{-1} B^T P A \quad (61)$$

Cette dernière conclusion expriment que le système bouclé par la commande

$$\hat{U}_t = - \begin{bmatrix} (B^T P B + R)^{-1} B^T P A \end{bmatrix} X_t \quad (62)$$

est asymptotiquement stable

enfin le coût optimal $J(\hat{U})$ est donné par :

$$J = \frac{1}{2} X_0^T P X_0 \quad (63)$$

Considérons maintenant l'état adjoint λ_t que l'on peut introduire pour résoudre le problème de minimisation avec contraintes défini par (57) et (58).

D'une part on a

$$\lambda_t = P_t X_t \quad (64)$$

et d'autre part l'ensemble des équations hamiltoniennes discrètes s'écrivent :

$$\begin{bmatrix} X_{t+1} \\ \lambda_t \end{bmatrix} = \begin{bmatrix} A & -BR^{-1}B^T \\ Q & A^T \end{bmatrix} \cdot \begin{bmatrix} X_t \\ \lambda_{t+1} \end{bmatrix} \quad (65)$$

Notons

$$M = \begin{bmatrix} A & -BR^{-1}B^T \\ Q & A^T \end{bmatrix} \quad (66)$$

$$B^T Q B + R = (B^T Q B + R)^T \quad (59c)$$

On sait $\begin{bmatrix} I \\ \end{bmatrix}$ que

$$P + P^T \geq 0 \quad P \text{ unique} \quad (60)$$

$$|\lambda_i(\tilde{A})| < 1 \quad \tilde{A} = A - B (B^T P B + R)^{-1} B^T P A \quad (61)$$

Cette dernière conclusion expriment que le système bouclé par la commande

$$\hat{U}_t = - \begin{bmatrix} (B^T P B + R)^{-1} B^T P A \end{bmatrix} X_t \quad (62)$$

est asymptotiquement stable

enfin le coût optimal $J(U)$ est donné par :

$$J = \frac{1}{2} X_0^T P X_0 \quad (63)$$

Considérons maintenant l'état adjoint λ_t que l'on peut introduire pour résoudre le problème de minimisation avec contraintes défini par (57) et (58).

D'une part on a

$$\lambda_t = P_t X_t \quad (64)$$

et d'autre part l'ensemble des équations hamiltoniennes discrètes s'écrivent :

$$\begin{bmatrix} X_{t+1} \\ \lambda_t \end{bmatrix} = \begin{bmatrix} A & -BR^{-1}B^T \\ Q & A^T \end{bmatrix} \cdot \begin{bmatrix} X_t \\ \lambda_{t+1} \end{bmatrix} \quad (65)$$

Notons

$$M = \begin{bmatrix} A & -BR^{-1}B^T \\ Q & A^T \end{bmatrix} \quad (66)$$

La matrice associée à ces équations, précisons le, n'est pas elle-même hamiltonienne, comme dans le cas continu, ni symplectique (annexe) bien que l'on puisse en déduire une forme symplectique à l'aide d'une transformation notée $\Psi(\cdot)$, que l'on introduira au paragraphe suivant.

3- PRODUIT SIMPLE ET PRODUIT " étoile " DE MATRICES

Tout un chacun connaît l'illustration classique du produit de matrices 2×2 à travers la combinaison en série de quadripôles (Fig I) :

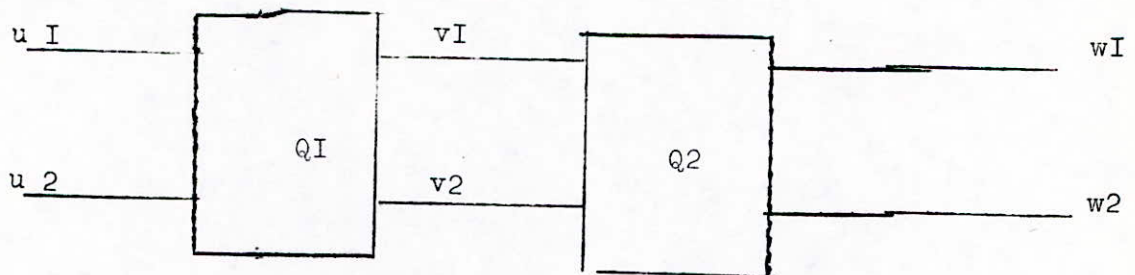


Figure I

que l'on formalise comme suit :

$$\begin{bmatrix} v1 \\ v2 \end{bmatrix} = \begin{bmatrix} a1 & b1 \\ c1 & d1 \end{bmatrix} \times \begin{bmatrix} u1 \\ u2 \end{bmatrix} \quad \Longrightarrow \quad v = Q_1 u \quad (67)$$

$$\begin{bmatrix} w1 \\ w2 \end{bmatrix} = \begin{bmatrix} a2 & b2 \\ c2 & d2 \end{bmatrix} \times \begin{bmatrix} v1 \\ v2 \end{bmatrix} \quad \Longrightarrow \quad w = Q_2 v \quad (68)$$

Le quadripôle équivalent Q est tel que

$$W = Q u \quad (69)$$

$$Q = Q_2 Q_1$$

d'où

$$Q = \begin{bmatrix} a2 a1 + b2 c1 & a2 b1 + b2 d1 \\ c2 a1 + d2 c1 & c2 b1 + d2 d1 \end{bmatrix} \quad (70)$$

La matrice associée à ces équations, précisons le, n'est pas elle-même hamiltonienne, comme dans le cas continu, ni symplectique (annexe) bien que l'on puisse en déduire une forme symplectique à l'aide d'une transformation notée $\Psi(\cdot)$, que l'on introduira au paragraphe suivant.

3- PRODUIT SIMPLE ET PRODUIT " étoile " DE MATRICES

Tout un chacun connaît l'illustration classique du produit de matrices 2X2 à travers la combinaison en série de quadripôles (Fig I) :

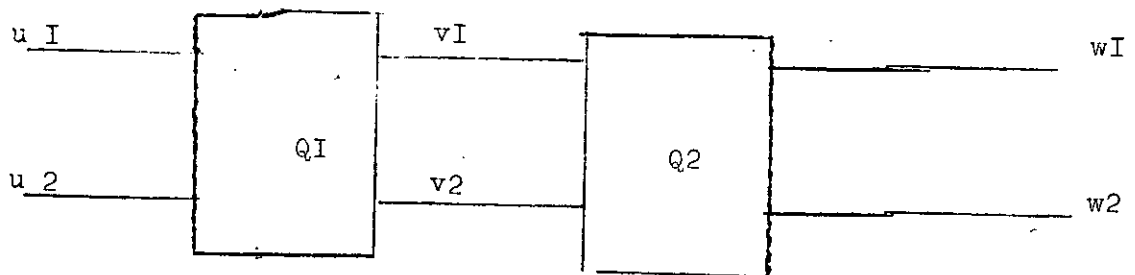


Figure I

que l'on formalise comme suit :

$$\begin{bmatrix} v1 \\ v2 \end{bmatrix} = \begin{bmatrix} a1 & b1 \\ c1 & d1 \end{bmatrix} \times \begin{bmatrix} u1 \\ u2 \end{bmatrix} \quad \Longrightarrow \quad v = Q_1 u \quad (67)$$

$$\begin{bmatrix} w1 \\ w2 \end{bmatrix} = \begin{bmatrix} a2 & b2 \\ c2 & d2 \end{bmatrix} \times \begin{bmatrix} v1 \\ v2 \end{bmatrix} \quad \Longrightarrow \quad w = Q_2 v \quad (68)$$

Le quadripôle équivalent Q est tel que

$$W = Q u \quad (69)$$

$$Q = Q_2 Q_1$$

d'où

$$Q = \begin{bmatrix} a2 a1 + b2 c1 & a2 b1 + b2 d1 \\ c2 a1 + d2 c1 & c2 b1 + d2 d1 \end{bmatrix} \quad (70)$$

Ce produit simple se généralise comme on le sait aux matrices de dimension quelconque.

IL est distributif par rapport à l'addition, associatif, possède un élément nul et un élément neutre.

Ainsi l'ensemble des matrices (n,n) muni de la loi d'addition et de ce produit forme un anneau et constitue une algèbre.

On va introduire maintenant un second type de produit matriciel basé lui aussi sur la composition en série de quadripôles et montre qu'il possède les mêmes propriétés que le produit simple.

Considérons donc le problème de propagation d'ondes suivant:

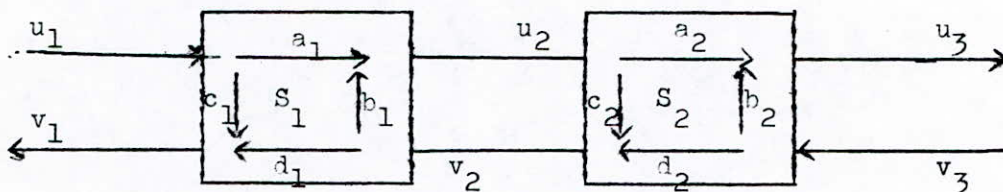


Figure 2

u_1, u_2 et u_3 sont des rayons incidents

v_1, v_2 et v_3 sont des rayons émergents

les milieux séparant les phases 1, 2 et 3 sont caractérisés par des Coefficients de réflexion de transmission désignés par $a, b, c,$ et d et qui constituent les matrices S_1 et S_2 (2.2). Les relations "entrées" "sorties" pour chacune des transitions s'écrivent:

$$\begin{bmatrix} u_2 \\ v_1 \end{bmatrix} = \begin{bmatrix} a_1 & b_1 \\ c_1 & d_1 \end{bmatrix} \begin{bmatrix} u_1 \\ v_2 \end{bmatrix} = S_1 \begin{bmatrix} u_1 \\ v_2 \end{bmatrix} \quad (71)$$

$$\begin{bmatrix} u_3 \\ v_2 \end{bmatrix} = \begin{bmatrix} a_2 & b_2 \\ c_2 & d_2 \end{bmatrix} \begin{bmatrix} u_2 \\ v_3 \end{bmatrix} = S_2 \begin{bmatrix} u_2 \\ v_3 \end{bmatrix} \quad (72)$$

Les matrices S sont ce qu'on appelle des matrices d'éparpillement. La matrice globale est définie implicitement par:

$$\begin{bmatrix} u_3 \\ v_1 \end{bmatrix} = S \begin{bmatrix} u_1 \\ v_3 \end{bmatrix} \quad (73)$$

Ce produit simple se généralise comme on le sait aux matrices de dimension quelconque.

IL est distributif par rapport à l'addition, associatif, possède un élément nul et un élément neutre.

Ainsi l'ensemble des matrices (n,n) muni de la loi d'addition et de ce produit forme un anneau et constitue une algèbre.

On va introduire maintenant un second type de produit matriciel basé lui aussi sur la composition en série de quadripôles et montre qu'il possède les mêmes propriétés que le produit simple.

Considérons donc le problème de propagation d'ondes suivant:

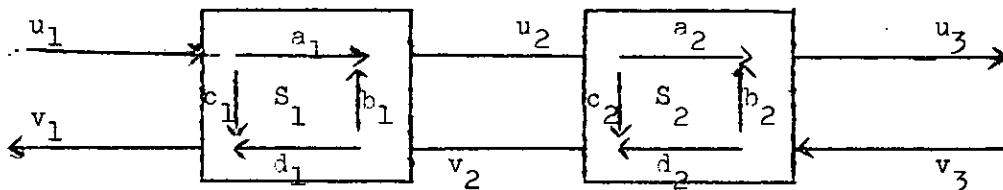


Figure 2

u_1, u_2 et u_3 sont des rayons incidents

v_1, v_2 et v_3 sont des rayons émergents

les milieux séparant les phases 1, 2 et 3 sont caractérisés par des Coefficients de réflexion de transmission désignés par $a, b, c,$ et d et qui constituent les matrices S_1 et S_2 (2.2). Les relations "entrées" "sorties" pour chacune des transitions s'écrivent:

$$\begin{bmatrix} u_2 \\ v_1 \end{bmatrix} = \begin{bmatrix} a_1 & b_1 \\ c_1 & d_1 \end{bmatrix} \begin{bmatrix} u_1 \\ v_2 \end{bmatrix} = S_1 \begin{bmatrix} u_1 \\ v_2 \end{bmatrix} \quad (71)$$

$$\begin{bmatrix} u_3 \\ v_2 \end{bmatrix} = \begin{bmatrix} a_2 & b_2 \\ c_2 & d_2 \end{bmatrix} \begin{bmatrix} u_2 \\ v_3 \end{bmatrix} = S_2 \begin{bmatrix} u_2 \\ v_3 \end{bmatrix} \quad (72)$$

Les matrices S sont ce qu'on appelle des matrices d'éparpillement. La matrice globale est définie implicitement par:

$$\begin{bmatrix} u_3 \\ v_1 \end{bmatrix} = S \begin{bmatrix} u_1 \\ v_3 \end{bmatrix} \quad (73)$$

ce qui donne

$$S = \begin{bmatrix} \frac{a_1 \cdot a_2}{1 - c_2 b_1} & b_2 + \frac{a_2 b_1 d_2}{1 - c_2 b_1} \\ c_1 + \frac{a_1 c_2 d_1}{1 - c_2 b_1} & \frac{d_1 d_2}{1 - c_2 b_1} \end{bmatrix} \quad (74)$$

Ainsi le produit étoile est défini explicitement par:

$$S = S_1 \star S_2 \quad (75)$$

Remarque:

S n'a de sens que si $c_2 b_1 \neq 1$, ce qui est évident en terme de réflexion et transmission.

Cas où les quantités u_i, v_i deviennent des n -vecteurs, Q et S étant alors des matrices $(2n \cdot 2n)$ constituées de $4(n \cdot n)$ matrices A, B, C, D .

Le produit simple donné par (74) se généralise par extension par:

$$Q = Q_2 Q_1^+ = \begin{bmatrix} A_2 A_1 + B_2 C_1 & A_2 B_1 + B_2 D_1 \\ C_2 A_1 + D_2 C_1 & C_2 B_1 + D_2 D_1 \end{bmatrix} \quad (76)$$

Pour le produit étoilé on procède de la même façon que pour le cas scalaire

Faisons:

$$S = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \quad (77)$$

ce qui donne

$$S = \begin{bmatrix} \frac{a_1 \cdot a_2}{1 - c_2 b_1} & b_2 + \frac{a_2 b_1 d_2}{1 - c_2 b_1} \\ c_1 + \frac{a_1 c_2 d_1}{1 - c_2 b_1} & \frac{d_1 d_2}{1 - c_2 b_1} \end{bmatrix} \quad (74)$$

Ainsi le produit étoile est défini explicitement par:

$$S = S_1 \star S_2 \quad (75)$$

Remarque:

S n'a de sens que si $c_2 b_1 \neq 1$, ce qui est évident en terme de réflexion et transmission.

Cas où les quantités u_i, v_i deviennent des n - vecteurs, Q et S étant alors des matrices $(2n \cdot 2n)$ constituées de $4(n \cdot n)$ matrices A, B, C, D .

Le produit simple donné par (60) se généralise par extension par:

$$Q = Q_2 Q_1^+ = \begin{bmatrix} A_2 A_1 + B_2 C_1 & A_2 B_1 + B_2 D_1 \\ C_2 A_1 + D_2 C_1 & C_2 B_1 + D_2 D_1 \end{bmatrix} \quad (76)$$

Pour le produit étoilé on procède de la même façon que pour le cas scalaire

Faisons:

$$S = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \quad (77)$$

D'après (71) et (72) on a :

$$V_1 = C_{11} U_1 + D_{12} V_2$$

$$V_2 = C_{22} U_2 + D_{23} V_3$$

$$U_2 = A_{11} U_1 + B_{12} V_2$$

$$U_3 = A_{22} U_2 + B_{23} V_3$$

Ce qui donne :

$$V_1 = \left[C_{11} + D_{12} (1 - C_{22} B_{12})^{-1} C_{22} A_{11} \right] U_1 + \left[D_{12} (1 - C_{22} B_{12})^{-1} D_{23} \right] V_3$$

$$U_3 = A_{22} \left[1 + B_{12} (1 - C_{22} B_{12})^{-1} C_{22} \right] A_{11} U_1 + \left[B_{23} + A_{22} B_{12} (1 - C_{22} B_{12})^{-1} D_{23} \right] V_3$$

Soit finalement :

$$\begin{bmatrix} A_1 & B_1 \\ C_1 & D_1 \end{bmatrix} * \begin{bmatrix} A_2 & B_2 \\ C_2 & D_2 \end{bmatrix} = \begin{bmatrix} A_2 \left[1 + B_{12} (1 - C_{22} B_{12})^{-1} C_{22} \right] A_{11} & B_{23} + A_{22} B_{12} (1 - C_{22} B_{12})^{-1} D_{23} \\ C_{11} + D_{12} (1 - C_{22} B_{12})^{-1} C_{22} A_{11} & D_{12} (1 - C_{22} B_{12})^{-1} D_{23} \end{bmatrix} \quad (78)$$

Résultat qui s'identifie au cas scalaire sauf pour le terme (1,1). On peut néanmoins y parvenir en opérant la transformation suivante.

à l'aide de l'identité matricielle :

$$(A + BCD)^{-1} = A^{-1} - A^{-1} B (C^{-1} + D A^{-1} B)^{-1} D A^{-1} \quad (79)$$

on peut vérifier que

$$I + B(I - CB)^{-1} C = (I - bc)^{-1} \quad (80)$$

En appliquant (80) il vient une seconde expression du produit étoile :

D'après (71) et (72) on a :

$$V_1 = C_{11} U_1 + D_{12} V_2$$

$$V_2 = C_{22} U_2 + D_{23} V_3$$

$$U_2 = A_{11} U_1 + B_{12} V_2$$

$$U_3 = A_{22} U_2 + B_{23} V_3$$

Ce qui donne :

$$V_1 = \left[C_{11} + D_{12} (1 - C_{22} B_{12})^{-1} C_{22} A_{11} \right] U_1 + \left[D_{12} (1 - C_{22} B_{12})^{-1} D_{23} \right] V_3$$

$$U_3 = A_{22} \left[1 + B_{12} (1 - C_{22} B_{12})^{-1} C_{22} A_{11} \right] U_1 + \left[B_{23} + A_{22} B_{12} (1 - C_{22} B_{12})^{-1} D_{23} \right] V_3$$

Soit finalement :

$$\begin{bmatrix} A_1 & B_1 \\ C_1 & D_1 \end{bmatrix} \times \begin{bmatrix} A_2 & B_2 \\ C_2 & D_2 \end{bmatrix} =$$

$$\begin{bmatrix} A_2 \left[1 + B_{12} (1 - C_{22} B_{12})^{-1} C_{22} A_{11} \right] & B_2 + A_{22} B_{12} (1 - C_{22} B_{12})^{-1} D_{23} \\ C_1 + D_{12} (1 - C_{22} B_{12})^{-1} C_{22} A_{11} & D_{12} (1 - C_{22} B_{12})^{-1} D_{23} \end{bmatrix} \quad (78)$$

Résultat qui s'identifie au cas scalaire sauf pour le terme (1,1). On peut néanmoins y parvenir en opérant la transformation suivante.

à l'aide de l'identité matricielle :

$$(A + BCD)^{-1} = A^{-1} - A^{-1} B (C^{-1} + D A^{-1} B)^{-1} D A^{-1} \quad (79)$$

on peut vérifier que

$$I + B(I - CB)^{-1} C = (I - BC)^{-1} \quad (80)$$

En appliquant (80) il vient une seconde expression du produit étoile :

$$\begin{bmatrix} A_1 & B_1 \\ C_1 & D_1 \end{bmatrix} \star \begin{bmatrix} A_2 & B_2 \\ C_2 & D_2 \end{bmatrix} = \begin{bmatrix} A_2(I-B_1C_2)^{-1}A_1 & B_2+A_2(I-B_1C_2)^{-1}B_1D_2 \\ C_1+D_1C_2(I-B_1C_2)^{-1}A_1 & D_1(I-C_2B_1)^{-1}D_2 \end{bmatrix} \quad (81)$$

qui n'est définie que si les matrices $(I-B_1C_2)$ et $(I-C_2B_1)$ sont régulières. Ceci étant admis, on rappellera que l'ensemble des matrices $(2n, 2n)$ muni de la loi d'addition classique et du produit $Q = Q_2 Q_1$ ou du produit étoile $S = S_1 \star S_2$ forme un anneau et constitue une algèbre.

On va maintenant aborder une seconde voie dont l'intérêt est de mettre en évidence un isomorphisme entre le produit \star et \cdot appelé transformation Ψ .

4-CONSTRUCTION DE L'ISOMORPHISME

Soit le problème élémentaire suivant de propagation:

$$\begin{bmatrix} U_2 \\ V_1 \end{bmatrix} = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} U_1 \\ V_2 \end{bmatrix} = S \begin{bmatrix} U_1 \\ V_2 \end{bmatrix} \quad (82)$$

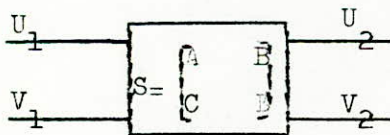


fig 3

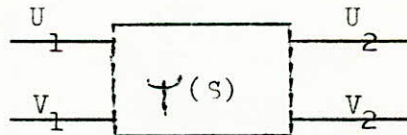


fig 4

Introduisons une transformation Ψ définie par

$$\begin{bmatrix} U_2 \\ V_2 \end{bmatrix} = \Psi(S) \begin{bmatrix} U_1 \\ V_1 \end{bmatrix} \quad (83)$$

Ce résultat est illustré par la figure 4

D'après (82) on a:

$$\begin{aligned} V_2 &= -D^{-1}CU_1 + D^{-1}V_2 \\ U_2 &= (A-BD^{-1}C)U_1 + BD^{-1}V_2 \end{aligned}$$

$$\begin{bmatrix} A_1 & B_1 \\ C_1 & D_1 \end{bmatrix} \times \begin{bmatrix} A_2 & B_2 \\ C_2 & D_2 \end{bmatrix} = \begin{bmatrix} A_2(I-B_1C_2)^{-1}A_1 & B_2+A_2(I-B_1C_2)^{-1}B_1D_2 \\ C_1+D_1C_2(I-B_1C_2)^{-1}A_1 & D_1(I-C_2B_1)^{-1}D_2 \end{bmatrix} \quad (81)$$

qui n'est définie que si les matrices $(I-B_1C_2)$ et $(I-C_2B_1)$ sont régulières. Ceci étant admis, on rappellera que l'ensemble des matrices $(2n, 2n)$ muni de la loi d'addition classique et du produit $Q = Q_2 Q_1$ ou du produit étoile $S = S_1 \times S_2$ forme un anneau et constitue une algèbre.

On va maintenant aborder une seconde voie dont l'intérêt est de mettre en évidence un isomorphisme entre le produit \times et \cdot appelé transformation Ψ .

4-CONSTRUCTION DE L'ISOMORPHISME

Soit le problème élémentaire suivant de propagation:

$$\begin{bmatrix} U_2 \\ V_1 \end{bmatrix} = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} U_1 \\ V_2 \end{bmatrix} = S \begin{bmatrix} U_1 \\ V_2 \end{bmatrix} \quad (82)$$

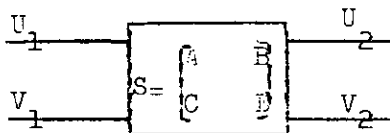


fig 3

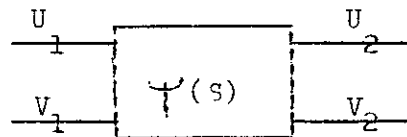


fig 4

Introduisons une transformation Ψ définie par

$$\begin{bmatrix} U_2 \\ V_2 \end{bmatrix} = \Psi(S) \begin{bmatrix} U_1 \\ V_1 \end{bmatrix} \quad (83)$$

Ce résultat est illustré par la figure 4

D'après (82) on a:

$$V_2 = -D^{-1}CU_1 + D^{-1}V_2$$

$$U_2 = (A-BD^{-1}C)U_1 + BD^{-1}V_2$$

d'où

$$\Psi(s) = \begin{bmatrix} A - BD^{-1}C & BD^{-1} \\ -D^{-1}C & D^{-1} \end{bmatrix} = \tilde{s} \quad (84)$$

Il est aisé de vérifier que

$$\Psi(\tilde{s}) = s \quad (85)$$

donc pourvu que le bloc D soit régulier, la transformation Ψ vérifie

$$\Psi[\Psi(s)] = s \implies \Psi^{-1}(s) = \Psi(s) \quad (86)$$

où Ψ^{-1} est la transformation inverse.

En appliquant ce résultat au problème défini à la figure 2, il vient:

$$\Psi(s) = \Psi(s_2) \cdot \Psi(s_1) \implies s = \Psi[\Psi(s_2) \cdot \Psi(s_1)] = s_1 * s_2 \quad (87)$$

où encore d'après (86):

$$\Psi(s_1 * s_2) = \Psi(s_2) \cdot \Psi(s_1) \quad (88)$$

Inversement en raisonnant à partir de la figure 1 et en passant par l'étape de la figure 2 on aura:

$$\Psi(q_1) * \Psi(q_2) = \Psi(q_2 \cdot q_1) \quad (89)$$

En résumé, la transformation involutive Ψ définit donc un isomorphisme entre le produit étoile et le produit classique de matrice.

Il faut toutefois noter que cette voie requiert l'hypothèse que D_1 et D_2 soient réguliers.

d'où,

$$\Psi'(s) = \begin{bmatrix} A - BD^{-1}C & BD^{-1} \\ -D^{-1}C & D^{-1} \end{bmatrix} = \tilde{s} \quad (84)$$

Il est aisé de vérifier que

$$\Psi(\tilde{s}) = s \quad (85)$$

donc pourvu que le bloc D soit régulier, la transformation Ψ vérifie

$$\Psi[\Psi'(s)] = s \implies \Psi^{-1}(s) = \Psi'(s) \quad (86)$$

où Ψ^{-1} est la transformation inverse.

En appliquant ce résultat au problème défini à la figure 2, il vient:

$$\Psi(s) = \Psi(s_2) \cdot \Psi(s_1) \implies s = \Psi[\Psi(s_2) \cdot \Psi(s_1)] = s_1 * s_2 \quad (87)$$

où encore d'après (86):

$$\Psi(s_1 * s_2) = \Psi(s_2) \cdot \Psi(s_1) \quad (88)$$

Inversement en raisonnant à partir de la figure 1 et en passant par l'étape de la figure 2 on aura:

$$\Psi(q_1) * \Psi(q_2) = \Psi(q_2 \cdot q_1) \quad (89)$$

En résumé, la transformation involutive Ψ définit donc un isomorphisme entre le produit étoile et le produit classique de matrice.

Il faut toutefois noter que cette voie requiert l'hypothèse que D_1 et D_2 soient réguliers.

