



UNIVERSITÉ  
University of Namur  
DE NAMUR

# Institutional Repository - Research Portal Dépôt Institutionnel - Portail de la Recherche

[researchportal.unamur.be](http://researchportal.unamur.be)

## THESIS / THÈSE

### MASTER EN SCIENCES MATHÉMATIQUES

#### La méthode de Ritz Galerkin en contrôle optimal : convergence et calcul de l'erreur

Toint, Philippe

*Award date:*  
1974

*Awarding institution:*  
Universite de Namur

[Link to publication](#)

#### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

#### Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

# La méthode de Ritz Galerkin en contrôle optimal

Convergence et calcul de l'erreur

Philippe Toint

1973-1974

La méthode de Ritz-Galerkin  
en contrôle optimal :

Convergence et calcul de l'erreur .

Philippe TOINT 1973-74

## Introduction

On sait que tout problème de physique mathématique conduit naturellement à la résolution d'une ou plusieurs équations fonctionnelles :

$$Ax = f \quad (\text{I})$$

où  $A$  opère d'un espace  $X$  dans un espace  $Y$ ,  $f$  est donné dans  $Y$  et  $x$  est donné dans  $X$  (exemples : équations différentielles ordinaires, équations intégrales, équations aux dérivées partielles)

En général, la solution de l'équation (I) est impossible à déterminer explicitement ou encore sa forme explicite est si compliquée qu'elle est inutilisable pratiquement, et on s'intéresse par conséquent à la résolution approchée de l'équation. L'idée est alors de remplacer les espaces  $X$  et  $Y$  par des espaces "plus simples"  $X_n$  et  $Y_n$  et d'associer à (I) une famille de problèmes approchés

$$A_n x_n = f_n$$

où  $A_n$  approche  $A$ ,  $f_n$  approche  $f$  dans  $Y_n$  et  $x_n$  approche  $x$

dans  $X_n$  (du moins on le souhaite). Ce procédé porte généralement le nom de Ritz-Galerkin.

Les problèmes qui se posent sont les suivants :

- 1) étude de l'équation exacte
- 2) étude des équations approchées
- 3) étude de la stabilité et de la convergence
- 4) étude de l'erreur ( $\|x - x_n\|$ ) pour un  $n$  fixé.

Au cours de ce mémoire, nous montrons, via la théorie des multiplicateurs de Lagrange, comment on peut appliquer les idées précédentes à des problèmes de contrôle optimal.

Rappelons, en termes volontairement vagues, que la théorie mathématique du contrôle optimal se fait à partir des données suivantes : 1) un contrôle  $u$  "à notre disposition" dans un ensemble  $U_{adm}$  ("l'ensemble des "contrôles admissibles"), 2) l'état  $x(u)$  du système à commander qui est donné, pour  $u$  choisi, par la résolution d'une équation

$$Bx(u) = \text{fonction donnée de } u$$

où  $B$  est l'opérateur (supposé connu) qui "représente" le système à contrôler ( $B$  est le modèle du système), 3) la fonction économique  $J[u]$ . On cherche

$$\inf \{J(u) \mid u \in U_{adm}\}$$

Les questions étudiées dans ce travail sont les suivantes.

Le procédé de Ritz-Galerkin est-il convergent ? Comment calculer les contrôles approchés, les états approchés et les valeurs de la fonction économique  $J$  ?

Le plan que nous adopterons est le suivant.

Dans un premier chapitre, nous rappellerons différents points de la théorie générale de l'optimisation pour pouvoir aborder, avec les outils nécessaires, le second chapitre qui posera, de manière précise, le problème de contrôle optimal. Un dernier chapitre sera alors consacré à l'application de la méthode de Ritz-Galerkin à quelques modèles plus précis.

## Chapitre I

### Questions générales d'optimisation

#### I Introduction.

Ce chapitre a pour but essentiel de rappeler quelques outils fréquemment utilisés dans la théorie de l'optimisation, et qui nous seront utiles dans la suite. Il n'est évidemment pas exhaustif, et beaucoup des cas qu'il traite peuvent l'être aussi dans un cadre beaucoup plus général.

Nous parlerons d'abord des problèmes de dualité, qui seront introduits par le biais des formes bilinéaires. Ensuite, nous aborderons la théorie de l'optimisation avec contraintes, globale et locale, vue sous l'angle des multiplicateurs de Lagrange.

#### II La dualité : espaces doux

61

Définissons d'abord le concept même de dualité dans le cadre des espaces vectoriels.