5-APPLICATION A L'EQUATION DE RICCATI DISCRETE

Considérons de nouveau les équations hamiltoniennes qui s'écrivent:

$$\begin{bmatrix} X_{t+1} \\ \lambda_{t+1} \end{bmatrix} = \begin{bmatrix} A & -BR^{-1}B^T \\ 0 & A^T \end{bmatrix} \begin{bmatrix} X_t \\ \lambda_t \end{bmatrix}$$

Appliquant la transformation Ψ tout en supposant pour l'instant que A est régulière, il vient d'après (84):

$$\begin{bmatrix} X_{t+1} \\ \lambda_{t+1} \end{bmatrix} = \begin{bmatrix} A + VA^{-T}Q & -VA^{-T} \\ -A^{-T}Q & A^{-T} \end{bmatrix} \begin{bmatrix} X_t \\ \lambda_t \end{bmatrix} = S \begin{bmatrix} X_t \\ \lambda_t \end{bmatrix} \quad (90)$$

où $V = BR^{-1}B^T$ et A^{-T} désigne la transposée de l'inverse

D'après VAUGHAN [IC], il vient:

$$S = \begin{bmatrix} W_1 & W_{12} \\ W_{21} & W_2 \end{bmatrix} \begin{bmatrix} J^{-1} & 0 \\ 0 & J \end{bmatrix} \begin{bmatrix} W_1 & W_{12} \\ W_{21} & W_2 \end{bmatrix}^{-1} \quad (91)$$

J : blocs de Jordan

avec

$$F = W_{21} \cdot W_1^{-1} \quad (92)$$

Si tant une matrice symplectique (annexe), on peut alors lui associer une matrice hamiltonienne M par transformation bilinéaire (annexe);

il vient d'après (91)

$$M = (S-1)(S+1) = W \begin{bmatrix} (J^{-1}-1)^{-1}(J^{-1}+1) & 0 \\ 0 & (J-1)^{-1}(J+1) \end{bmatrix} W^{-1}$$

5-AFFLICTION A L'EQUATION DE RICCATI DISCRETE

Considérons de nouveau les équations hamiltoniennes qui s'écrivent:

$$\begin{bmatrix} X_{t+1} \\ \lambda_t \end{bmatrix} = \begin{bmatrix} A & -BR^{-1}B^T \\ 0 & A^T \end{bmatrix} \begin{bmatrix} X_t \\ \lambda_{t+1} \end{bmatrix}$$

Appliquant la transformation S tout en supposant pour l'instant que A est régulière, il vient d'après (84):

$$\begin{bmatrix} X_{t+1} \\ \lambda_{t+1} \end{bmatrix} = \begin{bmatrix} A + VA^{-T}Q & -VA^{-T} \\ -A^{-T}Q & A^{-T} \end{bmatrix} \begin{bmatrix} X_t \\ \lambda_t \end{bmatrix} = S \begin{bmatrix} X_t \\ \lambda_t \end{bmatrix} \quad (90)$$

où $V = BR^{-1}B^T$ et A^{-T} désigne la transposée de l'inverse

D'après VAUGHAN [IC], il vient:

$$S = \begin{bmatrix} W_1 & W_{12} \\ W_{21} & W_2 \end{bmatrix} \begin{bmatrix} J^{-1} & 0 \\ 0 & J \end{bmatrix} \begin{bmatrix} W_1 & W_{12} \\ W_{21} & W_2 \end{bmatrix}^{-1} \quad (91)$$

J : blocs de Jordan

avec

$$F = W_{21} \cdot W_1^{-1} \quad (92)$$

Étant une matrice symplectique (annexe), on peut alors lui associer une matrice hamiltonienne M par transformation bilinéaire (annexe);

il vient d'après (91)

$$M = (S-1)(S+1) = W \begin{bmatrix} (J^{-1}-1)^{-1}(J^{-1}+1) & 0 \\ 0 & (J-1)^{-1}(J+1) \end{bmatrix} W^{-1}$$

$$= W \begin{bmatrix} -\wedge & 0 \\ 0 & \wedge \end{bmatrix} W^{-1} \quad (93)$$

avec $\wedge = (J-1)^{-1}(J+1)$, réel $[\wedge(\wedge)] > 0$ (94)

La matrice signe de H s'écrit alors:

$$\text{signe}(F) = Z \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \text{signe}(\lambda_i) & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} Z^{-1} ; M = ZJZ^{-1} \quad (95)$$

J étant la forme de Jordan

$$\text{signe}(\lambda) = \begin{cases} +1 & \text{si } \text{réel}(\lambda) > 0 \\ -1 & \text{si } \text{réel}(\lambda) < 0 \end{cases} \quad (96)$$

tenant compte des relations (93) et (94), on a alors:

$$\text{signe}(F) = W \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} W^{-1} \quad (97)$$

Introduisons alors la matrice:

$$F = \text{signe}(F) + \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} = W \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} W^{-1} + \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}$$

d'où

$$F = 2 \begin{bmatrix} W_1 & W_{12} \\ W_{21} & W_2 \end{bmatrix} \begin{bmatrix} -V_1 & 0 \\ 0 & V_2 \end{bmatrix} = \begin{bmatrix} F_1 & F_{12} \\ F_{21} & F_2 \end{bmatrix} \quad \text{avec } W^{-1} = \begin{bmatrix} V_1 & V_{12} \\ V_{21} & V_2 \end{bmatrix}$$

soit finalement d'après (92)

$$F_{21} F_1^{-1} = W_{21} W_1^{-1} = F$$

Un point important est qu'en définitive, on peut s'affranchir de l'hypothèse que l'on avait faite: \wedge régulière. En effet la matrice symplectique S peut s'écrire:

$$= W \begin{bmatrix} -\wedge & 0 \\ 0 & \wedge \end{bmatrix} W^{-1} \quad (93)$$

$$\text{avec } \wedge = (J-1)^{-1}(J+1), \quad \text{réel } [\lambda(\wedge)] > 0 \quad (94)$$

La matrice signe de H s'écrit alors:

$$\text{signe}(F) = Z \begin{bmatrix} \cdot & 0 \\ \text{signe}(\lambda_i) & \cdot \end{bmatrix} Z^{-1}; \quad M = ZJZ^{-1} \quad (95)$$

J étant la forme de Jordan

$$\text{signe}(\lambda) = \begin{cases} +1 & \text{si } \text{réel}(\lambda) > 0 \\ -1 & \text{si } \text{réel}(\lambda) < 0 \end{cases} \quad (96)$$

tenant compte des relations (93) et (94), on a alors:

$$\text{signe}(H) = W \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} W^{-1} \quad (97)$$

Introduisons alors la matrice:

$$F = \text{signe}(F) + \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} = W \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix} W^{-1} + \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}$$

d'où

$$F = 2 \begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{bmatrix} \begin{bmatrix} -V_1 & 0 \\ 0 & V_2 \end{bmatrix} = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix} \quad \text{avec } W^{-1} = \begin{bmatrix} V_1 & V_{12} \\ V_{21} & V_2 \end{bmatrix}$$

soit finalement d'après (92)

$$F_{21} F_{11}^{-1} = W_{21} W_{11}^{-1} = F$$

Un point important est qu'en définitive, on peut s'affranchir de l'hypothèse que l'on avait faite: Δ régulière. En effet la matrice symplectique Δ peut s'écrire:

$$S = \begin{bmatrix} I & BR^{-1}B^T \\ 0 & A^T \end{bmatrix} \begin{bmatrix} A & 0 \\ -Q & I \end{bmatrix} = U^{-1} \cdot L$$

d'où

$$H = (S-1)^{-1}(S+1) = (U^{-1}L-1)^{-1}(U^{-1}L+1) = (L-U)^{-1}(L+U)$$

On remarque dans cette expression que A^{-1} n'apparaît plus. Par ailleurs H étant hamiltonienne sans valeurs propres sur l'axe imaginaire, donc $L-U$ est régulière ainsi que $L+U$.

En résumé, l'algorithme que l'on se propose est défini par:

$$F = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix} = \begin{bmatrix} \text{signe}(L-U)^{-1}(L+U) \end{bmatrix} + \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \quad (98)$$

6-COMPARAISON AVEC D'AUTRES METHODES

L'étude comparative que l'on développe maintenant est décomposée en deux parties respectivement consacrées aux méthodes à caractère implicite et celles à formulation itérative explicite.

6-1-Méthodes à caractère implicite

Les deux méthodes qui font l'objet de ce paragraphe reposent, tout comme la méthode de la fonction signe de matrice, sur les propriétés de la matrice symplectique (9C) associée à l'équation de Riccati. Cependant les deux méthodes supposent que la matrice A est régulière. De plus, comme pour la méthode de la fonction signe de matrice, aucune ne permet de garantir le caractère symétrique non négatif de P , par une approche du type "square root".

6-1-1-Méthode de VAUGHAN

Nous considérons d'abord la méthode de VAUGHAN ~~/TC/~~; sans doute la plus classique, qui peut se résumer dans la relation:

$$S = \begin{bmatrix} I & BR^{-1}B^T \\ 0 & A^T \end{bmatrix} \begin{bmatrix} A & 0 \\ -Q & I \end{bmatrix} = U^{-1} \cdot L$$

d'où

$$H = (S-1)^{-1}(s+1) = (U^{-1}L-1)^{-1}(U^{-1}L+1) = (L-U)^{-1}(L+U)$$

On remarque dans cette expression que A^{-1} n'apparaît plus. Par ailleurs H étant hamiltonienne sans valeurs propres sur l'axe imaginaire, donc $L-U$ est régulière ainsi que $L+U$.

En résumé, l'algorithme que l'on se propose est défini par:

$$F = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix} = \begin{bmatrix} \text{signe}(L-U)^{-1}(L+U) \end{bmatrix} + \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \quad (98)$$

6-COMPARAISON AVEC D'AUTRES METHODES

L'étude comparative que l'on développe maintenant est décomposée en deux parties respectivement consacrées aux méthodes à caractère implicite et celles à formulation itérative explicite.

6-1-Méthodes à caractère implicite

Les deux méthodes qui font l'objet de ce paragraphe reposent, tout comme la méthode de la fonction signe de matrice, sur les propriétés de la matrice symplectique (9C) associée à l'équation de Riccati. Cependant les deux méthodes supposent que la matrice A est régulière. De plus, comme pour la méthode de la fonction signe de matrice, aucune ne permet de garantir le caractère symétrique non négatif de P , par une approche du type "square root".

6-1-1-Méthode de VAUGHAN

Nous considérons d'abord la méthode de VAUGHAN ~~/TC/~~; sans doute la plus classique, qui peut se résumer dans la relation:

$$P = W_{21} W_1^{-1}$$

(99)

avec (90) on a :

$$S = \begin{bmatrix} A + VA^{-T}Q & -VA^{-T} \\ -A^{-T}Q & A^{-T} \end{bmatrix} = \begin{bmatrix} W_1 & W_{12} \\ W_{21} & W_2 \end{bmatrix} \begin{bmatrix} J^{-1} & 0 \\ 0 & J \end{bmatrix} \begin{bmatrix} W_1 & W_{12} \\ W_{21} & W_2 \end{bmatrix}^{-1} \quad (100)$$

et

$$|\lambda(J)| > 1 \quad J \text{ blocs de Jordan}$$

En fait, dans la mesure où la décomposition sous forme de Jordan est impossible, ce résultat ne peut constituer la base d'un algorithme que si S est strictement diagonalisable.

Il n'en reste pas moins vrai que peuvent toujours subsister les problèmes liés au mauvais conditionnement (nombre de condition de W élevé) ou encore la grande sensibilité de certaines valeurs propres voisines entre elles ou/et de 1 (en module).

6-1-2-Méthode de LAUB

L'alternative proposée récemment par Laub [11], constitue un progrès considérable par rapport aux résultats précédents dans la mesure où elle supprime l'essentiel des difficultés numériques.

Sa méthode repose sur le fait qu'une relation du type (99) peut être obtenue sans imposer que S soit mis sous forme de Jordan, mais seulement sous forme de Schur réelle :

$$S = \begin{bmatrix} V_1 & V_{12} \\ V_{21} & V_2 \end{bmatrix} \begin{bmatrix} \Lambda_1 & \Lambda_{12} \\ 0 & \Lambda_2 \end{bmatrix} \begin{bmatrix} V_1 & V_{12} \\ V_{21} & V_2 \end{bmatrix}^T \quad (101a)$$

$$|\lambda(\Lambda_1)| < 1 \quad |\lambda(\Lambda_2)| > 1$$

Λ_1 et Λ_2 ne sont plus des blocs de Jordan inverses l'un de l'autre, mais

$$P = W_{21} W_1^{-1}$$

(99)

avec (90) on a :

$$S = \begin{bmatrix} A + VA^{-T}Q & -VA^{-T} \\ -A^{-T}Q & A^{-T} \end{bmatrix} = \begin{bmatrix} W_1 & W_{12} \\ W_{21} & W_2 \end{bmatrix} \begin{bmatrix} J^{-1} & 0 \\ 0 & J \end{bmatrix} \begin{bmatrix} W_1 & W_{12} \\ W_{21} & W_2 \end{bmatrix}^{-1} \quad (100)$$

et

$$|\lambda(J)| > 1 \quad J \text{ blocs de Jordan}$$

En fait, dans la mesure où la décomposition sous forme de Jordan est impossible, ce résultat ne peut constituer la base d'un algorithme que si S est strictement diagonalisable.

Il n'en reste pas moins vrai que peuvent toujours subsister les problèmes liés au mauvais conditionnement (nombre de condition de W élevé) ou encore la grande sensibilité de certaines valeurs propres voisines entre elles ou/et de 1 (en module).

6-1-2-Méthode de LAUB

L'alternative proposée récemment par Laub [117], constitue un progrès considérable par rapport aux résultats précédents dans la mesure où elle supprime l'essentiel des difficultés numériques.

Sa méthode repose sur le fait qu'une relation du type (99) peut être obtenue sans imposer que S soit mis sous forme de Jordan, mais seulement sous forme de Schur réelle :

$$S = \begin{bmatrix} V_1 & V_{12} \\ V_{21} & V_2 \end{bmatrix} \begin{bmatrix} \Lambda_1 & \Lambda_{12} \\ 0 & \Lambda_2 \end{bmatrix} \begin{bmatrix} V_1 & V_{12} \\ V_{21} & V_2 \end{bmatrix}^T \quad (101a)$$

$$|\lambda(\Lambda_1)| < 1 \quad |\lambda(\Lambda_2)| > 1$$

Λ_1 et Λ_2 ne sont plus des blocs de Jordan inverses l'un de l'autre, mais

ont leur spectre réparti comme J^{-1} et J . De plus la matrice V est orthogonale donc parfaitement conditionnée. Dans ces conditions le calcul de P comme suit ne pose plus de problème particulier

$$P = V_{21} V_1^{-1} \quad (101b)$$

Quand à la forme de Schur, elle est obtenue à l'aide du package Eispack [12] et des modules de Stewart [13].

6-2-ALGORITHMES ITERATIFS

6-2-1-Méthode du produit étoile

La première méthode que l'on se propose d'aborder est en fait une application du produit étoile.

Compte tenu de (90), on a le résultat suivant:

Théorème:

$$P = \lim_{K \rightarrow \infty} \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \end{array} (Z_K) \right]_{2,1} \quad (102)$$

$$\text{avec } Z_{K+1} = Z_K \cdot S \quad Z_0 = \begin{bmatrix} 1 & 0 \\ -P_0 & 1 \end{bmatrix} \quad (103)$$

$$\forall P_0 = P_0^T \geq 0$$

où P est solution de l'équation de Riccati et où $\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \end{array} \right]_{2,1}$ désigne le bloc (2,1) de dimension $n.n$ des matrices partitionnées $2n.2n$.

Ceci étant, on peut maintenant en déduire une seconde relation en terme de produit étoile. Rappelons que (90) s'écrit

$$S = \psi(M)$$

Posons alors

$$Y_K = \psi(Z_K) \implies Z_K = \psi(Y_K)$$

$$\psi(Y_{K+1}) = Z_{K+1} = \psi(Y_K) \cdot \psi(M)$$

ont leur spectre réparti comme J^{-1} et J . De plus la matrice V est orthogonale donc parfaitement conditionnée. Dans ces conditions le calcul de Γ comme suit ne pose plus de problème particulier

$$P = V_{21} V_1^{-1} \quad (101b)$$

Quand à la forme de Schur, elle est obtenue à l'aide du package Eispack [12] et des modules de Stewart [13].

6-2-ALGORITHMES ITERATIFS

6-2-1-Méthode du produit étoile

La première méthode que l'on se propose d'aborder est en fait une application du produit étoile.

Compte tenu de (90), on a le résultat suivant:

Théorème:

$$P = \lim_{k \rightarrow \infty} \left[\begin{array}{c} \Psi(Z_k) \\ \cdot \\ \cdot \end{array} \right]_{2,1} \quad (102)$$

$$\text{avec } Z_{K+1} = Z_K \cdot S \quad Z_0 = \begin{bmatrix} 1 & 0 \\ -P_0 & 1 \end{bmatrix} \quad (103)$$

$$\forall P_0 = P_0^T \geq 0$$

où P est solution de l'équation de Riccati et où $\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \end{array} \right]_{2,1}$ désigne le bloc (2,1) de dimension $n.n$ des matrices partitionnées $2n.2n$.

Ceci étant, on peut maintenant en déduire une seconde relation en terme de produit étoile. Rappelons que (90) s'écrit

$$S = \Psi(M)$$

Posons alors

$$Y_K = \Psi(Z_Y) \implies Z_Y = \Psi(Y_K)$$

$$\Psi(Y_{K+1}) = Z_{K+1} = \Psi(Y_K) \cdot \Psi(M)$$

soit

$$Y_{K+1} = M * Y_K \quad Y_0 = \begin{bmatrix} 1 & 0 \\ P_0 & 1 \end{bmatrix}$$

de plus

$$Y_K = \Psi(Z_K) = \begin{bmatrix} * & * \\ \Gamma_K & * \end{bmatrix}$$

donc

$$P = \lim_{K \rightarrow \infty} \begin{bmatrix} Y_K \end{bmatrix}_{2,1}$$

$$Y_{K+1} = M * Y_K \quad Y_0 = \begin{bmatrix} 1 & 0 \\ P_0 & 1 \end{bmatrix} \quad \forall P_0 = P_0^T \gg 0$$

Ce formalisme équivalent à (102)-(103) présente l'avantage de lever la contrainte A régulière.

Un algorithme dit de doublement, dont la convergence est d'ordre 2 s'obtient alors en posant:

$$M_1 = M = \begin{bmatrix} A & -BP^{-1}B^T \\ Q & A^T \end{bmatrix}$$

$$M_2 = M_1 * M_1$$

$$M_{2^k} = M_{2^{k-1}} * M_{2^{k-1}} \quad \text{avec} \quad M_{2^k} = \begin{bmatrix} U_k & -V_k \\ W_k & U_k^T \end{bmatrix}$$

d'où tout en simplifiant les notations

$$U_{K+1} = U_+ = U(I+VW)^{-1}U$$

$$U_0 = A \quad (a)$$

$$V_{K+1} = V_+ = V + U(I+VW)^{-1}VU^T$$

$$V_0 = BP^{-1}B^T \quad (b)$$

$$W_{K+1} = W_+ = W + U^T W(I+VW)^{-1}U$$

$$W_0 = Q \quad (c)$$

$$P = \lim_{K \rightarrow \infty} W_K = P_{2^k}$$

$$(d)$$

(104)

soit

$$Y_{K+1} = M * Y_K \quad Y_0 = \begin{bmatrix} 1 & 0 \\ P_0 & 1 \end{bmatrix}$$

de plus

$$Y_K = \Psi(Z_K) = \begin{bmatrix} * & * \\ \Gamma_K & * \end{bmatrix}$$

donc

$$P = \lim_{K \rightarrow \infty} \begin{bmatrix} Y_K \\ Y_K \end{bmatrix}_{2,1}$$

$$Y_{K+1} = M * Y_K \quad Y_0 = \begin{bmatrix} 1 & 0 \\ P_0 & 1 \end{bmatrix} \quad \forall P_0 = P_0^T \gg 0$$

Ce formalisme équivalent à (102)-(103) présente l'avantage de lever la contrainte A régulière.

Un algorithme dit de doublement, dont la convergence est d'ordre 2 s'obtient alors en posant:

$$M_1 = M = \begin{bmatrix} A & -BF^{-1}B^T \\ Q & A^T \end{bmatrix}$$

$$M_2 = M_1 * M_1$$

$$M_2^K = M_2^{K-1} * M_2^{K-1} \quad \text{avec} \quad M_2^K = \begin{bmatrix} U_K & -V_K \\ W_K & U_K^T \end{bmatrix}$$

d'où tout en simplifiant les notations

$$U_{K+1} = U_+ = U(I+VW)^{-1}U$$

$$U_0 = A \quad (a)$$

$$V_{K+1} = V_+ = V + U(I+VW)^{-1}VU^T$$

$$V_0 = BF^{-1}B^T \quad (b)$$

$$W_{K+1} = W_+ = W + U^T W(I+VW)^{-1}U$$

$$W_0 = Q \quad (c)$$

$$P = \lim_{K \rightarrow \infty} W_K = P_{2^K}$$

$$(d)$$

(104)

(b) et (c) sont des équations de Riccati écrites sous leur forme équivalente (annexe). On pose alors

$$V = X^T \cdot X \quad \text{et} \quad W = Y^T \cdot Y$$

Algorithme: calculer $D = X \cdot U^T$, $E = X \cdot Y^T$, $F = Y \cdot U$

1^o/ réaliser la triangularisation

$$\begin{bmatrix} 1 & 0 \\ E & D \\ 0 & X \end{bmatrix} \begin{bmatrix} K & L \\ 0 & X_+ \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot \end{bmatrix} \quad (105)$$

au moyen de la matrice orthogonale $\underline{\tau}$, d'où X_+ triangulaire supérieure

2^o/ résoudre le système triangulaire inférieur

$$K^T \cdot M = F$$

réaliser la triangularisation "orthogonale":

$$\tau \begin{bmatrix} Y \\ M \end{bmatrix} = \begin{bmatrix} Y_+ \\ C \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot \end{bmatrix}$$

dont on tire Y_+ triangulaire supérieur.

3^o/ (a) Y de rang maximal (W régulier)

-résoudre les systèmes triangulaires inférieurs:

$$Y^T \cdot H = U^T \quad \text{puis} \quad K^T \cdot G = H$$

-calculer

$$U_+ = G^T \cdot M$$

(b) Y quelconque (W singulier):

-calculer

$$G = I + X^T E Y$$

-résoudre le système linéaire $G \cdot F = U$

-calculer

$$U_+ = U \cdot F$$

(106)

Le coût total (en multiplications) de cet algorithme est de $7n^3$ approximativement.

(b) et (c) sont des équations de Riccati écrites sous leur forme équivalente (annexe). On pose alors

$$V = X^T \cdot X \quad \text{et} \quad W = Y^T \cdot Y$$

Algorithme: calculer $D = X \cdot U^T$, $E = X \cdot Y^T$, $F = Y \cdot U$

1^o/ réaliser la triangularisation

$$\begin{bmatrix} I & O \\ E & D \\ O & X \end{bmatrix} \begin{bmatrix} K & L \\ O & X_+ \\ O & O \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ O & \cdot & \cdot & \cdot \end{bmatrix} \quad (105)$$

au moyen de la matrice orthogonale \mathcal{U} , d'où X_+ triangulaire supérieure

2^o/ résoudre le système triangulaire inférieur

$$K^T \cdot M = F$$

réaliser la triangularisation "orthogonale":

$$\mathcal{U} \begin{bmatrix} Y \\ M \end{bmatrix} = \begin{bmatrix} Y_+ \\ O \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ O & \cdot & \cdot & \cdot \end{bmatrix}$$

dont on tire Y_+ triangulaire supérieur.

3^o/ (a) Y de rang maximal (W régulier)

-résoudre les systèmes triangulaires inférieurs:

$$Y^T \cdot H = U^T \quad \text{puis} \quad K^T \cdot G = H$$

-calculer

$$U_+ = G^T \cdot M$$

(b) Y quelconque (W singulier):

-calculer

$$G = I + X^T E Y$$

-résoudre le système linéaire $G \cdot H = U$

-calculer

$$U_+ = U \cdot H$$

(106)

Le coût total (en multiplications) de cet algorithme est de $7n^3$ approximativement.

6-2-2-Méthode de HEWER

La deuxième alternative que l'on se propose de rappeler est due à Hewer [77].

Le principe de l'algorithme consiste à résoudre l'équation de Riccati par la méthode de Newton; ce qui conduit au processus itératif suivant:

$$M_j - \tilde{A}_j^T M_j \tilde{A}_j = G_j^T R G_j + Q \quad (107a)$$

$$G_j = (B^T M_{j-1} B + R)^{-1} B^T M_{j-1} A \quad (107b)$$

$$\tilde{A}_j = A - B G_j \quad (107c)$$

où apparaît donc à chaque pas la résolution d'une équation de Lyapunov (107a). On sait alors que la suite M_j présente une convergence monotone décroissante d'ordre 2 vers la solution P cherchée quelque soit l'initialisation G_0 tel que

$$\left| \lambda(\tilde{A}_0 = A - B G_0) \right| < 1$$

avec

$$P_{K+1} = A^T P_K A - A^T P_K B (B^T P_K B + R)^{-1} B^T P_K A + Q \quad P_0 = C \quad (108)$$

Sur le plan coûtcalcul, chaque itération est dominée par la résolution de l'équation de Lyapunov (107a). Le coût global de (107) s'élève en moyenne à environ [78] : $20n^3$ ($m < n$, $B: m.n$).

En conséquence, cet algorithme est a priori coûteux, chose à laquelle s'ajoute la non garantie sur le caractère symétrique non négatif de la solution obtenue et un logiciel beaucoup plus lourd (ég de Lyapunov).

6-2-3-Méthode classique stabilisée

Une seconde formulation itérative est directement fournie par (108) ou sa version dite stabilisée (annexe):

$$P_{K+1} = \tilde{A}_K^T P_K A + G_K^T R G_K + Q$$

6-2-2-Méthode de HEWER

La deuxième alternative que l'on se propose de rappeler est due à Hewer [77].

Le principe de l'algorithme consiste à résoudre l'équation de Riccati par la méthode de Newton; ce qui conduit au processus itératif suivant:

$$M_j - \tilde{A}_j^T M_j \tilde{A}_j = G_j^T R G_j + Q \quad (107a)$$

$$G_j = (B^T M_{j-1} B + R)^{-1} B^T M_{j-1} A \quad (107b)$$

$$\tilde{A}_j = A - B G_j \quad (107c)$$

où apparaît donc à chaque pas la résolution d'une équation de Lyapunov (107a). On sait alors que la suite M_j présente une convergence monotone décroissante d'ordre 2 vers la solution P cherchée quelque soit l'initialisation G_0 tel que

$$\left| \lambda(\tilde{A}_0 = A - B G_0) \right| < 1$$

avec

$$P_{Y+1} = A^T P_K A - A^T P_K B (B^T P_K B + R)^{-1} B^T P_K A + Q \quad P_0 = C \quad (108)$$

Sur le plan coûtcalcul, chaque itération est dominée par la résolution de l'équation de Lyapunov (107a). Le coût global de (107) s'élève en moyenne à environ [78] : $20n^3$ ($m < n$, $B: m.n$).

En conséquence, cet algorithme est a priori coûteux, chose à laquelle s'ajoute la non garantie sur le caractère symétrique non négatif de la solution obtenue et un logiciel beaucoup plus lourd (ég de Lyapunov).

6-2-3-Méthode classique stabilisée

Une seconde formulation itérative est directement fournie par (108) ou sa version dite stabilisée (annexe):

$$P_{K+1} = \tilde{A}_K^T P_K A + G_K^T R G_K + Q$$

$$\begin{aligned}
 G_K &= (R^T \Gamma_K P + R)^{-1} R^T P A \\
 \tilde{A}_K &= A - R G_K
 \end{aligned}
 \tag{109}$$

qu'il ne faut pas confondre, malgré sa ressemblance avec le processus itératif de Hewer(107).

6-2-4-Square root d'ordre 1

A la différence de la méthode de Hewer, on peut réaliser une formulation du type square root de (108) garantissant le caractère symétrique non négatif de la solution cherchée.

$$\begin{aligned}
 R &= X^T \cdot X & Q &= V^T \cdot V \\
 \Gamma_K &= S_K^T \cdot S_K & S_K &\text{ triangulaire supérieure}
 \end{aligned}
 \tag{110}$$

itération K:

-calculer $D = S_K \cdot B \quad E = S_K \cdot A$ (111)

-réaliser la triangularisation orthogonale

$$\tau \begin{bmatrix} X & C \\ D & E \\ O & V \end{bmatrix} = \begin{bmatrix} L & K \\ O & S_{K+1} \\ O & O \end{bmatrix}
 \tag{112}$$

Le coût calcul pour $m < n$ ($B: n \cdot m$), reste inférieur à $3n^3$. Néanmoins, il apparaît clairement dans le tableau suivant, que le coût relatif à cette méthode est assez grand lorsque le régime permanent est suffisamment long à atteindre, c'est à dire lorsque le temps de réponse du système bouclé est important.

Tableau III

$N = 2^K$	$C = N \cdot 3n^3$
2	$6n^3$
4	$12n^3$
8	$24n^3$
16	$48n^3$
32	$96n^3$
64	$192n^3$

N: nombre d'itérations

C: coût calcul

$$\left. \begin{aligned} G_K &= (R^T P_K P + R)^{-1} R^T P_K A \\ \tilde{A}_K &= A - P G_K \end{aligned} \right\} \quad (109)$$

qu'il ne faut pas confondre, malgré sa ressemblance avec le processus itératif de Hewer(107).

6-2-4-Square root d'ordre 1

A la différence de la méthode de Hewer, on peut réaliser une formulation du type square root de (108) garantissant le caractère symétrique non négatif de la solution cherchée.

$$\left. \begin{aligned} R &= X^T \cdot X & Q &= V^T \cdot V \\ P_K &= S_K^T \cdot S_K & S_K &\text{ triangulaire supérieure} \end{aligned} \right\} \quad (110)$$

itération K:

-calculer $D = S_K \cdot B \quad E = S_K \cdot A.$ (111)

-réaliser la triangularisation orthogonale

$$\tau \begin{bmatrix} X & 0 \\ D & E \\ 0 & V \end{bmatrix} = \begin{bmatrix} L & K \\ 0 & S_{K+1} \\ 0 & 0 \end{bmatrix} \quad (112)$$

Le coût calcul pour $m < n$ ($B: n \cdot m$), reste inférieur à $3n^3$. Néanmoins, il apparaît clairement dans le tableau suivant, que le coût relatif à cette méthode est assez grand lorsque le régime permanent est suffisamment long à atteindre, c'est à dire lorsque le temps de réponse du système bouclé est important.

Tableau III

$N = 2^K$	$C = N \cdot 3n^3$
2	$6n^3$
4	$12n^3$
8	$24n^3$
16	$48n^3$
32	$96n^3$
64	$192n^3$

N: nombre d'itérations

C: coût calcul

7-CONCLUSIONS

Il est clair que quels qu'en soient les fondements théoriques, toutes les techniques numériques permettant de calculer P sont essentiellement itératives, dès lors on peut les classer par exemple, suivant leur complexité de mise en oeuvre, leur stabilité numérique, leur coût calcul par itération leur vitesse de convergence...

A ceci on peut ajouter leur domaine d'application: A régulière par exemple ou encore leur garantie concernant le caractère symétrique non négatif de P. C'est dans cet esprit que l'on va dégager parmi les huit méthodes, celles qui sont susceptibles de constituer les éléments de base d'un logiciel pour la résolution des équations de Riccati.