Nous dirons que

Deux espaces vectoriels  $F$  et  $G$  sont mis en dualité par la forme  $B$  si

- 1)  $B(x,y)$  est bilinéaire sur  $F \times G$
- 2)  $\forall_{Fx \neq 0} \exists_G y \quad B(x,y) \neq 0$
- 3)  $\forall_G y \neq 0 \exists_F x \quad B(x,y) \neq 0$

Exemples fondamentaux :

a) Considérons d'abord  $B(x,y) = \langle x, y \rangle = y(x)$  où  $G$  est l'espace vectoriel des formes linéaires sur  $F$ .

Nous retrouvons alors la notion de dual algébrique ( $F^*$ ). La propriété 2) est vraie puisque  $y$  est une forme linéaire et, pour vérifier 3), on s'appuie sur le fait que  $\forall x \neq 0$  dans  $E$ , on peut trouver une forme linéaire sur  $E$  telle que  $y(x) \neq 0$ .

b) Choisissons maintenant pour  $G$  l'espace des formes linéaires continues sur  $F$ , soit  $F'$ . On a donc  $F'$  est un sous espace de  $F^*$ . On prend aussi  $B(x,y) = \langle x, y \rangle$  qui vérifie de la même manière 2). Pour vérifier 3) il faut utiliser un corollaire du théorème de Hahn-Banach pour affirmer la non-trivialité de  $F'$  lorsque  $F$  est de dimension infinie. On retrouve ici le dual topologique  $F'$ .  
Remarque: il est clair que  $F'$  dépend de la topologie de  $F$ .

## §2 Topologies sur le dual topologique

Soit  $X$  un espace linéaire normé et  $Y = \mathbb{R}$  (ou  $\mathbb{C}$ ) (muni de la topologie usuelle)

$\mathcal{L}(X, Y) = \{ \text{formes linéaires continues de } X \text{ dans } Y \}$  peut donc être considéré comme le dual topologique  $X'$  de  $X$ .

La topologie de la convergence uniforme sur  $X'$  s'appelle la topologie forte de  $X'$

La topologie de la convergence simple dans  $Y$  est appelée la topologie faible\* de  $X'$

Rappelons aussi que  $X'$  muni de la norme  $\|f\| = \sup_{\|x\|=1} |f(x)|$  est un espace de Banach.

Considérons maintenant  $X''$  le bidual de  $X$ , i.e. le dual topologique de  $X'$  muni de la topologie forte

Une suite  $(x'_n)_{n \in \mathbb{N}} \in X'^{\mathbb{N}}$  converge faiblement vers  $x' \in X'$  ssi  $\forall x'' \in X'' \quad \langle x'_n, x'' \rangle \rightarrow \langle x', x'' \rangle$   
La topologie associée est la topologie faible de  $X'$

## §3 Prolongement de $X$ dans $X''$

Pour chaque  $x_0 \in X$ , on peut définir une forme linéaire continue  $f_0(x')$  sur  $X'$  par  $f_0(x') = \langle x_0, x' \rangle$ . Nous avons donc construit

$$J: X \rightarrow X''; \quad x_0 \mapsto Jx_0 = f_0$$

On peut alors démontrer que  $J$  est un isomorphisme de  $X$  sur  $JX \subseteq X''$

On dit qu'un espace de Banach est réflexif si  $X$  peut être identifié à  $X''$  (algébriquement et topologiquement) en utilisant  $J$ .

Il est immédiat, par le théorème de Riesz, que tout espace de Hilbert est réflexif.

De plus, comme remarqué dans Yosida,  $X''$  est un espace de Banach (cf supra) et donc tout espace réflexif normé linéaire est au moins un espace de Banach.

#### § 4

#### Remarque pour les espaces de Hilbert.

Par le théorème de Riesz, on peut identifier tout espace de Hilbert avec son dual. La forme de dualité définie sur  $X \times X'$  devient donc définie sur  $X \times X$  par  $(x, y)$  où  $(\cdot, \cdot)$  est le produit scalaire.

### III La dualité : opérateurs duaux.

#### § 1 Opérateurs duaux

Nous allons étendre la notion de "matrice transposée" à celle d'opérateur dual.

Théorème I.1

Soient  $X, Y$  deux espaces linéaires normés et  $X'$  et  $Y'$  leurs espaces duals topologiques forts.

Soit  $T : \text{dom}(T) \subset X \rightarrow Y$  linéaire.

Considérons les points  $(x', y') \in X' \times Y'$  qui satisfont

$$\langle Tx, y' \rangle = \langle x, x' \rangle \quad \forall x \in \text{dom}(T)$$

Alors  $x'$  est déterminé uniquement par  $y'$   
 ssi  $\overline{\text{dom}(T)} = X$

En effet, par la binarité du problème, il suffit de considérer l'équation

$$\langle x, x' \rangle = 0 \quad \forall x \in \text{dom}(T) \Rightarrow x' = 0$$

La condition suffisante découle alors de la continuité de  $x'$ .