- 1-Méthode classique: MC équation (108)
- 2-Méthode classique stabilisée: MCS équation (109)
- 3-Square root d'ordre 1: SR1 équations (111)-(112)
- 4-Square root d'ordre 2:(produit *): SR2 équations (105)-(106)
- 5-Newton: N équation (107)
- 6-Vaughan: V équations (99)-(100)
- 7-Laub: L équations (101a)-(101b)
- 8-Fonction signe: FS équation (98)

On peut déjà noter que les cinq dernières présentent une convergence d'ordre 2 dont seule SR2 garantit le caractère symétrique non négatif de P.

Par ailleurs V et L supposent A régulière.

Sur le plan de la stabilité numérique, V peut être éliminée d'emblée (mise sous forme de Jordan, conditionnement) comparativement à L et FS qui opèrent sur les matrices symplectiques. Il en va de même des méthodes d'ordre 1: MC et MCS malgré son appellation stabilisée.

Deux points restent à considérer la mise en oeuvre et le coût calcul

7-CONCLUSIONS

Il est clair que quels qu'en soient les fondements théoriques, toutes les techniques numériques permettant de calculer F sont essentiellement itératives, dès lors on peut les classer par exemple, suivant leur complexité de mise en oeuvre, leur stabilité numérique, leur coût calcul par itération leur vitesse de convergence...

A ceci on peut ajouter leur domaine d'application: A régulière par exemple ou encore leur garantie concernant le caractère symétrique non négatif de P . C'est dans cet esprit que l'on va dégager parmi les huit méthodes, celles qui sont susceptibles de constituer les éléments de base d'un logiciel pour la résolution des équations de Riccati.

- 1-Méthode classique: MC équation (108)
- 2-Méthode classique stabilisée: MCS équation (109)
- 3-Square root d'ordre 1: SR1 équations (111)-(112)
- 4-Square root d'ordre 2:(produit $*$): SR2 équations (105)-(106)
- 5-Newton: N équation (107)
- 6-Vaughan: V équations (99)-(100)
- 7-Laub: L équations (101a)-(101b)
- 8-Fonction signe: FS équation (98)

On peut déjà noter que les cinq dernières présentent une convergence d'ordre 2 dont seule SR2 garantit le caractère symétrique non négatif de P .

Par ailleurs V et L supposent A régulière.

Sur le plan de la stabilité numérique, V peut être éliminée d'emblée

(mise sous forme de Jordan, conditionnement) comparativement à L et FS qui opèrent sur les matrices symplectiques. Il en va de même des méthodes d'ordre 1: MC et MCS malgré son appellation stabilisée.

Deux points restent à considérer la mise en oeuvre et le coût calcul

qu'on peut apprécier à l'aide du tableau 3 et ce pour chaque méthode. Ainsi deux méthodes nous semblent particulièrement intéressantes. Tout d'abord la version SR2 qui est peut être potentiellement la plus séduisante, car pour un codage de complexité moyenne, elle offre une convergence d'ordre 2 tout en garantissant numériquement le caractère symétrique non négatif de la solution cherchée. Et la méthode F.S que l'on a développée et qui semble être très compétitive sur le plan rapidité, est également la plus générale.

En effet le même code peut servir pour résoudre les équations de Riccati relatives aux systèmes discrets que celles relatives aux systèmes continus. Seuls les codages de l'hamiltonien et de P changent.

$$\begin{array}{l}
 \text{RICCATI} \\
 \text{continu:} \quad PA + A^T P - PVP + Q = 0 \\
 \text{discret:} \quad A^T P (I + VP)^{-1} A - P + Q = 0 \\
 \text{avec} \quad V = BR^{-1}B^T
 \end{array}$$

Enfin si l'on pose $V=0$, le même algorithme permet de résoudre les équations de Lyapunov correspondantes

$$\begin{array}{l}
 \text{LYAPUNOV} \\
 \text{continu:} \quad PA + A^T P + Q = 0 \\
 \text{discret:} \quad A^T P A - P + Q = 0
 \end{array}$$

qu'on peut apprécier à l'aide du tableau 3 et ce pour chaque méthode. Ainsi deux méthodes nous semblent particulièrement intéressantes. Tout d'abord la version SR2 qui est peut être potentiellement la plus séduisante, car pour un codage de complexité moyenne, elle offre une convergence d'ordre 2 tout en garantissant numériquement le caractère symétrique non négatif de la solution cherchée. Et la méthode F.S que l'on a développée et qui semble être très compétitive sur le plan rapidité, est également la plus générale.

En effet le même code peut servir pour résoudre les équations de Riccati relatives aux systèmes discrets que celles relatives aux systèmes continus. Seuls les codages de l'hamiltonien et de P changent.

$$\begin{array}{l}
 \text{RICCATI} \\
 \text{continu:} \quad PA + A^T P - P V F + Q = 0 \\
 \text{discret:} \quad A^T P (I + V F)^{-1} A - P + Q = 0 \\
 \text{avec} \quad V = B R^{-1} B^T
 \end{array}$$

Enfin si l'on pose $V=0$, le même algorithme permet de résoudre les équations de Lyapunov correspondantes

$$\begin{array}{l}
 \text{LYAPUNOV} \\
 \text{continu:} \quad PA + A^T P + Q = 0 \\
 \text{discret:} \quad A^T P A - P + Q = 0
 \end{array}$$

TABLEAU 3

	MC	MCS	SR1	SR2	N	V	L	F.S
Ordre de convergence	1	1	1	2	2	2	2	2
A régulière	NON	NON	NON	NON	NON	OUI	OUI	NON
$P \geq 0$	NON	NON	OUI	OUI	NON	NON	NON	NON
stabilité numérique	NON	NON	OUI	OUI	OUI	NON	OUI	OUI
mise en oeuvre	S	S	S	M	I	M	I	S
coût CALCUL	$\# 5 n^3$ ($m = \frac{1}{2}$)	$\# 5 n^3$ ($m = \frac{1}{2}$)	$\# 2 n^3$ ($m = \frac{1}{2}$)	$7 n^3$ $\forall m$	$3 n^3$ ($m = \frac{1}{2}$) $+ 20 n^3$ (eq. Lyg)	$70 n^3$ (*)	$75 n^3$ (*)	$5 n^3 \forall m$ (20 à $30 n^3$) (*)

* coût total pour obtenir P

S simple

M moyenne

I importante

TABLEAU 3

	MC	MCS	SR1	SR2	N	V	L	F.S
Ordre de convergence	1	1	1	2	2	2	2	2
A régulière	NON	NON	NON	NON	NON	OUI	OUI	NON
$P \geq 0$	NON	NON	OUI	OUI	NON	NON	NON	NON
stabilité numérique	NON	NON	OUI	OUI	OUI	NON	OUI	OUI
mise en oeuvre	S	S	S	M	I	M	I	S
coût calcul	$\# 5 n^3$ ($m = n/2$)	$\# 5 n^3$ ($m = n/2$)	$\# 2 n^3$ ($m = n/2$)	$7 n^3$ γm	$3 n^3$ ($m = n/2$) $+ 20 n^3$ (eq. Lyg)	$70 n^3$ (*)	$75 n^3$ (*)	$5 n^3 \gamma m$ (20 à $30 n^3$) (*)

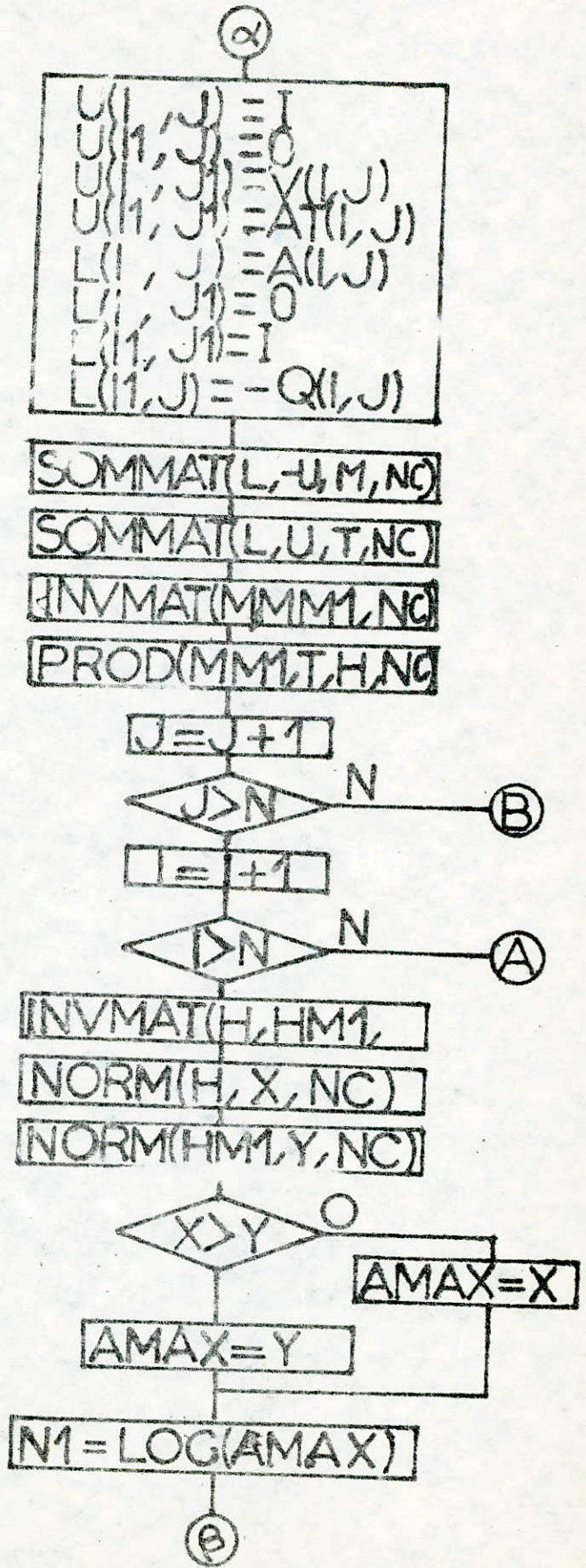
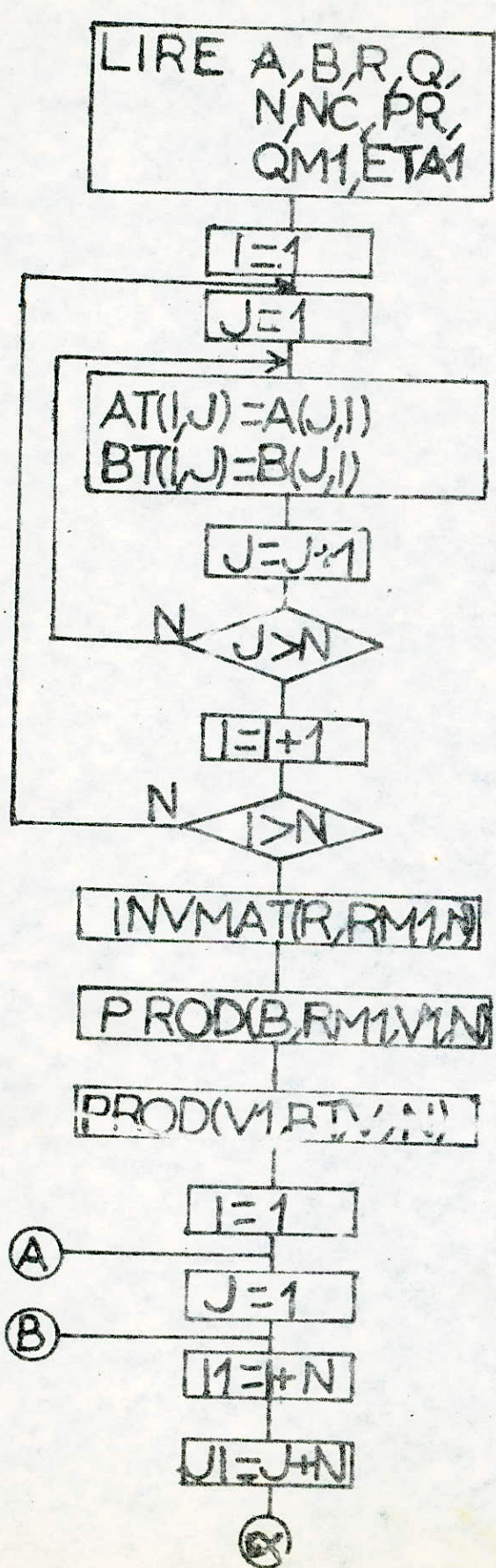
* coût total pour obtenir P

S simple

M moyenne

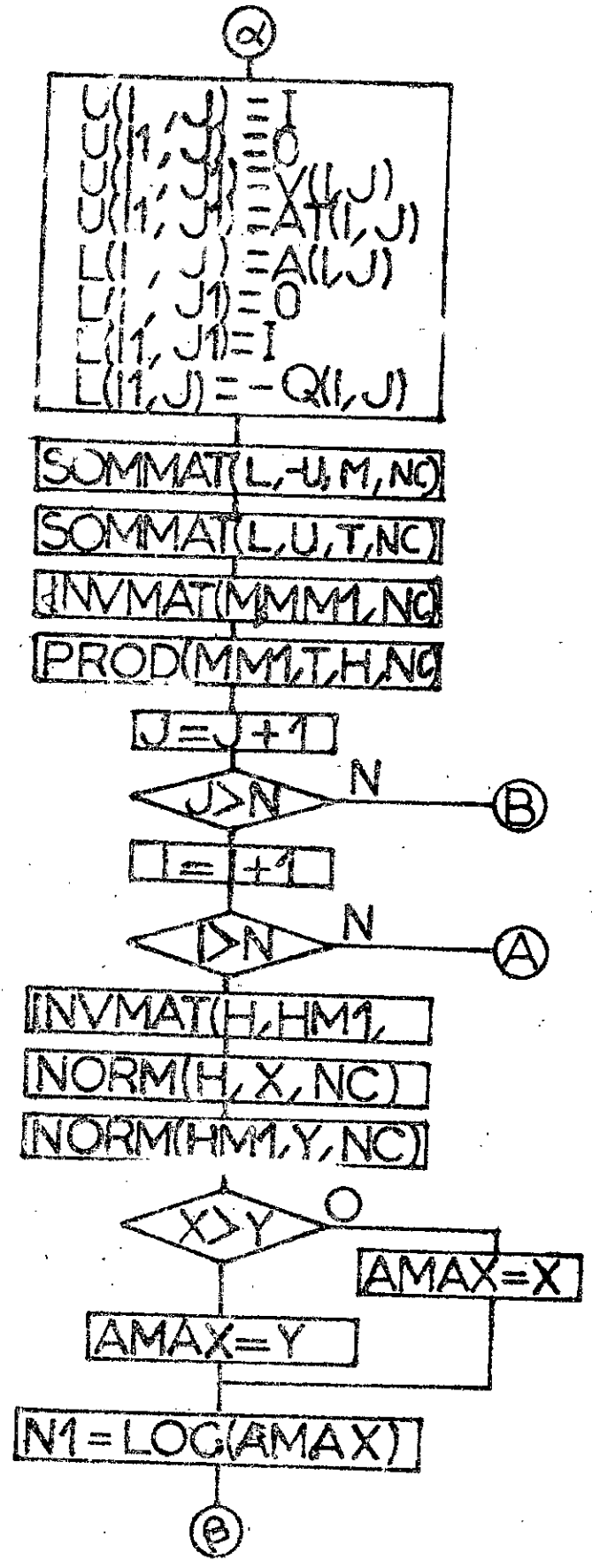
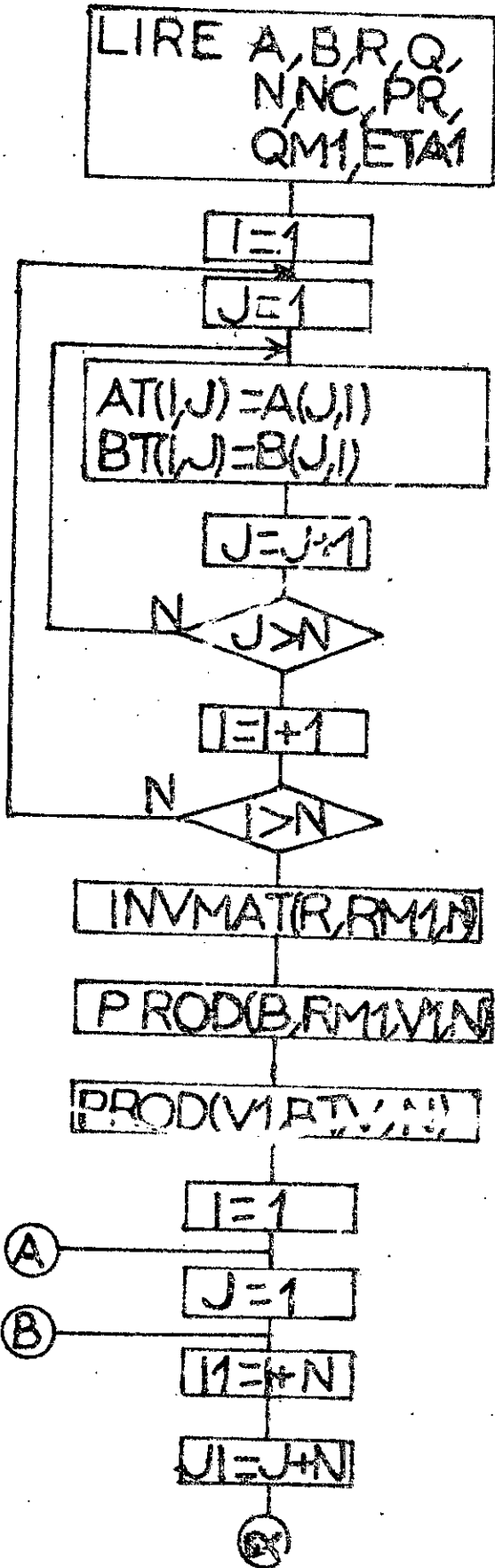
I importante

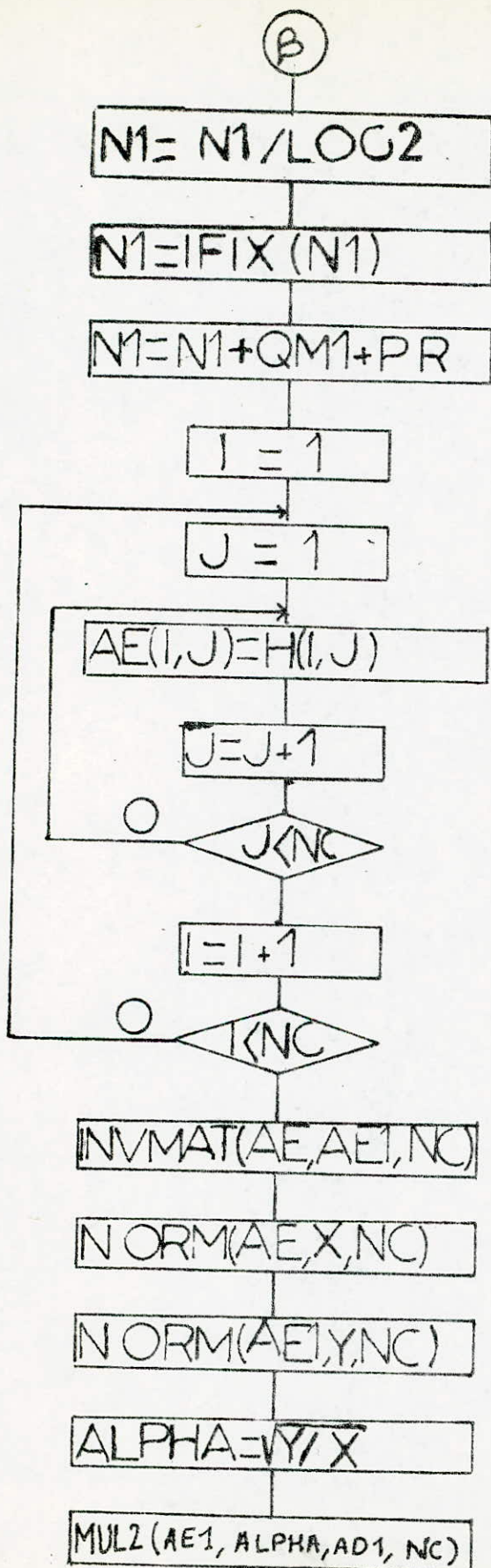
ORGANIGRAMME PRINCIPAL



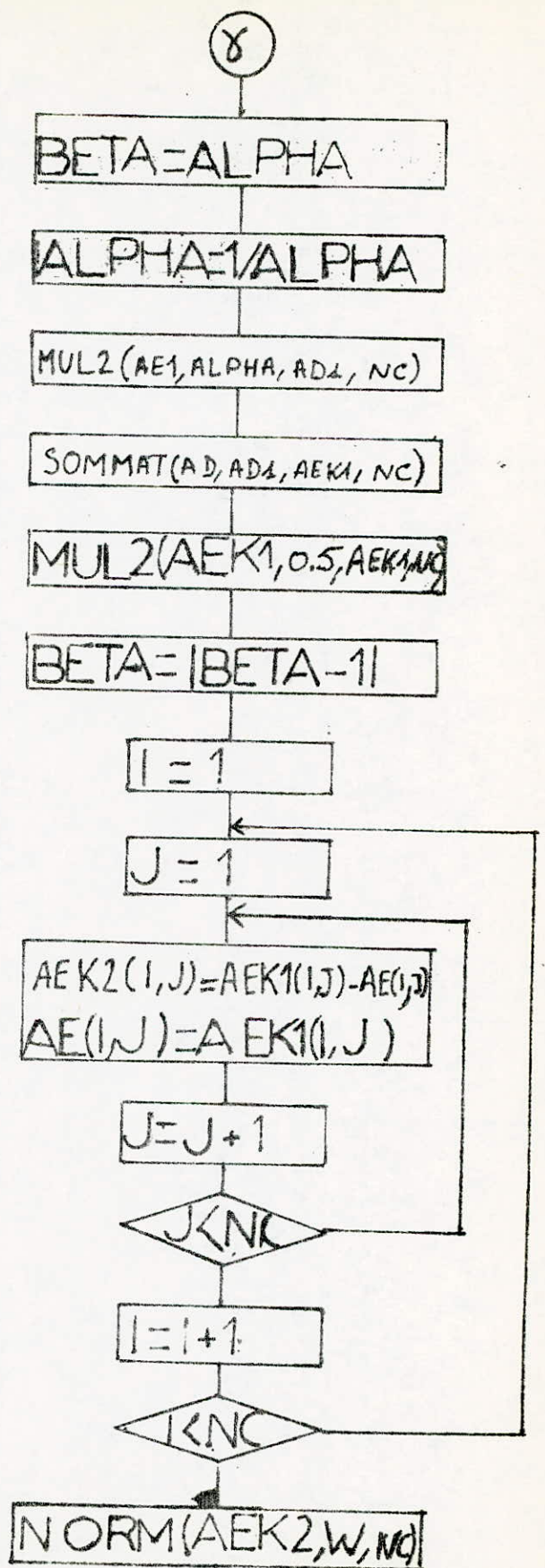
ORGANIGRAMME

PRINCIPAL

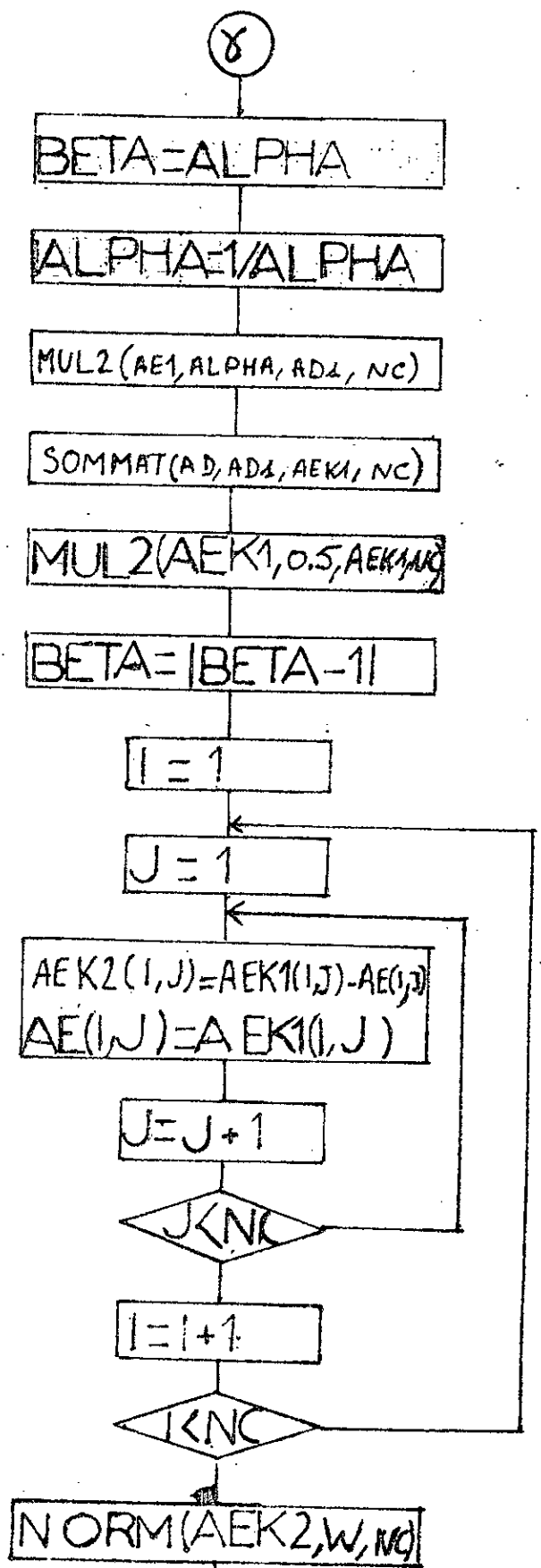
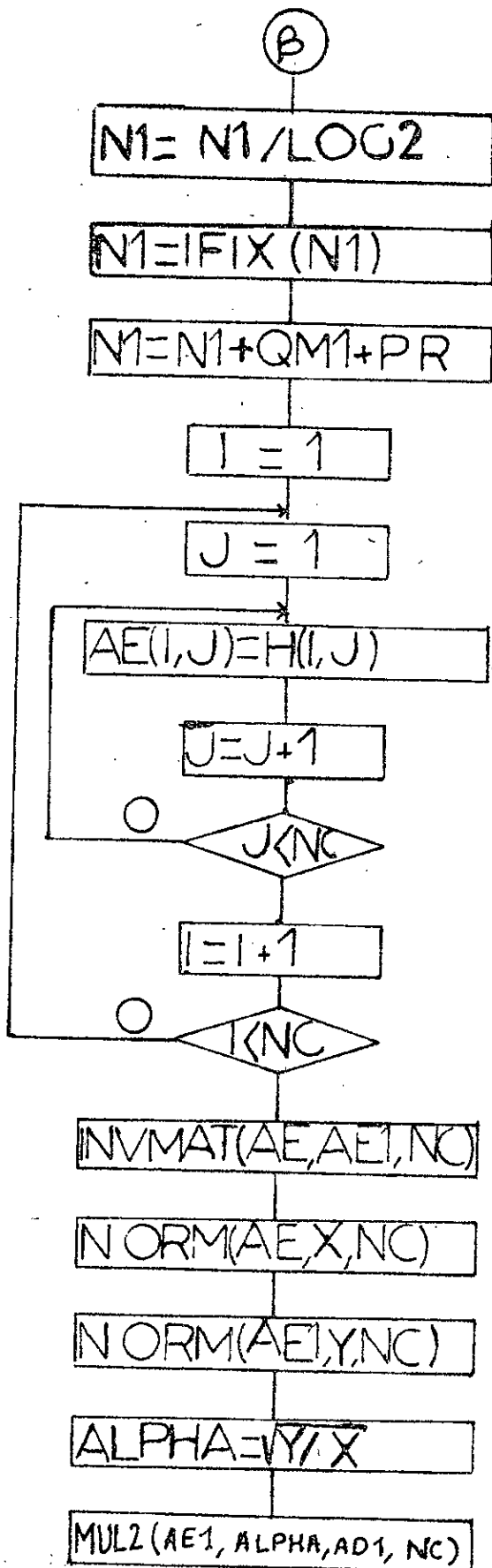