Supposons que  $\overline{\text{dom}(T)} \neq X$ . Alors, par le théorème de Hahn-Banach, il existe un  $x_0 \neq 0$  tel que  $\langle x, x_0 \rangle = 0$  pour tout  $x \in \text{dom}(T)$ . Ce qui prouve alors la condition nécessaire. ■

Nous pouvons donc définir l'opérateur  $T'$  dual de  $T$

Si  $\overline{\text{dom}(T)} = X$ , il existe  $T' : \text{dom}(T') \subseteq Y' \rightarrow X'$  et qui associe à  $y'$  l'élément  $Ty' = x'$  suivant la condition du théorème 4 (de manière unique)

Il découle immédiatement de la binarité du crochet de dualité que  $T'$  est un opérateur linéaire tel que

$$\langle Tx, y' \rangle = \langle x, T'y' \rangle$$

pour tout  $x \in \text{dom}(T)$  et tout  $y' \in \text{dom}(T')$ , où  $\text{dom}(T')$  est l'ensemble des  $y'$  tel qu'il existe  $x'$  pour satisfaire la condition du th. 4.

Théorème I.2

Si  $\text{dom}(T) = X$  et  $T \in \mathcal{L}(X, Y)$  (i.e.  $T$  lin.  $C^0$ ) alors  $T' \in \mathcal{L}(Y', X')$

En effet, nous définissons l'opérateur dual  $T'$  comme plus haut

Soit  $B$  un ensemble borné de  $X$ . Puisque  $T$  est continue,  $T(B) = \{Tx \mid x \in B\}$  est un ensemble borné de  $Y$ . Mais  $\langle Tx, y' \rangle = \langle x, x' \rangle$ .  
 Donc, si  $y' \rightarrow 0$  dans  $Y'$ ,  $x' \rightarrow 0$  dans  $X'$ . Ce qui implique  
 la continuité de  $T'$  de  $Y'$  dans  $X'$ .

Exemple: Soient  $X = Y = \mathbb{R}^n$  normé par la norme  $\ell^2$ .

Si  $T \in \mathcal{L}(X, X)$ , soit  $Tx = y$  où  $x = (x_1, \dots, x_n)$  et  $y = (y_1, \dots, y_n)$

$$\text{Alors } y_i = \sum_{j=1}^n t_{ij} x_j$$

Donc

$$\langle Tx, z \rangle = \langle y, z \rangle = \sum y_j z_j = \sum_i \left( \sum_j t_{ij} x_j \right) z_i = \sum x_j \left( \sum_i t_{ij} z_i \right)$$

et donc  $\langle Tx, z \rangle = \langle x, T'z \rangle$  et  $T'$  est la matrice  $T$  transposée.  
 L'opérateur dual généralise donc bien la transposée.

De plus, dans le cas continu, on peut énoncer le

### Théorème I.3

Soient  $X$  et  $Y$  deux espaces linéaires normés.

Si  $T \in \mathcal{L}(X, Y)$  on sait que  $T' \in \mathcal{L}(Y', X')$

De plus

$$\|T\| = \|T'\|$$

Démonstration : Par la relation de définition  $\langle Tx, y' \rangle = \langle x, x' \rangle$ ,  
 on obtient  $\|T'y'\| = \|x'\| = \sup_{\|x\| \leq 1} |\langle x, x' \rangle| = \sup_{\|x\| \leq 1} |\langle Tx, y' \rangle|$

$$\leq \|y'\| \sup_{\|x\| \leq 1} \|Tx\| \leq \|y'\| \cdot \|T\|$$

et donc  $\|T'\| \leq \|T\|$ .

D'autre part,  $\forall x_0 \in X \exists f_0 \in Y' \quad \|f_0\|=1$  et  $f_0(Tx_0) = \langle Tx_0, f_0 \rangle$

$\|Tz_0\| = \|Tf_0\|$ . Donc  $f'_0 = Tf_0$  satisfait  $\langle z_0, f'_0 \rangle = \|Tz_0\|$  et donc

$$\|Tz_0\| = \langle z_0, Tf'_0 \rangle \leq \|T'\| \|f'_0\| \|z_0\| = \|T'\| \|z_0\|$$

ce qui implique que  $\|T\| \leq \|T'\|$

On peut de plus démontrer que

- 1) Si  $T$  et  $S$  appartiennent à  $\mathcal{L}(X, Y)$ , alors  $(\alpha T + \beta S)' = \alpha T' + \beta S'$
- 2) Soient  $T$  et  $S$  linéaires tels que  $\text{dom}(T), \text{dom}(S), R(T)$  et  $R(S)$  sont tous contenus dans  $X$ .  
Si  $S \in \mathcal{L}(X, X)$  et  $\overline{\text{dom}(T)} = X$ , alors  $(ST)' = T'S'$   
De plus, si  $\overline{\text{dom}(TS)} = X$ ,  $(TS)'$  est une extension de  $S'T'$