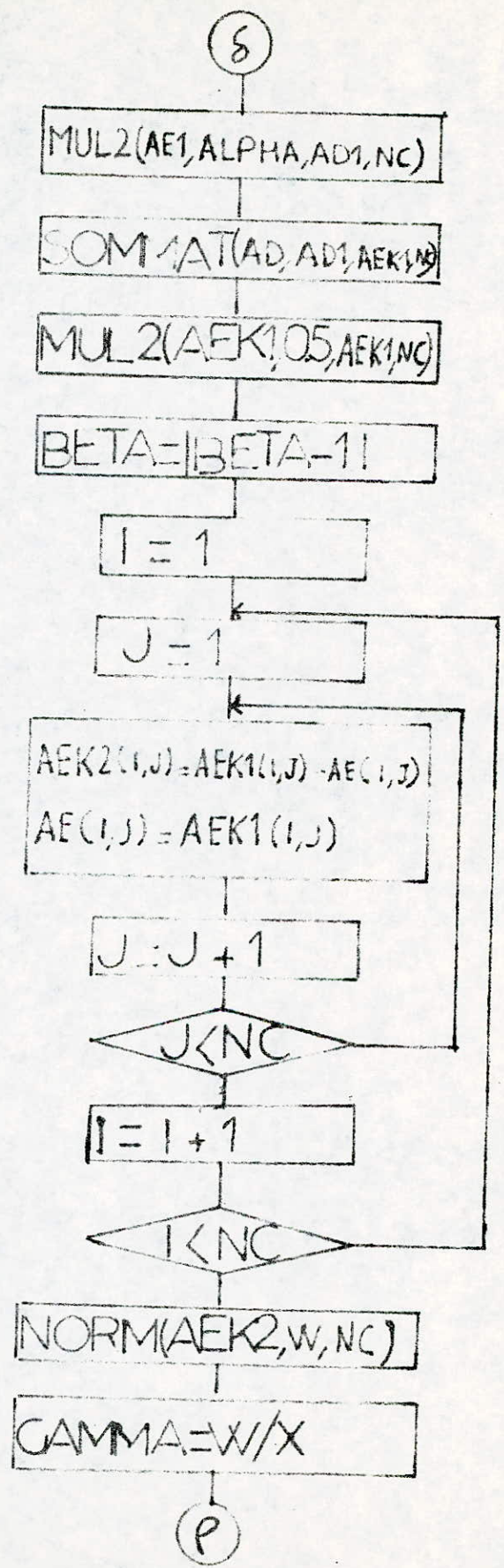
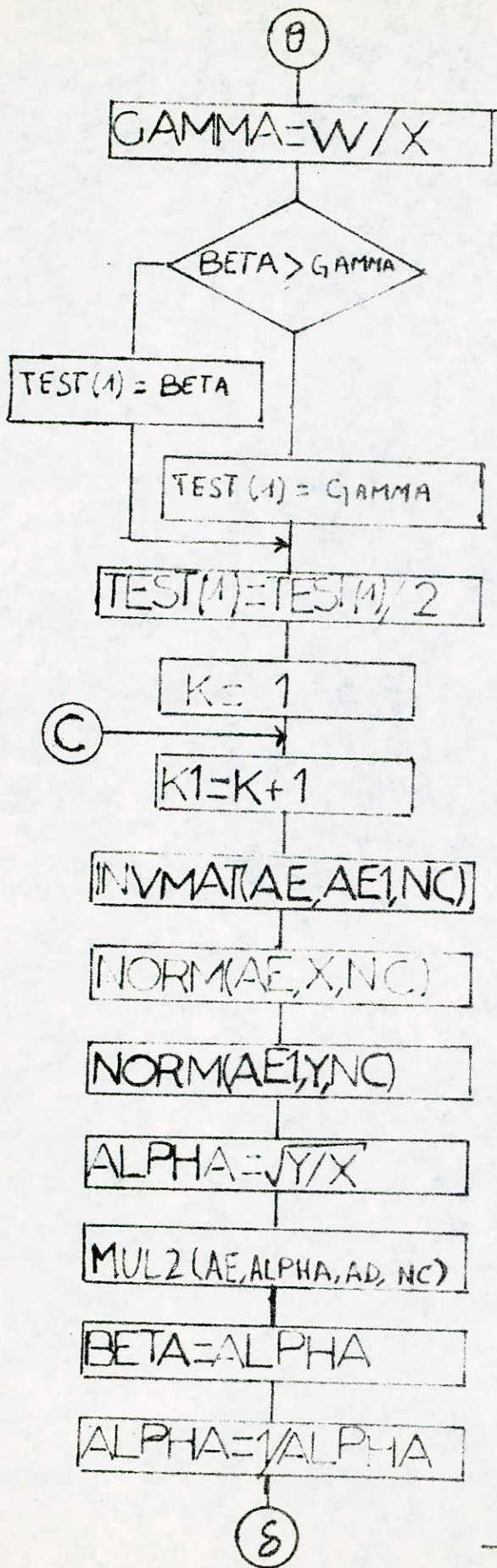


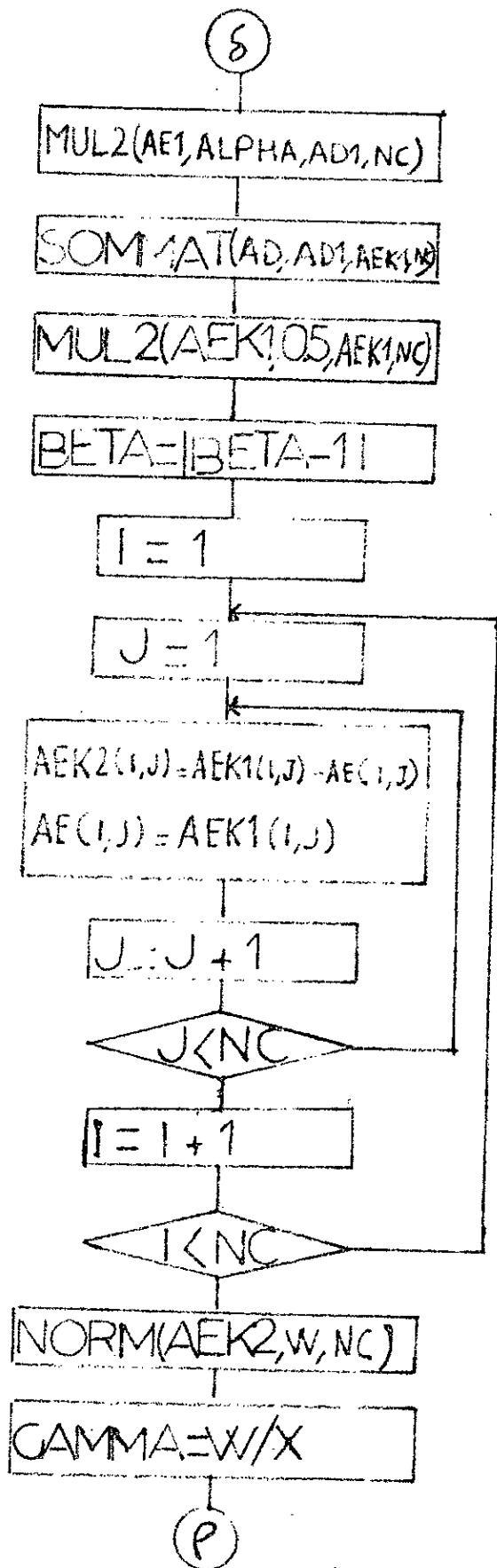
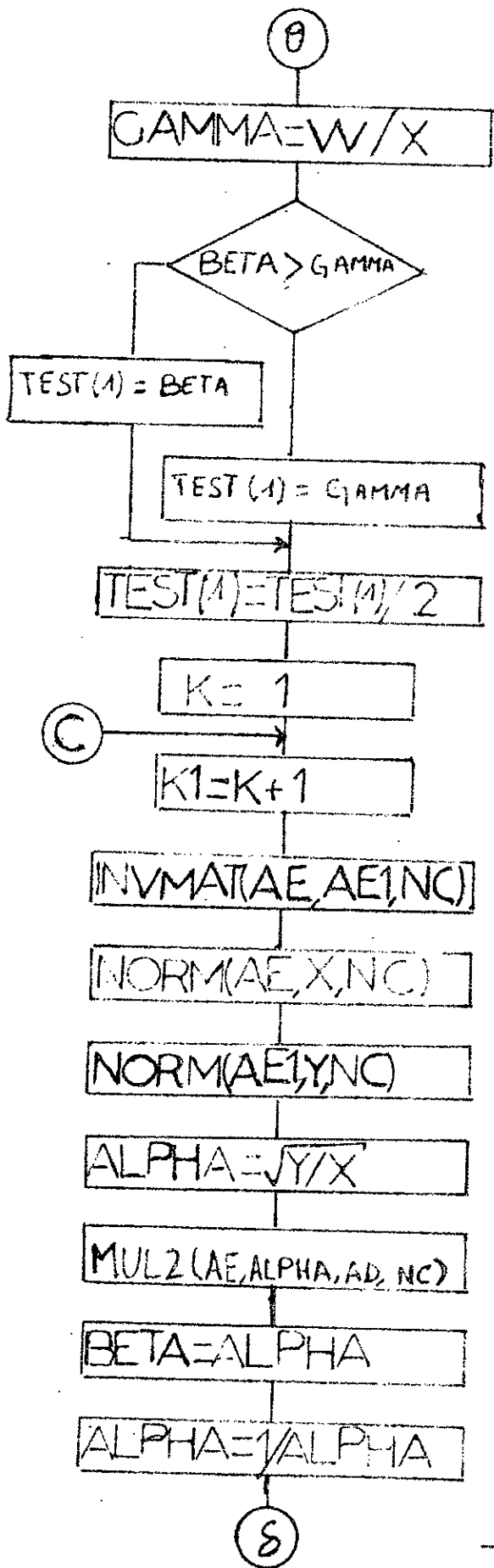
⊙ γ

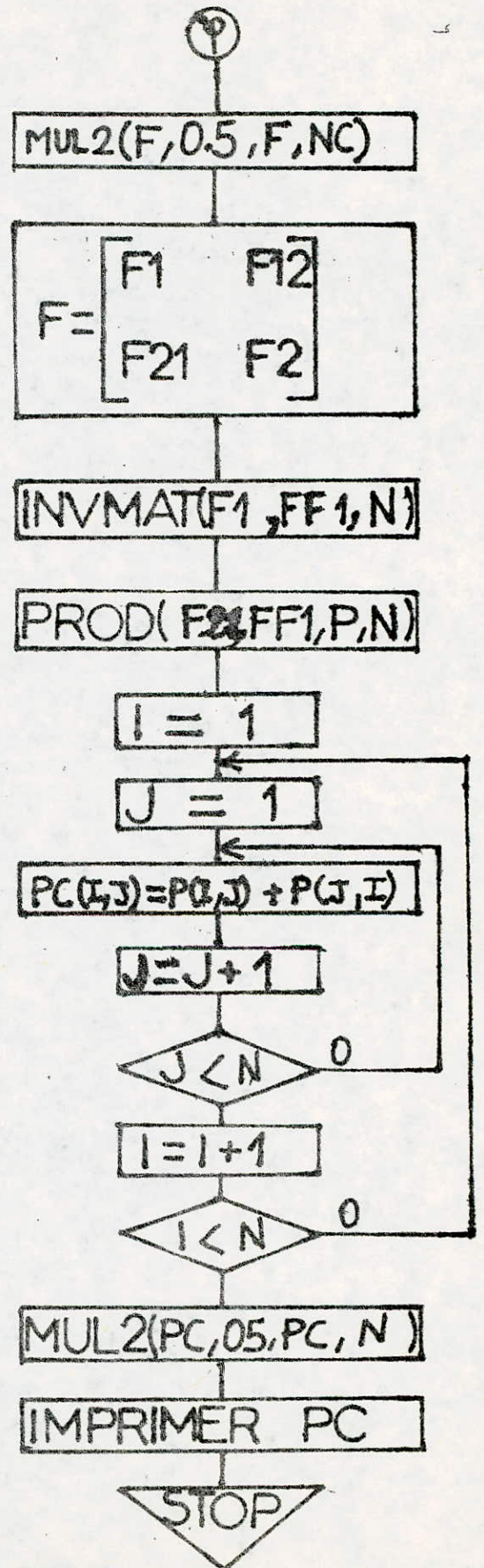
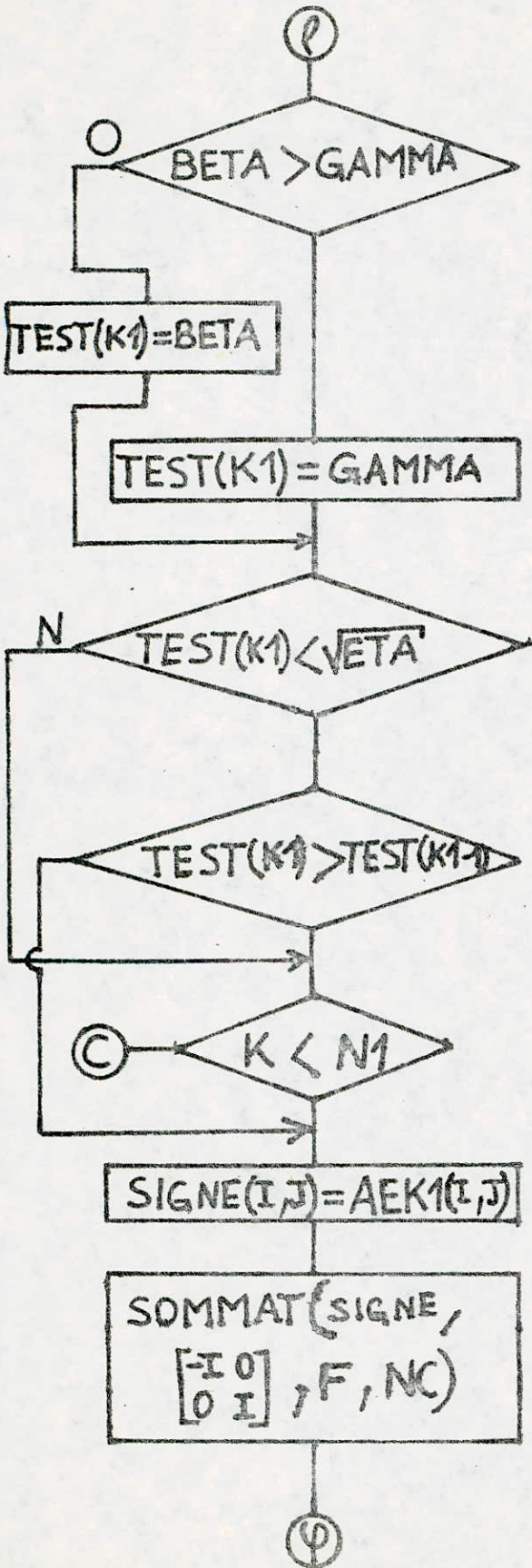


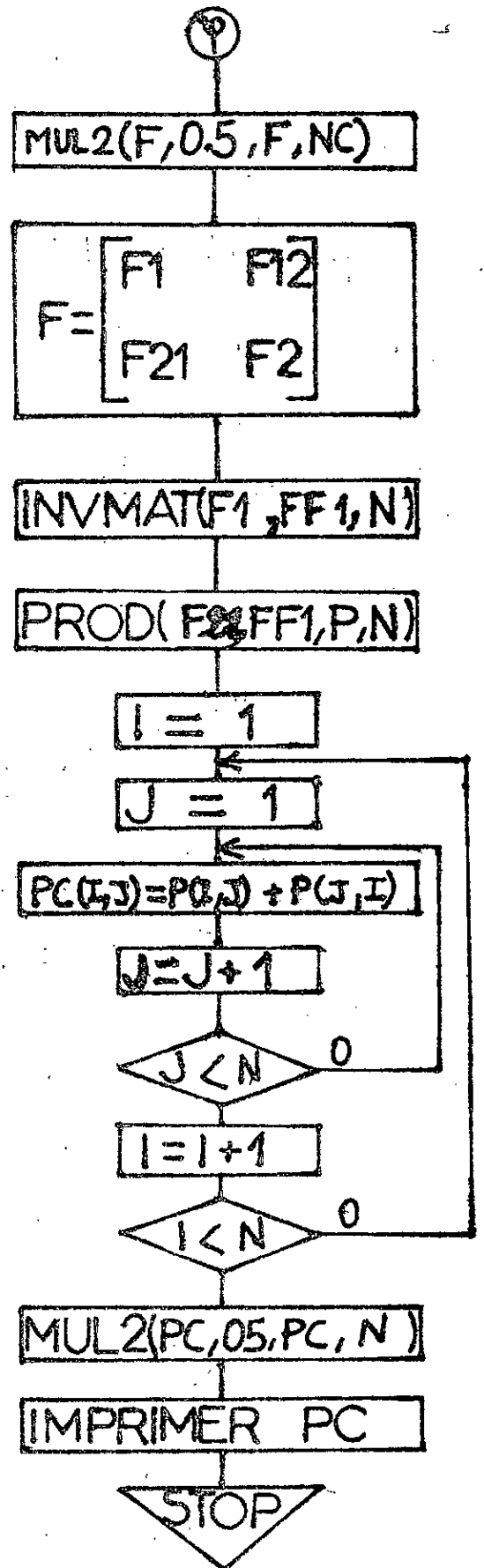
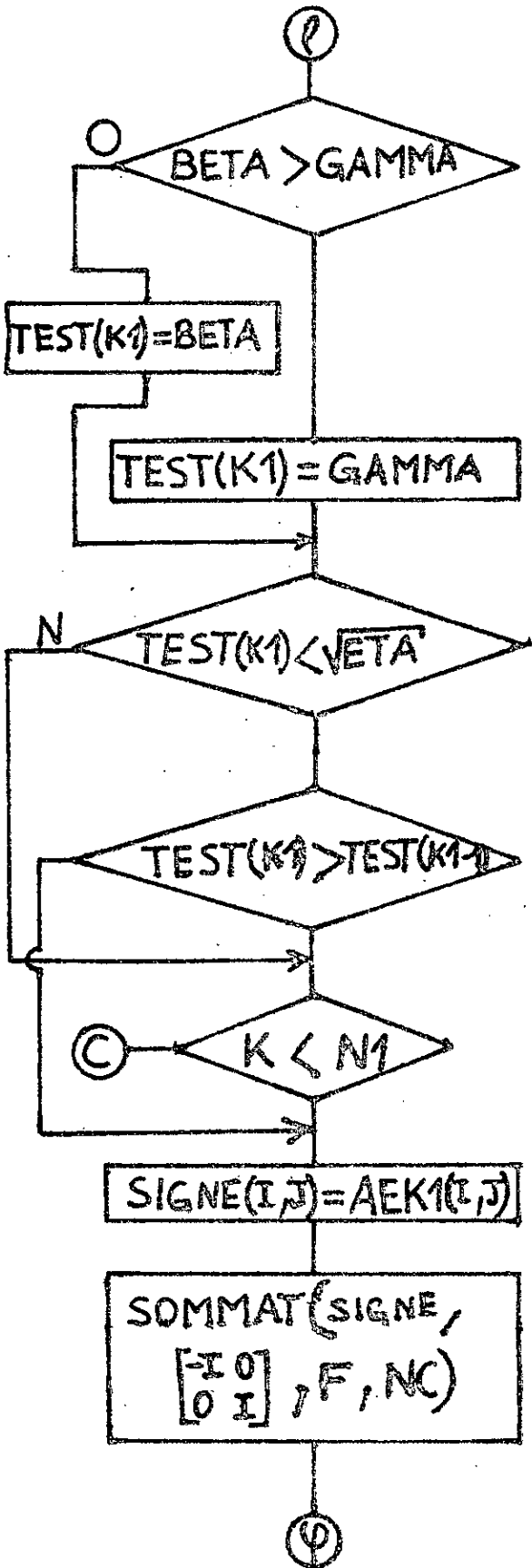
⊙ θ











CHAPITRE III/

APPLICATIONS NUMERIQUES

CHAPITRE III/

APPLICATIONS NUMERIQUES

APPLICATIONS

1-DESCRIPTION DES PROGRAMMES

1-1-Objet des programmes

On considère ici la résolution des équations de Riccati sous la forme:

$$PA + A^T P - PVP + Q = 0 \quad \text{cas continu}$$

$$A^T P (I + VP)^{-1} A - P + Q = 0 \quad \text{cas discret}$$

Ce programme peut être étendu à la résolution des équations de Lyapounov ($V=C$)

$$PA + A^T P + Q = C \quad \text{Cas continu}$$

$$A^T P - P + Q = C \quad \text{Cas discret}$$

I-2-Algorithmie

IL consiste à calculer la matrice signe de H , ($2n \cdot 2n$), associée aux équations hamiltoniennes système. Cette matrice doit être elle-même hamiltonienne, pour les systèmes continus, la matrice H est par nature hamiltonienne tandis que pour le cas discret, elle doit subir deux transformations; une première pour la rendre symplectique, une seconde transformation bilinéaire pour la rendre hamiltonienne.

Les principales étapes du programme se résument de la manière suivante

A/ Calculs préliminaires:

- Déterminer la matrice V telle que $V = BR^{-1}B^T$
- calcul de H
pour le cas continu.

$$H = \begin{bmatrix} A & V \\ -Q & -A \end{bmatrix}$$

Pour le cas discret, on doit calculer d'abord U et L telles que

$$U = \begin{bmatrix} I \\ 0 \end{bmatrix} \quad V = \begin{bmatrix} V \\ A^T \end{bmatrix} \quad L = \begin{bmatrix} A & 0 \\ -Q & I \end{bmatrix}$$

APPLICATIONS

1-DESCRIPTION DES PROGRAMMES

1-1-Objet des programmes

On considère ici la résolution des équations de Riccati sous la forme:

$$PA + A^T P - PVP + Q = 0 \quad \text{cas continu}$$

$$A^T P (I + VP)^{-1} A - P + Q = 0 \quad \text{cas discret}$$

Ce programme peut être étendu à la résolution des équations de Lyapunov ($V=0$)

$$PA + A^T P + Q = C \quad \text{Cas continu}$$

$$A^T P - P + Q = C \quad \text{Cas discret}$$

I-2-Algorithmes

IL consiste à calculer la matrice signe de H , $(2n \cdot 2n)$, associée aux équations hamiltoniennes système. Cette matrice doit être elle-même hamiltonienne, pour les systèmes continus, la matrice H est par nature hamiltonienne tandis que pour le cas discret, elle doit subir deux transformations; une première pour la rendre symplectique, une seconde transformation bilinéaire pour la rendre hamiltonienne.

Les principales étapes du programme se résument de la manière suivante

A/ Calculs préliminaires:

- Déterminer la matrice V telle que $V = BR^{-1}B^T$
- calcul de H
pour le cas continu.

$$H = \begin{bmatrix} A & V \\ -Q & -A^T \end{bmatrix}$$

Pour le cas discret, on doit calculer d'abord U et L telles que

$$U = \begin{bmatrix} I & \\ & V \end{bmatrix} \quad L = \begin{bmatrix} A & C \\ -Q & I \end{bmatrix}$$

Poser H telle que

$$H = (L - U)^{-1} (L + U)$$

- calculer le nombre d'itérations total tel que

$$NI = N + QMI + FR$$

$$\text{avec } N = \lfloor \log_2 \max(\|H\|, \|H^{-1}\|) \rfloor$$

$$QMI(\eta) = 1 + \lfloor \log_2 (0,64/\eta) \rfloor$$

FR(η) = nombre d'itérations supplémentaires en fonction
de la précision demandée

B/ calculs des itérés

On a posé à priori $AE = A_c^*$

puis on a calculé α comme étant

$$\alpha = \sqrt{\frac{\|AE^{-1}\|}{\|AE\|}}$$

Puis on a calculé le premier itéré à savoir

$$AEI = \frac{1}{2} \left(\alpha AE + \frac{1}{\alpha} AE^{-1} \right)$$

Après cela on réalise la boucle suivante:

I- Calculer

$$\alpha_{k-1} = \sqrt{\frac{\|AE_{k-1}^{-1}\|}{\|AE_{k-1}\|}}$$

$$AE_k = \frac{1}{2} \left(\alpha_{k-1} \cdot AE_{k-1} + \frac{1}{\alpha_{k-1}} \cdot AE_{k-1}^{-1} \right)$$

si $TEST(k) < \sqrt{\eta} = \text{ETA1}$ on passe à l'étape 3

2- On passe à l'étape 4

3- Si $test(k) > (1/2) \cdot test(k-1)$ on passe à l'étape 5

4- Si $k < N+1$ on retourne vers 1

5- La matrice signe est donnée par

$$SIGNE(I, J) = AE_k(I, J)$$

Poser H telle que

$$H = (L - U)^{-1} (L + U)$$

- calculer le nombre d'itérations total tel que

$$NI = N + QMI + FR$$

$$\text{avec } N = \lfloor \log_2 \max(\|H\|, \|H^{-1}\|) \rfloor$$

$$QMI(\eta) = 1 + \lfloor \log_2 (0,64/\eta) \rfloor$$

FR(η) = nombre d'itérations supplémentaires en fonction
de la précision demandée

B/ calculs des itérés

On a posé à priori $AE = A_c^*$

puis on a calculé α comme étant

$$\alpha = \sqrt{\frac{\|AE^{-1}\|}{\|AE\|}}$$

Puis on a calculé le premier itéré à savoir

$$AEI = \frac{1}{2} (\alpha AE + \frac{1}{\alpha} AE^{-1})$$

Après cela on réalise la boucle suivante:

I- Calculer

$$\alpha_{k-1} = \sqrt{\frac{\|AE_{k-1}^{-1}\|}{\|AE_{k-1}\|}}$$

$$AE_k = \frac{1}{2} (\alpha_{k-1} \cdot AE_{k-1} + \frac{1}{\alpha_{k-1}} \cdot AE_{k-1}^{-1})$$

si $TEST(k) < \sqrt{\eta} = \text{ETA1}$ on passe à l'étape 3

2- On passe à l'étape 4

3- Si $test(k) > (1/2) \cdot test(k-1)$ on passe à l'étape 5

4- Si $k < N+1$ on retourne vers 1

5- La matrice signe est donnée par

$$SIGNE(I, J) = AE_k(I, J)$$

C/Calcul de la matrice P: solution de l'équation de Riccati

1-cas continu:

On calcule F comme étant

$$F = \frac{1}{2}(I + \text{SIGNE}) = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix}$$

La matrice P est alors donnée par

$$P = -F_{12}^{-1} \cdot F_{21}$$

2-cas discret

La matrice F est calculée comme suit

$$F = \text{SIGNE} + \begin{bmatrix} -I & C \\ 0 & I \end{bmatrix} = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix}$$

d'où

$$P = F_{21} \cdot F_{11}^{-1}$$

1-3-Définition des arguments

Soit le système

$$\dot{X} = A \cdot X + B \cdot U$$

$$Y = C \cdot X$$

et le coût associé

$$J = \int_0^{\infty} (X^T Q X + U^T R U) dt$$

On définit alors les arguments suivants:

A: matrice réelle(n.n)

B: matrice réelle(n.n)

C: matrice réelle(n.n)

R: matrice réelle(n.n)

N: entier égal à n(dimension de A, B, C, R)

NC: entier égal à 2n(dimension de F)

C/Calcul de la matrice P: solution de l'équation de Riccati.

1-cas continu:

On calcule F comme étant

$$F = \frac{1}{2}(I + \text{SIGNE}) = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix}$$

La matrice P est alors donnée par

$$P = -F_{12}^{-1} \cdot F_{21}$$

2-cas discret

La matrice F est calculée comme suit

$$F = \text{SIGNE} + \begin{bmatrix} -I & C \\ 0 & I \end{bmatrix} = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix}$$

d'où

$$P = F_{21} \cdot F_{11}^{-1}$$

1-3-Définition des arguments

Soit le système

$$\dot{X} = A \cdot X + B \cdot U$$

$$Y = C \cdot X$$

et le coût associé

$$J = \int_0^{\infty} (X^T Q X + U^T R U) dt$$

On définit alors les arguments suivants:

A: matrice réelle(n.n)

B: matrice réelle(n.n)

C: matrice réelle(n.n)

R: matrice réelle(n.n)

N: entier égal à n(dimension de A, B, C, R)

NC:entier égal à 2n(dimension de F)

η : précision machine (simple précision $\eta = 10^{-16}$)

ETA1: $\sqrt{\eta}$

PR: nombre d'itérations supplémentaires pour obtenir une précision donnée (on a pris PR=5)

QMI: nombre d'itérations relatif à la convergence de l'argument e_k
Pour le cas réel QMI=0

NI: nombre d'itérations total

TEST: variable logique

SIGNE: matrice signe de H (2n.2n)

P: solution de l'équation de Riccati

PC: matrice P corrigée

1-4- Structure du programme

Le programme fait appel à 5 sous programmes et utilise les fonctions mathématiques et utilitaires suivantes:

SQRT(X): racine carrée de X

ALOG(X): logarithme népérien de X

ABS(X): valeur absolue de X

IFIX : conversion réel- entier

sous programmes

PROD(M1, M2, M3, N): fait le produit de la matrice M1 par M2 toute deux de dimension N et met le résultat dans M3
 $M1.M2=M3$

MUL2(M, E, MP, N): fait le produit d'une matrice (N.N) par un scalaire E et met le résultat dans MP
 $E.M=MP$

SOMMAT(D, E, F, N): fait la somme de deux matrice DetE de dimension (N.N) et met le résultat dans F: $(D+E)=F$

η : précision machine (simple précision $\eta = 10^{-16}$)

ETA1: $\sqrt{\eta}$

PR: nombre d'itérations supplémentaires pour obtenir une précision donnée (on a pris PR=5)

QMI: nombre d'itérations relatif à la convergence de l'argument ϵ_k
Pour le cas réel QMI=0

NI: nombre d'itérations total

TEST: variable logique

SIGNE: matrice signe de H (2n.2n)

P: solution de l'équation de Riccati

PC: matrice P corrigée

1-4- Structure du programme

Le programme fait appel à 5 sous programmes et utilise les fonctions mathématiques et utilitaires suivantes:

SQRT(X): racine carrée de X

ALOG(X): logarithme népérien de X

ABS(X): valeur absolue de X

IFIX : conversion réel- entier

sous programmes

PROD(M1, M2, M3, N): fait le produit de la matrice M1 par M2 toute deux de dimension N et met le resultat dans M3.

$$M1.M2=M3.$$

MUL2(M, E, MP, N): fait le produit d'une matrice (N.N) par un scalaire E et met le resultat dans MP

$$E.M=MP$$

SOMMAT(D, E, F, N): fait la somme de deux matrice DetE de dimension (N,N) et met le resultat dans F: (D+E)=F

NORM(A,X,N): calcule la norme de deux façons et choisit la plus petite

$$X = \min(\|A\|_1, \|A\|_\infty)$$

MRINV(A,B,KCD,DET,EPS,IL,IC): inverse la matrice A et la met dans B

$$B = A^{-1}$$

(on a fait appel à ce sous programme directement

de la bibliothèque du M I T R A 125)

2-APPLICATIONS NUMERIQUES

Notre but dans ce chapitre est d'illustrer notre étude avec des exemples précis tout en mettant en évidence quelques caractéristiques de nos programmes

2-1-RESOLUTION DE L'EQUATION DE RICCATI

2-1-1 Cas continu

Dans ce paragraphe nous allons considérer deux cas de variations du paramètre N1 (nombre d'itérations total) par rapport:

- aux pôles
- à la dimension de notre système

a/Variation de N1 par rapport aux pôles

Pour cela nous allons considérer deux systèmes de même dimension l'un stable l'autre instable.

→ Soit un système dynamique continu dont la matrice de gain en boucle ouverte est:

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

et le coût quadratique: $J = \int_0^{\infty} (x^T Q x + u^T R u) dt$

avec $Q = \begin{bmatrix} 4 & 0 \\ 0 & L \end{bmatrix}$ et $V = \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix}$

NORM(A,X,N): calcule la norme de deux façons et choisit la plus petite

$$X = \min(\|A\|_1, \|A\|_\infty)$$

MRINV(A,B,KCD,DET,EPS,IL,IC): inverse la matrice A et la met dans B

$$B = A^{-1}$$

(on a fait appel à ce sous programme directement

de la bibliothèque du M I T R A 125)

2-APPLICATIONS NUMERIQUES

Notre but dans ce chapitre est d'illustrer notre étude avec des exemples précis tout en mettant en évidence quelques caractéristiques de nos programmes

2-1-RESOLUTION DE L'EQUATION DE RICCATI

2-1-1 Cas continu

Dans ce paragraphe nous allons considérer deux cas de variations du paramètre N1 (nombre d'itérations total) par rapport:

- aux pôles
- à la dimension de notre système

a/Variation de N1 par rapport aux pôles

Pour cela nous allons considérer deux systèmes de même dimension l'un stable l'autre instable.

→ Soit un système dynamique continu dont la matrice de gain en boucle ouverte est:

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

et le coût quadratique: $J = \int_0^{\infty} (x^T Q x + u^T R u) dt$

avec $Q = \begin{bmatrix} 4 & 0 \\ 0 & L \end{bmatrix}$ et $V = \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix}$

L'équation de Riccati à résoudre sera donc:

$$\begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} P + P \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} + P \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} P + \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

résultats:

-nombre d'itérations $N1 = 7$

-matrice P:

$$P = \begin{bmatrix} 5,999988 & 1,999993 \\ 2,0 & 4,99994 \end{bmatrix}$$

la dissymétrie est de l'ordre de 7.10^{-6}

Après correction la matrice P devient:

$$P = \begin{bmatrix} 5,999988 & 1,99996 \\ 1,99996 & 4,99994 \end{bmatrix} \approx \begin{bmatrix} 6 & 2 \\ 2 & 5 \end{bmatrix}$$

Le temps d'exécution était de 3" 1

- le deuxième système est tel que le gain en boucle ouverte soit:

$$A = \begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix}$$

Q et V sont les mêmes que pour le cas précédent

l'équation de Riccati à résoudre serait donc

$$\begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix} P + P \begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix} + P \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} P + \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

- nombre d'itérations: $N1 = 7$

-matrice P(solution):

$$P = \begin{bmatrix} 5,999986 & 1,999987 \\ 2,000002 & 0,999999 \end{bmatrix}$$

L'équation de Riccati à résoudre sera donc:

$$\begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} P + P \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} + P \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} P + \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

résultats:

-nombre d'itérations N1 = 7

-matrice P:

$$P = \begin{bmatrix} 5,999988 & 1,999993 \\ 2,0 & 4,99994 \end{bmatrix}$$

la dissymétrie est de l'ordre de 7.10^{-6}

Après correction la matrice P devient:

$$P = \begin{bmatrix} 5,999988 & 1,99996 \\ 1,99996 & 4,99994 \end{bmatrix} \approx \begin{bmatrix} 6 & 2 \\ 2 & 5 \end{bmatrix}$$

Le temps d'exécution était de 3" 1

- le deuxième système est tel que le gain en boucle ouverte soit:

$$A = \begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix}$$

Q et V sont les mêmes que pour le cas précédent

l'équation de Riccati à résoudre serait donc

$$\begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix} P + P \begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix} + P \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} P + \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

- nombre d'itérations: N1 = 7

-matrice P(solution):

$$P = \begin{bmatrix} 5,999986 & 1,999987 \\ 2,000002 & 0,999999 \end{bmatrix}$$

-dissymétrie est d'environ 10^{-5}

-après correction la matrice P devient

$$P = \begin{bmatrix} 5,999986 & 1,999995 \\ 1,999995 & 0,999999 \end{bmatrix} \approx \begin{bmatrix} 6 & 2 \\ 2 & 1 \end{bmatrix}$$

- temps d'exécution 2" 2

b/ Variation de N1 par rapport à la dimension

Soit le système dynamique continu dont la matrice gain en boucle ouverte est:

$$A = \begin{bmatrix} -1 & 0 & 0 \\ -1 & 0 & -2 \\ 0 & 1 & -1 \end{bmatrix}$$

le coût quadratique J est tel que

$$Q = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} \quad \text{et} \quad V = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

L'équation de Riccati à résoudre serait donc:

$$\begin{bmatrix} -1 & -1 & 0 \\ 0 & 0 & 1 \\ 0 & -2 & -1 \end{bmatrix} P + P \begin{bmatrix} -1 & 0 & 0 \\ -1 & 0 & -2 \\ 0 & 1 & -1 \end{bmatrix} + P \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} P + \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} = 0$$

Résultats

-nombre d'itérations: N1 = 7

-matrice P(solution):

$$\begin{bmatrix} 0,5944349 & -0,3234051 & 0,3047468 \\ -0,3234048 & 1,2125367 & -0,2125883 \\ 0,3047490 & -0,2126026 & 1,8561363 \end{bmatrix}$$

-dissymétrie est d'environ 10^{-5}

-après correction la matrice P devient

$$P = \begin{bmatrix} 5,999986 & 1,999995 \\ 1,999995 & 0,999999 \end{bmatrix} \approx \begin{bmatrix} 6 & 2 \\ 2 & 1 \end{bmatrix}$$

- temps d'exécution 2" 2

b/ Variation de N1 par rapport à la dimension

Soit le système dynamique continu dont la matrice gain en boucle ouverte est:

$$A = \begin{bmatrix} -1 & 0 & 0 \\ -1 & 0 & -2 \\ 0 & 1 & -1 \end{bmatrix}$$

le coût quadratique J est tel que

$$Q = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} \quad \text{et} \quad V = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

L'équation de Riccati à résoudre serait donc:

$$\begin{bmatrix} -1 & -1 & 0 \\ 0 & 0 & 1 \\ 0 & -2 & -1 \end{bmatrix} P + P \begin{bmatrix} -1 & 0 & 0 \\ -1 & 0 & -2 \\ 0 & 1 & -1 \end{bmatrix} + P \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} P + \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} = 0$$

Résultats

-nombre d'itérations: N1 = 7

-matrice P(solution):

$$\begin{bmatrix} 0,5944349 & -0,3234051 & 0,3047468 \\ -0,3234048 & 1,2125367 & -0,2125883 \\ 0,3047490 & -0,2126026 & 1,8561363 \end{bmatrix}$$

- la dissymétrie est de l'ordre de 10^{-5} à 10^{-6}

- la matrice P corrigée est

$$P = \begin{bmatrix} 0,5944347 & -0,3234050 & 0,3047479 \\ -0,3234050 & 1,2125367 & -0,2125955 \\ 0,3047479 & -0,2125955 & 1,8561363 \end{bmatrix}$$

- le temps d'exécution est de 3" 8

Conclusions

D'après les résultats obtenus, on vient de confirmer quelques caractéristiques de cette méthode (que l'on avait citées dans la partie théorique) à savoir que:

- l'on peut trouver la commande optimale d'un système qu'il soit stable ou instable, de plus l'instabilité n'influe en rien sur le nombre d'itérations total, c'est à dire, sur la rapidité de la convergence de notre algorithme.
- la dimension de notre système n'influe pas ou du moins très peu sur le nombre d'itérations que doit effectuer notre algorithme pour converger.

Remarque

En ce qui concerne la variation du nombre d'itérations par rapport au spectre de la matrice, cela sera mis en évidence dans la résolution de l'équation de Lyapunov et les conclusions s'y rapportant s'appliquent intégralement au cas de la résolution de l'équation de Riccati.

2-1-2/ Cas discret

Pour le cas discret les conclusions sont identiques à celles du cas continu vu que l'algorithme est presque le même.

Néanmoins, on a pris un exemple particulier où l'on doit faire d'abord des calculs préliminaires (calcul de V et de Q), c'est à dire en partant uniquement du système d'équations d'état.

- la dissymétrie est de l'ordre de 10^{-5} à 10^{-6}

- la matrice P corrigée est

$$P = \begin{bmatrix} 0,5944347 & -0,3234050 & 0,3047479 \\ -0,3234050 & 1,2125367 & -0,2125955 \\ 0,3047479 & -0,2125955 & 1,8561363 \end{bmatrix}$$

- le temps d'exécution est de 3" 8

Conclusions

D'après les résultats obtenus, on vient de confirmer quelques caractéristiques de cette méthode (que l'on avait citées dans la partie théorique) à savoir que:

- l'on peut trouver la commande optimale d'un système qu'il soit stable ou instable, de plus l'instabilité n'influe en rien sur le nombre d'itérations total, c'est à dire, sur la rapidité de la convergence de notre algorithme.
- la dimension de notre système n'influe pas ou du moins très peu sur le nombre d'itérations que doit effectuer notre algorithme pour converger.

Remarque

En ce qui concerne la variation du nombre d'itérations par rapport au spectre de la matrice, cela sera mis en évidence dans la résolution de l'équation de Lyapunov et les conclusions s'y rapportant s'appliquent intégralement au cas de la résolution de l'équation de Riccati.

2-1-2/ Cas discret

Pour le cas discret les conclusions sont identiques à celles du cas continu vu que l'algorithme est presque le même.

Néanmoins, on a pris un exemple particulier où l'on doit faire d'abord des calculs préliminaires (calcul de V et de Q), c'est à dire en partant uniquement du système d'équations d'état.

Soit donc un système discret de dimension 5 dont l'équation d'état est

$$X_{K+1} = A \cdot X_K + B \cdot U_K$$

et le coût

$$J = \sum_{K=0}^{\infty} (X_K^T \cdot Q \cdot X_K + U_K^T \cdot R \cdot U_K)$$

avec

$$A = \begin{bmatrix} 0,75 & 0,009 & 0 & 0 & 0 \\ -1,74 & 0,91 & 0 & 0 & 0 \\ -0,3 & -0,0015 & 0,95 & 0 & 0 \\ 0 & 0 & 0 & 0,55 & 0 \\ -0,15 & -0,008 & 0 & 0 & 0,905 \end{bmatrix}$$

$$B = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} \quad H = \begin{bmatrix} 0 & 0 & 24,64 & 0 & 0 \\ 0 & 0 & 0 & 0,835 & 0 \\ 0 & 0 & 0 & 0 & 1,83 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

les matrices Q et V sont calculées comme suit

$$Q = H^T H \quad V = B R^{-1} B^T$$

L'équation de Riccati étant:

$$A^T P (I + V P)^{-1} A - P + Q = 0$$

on a aboutit aux résultats suivants:

- matrice P (solution):

$$P = \begin{bmatrix} 72,30322 & 2,65930 & -200,53321 & -0,37338 & -9,98702 \\ 2,65924 & 1,14270 & -6,12289 & -0,31483 & -2,33572 \\ -200,54036 & -6,12382 & 1212,88950 & 1,29670 & 13,77488 \\ -0,37376 & -0,31487 & -1,29774 & -0,93206 & 0,46769 \\ -9,98926 & -2,33575 & 13,77612 & -0,46759 & 11,74435 \end{bmatrix}$$

Soit donc un système discret de dimension 5 dont l'équation d'état est

$$X_{K+1} = A \cdot X_K + B \cdot U_K$$

et le coût

$$J = \sum_{K=0}^{\infty} (X_K^T \cdot Q \cdot X_K + U_K^T \cdot R \cdot U_K)$$

avec

$$A = \begin{bmatrix} 0,75 & 0,009 & 0 & 0 & 0 \\ -1,74 & 0,91 & 0 & 0 & 0 \\ -0,3 & -0,0015 & 0,95 & 0 & 0 \\ 0 & 0 & 0 & 0,55 & 0 \\ -0,15 & -0,008 & 0 & 0 & 0,905 \end{bmatrix}$$

$$B = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} \quad H = \begin{bmatrix} 0 & 0 & 24,64 & 0 & 0 \\ 0 & 0 & 0 & 0,835 & 0 \\ 0 & 0 & 0 & 0 & 1,83 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

les matrices Q et V sont calculées comme suit

$$Q = H^T H \quad V = B R^{-1} B^T$$

L'équation de Riccati étant:

$$A^T P (I + V P)^{-1} A - P + Q = 0$$

on a aboutit aux résultats suivants:

- matrice P (solution):

$$P = \begin{bmatrix} 72,30322 & 2,65930 & -200,53321 & -0,37338 & -9,98702 \\ 2,65924 & 1,14270 & -6,12289 & -0,31483 & -2,33572 \\ -200,54036 & -6,12382 & 1212,88950 & 1,29670 & 13,77488 \\ -0,37376 & -0,31487 & -1,29774 & -0,93206 & 0,46769 \\ -9,98926 & -2,33575 & 13,77612 & -0,46759 & 11,74435 \end{bmatrix}$$

- nombre d'itérations $N1 = 19$
- la dissymétrie est de l'ordre de 10^{-3} à 10^{-5}
- après correction la matrice P devient

$$P = \begin{bmatrix} 72,30322 & 2,65927 & -200,53685 & -0,37357 & -9,98814 \\ 2,65927 & 1,14270 & -6,12336 & -0,31485 & -2,33573 \\ -200,53685 & -6,12336 & 1212,88950 & 1,29722 & 13,77550 \\ -0,37357 & -0,31485 & 1,29722 & 0,93206 & 0,46764 \\ -9,98814 & -2,33573 & 13,77550 & 0,46764 & 11,74434 \end{bmatrix}$$

- le temps d'exécution était de 17" 6

- nombre d'itérations $N1 = 19$
- la dissymétrie est de l'ordre de 10^{-3} à 10^{-5}
- après correction la matrice P devient

$$P = \begin{bmatrix} 72,30322 & 2,65927 & -200,53685 & -0,37357 & -9,98814 \\ 2,65927 & 1,14270 & -6,12336 & -0,31485 & -2,33573 \\ -200,53685 & -6,12336 & 1212,88950 & 1,29722 & 13,77550 \\ -0,37357 & -0,31485 & 1,29722 & 0,93206 & 0,46764 \\ -9,98814 & -2,33573 & 13,77550 & 0,46764 & 11,74434 \end{bmatrix}$$

- le temps d'exécution était de 17" 6

3- APPLICATION A LA STABILITE DES SYSTEMES

3-1- Rappel sur la deuxième methode de Lyapunov

Nous allons étudier le cas d'un système à trois variables dans la représentation d'état. Ce système étant linéaire et son point d'équilibre se trouve à l'origine.

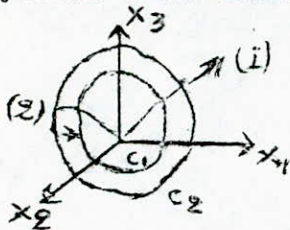
$$\dot{x}_1 = X_1(x_1, x_2, x_3)$$

$$\dot{x}_2 = X_2(x_1, x_2, x_3)$$

$$\dot{x}_3 = X_3(x_1, x_2, x_3)$$

Supposons que l'on puisse mettre en évidence dans l'espace (x_1, x_2, x_3) une famille de surfaces fermées entourant l'origine, telles que par chaque point de l'espace passe une surface unique de la famille. Si dans toute une région entourant l'origine le comportement du système est tel que la vitesse du point (x_1, x_2, x_3) doit toujours être dirigée vers l'intérieur de la surface de la famille passant par ce point; le point représentatif finira par arriver à l'origine et par conséquent le système sera stable.

Inversement on conçoit que si la vitesse est dirigée vers l'extérieur le système sera instable.



(1) instable

(2) stable

$$V(x) = \text{Cste}$$

Les équations des surfaces fermées sont du type

$$V(x_1, x_2, x_3) = \text{Cste}$$

D'où plus précisément pour l'application de la méthode directe de Lyapunov, on considère des fonctions auxiliaires $V(x_1, \dots, x_n)$

3- APPLICATION A LA STABILITE DES SYSTEMES

3-1- Rappel sur la deuxième methode de Lyapunov

Nous allons étudier le cas d'un système à trois variables dans la représentation d'état. Ce système étant linéaire et son point d'équilibre se trouve à l'origine.

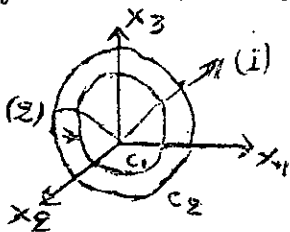
$$\dot{x}_1 = X_1(x_1, x_2, x_3)$$

$$\dot{x}_2 = X_2(x_1, x_2, x_3)$$

$$\dot{x}_3 = X_3(x_1, x_2, x_3)$$

Supposons que l'on puisse mettre en évidence dans l'espace (x_1, x_2, x_3) une famille de surfaces fermées entourant l'origine, telles que par chaque point de l'espace passe une surface unique de la famille. Si dans toute une région entourant l'origine le comportement du système est tel que la vitesse du point (x_1, x_2, x_3) doit toujours être dirigée vers l'intérieur de la surface de la famille passant par ce point; le point représentatif finira par arriver à l'origine et par conséquent le système sera stable.

Inversement on conçoit que si la vitesse est dirigée vers l'extérieur le système sera instable.



(1) instable

(2) stable

$$V(x) = \text{Cste}$$

Les équations des surfaces fermées sont du type

$$V(x_1, x_2, x_3) = \text{Cste}$$

D'où plus précisément pour l'application de la méthode directe de Lyapunov, on considère des fonctions auxiliaires $V(x_1, \dots, x_n)$

définies positives vérifiant les conditions suivantes:

a/V est nulle à l'origine $V(0, \dots, 0) = 0$

b/V est infinie pour des variables infinies $V(x_1, \dots, x_n) \longrightarrow \infty$
 $x_i \longrightarrow \infty$

c/V reste finie lorsque les variables sont finies

Les fonctions V sont dites fonction de Lyapunov.

Nous sommes maintenant en mesure d'énoncer les deux théorèmes fondamentaux de Lyapunov relatifs respectivement à la stabilité asymptotique et la stabilité simple.

théorème 1

S'il est possible de trouver une fonction V de signe défini (dans un domaine G comprenant la position d'équilibre) et dont la dérivée totale par rapport au temps dV/dt soit définie et de signe opposé dans le même domaine, l'équilibre sera asymptotiquement stable dans ce domaine G.

théorème 2

S'il est possible de trouver une fonction V de signe défini (dans un domaine G comprenant la position d'équilibre) et dont la dérivée totale par rapport au temps dV/dt soit semi définie et de signe opposé dans le même domaine l'équilibre est stable dans ce domaine

Rappels

1- Une fonction sera dite semi définie dans tout un domaine G si elle conserve le même signe en tout point de ce domaine mais s'annule aussi en d'autres points que l'origine.

2- Une fonction sera dite indéfinie dans le domaine G si elle prend des signes opposés en différents points de ce domaine.

définies positives vérifiant les conditions suivantes:

a/V est nulle à l'origine $V(0, \dots, 0) = 0$

b/V est infinie pour des variables infinies $V(x_1, \dots, x_n) \longrightarrow \infty$
 $x_i \longrightarrow \infty$

c/V reste finie lorsque les variables sont finies

Les fonctions V sont dites fonction de Lyapunov.

Nous sommes maintenant en mesure d'énoncer les deux théorèmes fondamentaux de Lyapunov relatifs respectivement à la stabilité asymptotique et la stabilité simple.

théorème 1

S'il est possible de trouver une fonction V de signe défini (dans un domaine G comprenant la position d'équilibre) et dont la dérivée totale par rapport au temps dV/dt soit définie et de signe opposé dans le même domaine, l'équilibre sera asymptotiquement stable dans ce domaine G.

théorème 2

S'il est possible de trouver une fonction V de signe défini (dans un domaine G comprenant la position d'équilibre) et dont la dérivée totale par rapport au temps dV/dt soit semi définie et de signe opposé dans le même domaine l'équilibre est stable dans ce domaine

Rappels

1- Une fonction sera dite semi définie dans tout un domaine G si elle conserve le même signe en tout point de ce domaine mais s'annule aussi en d'autres points que l'origine.

2- Une fonction sera dite indéfinie dans le domaine G si elle prend des signes opposés en différents points de ce domaine.

Remarque importante

Lorsque le système étudié est linéaire on peut montrer que les théorèmes 1 et 2 deviennent des conditions nécessaires et suffisantes de stabilité. En effet le signe de dV/dt est lié au signe du produit de la vitesse du point sur la trajectoire et de ∇V

3-2- Etablissement de l'équation de Lyapunov.

Soit le système linéaire autonome représenté à l'aide de la représentation d'état par l'équation:

$$\dot{X} = A X \text{ et } V(X) \text{ une forme quadratique du type } V(X) = X^T P X$$

où P est une matrice symétrique, définie positive. Donc $V(X)$ est définie positive au sens de Lyapunov.

pour que le système soit stable il faut et il suffit que:

$$(dV/dt) < 0 \implies \dot{X}^T P X + X^T P \dot{X} = (AX)^T P X + X^T P A X.$$

$$= X^T A^T P X + X^T P A X = X^T (A^T P + P A) X < 0$$

Ce qui revient à résoudre l'équation :

$$A^T P + P A = -Q \text{ où } Q \text{ est définie positive}$$

Conclusion

Ainsi des mêmes algorithmes résolvant les équations de Riccati nous permettent de résoudre celles de Lyapunov relatives à la stabilité. En effet l'équation de Riccati étant:

$$A^T P + P A + P D P + Q = 0$$

D'où en mettant $D=0$, on retrouve l'équation de Lyapunov.

3-3- Applications numériques

On a vu ultérieurement que, pratiquement le même algorithme ayant servi pour la résolution de l'équation de Riccati, peut servir à résoudre l'équation de Lyapunov (en prenant $V=0$).

Remarque importante

Lorsque le système étudié est linéaire on peut montrer que les théorèmes 1 et 2 deviennent des conditions nécessaires et suffisantes de stabilité. En effet le signe de dV/dt est lié au signe du produit de la vitesse du point sur la trajectoire et de ∇V

3-2- Etablissement de l'équation de Lyapunov.

Soit le système linéaire autonome représenté à l'aide de la représentation d'état par l'équation:

$$\dot{X} = A X \text{ et } V(X) \text{ une forme quadratique du type } V(X) = X^T P X$$

où P est une matrice symétrique, définie positive. Donc $V(X)$ est définie positive au sens de Lyapunov.

pour que le système soit stable il faut et il suffit que:

$$(dV/dt) < 0 \implies \dot{X}^T P X + X^T P \dot{X} = (AX)^T P X + X^T P A X.$$

$$= X^T A^T P X + X^T P A X = X^T (A^T P + P A) X < 0$$

Ce qui revient à résoudre l'équation :

$$A^T P + P A = -Q \text{ où } Q \text{ est définie positive}$$

Conclusion

Ainsi des mêmes algorithmes résolvant les équations de Riccati nous permettent de résoudre celles de Lyapunov relatives à la stabilité. En effet l'équation de Riccati étant:

$$A^T P + P A + P D P + Q = 0$$

D'où en mettant $D=0$, on retrouve l'équation de Lyapunov.

3-3- Applications numériques

On a vu ultérieurement que, pratiquement le même algorithme ayant servi pour la résolution de l'équation de Riccati, peut servir à résoudre l'équation de Lyapunov (en prenant $V=0$).

Il nous a paru utile d'étudier dans ce cas, la variation du nombre d'itérations total en fonction du spectre.

L'équation de Lyapunov dans le cas continu étant

$$A^T P + P A + Q = 0$$

on va considérer dans ce qui suit trois cas d'applications dans lesquels la matrice A reste la même.

Quand à la matrice Q , elle est mise sous forme diagonale

$$Q = M.D.M^{-1}$$

où

M est la matrice des vecteurs propres, fixe pour les trois exemples.

D est la matrice diagonale de Q . Elle varie d'un exemple à l'autre.

$$A = \begin{bmatrix} 4 & 1 & -1 \\ 1 & 3 & 0 \\ -1 & 0 & 5 \end{bmatrix}$$

$$M = \begin{bmatrix} 1 & -1 & 1 \\ 2 & 0 & 1 \\ 1 & 2 & -1 \end{bmatrix}$$

Premier cas

$$Q = M.D_1.M^{-1}$$

avec

$$D_1 = \text{diag}(-1,1 ; -1,01 ; -0,95)$$

résultats

- nombre d'itérations total $N1 = 13$
- la matrice P après correction est

Il nous a paru utile d'étudier dans ce cas, la variation du nombre d'itérations total en fonction du spectre.

L'équation de Lyapunov dans le cas continu étant

$$A^T P + P A + Q = 0$$

on va considérer dans ce qui suit trois cas d'applications dans lesquels la matrice A reste la même.

Quand à la matrice Q , elle est mise sous forme diagonale

$$Q = M.D.M^{-1}$$

où

M est la matrice des vecteurs propres, fixe pour les trois exemples.

D est la matrice diagonale de Q . Elle varie d'un exemple à l'autre.

$$A = \begin{bmatrix} 4 & 1 & -1 \\ 1 & 3 & 0 \\ -1 & 0 & 5 \end{bmatrix}$$

$$M = \begin{bmatrix} 1 & -1 & 1 \\ 2 & 0 & 1 \\ 1 & 2 & -1 \end{bmatrix}$$

Premier cas

$$Q = M.D_1.M^{-1}$$

avec

$$D_1 = \text{diag}(-1,1 ; -1,01 ; -0,95)$$

résultats

- nombre d'itérations total $N_I = 13$
- la matrice P après correction est

$$P = \begin{bmatrix} -9,0765685 & 2,9470474 & -2,0365685 \\ 2,9470474 & -2,6469392 & 3,8338310 \\ -2,0365685 & 3,8338310 & -4,4555342 \end{bmatrix}$$

- le temps d'exécution était de 4" 4

Deuxième cas

$$Q = M \cdot D_2 \cdot M^{-1} \quad \text{avec} \quad D_2 = \text{diag}(-1 ; -2 ; -3)$$

résultats

- nombre d'itérations total $N_1 = 16$

- la matrice P après correction est

$$P = \begin{bmatrix} -0,8173513 & 1,4871061 & 0,6029091 \\ 1,4871061 & -0,7418242 & 1,6044044 \\ 0,6029091 & 1,6044044 & -1,5638620 \end{bmatrix}$$

- le temps d'exécution étant 5" 6

Troisième cas

$$Q = M \cdot D_3 \cdot M^{-1} \quad \text{avec} \quad D_3 = \text{diag}(-10 ; -1 ; -1000)$$

résultats

- nombre d'itérations total $N_1 = 33$

- la matrice P après correction est

$$\begin{bmatrix} -0,3189621 & -0,0882767 & -0,2126414 \\ -0,0882767 & -0,0295951 & -0,0537453 \\ -0,2126414 & -0,0537453 & 0,0140174 \end{bmatrix}$$

- le temps d'exécution est de 8" 3

Conclusion

On constate que, d'après ces résultats, le nombre d'itérations varie en fonction du spectre. Ceci était prévisible car à chaque itération, il fallait équilibrer la matrice H dont on calcula la matrice signe jusqu'à avoir toutes les valeurs propres égales ou très proches les unes des autres. Cet équilibrage est donc essentiellement fonction du spectre de H. Ce qui augmente N_1 lorsque les valeurs propres sont fortement distinctes.

$$P = \begin{bmatrix} -9,0755685 & 2,9470474 & -2,0365685 \\ 2,9470474 & -2,6469392 & 3,8338310 \\ -2,0365685 & 3,8338310 & -4,4555342 \end{bmatrix}$$

- le temps d'exécution était de 4" 4

Deuxième cas

$$Q = M \cdot D_2 \cdot M^{-1} \quad \text{avec} \quad D_2 = \text{diag}(-1 ; -2 ; -3)$$

résultats

- nombre d'itérations total $N_1 = 16$

- la matrice P après correction est

$$P = \begin{bmatrix} -0,8173513 & 1,4871061 & 0,6029091 \\ 1,4871061 & -0,7418242 & 1,6044044 \\ 0,6029091 & 1,6044044 & -1,5638620 \end{bmatrix}$$

- le temps d'exécution étant 5" 6

Troisième cas

$$Q = M \cdot D_3 \cdot M^{-1} \quad \text{avec} \quad D_3 = \text{diag}(-10 ; -1 ; -1000)$$

résultats

- nombre d'itérations total $N_1 = 33$

- la matrice P après correction est

$$\begin{bmatrix} -0,3189621 & -0,0882767 & -0,2126414 \\ -0,0882767 & -0,0295951 & -0,0537453 \\ -0,2126414 & -0,0537453 & 0,0140174 \end{bmatrix}$$

- le temps d'exécution est de 8" 3

Conclusion

On constate que, d'après ces résultats, le nombre d'itérations varie en fonction du spectre. Ceci était prévisible car à chaque itération, il fallait équilibrer la matrice H dont on calcule la matrice sirne jusqu'à avoir toutes les valeurs propres égales ou très proches les unes des autres. Cet équilibrage est donc essentiellement fonction du spectre de H. Ce qui augmente N_1 lorsque les valeurs propres sont fortement distinctes.