On trouvera une preuve de ce théorème dans Ref (Y-1) p 195

## §2 Opérateurs adjoints.

Plaçons nous maintenant dans le cas hilbertien. Soient  $X$  et  $Y$  deux espaces de Hilbert. Soit  $T$  un opérateur linéaire défini sur  $\text{dom}(T) \subseteq X$  et à valeurs dans  $Y$ . Supposons de plus que  $\text{dom}(T)$  est dense dans  $X$  et définissons alors  $T'$  le dual de  $T$  comme précédemment. Donc  $\langle Tx, y' \rangle = \langle x, T'y' \rangle$  pour  $x \in \text{dom}(T)$  et  $y' \in \text{dom}(T')$ . Appelons  $J_X^*$  la bijection (isomorphisme) de  $X$  dans son dual  $X'$ , qui est la conséquence du théorème de Riesz.

$$\text{Alors } \langle Tx, y' \rangle = y'(Tx) = (Tx, J_Y y')$$

$$\langle x, T'y' \rangle = (T'y')(x) = (x, J_X T'y')$$

$$\text{Comme } \langle Tx, y' \rangle = \langle x, T'y' \rangle, \text{ on obtient } (Tx, J_Y y') = (x, J_X T'y')$$

$$\text{C'est à dire } (Tx, y) = (x, J_X T' J_Y^{-1} y)$$

Dans le cas où  $X = Y$ , on obtient

$$T^* = J_X T' J_X^{-1}$$

qui est appelé l'opérateur adjoint de  $T$

Exactement comme dans le cas des opérateurs duals, on démontre

Théorème I 4

$T^*$  existe ssi  $\overline{\text{dom}(T)} = X$

Il est défini de la manière suivante :

sait  $y \in X$  et  $\exists y \in \text{dom}(T) \Leftrightarrow \exists y^* \in X \quad (T_x y) = (x, y^*) \quad \forall x \in \text{dom}(T)$

alors  $T^*y = y^*$

Théorème I 5

Si  $\text{dom}(T) = X$

alors  $T \in \mathcal{L}(X, X) \Rightarrow T^* \in \mathcal{L}(X, X)$

et  $\|T\| = \|T^*\|$

De nombreux, de plus amples renseignements peuvent être trouvés  
dans Ref (Y-1) p 196 sq.

## IV Optimisation avec contraintes : version globale.

Après avoir examiné et s'être ravis en minimaire la théorie élémentaire de la dualité, nous allons maintenant aborder un sujet qui a plus directement trait à l'optimisation au sens usuel du terme : la théorie des multiplicateurs de Lagrange.

Dans cette première partie, nous nous placrons dans l'espace des contraintes où nous interpréterons le multiplicateur de Lagrange comme un hyperplan. Dans la seconde partie, nous considérerons la théorie locale dans l'espace primal  $X$ .

### §1 Multiplicateurs de Lagrange.

Nous allons considérer, dans les quelques sections suivantes, le problème suivant :

$$\begin{cases} \text{minimiser } f(x) \\ \text{sous la contrainte } x \in \Omega, g(x) \leq \Theta \end{cases} \quad (1)$$

où  $\Omega$  est un sous ensemble convexe d'un espace linéaire  $X$ ,  $f$  est une forme convexe de  $\Omega$  dans  $\mathbb{R}$ , et où  $g$  est une application convexe de  $\Omega$  dans un espace linéaire normé  $Z$  possédant un cône positif  $P$  et  $\Theta$  comme origine.

Nous examinerons le problème (1) et développerons la théorie des multiplicateurs de Lagrange en considérant la classe de problèmes

$$\begin{cases} \text{minimiser } f(x) \\ \text{sous la contrainte } x \in \Omega, g(x) \leq z \end{cases}$$

où  $z$  est un vecteur arbitraire de  $Z$ . La solution de ce problème dépend évidemment de  $z$  (problème partiel)

Définissons

$$\Gamma = \{z \mid \exists_{\Omega} x \quad G(x) \leq z\}$$

Nous avons immédiatement que  $\Gamma$  est convexe. En effet  $z_1, z_2 \in \Gamma$  implique l'existence de  $x_1, x_2 \in \Omega$  tels que  $G(x_1) \leq z_1$  et  $G(x_2) \leq z_2$ . Donc, pour  $\alpha \in ]0,1[$ ,  $G(\alpha x_1 + (1-\alpha)x_2) \leq \alpha z_1 + (1-\alpha)z_2$  par la convexité de  $G$ , ce qui implique que  $\alpha z_1 + (1-\alpha)z_2 \in \Gamma$ .