-oOo- C O N C L U S I O N -oOo-

On a développé dans cette étude un concept (la fonction signe de matrice), et ses applications à la commande optimale et à la stabilité des systèmes. D'où l'on a été amené à établir des algorithmes nécessaires à la résolution des équations de RICCATI et de LYAPUNOV dans les cas continu et discret. Ces algorithmes tiennent compte des contraintes sur le temps calcul, la précision, et les performances techniques du calculateur MITRA 125 se trouvant au centre de calcul du C E N.

On a montré sur le plan de la convergence comment évoluaient exactement les itérés engendrés par l'algorithme de Newton, introduit pour le calcul de la fonction signe. Les résultats obtenus ont débouché sur deux nouveaux algorithmes, l'un fini pour les matrices à spectre réel, l'autre relatif au cas général (spectre complexe), constituant ce qu'on a appelé un algorithme de Newton accéléré.

Du point de vue complexité de mise en oeuvre, stabilité numérique, coût calcul par itération et vitesse de convergence, la méthode de la fonction signe que l'on avait développée (au niveau de ses applications) nous semble être très performante. De plus elle est plus générale.

En effet le même code peut servir pour résoudre les équations de Riccati relatives aux systèmes continus et discrets. Seule la détermination de l'hamiltonien et le calcul de P changent.

-oOo- C O N C L U S I O N -oOo-

On a développé dans cette étude un concept (la fonction signe de matrice), et ses applications à la commande optimale et à la stabilité des systèmes. D'où l'on a été amené à établir des algorithmes nécessaires à la résolution des équations de RICCATI et de LYAPUNOV dans les cas continu et discret. Ces algorithmes tiennent compte des contraintes sur le temps calcul, la précision, et les performances techniques du calculateur MITRA 125 se trouvant au centre de calcul du C E N.

On a montré sur le plan de la convergence comment évoluaient exactement les itérés engendrés par l'algorithme de Newton, introduit pour le calcul de la fonction signe. Les résultats obtenus ont débouché sur deux nouveaux algorithmes, l'un fini pour les matrices à spectre réel, l'autre relatif au cas général (spectre complexe), constituant ce qu'on a appelé un algorithme de Newton accéléré.

Du point de vue complexité de mise en oeuvre, stabilité numérique, coût calcul par itération et vitesse de convergence, la méthode de la fonction signe que l'on avait développée (au niveau de ses applications) nous semble être très performante. De plus elle est plus générale.

En effet le même code peut servir pour résoudre les équations de Riccati relatives aux systèmes continus et discrets. Seule la détermination de l'hamiltonien et le calcul de P changent.

Enfin si l'on pose $V=0$, le même algorithme permet de résoudre les équations de Lyapunov.

Notre expérience dans le domaine numérique montre que les résultats obtenus (sauf mauvais conditionnement) comparés à ceux de Barraud et Beavers, sont acceptables.

Du point de vue symétrie, les résultats obtenus (avant correction) ont une dissymétrie de l'ordre de 10^{-5} à 10^{-6} .

Quant à la convergence, elle est du même ordre que la méthode SR2 seulement notre algorithme ne demande aucun calcul préliminaire, d'où une facilité de mise en œuvre, une diminution du coût calcul et de l'occupation mémoire.

Comme continuité à notre travail nous pensons que les remarques suivantes sont utiles:

- on pourrait développer d'autres algorithmes de la forme

$$A_{K+1} = \alpha A_K + \beta A_K^{-1}$$

où il s'agirait de calculer α et β de telle sorte que l'algorithme converge plus rapidement que celui dit de Newton accéléré (où $\alpha\beta = \frac{1}{4}$)

- suivant la nature des problèmes de contrôle auxquels on s'intéresse, on pourrait penser que si des calculs relativement longs du type conditionnement de la matrice H étaient faits hors ligne, on pourrait améliorer la rapidité de convergence des algorithmes pour une éventuelle application en ligne.

Enfin si l'on pose $V=0$, le même algorithme permet de résoudre les équations de Lyapunov.

Notre expérience dans le domaine numérique montre que les résultats obtenus (sauf mauvais conditionnement) comparés à ceux de Barraud et Beavers, sont acceptables.

Du point de vue symétrie, les résultats obtenus (avant correction) ont une dissymétrie de l'ordre de 10^{-5} à 10^{-6} .

Quant à la convergence, elle est du même ordre que la méthode SR2 seulement notre algorithme ne demande aucun calcul préliminaire, d'où une facilité de mise en œuvre, une diminution du coût calcul et de l'occupation mémoire.

Comme continuité à notre travail nous pensons que les remarques suivantes sont utiles:

- on pourrait développer d'autres algorithmes de la forme

$$A_{K+1} = \alpha A_K + \beta A_K^{-1}$$

où il s'agirait de calculer α et β de telle sorte que l'algorithme converge plus rapidement que celui dit de Newton accéléré (où $\alpha\beta = \frac{1}{4}$)

- suivant la nature des problèmes de contrôle auxquels on s'intéresse, on pourrait penser que si des calculs relativement longs du type conditionnement de la matrice H étaient faits hors ligne, on pourrait améliorer la rapidité de convergence des algorithmes pour une éventuelle application en ligne.

ANNEXE I

RAPPELS MATHÉMATIQUES

1-Valeurs propres d'une matrice

On appelle valeur propre d'une matrice toute racine λ de son équation caractéristique

$$\det(\lambda I - A) = 0$$

Si $A(n.n)$, on aura n valeurs propres réelles ou complexes, distinctes ou multiples.

2-Vecteurs propres d'une matrice

On appelle vecteur propre m associé à une valeur propre, tout vecteur solution de l'équation

$$[A - \lambda I] m = 0$$

donc à une valeur propre simple correspond un vecteur propre unique (à un coefficient de proportionnalité près). A une valeur propre multiple d'ordre μ , peut correspondre 1 à μ vecteurs propres selon le nombre de solutions indépendantes du système d'équations.

3-Forme diagonale d'une matrice

Si toutes les racines de l'équation caractéristique sont distinctes, il est possible de mettre la matrice A sous forme diagonale

$$\tilde{A} = \text{diag} [\lambda_1, \lambda_2, \dots, \lambda_n] \quad \lambda_i \neq \lambda_j \text{ si } i \neq j$$

La matrice de transformation M est la matrice des vecteurs propres

$$M = \begin{bmatrix} m_1 & m_2 & \dots & m_n \end{bmatrix} \quad \text{avec} \quad A m_1 = \lambda_1 m_1$$

$$\tilde{A} = M^{-1} A M \iff M A = \tilde{A} M \iff M \tilde{A} M^{-1} = A$$

ANNEXE I

RAPPELS MATHÉMATIQUES

1-Valeurs propres d'une matrice

On appelle valeur propre d'une matrice toute racine λ de son équation caractéristique

$$\det(\lambda I - A) = 0$$

Si $A(n,n)$, on aura n valeurs propres réelles ou complexes, distinctes ou multiples.

2-Vecteurs propres d'une matrice

On appelle vecteur propre r associé à une valeur propre, tout vecteur solution de l'équation

$$[A - \lambda I] m = 0$$

donc à une valeur propre simple correspond un vecteur propre unique (à un coefficient de proportionalité près). A une valeur propre multiple d'ordre μ , peut correspondre 1 à μ vecteurs propres selon le nombre de solutions indépendantes du système d'équations.

3-Forme diagonale d'une matrice

Si toutes les racines de l'équation caractéristique sont distinctes, il est possible de mettre la matrice A sous forme diagonale

$$\tilde{A} = \text{diag} [\lambda_1, \lambda_2, \dots, \lambda_n] \quad \lambda_i \neq \lambda_j \text{ si } i \neq j$$

La matrice de transformation M est la matrice des vecteurs propres

$$M = [m_1 m_2 \dots m_n] \quad \text{avec} \quad A m_1 = \lambda_1 m_1$$

$$\tilde{A} = M^{-1} A M \iff M A = A M \iff M \tilde{A} M^{-1} = A$$

donc, on peut dire que si on a une matrice A telle que $\lambda_i \neq \lambda_j$ si $i \neq j$

$$\implies \exists \tilde{A} = \text{diag} [\lambda_1 \dots \dots \lambda_n] \text{ telle que } A = \tilde{M} \tilde{A} \tilde{M}^{-1}$$

4- Forme quasi diagonale de Jordan

Appelons matrice élémentaire, une matrice $D_k(\lambda)$ ($k \times k$), dont la diagonale principale est formée des k même nombres complexes λ , et dont les éléments situés immédiatement à droite des éléments diagonaux sont égaux à 1 et tous les autres éléments nuls:

$$D_k(\lambda) = \begin{bmatrix} \lambda & 1 & & 0 \\ & \lambda & \ddots & \\ & & \ddots & 1 \\ 0 & & & \lambda \end{bmatrix}$$

Toute matrice $A(n \times n)$ est semblable à une matrice de la forme

$$\begin{bmatrix} D_{k_1}(\lambda_1) & & 0 \\ & D_{k_2}(\lambda_2) & \\ & & \ddots \\ 0 & & & 0 \end{bmatrix}$$

dans laquelle les éléments situés au voisinage de la diagonale principale se distribuent dans une suite de matrices élémentaires de rang k_1, k_2, \dots, k_n

avec $k_1 + k_2 + \dots + k_n = n$; tous les autres éléments sont nuls.

Les éléments diagonaux $\lambda_1, \lambda_2, \dots$ correspondant sont manifestement les valeurs propres de la matrice A .

Si $k_i > 1$, λ_i est valeur propre multiple d'ordre au moins égal à k_i .

Cet ordre vaut exactement k_i si λ_i sont distinctes. Mais il peut arriver

par exemple que $\lambda_2 = \lambda_1$, λ_1 est alors multiple d'ordre $k_1 + k_2$ même si $k_1 = 1$.

Si tous les k_i sont égaux à 1, A est semblable à une matrice diagonale;

donc, on peut dire que si on a une matrice A telle que $\lambda_i \neq \lambda_j$ si $i \neq j$

$$\implies \exists \tilde{A} = \text{diag} [\lambda_1 \dots \dots \lambda_n] \text{ telle que } A = M \tilde{A} M^{-1}$$

4- Forme quasi diagonale de Jordan

Appelons matrice élémentaire, une matrice $D_k(\lambda)$ ($k \times k$), dont la diagonale principale est formée des k même nombres complexes λ , et dont les éléments situés immédiatement à droite des éléments diagonaux sont égaux à 1 et tous les autres éléments nuls:

$$D_k(\lambda) = \begin{bmatrix} \lambda & 1 & & 0 \\ & \lambda & \ddots & \\ & & \ddots & 1 \\ 0 & & & \lambda \end{bmatrix}$$

Toute matrice A ($n \times n$) est semblable à une matrice de la forme

$$\begin{bmatrix} D_{k_1}(\lambda_1) & & 0 \\ & D_{k_2}(\lambda_2) & \\ & & \ddots \\ 0 & & & 0 \end{bmatrix}$$

dans laquelle les éléments situés au voisinage de la diagonale principale se distribuent dans une suite de matrices élémentaires de rang k_1, k_2, \dots, k_n

avec $k_1 + k_2 + \dots + k_n = n$; tous les autres éléments sont nuls.

Les éléments diagonaux $\lambda_1, \lambda_2, \dots$ correspondant sont manifestement les valeurs propres de la matrice A .

Si $k_i > 1$, λ_i est valeur propre multiple d'ordre au moins égal à k_i .

Cet ordre vaut exactement k_i si λ_i sont distinctes. Mais il peut arriver

par exemple que $\lambda_2 = \lambda_1$, λ_1 est alors multiple d'ordre $k_1 + k_2$ même si $k_1 = 1$.

Si tous les k_i sont égaux à 1, A est semblable à une matrice diagonale;

chacune des matrices élémentaires est alors de rang 1.

Donc en conclusion on peut dire que toute matrice A peut

se mettre sous la forme de Jordan $A = M J M^{-1}$

J étant la matrice formée de blocs de Jordan, M étant la matrice formée des vecteurs propres

5-Théorème concernant l'algorithme de Newton

Soit $f(x)$ une fonction de la variable réelle x , de classe C^2

tel que $f(\hat{x}) = 0$ et $f'(\hat{x}) \neq 0$ (\hat{x} racine simple). Alors il existe un intervalle ouvert $N(\hat{x})$ contenant \hat{x} tel que si $x_1 \in N(\hat{x})$ l'algorithme de Newton

$$x_{k+1} = x_k + f(x_k)/f'(x_k)$$

engendre une suite x_k qui converge vers \hat{x} , de plus

$$\lim_{k \rightarrow \infty} \frac{x_{k+1} - \hat{x}}{(x_k - \hat{x})^2} = \frac{f''(\hat{x})}{2f'(\hat{x})}$$

résultat qui traduit une convergence finale d'ordre 2.

Remarque

Ce théorème peut s'appliquer aux fonctions de variables complexes analytiques.

Démonstration de la formule (9)

considérons la relation suivante

$$\frac{x_k}{2} < x_{k+1} < \frac{1}{2}(1+x_k)$$

soit

$$\frac{x_{k-1}}{2} < x_k \quad \Rightarrow \quad \frac{x_{k-1}}{2^2} < x_{k+1}$$

chacune des matrices élémentaires est alors de rang 1.

Donc en conclusion on peut dire que toute matrice A peut se mettre sous la forme de Jordan $A = M J M^{-1}$

J étant la matrice formée de blocs de Jordan, M étant la matrice formée des vecteurs propres

5-Théorème concernant l'algorithme de Newton

Soit $f(x)$ une fonction de la variable réelle x , de classe C^2

tel que $f(\hat{x}) = 0$ et $f'(\hat{x}) \neq 0$ (\hat{x} racine simple). Alors il existe un intervalle ouvert $N(\hat{x})$ contenant \hat{x} tel que si $x_1 \in N(\hat{x})$ l'algorithme de Newton

$$x_{k+1} = x_k - f(x_k)/f'(x_k)$$

engendre une suite x_k qui converge vers \hat{x} , de plus

$$\lim_{k \rightarrow \infty} \frac{x_{k+1} - \hat{x}}{(x_k - \hat{x})^2} = \frac{f''(\hat{x})}{2f'(\hat{x})}$$

résultat qui traduit une convergence finale d'ordre 2.

Remarque

Ce théorème peut s'appliquer aux fonctions de variables complexes analytiques.

Démonstration de la formule (9)

considérons la relation suivante

$$\frac{x_k}{2} < x_{k+1} < \frac{1}{2}(1+x_k)$$

soit

$$\frac{x_{k-1}}{2} < x_k \implies \frac{x_{k-1}}{2^2} < x_{k+1}$$

$$\frac{x_{k-2}}{2} < x_{k-1} \implies \frac{x_{k-2}}{2^2} < x_{k+1}$$

⋮

$$\frac{x_0}{2} < x_1 \implies \frac{x_0}{2^n} < x_n \quad \text{en prenant } k+1=n$$

Prenons maintenant l'inégalité de droite

$$x_k < \frac{1}{2}(1+x_{k-1}) \implies x_{k+1} < \frac{1}{2} \left[1 + \frac{1}{2}(1+x_{k-1}) \right]$$

$$x_{k-1} < \frac{1}{2}(1+x_{k-2}) \implies x_{k+1} < \frac{1}{2} \left[1 + \frac{1}{2} \left(1 + \frac{1}{2}(1+x_{k-2}) \right) \right]$$

⋮

$$x_1 < \frac{1}{2}(1+x_0) \implies x_{k+1} < \frac{1}{2} \left[1 + \frac{1}{2} \left(1 + \frac{1}{2} (\dots (1+x_0)) \right) \right]$$

d'où

$$x_{k+1} < \frac{1}{2} + \frac{1}{2^2} + \dots + \frac{1}{2^n} + \frac{1}{2^n} x_0$$

le deuxième terme de l'inégalité est une suite géométrique de raison 1/2

dont la somme est égale à

$$S = (1 - 1/2^n)$$

d'où l'on peut écrire

$$\frac{x_0}{2^n} < x_n < (1 - 1/2^n) + x_0/2^n$$

Démonstration de la formule()

d'après (11) on a $\rho_{k+1}^2 = \frac{1}{4}\rho_k^2 + \frac{1}{4}\frac{1}{\rho_k^2} - 1/2 + \cos^2 e_k$

comme $\rho_k > 1 \implies 1/4 \cdot \rho_k^2 \ll 1 \implies \rho_{k+1}^2 > \frac{1}{4}\rho_k^2 - \frac{1}{2}$

maintenant si $1/\rho_k^2 \approx 1$ et $\cos^2 e_k \approx 1$ on aura $\rho_{k+1}^2 < \frac{1}{4}\rho_k^2 + 3/4$

$$\frac{x_{k-2}}{2} < x_{k-1} \implies \frac{x_{k-2}}{2^2} < x_{k+1}$$

⋮

$$\frac{x_0}{2} < x_1 \implies \frac{x_0}{2^n} < x_n \quad \text{en prenant } k+1=n$$

Prenons maintenant l'inégalité de droite

$$x_k < \frac{1}{2}(1+x_{k-1}) \implies x_{k+1} < \frac{1}{2} \left[1 + \frac{1}{2}(1+x_{k-1}) \right]$$

$$x_{k-1} < \frac{1}{2}(1+x_{k-2}) \implies x_{k+1} < \frac{1}{2} \left[1 + \frac{1}{2} \left(1 + \frac{1}{2}(1+x_{k-2}) \right) \right]$$

⋮

$$x_1 < \frac{1}{2}(1+x_0) \implies x_{k+1} < \frac{1}{2} \left[1 + \frac{1}{2} \left(1 + \frac{1}{2} (\dots (1+x_0)) \right) \right]$$

d'où

$$x_{k+1} < \frac{1}{2} + \frac{1}{2^2} + \dots + \frac{1}{2^n} + \frac{1}{2^n} x_0$$

le deuxième terme de l'inégalité est une suite géométrique de raison 1/2

dont la somme est égale à

$$S = (1 - 1/2^n)$$

d'où l'on peut écrire

$$\frac{x_0}{2^n} < x_n < (1 - 1/2^n) + x_0/2^n$$

Démonstration de la formule()

d'après (1f) on a $\rho_{k+1}^2 = \frac{1}{4}\rho_k^2 + \frac{1}{4}\frac{1}{\rho_k^2} - 1/2 + \cos^2 e_k$

comme $\rho_k > 1 \implies 1/4 \cdot \rho_k^2 \ll 1 \implies \rho_{k+1}^2 > \frac{1}{4}\rho_k^2 - \frac{1}{2}$

maintenant si $1/\rho_k^2 \approx 1$ et $\cos^2 e_k \approx 1$ on aura $\rho_{k+1}^2 < \frac{1}{4}\rho_k^2 + 3/4$

DEFINITIONS ET PROPRIETES IMPORTANTES

COMMANDABILITE

la paire (A,B) est dite commande si et seulement si tout état $x_{t_0} \neq 0$ peut être ramené à zéro en un temps fini.

$$(A,B) \text{ commandable} \iff \text{rang}(B, AB, \dots, A^{n-1}B) = n = \dim(A)$$

RECONSTRUCTIBILITE

Etant donné une sortie $y = cx$, et $\dot{x} = Ax$, la paire (C,A) est reconstructible si et seulement si tout état $x_{t_1} \neq 0$ peut être déterminé de façon unique à partir de la sortie passée $y(t), t \in [t_0, t_1]$ sur un horizon $(t_1 - t_0)$ fini:

$$(C,A) \text{ reconstructible} \implies \text{rang} \begin{bmatrix} C^T, A^T C^T, \dots, (A^T)^{n-1} C^T \end{bmatrix} = n$$

STABILISABILITE

La paire (A,B) est stabilisable si et seulement si:

$$\exists L: \text{réel } \lambda_i [A - BL] < 0$$

DETECTIBILITE

La paire (C,A) est détectable si et seulement si:

$$\exists M: \text{réel } \lambda_i [A + MC] < 0$$

PROPRIETES

-étant donnée une paire (A,B) non commandable, la paire(A,B) est stabilisable si et seulement si les états non commandables sont asymptotiquement stables.

-étant donnée une paire (C,A) non reconstructible, la paire(C,A) est détectable si et seulement si les états non reconstructibles sont asymptotiquement stables.

DEFINITIONS ET PROPRIETES IMPORTANTES

COMMANDABILITE

la paire (A,B) est dite commande si et seulement si tout état $x_{t_0} \neq 0$ peut être ramené à zéro en un temps fini.

$$(A,B) \text{ commandable} \iff \text{rang}(B, AB, \dots, A^{n-1}B) = n = \dim(A)$$

RECONSTRUCTIBILITE

Etant donné une sortie $y = cx$, et $\dot{x} = Ax$, la paire (C,A) est reconstructible si et seulement si tout état $x_{t_1} \neq 0$ peut être déterminé de façon unique à partir de la sortie passée $y(t), t \in [t_0, t_1]$ sur un horizon $(t_1 - t_0)$ fini:

$$(C,A) \text{ reconstructible} \implies \text{rang} \begin{bmatrix} C^T, A^T C^T, \dots, (A^T)^{n-1} C^T \end{bmatrix} = n$$

STABILISABILITE

La paire (A,B) est stabilisable si et seulement si:

$$\exists L: \text{réel } \lambda_i [A - BL] < 0$$

DETECTIBILITE

La paire (C,A) est détectable si et seulement si:

$$\exists M: \text{réel } \lambda_i [A + MC] < 0$$

PROPRIETES

-étant donnée une paire (A,B) non commandable, la paire (A,B) est stabilisable si et seulement si les états non commandables sont asymptotiquement stables.

-étant donnée une paire (C,A) non reconstructible, la paire (C,A) est détectable si et seulement si les états non reconstructibles sont asymptotiquement stables.

ANNEXE III

Démonstration du théorème concernant l'équation (continue) de Riccati.

Etant donné le problème de commande (40)-(40a) et les hypothèses (42), on sait que la matrice du système bouclé \tilde{A} (44) est asymptotiquement stable, la matrice P étant solution de l'équation de Riccati. Ceci étant, on peut alors vérifier que

$$H = \begin{bmatrix} A & -BR^{-1}B^T \\ -Q & -A^T \end{bmatrix} = \begin{bmatrix} I & -V \\ P & I-PV \end{bmatrix} \begin{bmatrix} A & 0 \\ 0 & -A \end{bmatrix} \begin{bmatrix} I-VP & V \\ -P & I \end{bmatrix}$$

Il est clair que cette relation peut encore s'écrire

$$H = U \begin{bmatrix} A & 0 \\ 0 & -A \end{bmatrix} U^{-1}$$

ce qui démontre les propriétés 1 et 2 du théorème

Définissons maintenant la factorisation de Jordan de A par

$$A = W_{12}(-J)W_{12}^{-1}$$

et reportons ce résultat dans la première relation ci-dessus, il vient

$$H = \begin{bmatrix} W_{12} & -VW_{12} \\ FW_{12} & (I-FV)W_{12} \end{bmatrix} \begin{bmatrix} -J & 0 \\ 0 & J \end{bmatrix} \begin{bmatrix} W_{12}^{-1}(I-VP) & W_{12}^{-1}V \\ -W_{12}^{-1} & W_{12}^{-1} \end{bmatrix}$$

$$= \begin{bmatrix} -VW_{12} & W_{12} \\ (I-FV)W_{12} & FW_{12} \end{bmatrix} \begin{bmatrix} J & 0 \\ 0 & -J \end{bmatrix} \begin{bmatrix} -W_{12}^{-1}P & W_{12}^{-1} \\ W_{12}^{-1}(I-VP) & W_{12}^{-1}V \end{bmatrix}$$

En identifiant cette relation avec la troisième relation du théorème, on obtient alors la relation (4) du théorème.

ANNEXE III

Démonstration du théorème concernant l'équation (continue) de Riccati.

Etant donné le problème de commande (40)-(40a) et les hypothèses (42), on sait que la matrice du système bouclé \tilde{A} (44) est asymptotiquement stable, la matrice P étant solution de l'équation de Riccati. Ceci étant, on peut alors vérifier que

$$H = \begin{bmatrix} A & -BR^{-1}B^T \\ -Q & -A^T \end{bmatrix} = \begin{bmatrix} I & -V \\ P & I-PV \end{bmatrix} \begin{bmatrix} A & 0 \\ 0 & -A \end{bmatrix} \begin{bmatrix} I-VP & V \\ -P & I \end{bmatrix}$$

Il est clair que cette relation peut encore s'écrire

$$H = U \begin{bmatrix} A & 0 \\ 0 & -A \end{bmatrix} U^{-1}$$

ce qui démontre les propriétés 1 et 2 du théorème

Définissons maintenant la factorisation de Jordan de A par

$$A = W_{12}(-J)W_{12}^{-1}$$

et reportons ce résultat dans la première relation ci-dessus, il vient

$$H = \begin{bmatrix} W_{12} & -VW_{12} \\ FW_{12} & (I-FV)W_{12} \end{bmatrix} \begin{bmatrix} -J & Q \\ 0 & J \end{bmatrix} \begin{bmatrix} W_{12}^{-1}(I-VP) & W_{12}^{-1}V \\ -W_{12}^{-1} & W_{12}^{-1} \end{bmatrix}$$

$$= \begin{bmatrix} -VW_{12} & W_{12} \\ (1-FV)W_{12} & FW_{12} \end{bmatrix} \begin{bmatrix} J & 0 \\ 0 & -J \end{bmatrix} \begin{bmatrix} -W_{12}^{-1}P & W_{12}^{-1} \\ W_{12}^{-1}(I-VP) & W_{12}^{-1}V \end{bmatrix}$$

En identifiant cette relation avec la troisième relation du théorème, on obtient alors la relation (4) du théorème.

Quelques formes équivalentes à: $P = A^T P A - A^T P B (B^T P B + R)^{-1} B^T P A + Q$ (1)

Soit le gain G introduit implicitement:

$$G = (B^T P B + R)^{-1} B^T P A$$

de sorte que le système bouclé \tilde{A} s'écrit

$$\tilde{A} = A - B G$$

Dans ces conditions, on peut vérifier que (1) est équivalente à

$$P = \tilde{A}^T P \tilde{A} + G^T R G + Q$$

De plus si $R > 0$, on a

$$\begin{aligned} P &= A^T P A - A^T P B (B^T P B + R)^{-1} B^T P A + Q \\ &= A^T P \left[I - B (B^T P B + R)^{-1} B^T P \right] A + Q \\ &= A^T P (I + B R^{-1} B^T P)^{-1} A + Q \end{aligned}$$

ou encore

$$\begin{aligned} P &= A^T \left[I - P B (B^T P B + R)^{-1} B^T \right] P A + Q \\ &= A^T (I + P B R^{-1} B^T)^{-1} P A + Q \end{aligned}$$

l'on pose $P = K^T K$ (forme factorisée), on aura

$$\begin{aligned} K^T K &= A^T K^T K A - A^T K^T P B (B^T K^T K B + R)^{-1} B^T K^T K A + Q \\ K^T K &= A^T K^T \left[I - K B (B^T K^T K B + R)^{-1} B^T K^T \right] K A + Q \\ K^T K &= A^T K^T (I + K B R^{-1} B^T K^T)^{-1} K A + Q \end{aligned}$$

Remarque

Les différentes transformations ont été créées en utilisant l'identité matricielle suivante:

$$(A + B C D)^{-1} = A^{-1} - A^{-1} B (C^{-1} + D A^{-1} B)^{-1} D A^{-1}$$

Quelques formes équivalentes à: $P = A^T P A - A^T P B (B^T P B + R)^{-1} B^T P A + Q$ (1)

Soit le gain G introduit implicitement:

$$G = (B^T P B + R)^{-1} B^T P A$$

de sorte que le système bouclé \tilde{A} s'écrit

$$\tilde{A} = A - B G$$

Dans ces conditions, on peut vérifier que (1) est équivalente à

$$P = \tilde{A}^T P \tilde{A} + G^T R G + Q$$

De plus si $R > 0$, on a

$$\begin{aligned} P &= A^T P A - A^T P B (B^T P B + R)^{-1} B^T P A + Q \\ &= A^T P \left[I - B (B^T P B + R)^{-1} B^T P \right] A + Q \\ &= A^T P (I + B R^{-1} B^T P)^{-1} A + Q \end{aligned}$$

ou encore

$$\begin{aligned} P &= A^T \left[I - P B (B^T P B + R)^{-1} B^T \right] P A + Q \\ &= A^T (I + P B R^{-1} B^T)^{-1} P A + Q \end{aligned}$$

l'on pose $P = K^T K$ (forme factorisée), on aura

$$\begin{aligned} K^T K &= A^T K^T K A - A^T K^T P B (B^T K^T K B + R)^{-1} B^T K^T K A + Q \\ K^T K &= A^T K^T \left[I - K B (B^T K^T K B + R)^{-1} B^T K^T \right] K A + Q \\ K^T K &= A^T K^T (I + P B R^{-1} B^T K^T)^{-1} K A + Q \end{aligned}$$

remarque

Les différentes transformations ont été créées en utilisant l'identité matricielle suivante:

$$(A + B C D)^{-1} = A^{-1} - A^{-1} B (C^{-1} + D A^{-1} B)^{-1} D A^{-1}$$

Matrice symplectique

Soit une matrice réelle $2n \times 2n$. Introduisons l'opérateur $2n \times 2n$;

$$g = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \iff g^{-1} = -g = g^T$$

Définition; La matrice M est symplectique si et seulement si:

$$M^{-1} = g^{-1} M^T g$$

Propriété 1:

$$M = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \quad M^{-1} = \begin{bmatrix} D^T & -B^T \\ -C^T & A^T \end{bmatrix}$$

Propriété 2:

$$M^{-1}M = I \quad \left\{ \begin{array}{l} A^T D - C^T B = D^T A - B^T C = I \\ B^T D \text{ symétrique; si } \exists D^{-1}: BD^{-1} \text{ symétrique} \\ A^T C \text{ symétrique; si } \exists A^{-1}: CA^{-1} \text{ symétrique} \end{array} \right.$$

Propriété 3:

$$MM^{-1} = I \quad \left\{ \begin{array}{l} AD^T - BC^T = DA^T - CB^T = I \\ AB^T \text{ symétrique; si } \exists A^{-1}: A^{-1}B \text{ symétrique} \\ CD^T \text{ symétrique; si } \exists D^{-1}: D^{-1}C \text{ symétrique} \end{array} \right.$$

Propriété 4:

Si λ est valeur propre de M alors $1/\lambda$ l'est aussi, d'où $\det M = 1$

Propriété 5:

Soit $v = \begin{bmatrix} x \\ y \end{bmatrix}$ le vecteur propre associé à λ , alors $(y^T, -x^T)$

est le vecteur propre gauche associé à $1/\lambda$.

Matrice symplectique

Soit une matrice réelle $2n \times 2n$. Introduisons l'opérateur $2n \times 2n$;

$$g = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \iff g^{-1} = -g = g^T$$

Définition; La matrice M est symplectique si et seulement si:

$$M^{-1} = g^{-1} M^T g$$

Propriété 1:

$$M = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \quad M^{-1} = \begin{bmatrix} D^T & -B^T \\ -C^T & A^T \end{bmatrix}$$

Propriété 2:

$$M^{-1}M = I \quad \left\{ \begin{array}{l} A^T D - C^T B = D^T A - B^T C = I \\ B^T D \text{ symétrique; si } \exists D^{-1}: BD^{-1} \text{ symétrique} \\ A^T C \text{ symétrique; si } \exists A^{-1}: CA^{-1} \text{ symétrique} \end{array} \right.$$

Propriété 3:

$$MM^{-1} = I \quad \left\{ \begin{array}{l} AD^T - BC^T = DA^T - CB^T = I \\ AB^T \text{ symétrique; si } \exists A^{-1}: A^{-1}B \text{ symétrique} \\ CD^T \text{ symétrique; si } \exists D^{-1}: D^{-1}C \text{ symétrique} \end{array} \right.$$

Propriété 4:

Si λ est valeur propre de M alors $1/\lambda$ l'est aussi, d'où $\det M = 1$

Propriété 5:

Soit $V = \begin{bmatrix} x \\ y \end{bmatrix}$ le vecteur propre associé à λ , alors $(y^T, -x^T)$

est le vecteur propre gauche associé à $1/\lambda$.

Matrice hamiltonienne

Soit H une matrice réelle $2n \times 2n$

Définition: la matrice H est hamiltonienne si et seulement si:

$$-H = J^{-1} H^T J$$

Propriété 1: forme générique

$$H = \begin{bmatrix} A & B^T B \\ C^T C & -A^T \end{bmatrix}$$

Propriété 2:

Si λ est valeur propre de H , alors $-\lambda$ est aussi valeur propre de H .

Propriété 3:

Soit $V = \begin{bmatrix} x \\ y \end{bmatrix}$ le vecteur propre associé à λ , alors $(y^T, -x^T)$ est le vecteur propre gauche associé à $-\lambda$.

Relation entre matrice hamiltonienne et sympléctique

Soit S une matrice sympléctique telle que $|\lambda(S)| \neq 1$. On peut alors lui associer une matrice hamiltonienne H , telle que réel $[\mu(H)] \neq 0$ à l'aide de la transformation bilinéaire

$$\mu = \frac{\lambda + 1}{\lambda - 1}$$

d'où
$$H = (S-1)^{-1}(S+1) = (S+1)(S-1)^{-1}$$

ou inversement

$$S = (H-1)^{-1}(H+1) = (H+1)(H-1)^{-1}$$

Matrice hamiltonienne

Soit H une matrice réelle $2n \times 2n$

Définition: la matrice H est hamiltonienne si et seulement si:

$$-H = J^{-1} H^T J$$

Propriété 1: forme générique

$$H = \begin{bmatrix} A & B^T B \\ C^T C & -A^T \end{bmatrix}$$

Propriété 2:

Si λ est valeur propre de H , alors $-\lambda$ est aussi valeur propre de H .

Propriété 3:

Soit $V = \begin{bmatrix} x \\ y \end{bmatrix}$ le vecteur propre associé à λ , alors $(y^T, -x^T)$ est le vecteur propre gauche associé à $-\lambda$.

Relation entre matrice hamiltonienne et sympléctique

Soit S une matrice sympléctique telle que $|\lambda(S)| \neq 1$. On peut alors lui associer une matrice hamiltonienne H , telle que réel $[\mu(H)] \neq 0$ à l'aide de la transformation bilinéaire

$$\mu = \frac{\lambda + 1}{\lambda - 1}$$

d'où
$$H = (S-1)^{-1}(S+1) = (S+1)(S-1)^{-1}$$

ou inversement

$$S = (H-1)^{-1}(H+1) = (H+1)(H-1)^{-1}$$

```

1      C      PROGRAMME PRINCIPAL
2      C      RESOLUTION DE L'EQUATION DE RICCATI DANS LE CAS CONTINU
3      C      VARIATIONS DU NOMBRE D'ITERATIONS/A LA DIMENSION
4      DIMENSION IL(36), IC(36), TEST(20)
5      DIMENSION A(3,3), Q(3,3), AT(3,3)
6      DIMENSION V(3,3), H(6,6), HML(6,6)
7      DIMENSION FL21(3,3), F1(3,3), P(3,3)
8      DIMENSION AE(6,6), SIGNE(6,6), F(6,6), FL2(3,3)
9      REAL ID(6,6)
10     DIMENSION AE1(6,6), AD(6,6), AD1(6,6)
11     DIMENSION AEK1(6,6), AEK2(6,6), AES(6,6)
12     DIMENSION PT(3,3), PC(3,3)
13     N=3
14     006     NC=6
15     00A     PR=5.
16     00E     QM1=0.
17     012     ETAL=0.00000001
18     012     C      MISE A ZERO
19     016     DO 1 I=1,N
20     01A     DO 1 J=1,N
21     01E     A(I,J)=0.
22     02E     V(I,J)=0.
23     03E     1     Q(I,J)=0.
24     062     A(1,1)=-1.
25     06C     A(2,1)=-1.
26     078     A(2,3)=-2.
27     084     A(3,2)=1.
28     08A     A(3,3)=-1.
29     096     V(1,1)=1.
30     09A     V(2,2)=1.
31     0A0     Q(1,1)=1.
32     0A4     Q(2,2)=2.
33     0AA     Q(3,3)=3.
34     0AA     C      CALCUL DE LA TRANSPOSEE DE A
35     0B0     DO 21 I=1,N
36     0B4     DO 5J=1:N
37     0B8     AT(I,J)=A(J,I)
38     0D8     5     CONTINUE
39     0E2     21    CONTINUE
40     0E2     C      REMPLISSAGE DE LA MATRICE H
41     0EC     DO 23 I=1,N
42     0FO     DO 6 J=1,N
43     0F4     H(I,J)=A(I,J)
44     114     6     CONTINUE
45     11E     23    CONTINUE
46     128     DO 25 I=1,N
47     12C     DO 14 J=1,N
48     130     J1=J+N
49     136     H(I,J1)=-V(I,J)
50     15C     14    CONTINUE
51     166     25    CONTINUE
52     170     DO 26 I=1,N
53     174     DO 13 J=1,N
54     178     J1=J+N
55     17E     I1=I+N

```

```

1  C PROGRAMME PRINCIPAL
2  C RESOLUTION DE L'EQUATION DE RICCATI DANS LE CAS CONTINU
3  C VARIATIONS DU NOMBRE D'ITERATIONS/A LA DIMENSION
4  DIMENSION IL(36), IC(36), TEST(20)
5  DIMENSION A(3,3), Q(3,3), AT(3,3)
6  DIMENSION V(3,3), H(6,6), HM1(6,6)
7  DIMENSION F121(3,3), F1(3,3), P(3,3)
8  DIMENSION AE(6,6), SIGNE(6,6), F(6,6), F12(3,3)
9  REAL ID(6,6)
10 DIMENSION AE1(6,6), AD(6,6), AD1(6,6)
11 DIMENSION AEK1(6,6), AEK2(6,6), AES(6,6)
12 DIMENSION PT(3,3), PC(3,3)
13 N=3
14 006 NC=6
15 00A PR=5.
16 00E QM1=0.
17 012 ETA1=0.00000001
18 012 C MISE A ZERO
19 016 DO 1 I=1,N
20 01A DO 1 J=1,N
21 01E A(I,J)=0.
22 02E V(I,J)=0.
23 03E 1 Q(I,J)=0.
24 062 A(1,1)=-1.
25 06C A(2,1)=-1.
26 078 A(2,3)=-2.
27 084 A(3,2)=1.
28 08A A(3,3)=-1.
29 096 V(1,1)=1.
30 09A V(2,2)=1.
31 0AO Q(1,1)=1.
32 0A4 Q(2,2)=2.
33 0AA Q(3,3)=3.
34 0AA C CALCUL DE LA TRANSPOSEE DE A
35 0BO DO 21 I=1,N
36 0B4 DO 5J=1,N
37 0B8 AT(I,J)=A(J,I)
38 0D8 5 CONTINUE
39 0E2 21 CONTINUE
40 0E2 C REMPLISSAGE DE LA MATRICE H
41 0EC DO 23 I=1,N
42 0FO DO 6 J=1,N
43 0F4 H(I,J)=A(I,J)
44 114 6 CONTINUE
45 11E 23 CONTINUE
46 128 DO 25 I=1,N
47 12C DO 14 J=1,N
48 130 J1=J+N
49 136 H(I,J1)=-V(I,J)
50 15C 14 CONTINUE
51 166 25 CONTINUE
52 170 DO 26 I=1,N
53 174 DO 13 J=1,N
54 178 J1=J+N
55 17E I1=I+N

```


56	184		H(IL, J1) = -AT(I, J)
57	1AA	13	CONTINUE
58	1B4	26	CONTINUE
59	1BE		DO 27 I=1, N
60	1CE		DO 17 J=1, N
61	1C6		IL=I+N
62	1CC		H(IL, J) = -Q(I, J)
63	1F2	17	CONTINUE
64	1FC	27	CONTINUE
65	206		PRINT 90
66	210		PRINT 71, ((H(I, J), J=1, NC), I=1, NC)
67	24A		CALL MRINV(H, HML, NC, KOD, DET, EPS, IL, IC)
68	24E		PRINT 60
69	258		PRINT 71, ((HML(I, J), J=1, NC), I=1, NC)
70	292		CALL NORM(H, X, NC)
71	296		PRINT 89, X
72	2A6		CALL NORM(HML, Y, NC)
73	2AA		PRINT 99, Y
74	2BA		AMAX=Y
75	2BE		IF(X.GT.Y) AMAX=X
76	2D6		GO TO 79
77	2D8	79	CONTINUE
78	2D8		PRINT 550, AMAX
79	2D8	C	CALCUL DU NOMBRE D'ITERATIONS TOTAL
80	2E8		N1=ALOG(AMAX)
81	2FO		N1=N1/ALOG(2)
82	302		N1=IFIX(N1)
83	30A		PRINT 80, N1
84	31A		N1=N1+QML+PR
85	326		N1=IFIX(N1)
86	32E		PRINT 110, N1
87	32E	C	CALCUL DE LA MATRICE SIGNE
88	33E		DO 9 I=1, NC
89	342		DO 9 J=1, NC
90	346	9	AE(I, J) = H(I, J)
91	372		CALL MRINV(AE, AEL, NC, KOD, DET, EPS, IL, IC)
92	376		CALL NORM(AE, X, NC)
93	37A		CALL NORM(AEL, Y, NC)
94	37E		ALPHA=SQRT(Y/X)
95	386		PRINT 82, ALPHA
96	396		BETA=ALPHA
97	39A		CALL MUL2(AE, ALPHA, AD, NC)
98	39E		ALPHA=1/ALPHA
99	3A6		CALL MUL2(AEL, ALPHA, AD1, NC)
100	3AA		CALL SOMMAT(AD, AD1, AES, NC)
101	3AE		CALL MUL2(AES, 0.5, AEK1, NC)
102	3B4		BETA=BETA-1
103	3C2		BETA=ABS(BETA)
104	3C8		DO 83 I=1, NC
105	3CE		DO 84 J=1, NC
106	3D6		AEK2(I, J) = AEK1(I, J) - AE(I, J)
107	3F8		AE(I, J) = AEK1(I, J)
108	410	84	CONTINUE
109	420	83	CONTINUE
110	42E		CALL NORM(AEK2, W, NC)

56	184		H(I1, J1)=-AT(I, J)
57	1AA	13	CONTINUE
58	1B4	26	CONTINUE
59	1BE		DO 27 I=1, N
60	1CE		DO 17 J=1, N
61	1C6		I1=I+N
62	1CC		H(I1, J)=-Q(I, J)
63	1F2	17	CONTINUE
64	1FC	27	CONTINUE
65	206		PRINT 90
66	210		PRINT 71, ((H(I, J), J=1, NC), I=1, NC)
67	24A		CALL MRINV(H, HM1, NC, KOD, DET, EPS, IL, IC)
68	24E		PRINT 60
69	258		PRINT 71, ((HM1(I, J), J=1, NC), I=1, NC)
70	292		CALL NORM(H, X, NC)
71	296		PRINT 89, X
72	2A6		CALL NORM(HM1, Y, NC)
73	2AA		PRINT 99, Y
74	2BA		AMAX=Y
75	2BE		IF(X.GT.Y) AMAX=X
76	2D6		GO TO 79
77	2D8	79	CONTINUE
78	2D8		PRINT 550, AMAX
79	2D8	C	CALCUL DU NOMBRE D'ITERATIONS TOTAL
80	2E8		N1=ALOG(AMAX)
81	2FO		N1=N1/ALOG(2)
82	302		N1=IFIX(N1)
83	30A		PRINT 80, N1
84	31A		N1=N1+QM1+PR
85	326		N1=IFIX(N1)
86	32E		PRINT 110, N1
87	32E	C	CALCUL DE LA MATRICE SIGNE
88	33E		DO 9 I=1, NC
89	342		DO 9 J=1, NC
90	346	9	AE(I, J)=H(I, J)
91	372		CALL MRINV(AE, AE1, NC, KOD, DET, EPS, IL, IC)
92	376		CALL NORM(AE, X, NC)
93	37A		CALL NORM(AE1, Y, NC)
94	37E		ALPHA=SQRT(Y/X)
95	386		PRINT 82, ALPHA
96	396		BETA=ALPHA
97	39A		CALL MUL2(AE, ALPHA, AD, NC)
98	39E		ALPHA=1/ALPHA
99	3A6		CALL MUL2(AE1, ALPHA, AD1, NC)
100	3AA		CALL SOMMAT(AD, AD1, AES, NC)
101	3AE		CALL MUL2(AES, 0.5, AEK1, NC)
102	3B4		BETA=BETA-1
103	3C2		BETA=ABS(BETA)
104	3C8		DO 83 I=1, NC
105	3CE		DO 84 J=1, NC
106	3D6		AEK2(I, J)=AEK1(I, J)-AE(I, J)
107	3F8		AE(I, J)=AEK1(I, J)
108	410	84	CONTINUE
109	420	83	CONTINUE
110	42E		CALL NORM(AEK2, W, NC)

```

111 432 GAMMA=W/X
112 43C IF(BETA.GT.GAMMA) TEST(1)=BETA
113 45C TEST(1)=GAMMA
114 462 TEST(1)=TEST(1)/2
115 472 DO 10 K=1,N1
116 478 K1=K+1
117 47E CALL NORM(AE,X,NC)
118 482 CALL MRINV(AE,AEL,NC,KOD,DET,EPS,IL,IC)
119 486 CALL NORM(AEL,Y,NC)
120 48A ALPHA=SQRT(Y/X)
121 494 PRINT 82,ALPHA
122 4A4 BETA=ALPHA
123 4A8 CALL MUL2(AE,ALPHA,NC)
124 4AE ALPHA=1/ALPHA
125 4B8 CALL MUL2(AEL,ALPHA,AD1,NC)
126 4BE CALL SOMMAT(AD,AD1,AES,NC)
127 4C2 CALL MUL2(AES,0.5,AEK1,NC)
128 4C6 BETA=BETA-1
129 4D4 BETA=ABS(BETA)
130 4DA DO 85 I=1,NC
131 4EO DO 86 J=1,NC
132 4E8 AEK2(I,J)=AEK1(I,J)-AE(I,J)
133 50A AE(I,J)=AEK1(I,J)
134 522 86 CONTINUE
135 532 85 CONTINUE
136 540 CALL NORM(AEK2,W,NC)
137 544 GAMMA=W/X
138 54E IF(BETA.GT.GAMMA) TEST(K1)=BETA
139 576 TEST(K1)=GAMMA
140 584 IF(TEST(K1).LT.ETA1) GOTO 33
141 5A4 GOTO 34
142 5A6 33 IF(TEST(K1).GT.TEST(K1-1)) GOTO 35
143 5D6 GOTO 34
144 5D8 34 TEST(K1)=TEST(K1)/2
145 5FA 10 CONTINUE
146 60A 35 CONTINUE
147 60A DO 11 I=1,NC
148 612 DO 11 J=1,NC
149 61A 11 SIGNE(I,J)=AEK1(I,J)
150 64E PRINT 100
151 658 PRINT 101,((SIGNE(I,J),J=1,NC),I=1,NC)
152 6A2 DO 7 I=1,NC
153 6AA DO 7 J=1,NC
154 6B2 ID(I,J)=0.
155 6C4 IF(I.EQ.J) ID(I,J)=1
156 6E8 7 CONTINUE
157 704 PRINT 150
158 70E PRINT 160,((ID(I,J),J=1,NC),I=1,NC)
159 70E C CALCUL DE LA MATRICE F
160 758 CALL SOMMAT(SIGNE,ID,F,NC)
161 75C PRINT 71,((F(I,J),J=1,NC),I=1,NC)
162 7A6 CALL MUL2(F,0.5,F,NC)
163 7AA DO 12 I=1,N
164 7B2 DO 12 J=1,N
165 7BA J1=J+N

```

```

111 432 GAMMA=W/X
112 43C IF(BETA.GT.GAMMA) TEST(1)=BETA
113 45C TEST(1)=GAMMA
114 462 TEST(1)=TEST(1)/2
115 472 DO 10 K=1,N1
116 478 K1=K+1
117 47E CALL NORM(AE,X,NC)
118 482 CALL MRINV(AE,AE1,NC,KOD,DET,EPS,IL,IC)
119 486 CALL NORM(AE1,Y,NC)
120 48A ALPHA=SQRT(Y/X)
121 494 PRINT 82,ALPHA
122 4A4 BETA=ALPHA
123 4A8 CALL MUL2(AE,ALPHA,NC)
124 4AE ALPHA=1/ALPHA
125 4B8 CALL MUL2(AE1,ALPHA,AD1,NC)
126 4BE CALL SOMMAT(AD,AD1,AES,NC)
127 4C2 CALL MUL2(AES,0.5,AEK1,NC)
128 4C6 BETA=BETA-1
129 4D4 BETA=ABS(BETA)
130 4DA DO 85 I=1,NC
131 4EO DO 86 J=1,NC
132 4E8 AEK2(I,J)=AEK1(I,J)-AE(I,J)
133 50A AE(I,J)=AEK1(I,J)
134 522 86 CONTINUE
135 532 85 CONTINUE
136 540 CALL NORM(AEK2,W,NC)
137 544 GAMMA=W/X
138 54E IF(BETA.GT.GAMMA) TEST(K1)=BETA
139 576 TEST(K1)=GAMMA
140 584 IF(TEST(K1).LT.ETA1) GOTO 33
141 5A4 GOTO 34
142 5A6 33 IF(TEST(K1).GT.TEST(K1-1)) GOTO 35
143 5D6 GOTO 34
144 5D8 34 TEST(K1)=TEST(K1)/2
145 5FA 10 CONTINUE
146 60A 35 CONTINUE
147 60A DO 11 I=1,NC
148 612 DO 11 J=1,NC
149 61A 11 SIGNE(I,J)=AEK1(I,J)
150 64E PRINT 100
151 658 PRINT 101,((SIGNE(I,J),J=1,NC),I=1,NC)
152 6A2 DO 7 I=1,NC
153 6AA DO 7 J=1,NC
154 6B2 ID(I,J)=0.
155 6C4 IF(I.EQ.J) ID(I,J)=1
156 6E8 7 CONTINUE
157 704 PRINT 150
158 70E PRINT 160,((ID(I,J),J=1,NC),I=1,NC)
159 70E C CALCUL DE LA MATRICE F
160 758 CALL SOMMAT(SIGNE,ID,F,NC)
161 75C PRINT 71,((F(I,J),J=1,NC),I=1,NC)
162 7A6 CALL MUL2(F,0.5,F,NC)
163 7AA DO 12 I=1,N
164 7B2 DO 12 J=1,N
165 7BA J1=J+N

```

```

166      7C2      F12(I,J)=F(I,J1)
167      7F2      12  F1(I,J)=F(I,J)
168      820      CALL MRINV(F12,F121,N,KOD,DET,EPS,IL,IC)
169      824      DO 8 I=1,N
170      82C      DO 8 J=1,N
171      834      8   F121(I,J)=-F121(I,J)
172      834      C   CALCUL DE LA MATRICE P
173      86E      CALL PROD(F121,F1,P,N)
174      872      PRINT 200
175      87C      PRINT 201,((P(I,J),J=1,N),I=1,N)
176      87C      C   CALCUL DE LA MATRICE P CORRIGEE
177      8C6      DO 56 I=1,N
178      8CE      DO 57 J=1,N
179      8D6      PT(I,J)=P(J,I)
180      8F8      57  CONTINUE
181      908      56  CONTINUE
182      916      CALL SOMMAT(P,PT,PC,N)
183      91A      CAMM MUL2(PC,0.5,PC,N)
184      91E      PRINT 300
185      928      PRINT 201,((PC(I,J),J=1,N),I=1,N)
186      972      90  FORMAT(////,10X,'MATRICE H')
187      972      150  FORMAT(////,10X,'MATRICE ID')
188      972      300  FORMAT(////,10X,'MATRICE PC')
189      972      17   FORMAT(////,20X,'NOMBRE D' ITERATIONS=',I2)
190      972      101  FORMAT(////,6(5X,F9.6))
191      972      100  FORMAT(////,10X,'MATRICE SIGNE. ')
192      972      200  FORMAT(////,10X,'MATRICE P')
193      972      82   FORMAT(////,10X,'ALPHA=',F12.5,/)
194      972      60   FORMAT(////,10X,'MATRICE HM1')
195      972      80   FORMAT(////,20X,'N1=',I2)
196      972      89   FORMAT(1H,////,20X,'NORM DE H=',F12.6)
197      972      99   FORMAT(////,10X,'NORM DE HM1=',F12.6)
198      972      550  FORMAT(////,10X,'AMAX=',F6.3,/)
199      972      201  FORMAT(////,3(4X,F10.7))
200      972      71   FORMAT(////,6(3X,F10.6))
201      972      160  FORMAT(////,6(3X,F2.0))
202      972      STOP
203      974      END

```

```

166      7C2      F12(I,J)=F(I,J1)
167      7F2      12  F1(I,J)=F(I,J)
168      820      CALL MRINV(F12,F121,N,KOD,DET,EPS,IL,IC)
169      824      DO 8 I=1,N
170      82C      DO 8 J=1,N
171      834      8   F121(I,J)=-F121(I,J)
172      834      C   CALCUL DE LA MATRICE P
173      86E      CALL PROD(F121,F1,P,N)
174      872      PRINT 200
175      87C      PRINT 201,((P(I,J),J=1,N),I=1,N)
176      87C      C   CALCUL DE LA MATRICE P CORRIGEE
177      8C6      DO 56 I=1,N
178      8CE      DO 57 J=1,N
179      8D6      PT(I,J)=P(J,I)
180      8F8      57  CONTINUE
181      908      56  CONTINUE
182      916      CALL SOMMAT(P,PT,PC,N)
183      91A      CAMM MUL2(PC,0.5,PC,N)
184      91E      PRINT 300
185      928      PRINT 201,((PC(I,J),J=1,N),I=1,N)
186      972      90  FORMAT(////,10X,'MATRICE H')
187      972      150 FORMAT(////,10X,'MATRICE ID')
188      972      300 FORMAT(////,10X,'MATRICE PC')
189      972      17  FORMAT(////,20X,'NOMBRE D' ITERATIONS=',I2)
190      972      101 FORMAT(////,6(5X,F9.6))
191      972      100 FORMAT(////,10X,'MATRICE SIGNE.')
```

```

192      972      200 FORMAT(////,10X,'MATRICE P')
193      972      82  FORMAT(////,10X,'ALPHA=',F12.5,/)
194      972      60  FORMAT(////,10X,'MATRICE HM1')
195      972      80  FORMAT(////,20X,'N1=',I2)
196      972      89  FORMAT(1H,////,20X,'NORM DE H=',F12.6)
197      972      99  FORMAT(////,10X,'NORME DE HM1=',F12.6)
198      972      550 FORMAT(////,10X,'AMAX=',F6.3,/)
199      972      201 FORMAT(////,3(4X,F10.7))
200      972      71  FORMAT(////,6(3X,F10.6))
201      972      160 FORMAT(////,6(3X,F2.0))
202      972      STOP
203      974      END
```

1		C	PROGRAMME PRINCIPAL
2		C	RESOLUTION DE L'EQUATION DE RICCATI DANS LE CAS DISCRET
3			DIMENSION IL (100), IC (100)
4			DIMENSION AEK1 (10,10), AEK2(10,10), AE1(10,10)
5			DIMENSION HT (5,5), SIGNE (10,10), F(10,10), AE(10,10)
6			DIMENSION A(5,5), B (5,5), R(5,5), H(5,5)
7			DIMENSION V(5,5), AT(5,5), BF (5,5), U (10,10)
8			DIMENSION VI (5,5), FT (5,5), PC (5,5)
9			DIMENSION Q (5,5), T (10,10), RM1 (5,5)
10			DIMENSION ID (5,5), U1 (10,10), S (10,10), SS1 (10,10)
11			DIMENSION F1 (5,5), P21 (5,5), PC(5,5)
12			DIMENSION FF1 (10,10), AD (10,10) AD1 (10,10)
13			REAL M (10,10), MM1 (10,10), L (10,10), IE (10,10)
14			DIMENSION TEST (20)
15			N = 5
16	006		NC =10
17	00A		QM1 = 0.
18	00E		PR = 5.
19	012		ETA1 = 0.00000001
20	016		DO 1 I=1,N
21	01A		DO 1J = 1,N
22	01E		A (I,J) = 0.
23	02E		R (I, J) = 0.
24	03E		Q (I,J) = 0.
25	04E	1	L (I,J) =0.
26	072		A (1,1) =-0.75
27	076		A (1,2) =-0.09
28	07C		A (2,1) =1.74
29	088		A (2,2) =-0.91
30	08E		A(3,1) =-0.3
31	09A		A (3,2) =-0.0015
32	0A6		A (3,3) =-0.95
33	0AC		A (4,4) =-0.55
34	0B2		A (5,1) =-0.15
35	0BE		A (5,2) =-0.008
36	0CA		A (5,5) =-0.905
37	0DO		PRINT 30
38	0DA		PRINT 31, (A(I,J)J=1,N), I=1,N)
39	114		H (1,3) = 24.64
40	11A		H (2,4) = 0.835
41	120		H (3,5) = 1.83
42	126		H (1,1) = 1.
43	12A		B (2,2) = 2.
44	130		B (4,2) = 1.
45	136		B (5,1) = 1.
46	13C		PRINT 40
47	146		PRINT 41, ((B (I,J),J =1,N), I=1,N)
48	180		DO 55 I =1, N
49	184		DO 4 J = 1, N
50	188		BT (I,J) = B (J,I)
51	1A8	4	CONTINUE
52	1B2	55	CONTINUE
53	1BC		DO 221 = 1, N
54	100		DO 15 J = 1,N
55	1C4	15	HT (I,J) =P (J,I)
56	1EE	22	CONTINUE