Sur  $\Gamma$  définissons maintenant la fonctionnelle primaire  $w$

$$w(z) = \inf \{f(x) \mid x \in \Omega \wedge G(x) \leq z\}$$

Cette fonctionnelle peut prendre des valeurs infinies.

Le problème (1) peut être vu sous l'angle suivant : déterminer  $w(\Theta)$ .

On a maintenant le

Théorème I.6

w est convexe

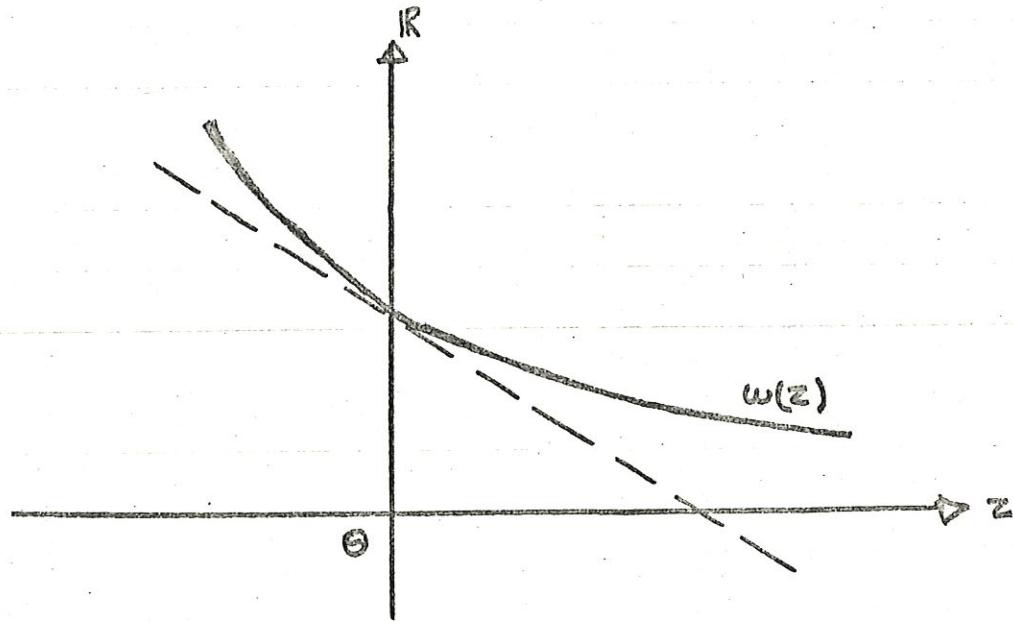
En effet, pour  $\alpha \in ]0,1[$

$$\begin{aligned} w(\alpha z_1 + (1-\alpha)z_2) &= \inf \{f(x) \mid x \in \Omega \wedge G(x) \leq \alpha z_1 + (1-\alpha)z_2\} \\ &\leq \inf \{f(x) \mid x = \alpha x_1 + (1-\alpha)x_2, x_1 \in \Omega, x_2 \in \Omega, \\ &\quad G(x_1) \leq z_1, G(x_2) \leq z_2\} \\ &\leq \alpha \inf \{f(x_1) \mid x_1 \in \Omega, G(x_1) \leq z_1\} \\ &\quad + (1-\alpha) \inf \{f(x_2) \mid x_2 \in \Omega, G(x_2) \leq z_2\} \\ &\leq \alpha w(z_1) + (1-\alpha)w(z_2) \quad \blacksquare \end{aligned}$$

Il est de plus immédiat que

w est décroissante

Typiquement, on peut représenter un  $w$  unidimensionnel comme suit



Conceptuellement, le théorème des multiplicateurs de Lagrange résulte du fait que si  $w$  est convexe, il existe un hyperplan tangent à  $w(z)$  en  $\theta$  et qui reste "en-dessous" de  $w(z)$  partout où  $w$  est définie. Si l'on considère l'hyperplan comme une nouvelle horizontale, on voit que  $w$  est minimisé en  $\theta$ , ce qui est équivalent à dire qu'en ajoutant à  $w(z)$  une fonctionnelle linéaire appropriée  $\langle z, z_0' \rangle$ , le résultat  $w(z) + \langle z, z_0' \rangle$  est minimum au point  $z = \theta$ . La fonctionnelle  $z_0'$  est le multiplicateur de Lagrange du problème.