1		C	PROGRAMME PRINCIPAL
2		C	RESOLUTION DE L'EQUATION DE RICCATI DANS LE CAS DISCRET
3			DIMENSION IL (100), IC (100)
4			DIMENSION AEK1 (10,10), AEK2(10,10), AE1(10,10)
5			DIMENSION HT (5,5), SIGNE (10,10), F(10,10), AE(10,10)
6			DIMENSION A(5,5), B (5,5), R(5,5), H(5,5)
7			DIMENSION V(5,5), AT(5,5), BF (5,5), U (10,10)
8			DIMENSION VI (5,5), FT (5,5), PC (5,5)
9			DIMENSION Q (5,5), T (10,10), RM1 (5,5)
10			DIMENSION ID (5,5), U1 (10,10), S (10,10), SS1 (10,10)
11			DIMENSION F1 (5,5), P21 (5,5), PC(5,5)
12			DIMENSION PF1 (10,10), AD (10,10) AD1 (10,10)
13			REAL M (10,10), MM1 (10,10), L (10,10), IE (10,10)
14			DIMENSION TEST (20)
15			N = 5
16	006		NC =10
17	00A		QM1 = 0.
18	00E		PR = 5.
19	012		ETA1 = 0.00000001
20	016		DO 1 I=1,N
21	01A		DO 1J = 1,N
22	01E		A (I,J) = 0.
23	02E		R (I, J) = 0.
24	03E		Q (I,J) = 0.
25	04E	1	L (I,J) =0.
26	072		A (1,1) =0.75
27	076		A (1,2) =0.09
28	07C		A (2,1) =1.74
29	088		A (2,2) =0.91
30	08E		A (3,1) =0.3
31	09A		A (3,2) =0.0015
32	0A6		A (3,3) =0.95
33	0AC		A (4,4) =0.55
34	0B2		A (5,1) =0.15
35	0BE		A (5,2) =0.008
36	0CA		A (5,5) =0.905
37	0DO		PRINT 30
38	0DA		PRINT 31, (A(I,J)J=1,N), I=1,N)
39	114		H (1,3) = 24.64
40	11A		H (2,4) = 0.835
41	120		H (3,5) = 1.83
42	126		H (1,1) = 1.
43	12A		B (2,2) = 2.
44	130		B (4,2) = 1.
45	136		B (5,1) = 1.
46	13C		PRINT 40
47	146		PRINT 41, ((3 (I,J),J =1,N), I=1,N)
48	180		DO 55 I =1, N
49	184		DO 4 J = 1, N
50	188		BT (I,J) = B (J,I)
51	1A8	4	CONTINUE
52	1B2	55	CONTINUE
53	1BC		DO 221 = 1, N
54	100		DO 15 J = 1,N
55	1C4	15	HT (I,J) =P (J,I)
56	1EE	22	CONTINUE

57	1F8		DO 21 I = 1,N
58	1FC		DO 5 J = 1,N
59	200	5	AT (I,J) = A(J,I)
60	22A	21	CONTINUE
61	234		CALL PROD (HT,H,Q,N)
62	238		PRINT 51, ((Q(I,J), J = 1,N), I = 1,N)
63	272		R (1,1) = 1.
64	276		R (2,2) = 1.
65	27C		R (3,3) = 1.
66	282		R (4,4) = 1.
67	288		R (5,5) = 1.
68	28E		PRINT 60
69	298		PRINT 41, ((R(I,J), = 1,N), I = 1,N)
70	2D2		CALL MRINV (R, RM1,N,KOD,DET,EPS,IL,1C)
71	2D8		PRINT 41, ((RM1 (I,J), J = 1,N), I = 1,N)
72	322		CALL PROD (B,RM1, V1, N)
73	326		CALL PROD (V1,BT,V?N)
74	32A		DO 23 I = 1,N
75	332		DO 6J = 1,N
76	33A		U(I,J) = 0
77	34C		IF (I.EQ.J)L (I,J) = 1.
78!	370		GOTO 6
79	372	6	CONTINUE
80	380	23	CONTINUE
81	38E		DO 25 I = 1,N
82	396		DO 14 J = 1,N
83	39E		J1 = J+N
84	3A8		U(I,J1) = V (I,J)
85	3CE	14	CONTINUE
86	3DE	25	CONTINUE
87	3EC		DO 26 I = 1,N
88	3F4		DO 13 J = 1,N
89	3FC		J1 = J + N
90	406		I1 = I + N
91	410		U (I1, J1) = AT (I,J)
92	436	13	CONTINUE
93	446	26	CONTINUE
94	454		DO 27 I = 1,N
95	45C		DO 17 J = 1,N
96	464		I 1 = I + 5
97	46E		U (I1,J) = 0.
98	484	17	CONTINUE
99	494	27	CONTINUE
100	4A2		PRINT 160
101	4AC		PRINT 71, (U(I,J)J=1;NC,I=1,NC)
102	4F6		DO 7 I=1,N
103	4FE		DO 7 J=1,N
104	506		I1 = 1,N
105	510		J1 = J+N
106	51A		L (I!,J1)=0
107	536		IF (I1.EQ.J1)L(I1,J1) = 1
108	564		GOTO 7
109	566	7	CONTINUE
110	582		DO 28 J = 1,N

57	1F8		DO 21 I = 1, N
58	1FC		DO 5 J = 1, N
59	200	5	AT (I, J) = A(J, I)
60	22A	21	CONTINUE
61	234		CALL PROD (HT, H, Q, N)
62	238		PRINT 51, ((Q(I, J), J = 1, N), I = 1, N)
63	272		R (1, 1) = 1.
64	276		R (2, 2) = 1.
65	27C		RR (3, 3) = 1.
66	282		R (4, 4) = 1.
67	288		R (5, 5) = 1.
68	28E		PRINT 60
69	298		PRINT 41, ((R(I, J), = 1, N), I = 1, N)
70	2D2		CALL MRINV (R, RM1, N, KOD, DET, EPS, IL, 1C)
71	2D8		PRINT 41, ((RM1 (I, J), J = 1, N), I = 1, N)
72	322		CALL PROD (B, RM1, V1, N)
73	326		CALL PROD (V1, BT, V?N)
74	32A		DO 23 I = 1, N
75	332		DO 6 J = 1, N
76	33A		U(I, J) = 0
77	34C		IF (I.EQ.J)L (I, J) = 1.
78!	370		GOTO 6
79	372	6	CONTINUE
80	380	23	CONTINUE
81	38E		DO 25 I = 1, N
82	396		DO 14 J = 1, N
83	39E		J1 = J + N
84	3A8		U(I, J1) = V (I, J)
85	3CE	14	CONTINUE
86	3DE	25	CONTINUE
87	3EC		DO 26 I = 1, N
88	3F4		DO 13 J = 1, N
89	3FC		J1 = J + N
90	406		I1 = I + N
91	410		U (I1, J1) = AT (I, J)
92	436	13	CONTINUE
93	446	26	CONTINUE
94	454		DO 27 I = 1, N
95	45C		DO 17 J = 1, N
96	464		I 1 = I + 5
97	46E		U (I1, J) = 0.
98	484	17	CONTINUE
99	494	27	CONTINUE
100	4A2		PRINT 160
101	4AC		PRINT 71, (U(I, J) J=1; NC, I=1, NC)
102	4F6		DO 7 I=1, N
103	4FE		DO 7 J=1, N
104	506		I1 = 1, N
105	510		J1 = J + N
106	51A		L (I1, J1) = 0
107	536		IF (I1.EQ.J1)L(I1, J1) = 1
108	564		GOTO 7
109	566	7	CONTINUE
110	582		DO 28 J = 1, N

111	58A		DO 16 J = 1,N
112	592		L (I,J)= A (I,J)
113	5B4	16	CONTINUE
114	5C4	28	CONTINUE
115	5D2		DO 24 J = 1,N
116	5DA		DO 18 J = 1,N
117	5E2		I1 = I+N
118	5EC		L(I1,J) = -Q(I,J)
119	618	18	CONTINUE
120	628	24	CONTINUE
121	636		DO 19 I = 1,N
122	63E		DO 19 J = 1,N
123	646		J1 = J+N
124	650		L(I,J1)=0
125	66C	19	CONTINUE
126	68A		PRINT 150
127	694		PRINT 71, ((L (I,J), J =1,NC), I=1,NC)
128	6DE		CALL SOMMAT (L,U,T,NC)
129	6E2		DO 81 = 1,NC
130	6EA		DO8 J= 1,NC
131	6F2		U1 (I,J)=-U(I,J)
132	734		CALL SOMMAT (L,U1,M,NC)
133	738		CALL MRINV (MM1,NC,KOD,DET,EPS,I1,IC)
134	73C		CALL PROD(MM1,T,S,NC)
135	740		PRINT 90
136	74A		PRINT 71, ((S(I,J), J=1,NC), I=1,NC)
137	794		CALL NORM(S,X,NC)
138	798		PRINT 100,X
139	7A8		CALL MRINV(S,SS1,NC,KOD,DET,EPS,IL,IC)
140	7AC		CALL NORM(SS1,Y,NC)
141	780		AMAX=Y
142	7B4		IF(X.GT.Y) AMAX=X
143	7CC		GOTO 79
144	7CE	79	CONTINUE
145	7CE		N1=ALOG(AMAX)
146	7D6		N1=N1/ALOG(2)
147	7FO		N1=IFIX(N1)
148	7F8		N1=N1+QM1+PR
149	808		PRINT 110,N1

111	58A		DO 16 J = 1,N
112	592		L (I,J)= A (I,J)
113	5B4	16	CONTINUE
114	5C4	28	CONTINUE
115	5D2		DO 24 J = 1,N
116	5DA		DO 18 J = 1,N
117	5E2		I1 = I+N
118	5EC		L(I1,J) = -Q(I,J)
119	618	18	CONTINUE
120	628	24	CONTINUE
121	636		DO 19 I = 1,N
122	63E		DO 19 J = 1,N
123	646		J1 = J+N
124	650		L(I,J1)=0
125	66C	19	CONTINUE
126	68A		PRINT 150
127	694		ERINT 71, ((L (I,J), J =1,NC), I=1,NC)
128	6DE		CALL SOMMAT (L,U,T,NC)
129	6E2		DO 81 = 1,NC
130	6EA		DO8 J= 1,NC
131	6F2		U1 (I,J)=-U(I,J)
132	734		CALL SOMMAT (L,U1,M,NC)
133	738		CALL MRINV (MM1,NC,KOD,DET,EPS,I1,IC)
134	73C		CALL PROD(MM1,T,S,NC)
135	740		PRINT 90
136	74A		PRINT 71, ((S(I,J), J=1,NC), I=1,NC)
137	794		CALL NORM(S,X,NC)
138	798		PRINT 100,X
139	7A8		CALL MRINV(S,SS1,NC,KOD,DET,EPS,IL,IC)
140	7AC		CALL NORM(SS1,Y,NC)
141	780		AMAX=Y
142	7B4		IF(X.GT.Y) AMAX=X
143	7CC		GOTO 79
144	7CE	79	CONTINUE
145	7CE		N1=ALOG(AMAX)
146	7D6		N1=N1/ALOG(2)
147	7FO		N1=IFIX(N1)
148	7F8		N1=N1+QM1+PR
149	808		PRINT 110,N1

```

150      818      DO      9I=1,NC
151      820      DO      9J=1,NC
152      828      9      AE(I,J)  =S(I,J)
153      864      CALL MRINV(AE,AE1,NC,KOD,DET,EPS,IL,IC)
154      868      CALL NORM(AE,X,NC)
155      86C      CALL NORM(AE1,Y,NC)
156      870      ALPHA=SQRT(Y/X)
157      878      PRINT 82,ALPHA
158      888      BETA=ALPHA
159      88C      CALL MUL2(AE,ALPHA,AD,NC)
160      890      ALPHA=1/ALPHA
161      898      CALL MUL2(AE1,ALPHA,AD1,NC)
162      89C      CALL SOMMAT4AD,AD1,AEK1,NC)
163      8A0      CALL MUL2(AEK1,0.5,AEK1,NC)
164      8A4      BETA=BETA-1
165      8B8      BETA=ABS(BETA)
166      8BE      DO 84 I=1,NC
167      8C6      DO 83 J=1 NC
168      8CE      AEK2(I,J)=AEK1(I,J)-AE(I,J)
169      8FC      AE(I,J)=AEK1(I,J)
170      91C      83      CONTINUE
171      92C      84      CONTINUE
172      93A      CALL NORM(AEK2,W,NC)
173      93E      GAMMA=W/X
174      944      IF(BETA.GT.GAMMA) TEST(1)=BETA
175      960      TEST(1)=GAMMA
176      966      TEST(1)=TEST(1)/2
177      976      DO 10 K=1,N1
178      97C      K1=K+1
179      982      CALL NORM(AE,X,NC)
180      986      CALL MRINV(AE,AE1,NC,KOD,DET,EPS,IL,IC)
181      98A      CALL NORM(AE1,Y,NC)
182      98E      ALPHA=SQRT(Y/X)
183      996      PRINT 82,ALPHA
184      9A6      BETA=ALPHA
185      9AA      CALL MUL2(AE,ALPHA,AD,NC)
186      9AE      ALPHA=1/ALPHA

```

```

150      818      DO      9I=1,NC
151      820      DO      9J=1,NC
152      828      9      AE(I,J)  =S(I,J)
153      864      CALL MRINV(AE,AE1,NC,KOD,DET,EPS,IL,IC)
154      86E      CALL NORM(AE,X,NC)
155      86C      CALL NORM(AE1,Y,NC)
156      870      ALPHA=SQRT(Y/X)
157      87E      PRINT 82,ALPHA
158      888      BETA=ALPHA
159      88C      CALL MUL2(AE,ALPHA,AD,NC)
160      890      ALPHA=1/ALPHA
161      898      CALL MUL2(AE1,ALPHA,AD1,NC)
162      89C      CALL SOMMAT4AD,AD1,AEK1,NC)
163      8A0      CALL MUL2(AEK1,0.5,AEK1,NC)
164      8A4      BETA=BETA-1
165      8B8      BETA=ABS(BETA)
166      8BE      DO 84 I=1,NC
167      8C6      DO 83 J=1 NC
168      8CE      AEK2(I,J)=AEK1(I,J)-AE(I,J)
169      8FC      AE(I,J)=AEK1(I,J)
170      91C      83      CONTINUE
171      92C      84      CONTINUE
172      93A      CALL NORM(AEK2,W,NC)
173      93E      GAMMA=W/X
174      944      IF(BETA.GT.GAMMA) TEST(1)=BETA
175      960      TEST(1)=GAMMA
176      966      TEST(1)=TEST(1)/2
177      976      DO 10 K=1,N1
178      97C      K1=K+1
179      982      CALL NORM(AE,X,NC)
180      986      CALL MRINV(AE,AE1,NC,KOD,DET,EPS,IL,IC)
181      98A      CALL NORM(AE1,Y,NC)
182      98E      ALPHA=SQRT(Y/X)
183      996      PRINT 82,ALPHA
184      9A6      BETA=ALPHA
185      9AA      CALL MUL2(AE,ALPHA,AD,NC)
186      9AE      ALPHA=1/ALPHA

```

```

187      9B6      CALL MUL2(AE1,ALPHA,AD1,NC)
188      9BA      CALL SOMNAT(AD,AD1,AEK1,NC)
189      9BE      CALL MUL2(AEK1,0.5,AEK1,NC)
190      9O2      BETA=BETA-1
191      9D6      BETA=ABS(BETA)
192      9DC      DO 85 I=1,NC
193      9E4      DO 86 J=1,NC
194      9EC      AEK2(I,J)=AEK1(I,J)-AE(I,J)
195      A1A      AE(I,J)=AEK1(I,J)
196      A3A      86. CONTINUE
197      A4A      85 CONTINUE
198      A53      CALL NORM(AEK2,W,NC)
199      A      GAMMA+W/X
200      A62      IF(BETA.GT.GAMMA) TEST(K1)=BETA
201      A86      TEST(K1)=GAMMA
202      A94      IF(TEST(K1).LT.ETA1) GOTO 33
203      AB4
204      AB6      33 IF(TEST(K1).GT.TEST(K1-1)) GOTE 35
205      AF6      GOTO 34
206      AEA      34 TEST(K1)=TEST(K1)/2
207      B10      10 CONTINUE
208      B1C      35 CONTINUE
209      B1E      DO 11 I=1,NC
210      B26      DO 11 J=1,NC
211      B2E      11 SIGNE(I,J)=AEK1(I,J)
212      B6A      PRINT 1101
213      B74      PRINT 71,((SIGNE(I,J),J=1,NC),I=1,NC)
214      BCO      DO 3 I=1,N
215      EC8      DO 3 J=1,N
216      BDO      I1=I+N
217      BDA      J1=J+N
218      BE4      IE(I1,J1)=0.
219      BF6      IF(I1.EQ.J1) IE(I1,J1)=1.
220      C1F      3 CONTINUE
221      C3A      DO 2 I=1,N
222      C42      DO 2 J=1,N

```

```

187      9B6 .      CALL MUL2(AE1,ALPHA,AD1,NC)
188      9BA      CALL SOMNAT(AD,AD1,AEK1,NC)
189      9BF      CALL MUL2(AEK1,0.5,AEK1,NC)
190      9O2      BETA=BETA-1
191      9D6 .      BETA=ABS(BETA)
192      9DC      DO 85 I=1,NC
193      9E4      DO 86 J=1,NC
194      9EC      AEK2(I,J)=AEK1(I,J)-AE(I,J)
195      A1A      AE(I,J)=AEK1(I,J)
196 .     A3A      86 . CONTINUE
197      A4A      85 . CONTINUE
198      A5B .     CALL NORM(AEK2,W,NC)
199      A .      GAMMA+W/X
200      A62      IF(BETA.GT.GAMMA) TEST(K1)=BETA
201      A86      TEST(K1)=GAMMA
202      A94      IF(TEST(K1).LT.ETA1) GOTO 33
203      AB4
204      AB6      33 IF(TEST(K1).GT.TEST(K1-1)) GOTE 35
205      AF6      GOTO 34
206      AEA      34 TEST(K1)=TEST(K1)/2
207      B1O      10 CONTINUE
208      B1C      35 CONTINUE
209      B1E      DO 11 I=1,NC
210      B26 .     DO 11 J=1,NC
211      B2E      11 SIGNE(I,J)=AEK1(I,J)
212      B6A      PRINT 1101
213      B74      PRINT 71,((SIGNE(I;J),J=1,NC),I=1,NC)
214      BCO      DO 3 I=1,N
215      EC8 .     DO 3 J=1,N
216      BDO      I1=I+N
217      BDA .     J1=J+N
218 .     BE4      IE(I1,J1)=0.
219      BF6 .     IF(I1.EQ:J1) IE(I1,J1)=1.
220      C1F      3 CONTINUE
221      C3A .     DO 2 I=1,N
222      C42      DO 2 J=1,N

```



```

223 C4A IE(I,J)=0.
224 C5C IF(I.EQ.J) IE(I,J)=-1.
225 C86 2 CONTINUE
226 CA2 DO 20 I=1,N
227 CAA DO 29 J=1,N
228 CB2 I1=I+N
229 CBC IE(I1,J)=0.
230 CD2 29 CONTINUE
231 CE2 20 CONTINUE
232 CFO DO 1500 I=1,N
233 CF8 DO 1501 J=1,N
234 DOO J1=J+N
235 DOA IF(I,J1)=0.
236 D1C 1501 CONTINUE
237 D2C 1500 CONTINUE
238 D3A CALL SOMMAT(SIGNE,IE,F,NC)
239 D3E PRINT 120
240 D48 PRINT 71,((F(I,J),J=1,NC),I=1,NC)
241 D94 DO 12 I=1,N
242 D9C DO 12 J=1,N
243 DA4 I1=I+N
244 DAE F1(I,J)=F(I,J)
245 DDO 12 F21(I,J)=F(I1,I)
246 E12 CALL MRINV(F1,FF1,N,KOD,DET,EPS,IL,IC)
247 E16 CALL PROD(F21,FF1,F,N)
248 E1A PRINT 140
249 E24 PRINT 81,((P(I,J),J=1,N),I=1,N)
250 E6E DO 56 I=1,N
251 E76 DO 57 J=1,N
252 E7F PT(I,J)=F(J,I)
253 E8O 57 CONTINUE
254 E8O 56 CONTINUE
255 E8E CALL SOMMAT(P,PT,PC,N)
256 EC2 CALL MUL 2 (PC,0.5,PC,N)
257 EC6 PRINT 700
258 EDO PRINT 81,((PC(I,J),j=1,N),I=1,N)
259 F1A 71 FORMAT (///,10(1x,F12.5))
260 F1A 81 FORMAT (///,5(4x F12.5))
261 F1A 90 FORMAT (////,10x,MATRICE S ' ;//)

```

```

223   C4A           IE(I,J)=0.
224   C5C           IF(I.EQ.J) IE(I,J)=-1.
225   C86           2   CONTINUE
226   CA2           DO 20 I=1,N
227   CAA           DO 29 J=1,N
228   CB2           I1=I+N
229   CBC           IE(I1,J)=0.
230   CD2           29  CONTINUE
231   CE2           20  CONTINUE
232   CFO           DO 1500 I=1,N
233   CF8           DO 1501 J=1,N
234   D00           J1=J+N
235   DOA           IF(I,J1)=0.
236   D1C           1501 CONTINUE
237   D2C           1500 CONTINUE
238   D3A           CALL SOMMAT(SIGNE,IE,F,NC)
239   D3E           PRINT 120
240   D48           PRINT 71,((F(I,J),J=1,NC),I=1,NC)
241   D94           DO 12 I=1,N
242   D9C           DO 12 J=1,N
243   DA4           I1=I+N
244   DAE           F1(I,J)=F(I,J)
245   DDO           12  F21(I,J)=F(I1,I)
246   E12           CALL MRINV(F1,FF1,N,KOD,DET,EPS,IL,IC)
247   E16           CALL PROD(F21,FF1,F,N)
248   E1A           PRINT 140
249   E24           PRINT 81,((P(I,J),J=1,N),I=1,N)
250   E6E           DO 56 I=1,N
251   E76           DO 57 J=1,N
252   E7F           PT(I,J)=F(J,I)
253   E80           57  CONTINUE
254   E8C           56  CONTINUE
255   E8E           CALL SOMMAT(P,PT,PC,N)
256   EC2           CALL MUL 2 (PC,0.5,PC,N)
257   EC6           PRINT 700
258   EDO           PRINT 81,((PC(I,J),j=1,N),I=1,N)
259   F1A           71  FORMAT (///,10(1x,F12.5))
260   F1A           81  FORMAT (///,5(4x F12.5))
261   F1A           90  FORMAT (////,10x,MATRICE S '7//)

```

262	F1A	1101	FORMAT (////,10x, MATRICE SIGNE.'//)
263	F1A	150	FORMAT (////,10x, MATRICE L')
264	F1A	160	FORMAT (////,10x, MATRICE L')
265	F1A	41	FORMAT (////,5 (10x,F2.0)
266	F1A	60	FORMAT (////,10x' MATRICE R. ' //)
267	F1A	51	FORMAT (////, 5 (10x,F6.3))
268	F1AA	40	FORMAT (////, 10x, MATRICE B. ' //)
269	F1A	30	FORMAT (////,10x'MATRICE A.'//)
270	F1A	31	FORMAT (////,5 (10x,F7.4)
271	F1A	700	FORMAT (1H,////,10x,' MATRICE PC.' //)
272	F1A	140	FORMAT (1H,////,10x,' MATRICE P.' //)
273	F1A	120	FORMAT (1H,////,10, MATRICE F.' //)
274	F1A	110	FORMAT (////,20 x ' NOMBRE D'ITERATIONS = ' ,12)
275	F1A	82	FORMAT (////,10x ' ALPHA = ', F.12.5,/)
276	F1	100	FORMAT (1 H, ////,20 x, NORME DE S = ', F12.5)
277	F1A		STOP
278	F1C		END

262	F1A	1101	FORMAT (////,10x, MATRICE SIGNE.'//)
263	F1A	150	FORMAT (////,10x, MATRICE L')
264	F1A	160	FORMAT (////,10x, MATRICE L')
265	F1A	41	FORMAT (////,5 (10x,F2.0)
266	F1A	60	FORMAT (////,10x, MATRICE R. ' //)
267	F1A	51	FORMAT (////, 5 (10x,F6.3))
268	F1AA	40	FORMAT (////, 10x, MATRICE B. ' //)
269	F1A	30	FORMAT (////,10x, MATRICE A.'//)
270	F1A	31	FORMAT (////,5 (10x,F7.4)
271	F1A	700	FORMAT (1H,////,10x, ' MATRICE PC.' //)
272	F1A	140	FORMAT (1H,////,10x, ' MATRICE P.' //)
273	F1A	120	FORMAT (1H,////,10, MATRICE F.' //)
274	F1A	110	FORMAT (////,20 x ' NOMBRE D'ITERATIONS = ' ,12)
275	F1A	82	FORMAT (////,10x ' ALPHA = ' , F.12.5,/)
276	F1	100	FORMAT (1 H, ////,20 x, NORME DE S = ' , F12.5)
277	F1A		STOP
278	F1C		END

BIBLIOGRAPHIE

- Cours d'automatique avancée :Optimalisation déterministe de la commande
PAUL LEFEVRE(école nationale supérieure des techniques avancées)
- The matrix sign function and computations in systems(EUGENE DENMAN and
ALEX N.BEAVERS, JR) Applied mathematics and computations 2,63-94(1976)
- 3 J.L.CASTI,Dynamical systems and their applications: linear théory,
Académic Fress 1977.
- 4 B.D.O.ANDERSON,Second ordr convergent algorithms for the steadystate
Riccati equation,Paper FA6-11:CO,1977 I.E.E.E
- 5 B.D.O.ANDERSON,Second ordr convergent algorithms for the steadystate
Riccati equation,Int.J.Contr.,1978,
- A.BARRAUD,Investigation autour de la fonction signe d'une matrice,appli-
cation à l'équation de Riccati,RAIRO automatique,n°4 1979,p335-368
- 7 G.A.HEWER,An Iterative Technique for the computations of the steadystate
Gain for the DiscreteOptimal Regulator,I.E.E.E vol.AC-16,n°4,1971
- 8 A.BARRAUD,Produit étoile et fonction signe de matrice,application à
l'équation de Riccati dans le cas discret,RAIRO automatique/Systems
Analysis and Control(vol.14,n°1,1980.p.55 à 85)
- 10 D.R.VAUGHAN.A Non Recursive Algebraic Solution for the Discrete Riccati
Equation I.E.E.E,vol,AC-15,n°5 1970.
- 11 A.J.LAUB,A Schur Method for Solving Algebraix Riccati Equations,I.E.E.E
1978
- 12 B.T.SMITH, Matrix Eigen-System Routines.EISPACK Guide,vol 6, 1976
- 13 G.W.STEWART,HQR and Exchng Fortran Subroutines vol.2,1976.
- 15 C.L.LAWSON et R.J.HANSON,Solving Least Squares Problems,Prentice Hall
1974.
- GILES.P , Systèmes asservis linéaires

BIBLIOGRAPHIE

- Cours d'automatique avancée :Optimalisation déterministe de la commande
PAUL LEFEVRE(école nationale supérieure des techniques avancées)
- The matrix sign function and computations in systems(EUGENE DENMAN and
ALEX N.BEAVERS, JR) Applied mathematics and computations 2,63-94(1976)
- 3 J.L.CASTI,Dynamical systems and their applications: linear théory,
Académic Press 1977.
- 4 B.D.O.ANDERSON,Second ordr convergent algorithms for the steadystate
Riccati equation,Paper FA6-11:CO,1977 I.E.E.E
- 5 B.D.O.ANDERSON,Second ordr convergent algorithms for the steadystate
Riccati equation,Int.J.Contr.,1978,
- A.BARRAUD,Investigation autour de la fonction signe d'une matrice,appli-
cation à l'équation de Riccati,RAIRO automatique,n°4 1979,p335-368
- 7 G.A.HEWER,An Iterative Technique for the computations of the steadystate
Gain for the DiscreteOptimal Regulator,I.E.E.E vol.AC-16,n°4,1971
- 8 A.BARRAUD,Produit étoile et fonction signe de matrice,application à
l'équation de Riccati dans le cas discret,RAIRO automatique/Systems
Analysis and Control(vol.14,n°1,1980.p.55 à 85)
- 10 D.R.VAUGHAN.A Non-Recursive Algebraic Solution for the Discrete Riccati
Equation I.E.E.E,vol,AC-15,n°5 1970.
- 11 A.J.LAUB,A Schur Method for Solving Algebraix Riccati Equations,I.E.E.E
1978
- 12 B.T.SMITH, Matrix Eigen-System Routines.EISPACK Guide,vol 6, 1976
- 13 G.W.STEWART,HQR and Exchng Fortran Subroutines vol.2,1976.
- 15 C.L.LAWSON et R.J.HANSON,Solving Least Squares Problems,Prentice Hall
1974.
- GILES.P , Systèmes asservis linéaires