### Théorème I.7

Soient  $X$  et  $Y$  respectivement un espace linéaire et un espace linéaire normé munis d'un cône positif  $P$ .

Si  $\Omega \subseteq X$  est convexe et si  $P \neq \emptyset$ ,

Soit  $f$  convexe de  $\Omega$  dans  $Y$  et  $g$  convexe de  $X$  dans  $Z$ .

Soit  $\exists x_0 \in \Omega \quad g(x_0) < 0$

Soit  $\mu_0 = \inf \{f(x) \mid \text{pour } x \in \Omega \text{ et } g(x) \leq 0\}$ . Alors

(L)

Alors, il existe  $z'_0 \geq 0$  dans  $Z'$  tel que

$$\mu_0 = \inf_{x \in \Omega} \{f(x) + \langle G(x), z'_0 \rangle\}$$

(3)

De plus, si l'infimum est atteint dans (2) en  $x_0 \in \Omega$ ,

$G(x_0) \leq 0$ , c'est aussi le cas dans (3) et

$$\langle G(x_0), z'_0 \rangle = 0$$

(4)

Remarque : Le cône positif de  $Z'$  est défini comme suit

$$P^+ = \{x' \in X' \mid \langle x, x' \rangle \geq 0 \quad \forall x \in P\}$$

Démonstration : Soit  $W = \mathbb{R} \times Z$ , où nous définissons

$$A = \{(r, z) \mid \exists x \in \Omega \quad r > f(x), z \geq G(x)\}$$

$$B = \{(r, z) \mid r \leq \mu_0, z \leq 0\}$$

Comme  $f$  et  $G$  sont convexes,  $A$  et  $B$  le sont aussi (D'ailleurs  $A$  représente la région convexe située au-dessus du graphe de la fonctionnelle primaire  $w$ ). La définition de  $\mu_0$  implique que  $A$  ne contient aucun point intérieur à  $B$ .

Soit  $N = -P$ . Par hypothèse, il contient un point dans son intérieur.

Donc l'intérieur de  $B$  n'est pas vide. On déduit alors du théorème de l'hyperplan séparant qu'il existe un élément non nul  $w'_0 = (r_0, z'_0)$  de  $W'$  tel que

$$\forall (r_1, z_1) \in A \quad \forall (r_2, z_2) \in B \quad r_0 r_1 + \langle z_1, z'_0 \rangle \geq r_0 r_2 + \langle z_2, z'_0 \rangle$$

De la nature de  $B$  on déduit immédiatement que  $w'_0 \geq 0$  i.e.  $r_0 \geq 0$  et  $z'_0 \geq 0$ . Mais on sait maintenant que  $r_0 > 0$ .

Le point  $(\mu_0, 0) \in B$ . Donc  $r_0 r + \langle z, z'_0 \rangle \geq r_0 \mu_0$  pour tout  $(r, z) \in A$ . Si  $r_0 = 0$ , on aurait en particulier  $\langle G(x_0), z'_0 \rangle \geq 0$  avec, de plus,  $z'_0 \neq 0$ . Cependant, puisque  $G(x_0) \in N$  et  $z'_0 \geq 0$ , on a  $\langle G(x_0), z'_0 \rangle < 0$ , ce qui est contradictoire.

Dès  $r_0 > 0$  et, sans perte de généralité, nous le prendrons égal à 1.

Puisque le point  $(\mu_0, \theta)$  est arbitrairement pris de A et de B, on obtient (avec  $r_0 = 1$ )

$$\mu_0 = \inf_{(x,z) \in A} [r + \langle z, z_0' \rangle] \leq \inf_{x \in \Omega} [\ell(x) + \langle b(x), z_0' \rangle]$$

$$\leq \inf_{x \in \Omega} \ell(x) = \mu_0$$

$$b(x) \leq 0$$

6. qui prouve la première partie du théorème.

De plus, s'il existe un  $x_0$  tel que  $b(x_0) \leq 0$ ,  $\mu_0 = \ell(x_0)$ , alors

$$\mu_0 \leq \ell(x_0) + \langle b(x_0), z_0' \rangle \leq \ell(x_0) = \mu_0$$

6. qui implique alors que  $\langle b(x_0), z_0' \rangle = 0$  ■

Donnons maintenant une version algébrique de ce théorème.

### Théorème I.7.1

Sous les hypothèses du th. 7.1, supposons que  $x_0$  fournit le minimum avec contrainte. Alors,  $\exists z' \geq 0$  tel que

$$L(x, z') = \ell(x) + \langle b(x), z' \rangle$$

possède un point de selle en  $x_0, z_0'$  i.e

$$\forall x, \forall z' \geq 0 \quad L(x_0, z') \leq L(x_0, z_0') \leq L(x, z')$$

En effet, choisissons  $z'$  comme dans le th. I.7.3) implique que  $L(x_0, z') \leq L(x, z')$

Mais (4) donne

$$L(x_0, z') - L(x_0, z'_0) = \langle G(x_0), z' \rangle - \langle G(x_0), z'_0 \rangle = \langle G(x_0), z' \rangle \leq 0$$

■

## §2 Sans les conditions de convexité.

Il est clair que sans les conditions de convexité et l'existence de points intérieurs, on ne peut généralement assurer l'existence d'un hyperplan séparant dans  $\mathbb{R} \times \mathbb{Z}$ . Cependant, si cet hyperplan existe tout de même, la technique des multiplicateurs de Lagrange est toujours d'application.

### Théorème I.8

Soit  $f: \Omega \rightarrow \mathbb{R}$  où  $\Omega \subseteq X$

Soit  $G: \Omega \rightarrow \mathbb{Z}$  borné, de cône positif  $P$

Si  $\exists z'_0 \in \mathbb{Z}', z'_0 \geq 0 \quad \exists x_0 \in \Omega$  tel que

$$\forall x \quad f(x_0) + \langle G(x_0), z'_0 \rangle \leq f(x) + \langle G(x), z'_0 \rangle$$

Alors  $x_0$  est la solution de

minimiser  $f(x)$

sous la contrainte  $G(x) \leq G(x_0)$ ,  $x \in \Omega$

Démonstration: Si  $\exists x_1 \in \Omega$  tel que  $f(x_1) < f(x_0)$  et  $G(x_1) \leq G(x_0)$ , et, comme  $z'_0 \geq 0$ , on déduit

$$\langle G(x_1), z'_0 \rangle \leq \langle G(x_0), z'_0 \rangle$$

et donc

$$f(x_1) + \langle G(x_1), z'_0 \rangle \leq f(x_0) + \langle G(x_0), z'_0 \rangle$$

ce qui contredit l'hypothèse

■

### Théorème I.9

Soient  $X, Z, \Omega, P, f$  et  $g$  comme supra.

Sait  $P$  fermé,  $L(x, z') = f(x) + \langle g(x), z' \rangle$

Si  $\exists z'_0 \in Z', z'_0 \geq 0$ ,  $\exists x_0 \in \Omega$  tel que

$$L(x_0, z') \leq L(x_0, z'_0) \leq L(x, z'_0)$$

pour  $\forall x \in \Omega$  et  $\forall z' \in Z', z' \geq 0$

Alors  $x_0$  est la solution de

$$\begin{cases} \text{minimiser } f(x) \\ \text{sous la contrainte } g(x) \leq 0, x \in \Omega \end{cases}$$

Démonstration : En particulierisant la condition de point de selle, on obtient

$$\langle g(x_0), z' \rangle \leq \langle g(x_0), z'_0 \rangle \quad \forall z' \in Z', z' \geq 0$$

Donc, pour  $z'_0 \geq 0$

$$\langle g(x_0), z'_0 + z'_1 \rangle \leq \langle g(x_0), z'_0 \rangle$$

ce qui implique  $\langle g(x_0), z'_1 \rangle \leq 0$

On peut en déduire que  $g(x_0) \leq 0$ .

La condition de point de selle implique alors  $\langle g(x_0), z'_0 \rangle = 0$

Supposons maintenant que  $x_0 \in \Omega$  et  $g(x_0) \leq 0$ . Alors,

$$f(x_0) = f(x_0) + \langle g(x_0), z'_0 \rangle \leq f(x_0) + \langle g(x_0), z'_0 \rangle \leq f(x_0)$$

Donc  $x_0$  est bien la solution cherchée ■

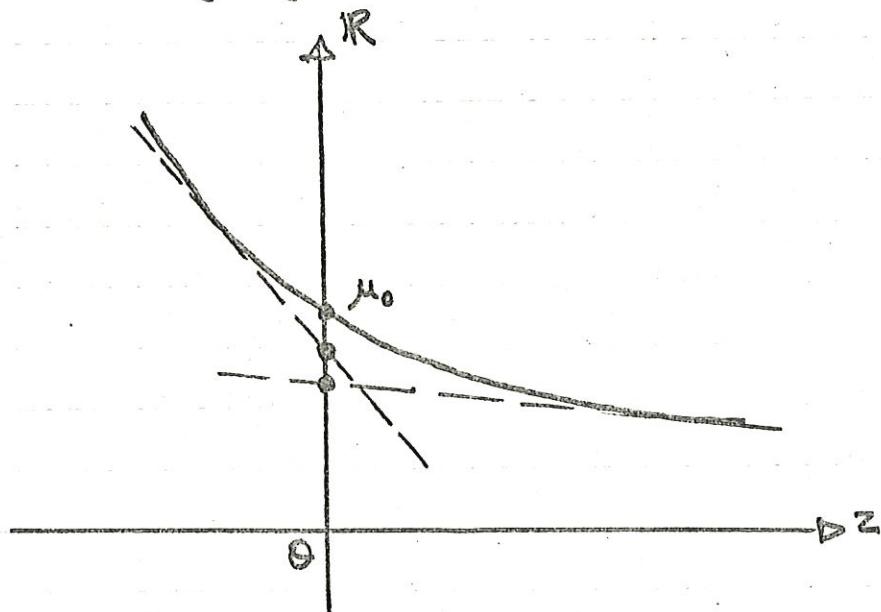
§3

### Dualité

Replongeons nous dans les hypothèses de convexité et considérons toujours

$$\begin{cases} \text{minimiser } f(x) \\ \text{sous la contrainte } g(x) \leq 0, x \in \Omega \end{cases}$$

Définissons comme supra  $w(z) = \inf \{ f(x) \mid g(x) \leq z, x \in \mathbb{L} \}$  et posons  $\mu_0 := w(0)$ . Le principe de dualité est basé sur le fait que  $\mu_0$  est l'intersection maximum de l'axe vertical avec les hyperplans qui sont en dessous de  $w(z)$ . Ce maximum est atteint pour l'hyperplan représentant le multiplicateur de Lagrange du problème.



Pour exprimer ce fait analytiquement, introduisons la fonctionnelle duale  $\varphi$  (correspondant à  $w$ ) définie sur le cône positif de  $\mathbb{Z}$  (soit  $P^+$ ).

$$\varphi(z') = \inf_{x \in \mathbb{L}} \{ f(x) + \langle g(x), z' \rangle \}$$

En général,  $\varphi$  n'est pas finie partout dans  $P^+$  mais la région où elle est finie est convexe. De plus, on voit aisément que  $\varphi$  elle-même est concave.

$$\text{En effet: } \varphi(\alpha z'_1 + (1-\alpha)z'_2) = \inf_x \{ f(x) + \langle g(x), \alpha z'_1 + (1-\alpha)z'_2 \rangle \}$$

$$\geq \inf_x \{ \alpha f(x) + \langle g(x), \alpha z'_1 \rangle + (1-\alpha) f(x) + \langle g(x), (1-\alpha)z'_2 \rangle \}$$

$$\geq \alpha \inf_x \{ f(x) + \langle g(x), z'_1 \rangle \} + (1-\alpha) \inf_x \{ f(x) + \langle g(x), z'_2 \rangle \}$$

Voyons de plus qui m'a le

Théorème I.10

$$\varphi(z') = \inf_{z \in \Gamma} \{ w(z) + \langle z, z' \rangle \}$$

où  $\Gamma$  est défini comme supra  $\Gamma = \{ z \mid \exists x \in \Omega \quad b(x) \leq z \}$ .

Démonstration Pour tout  $z' \geq 0$  et tout  $z \in \Gamma$ , on a

$$\begin{aligned} \varphi(z') &= \inf_{x \in \Omega} \{ f(x) + \langle b(x), z' \rangle \} \leq \inf_{x \in \Omega} \{ f(x) + \langle z, z' \rangle \mid b(x) \leq z, x \in \Omega \} \\ &= w(z) + \langle z, z' \rangle \end{aligned}$$

D'autre part, pour tout  $x \in \Omega$ , on obtient, en posant  $z_1 = b(x)$ ,

$$\begin{aligned} f(x) + \langle b(x), z' \rangle &\geq \inf \{ f(x) + \langle z_1, z' \rangle \mid b(x) \leq z_1, x \in \Omega \} \\ &= w(z) + \langle z_1, z' \rangle \end{aligned}$$

et donc  $\varphi(z') \geq \inf_{z \in \Gamma} \{ w(z) + \langle z, z' \rangle \}$

L'élément  $(1, z')$  de  $\mathbb{R} \times \mathbb{Z}'$  détermine une famille d'hyperplans dans  $\mathbb{R} \times \mathbb{Z}$ , chaque hyperplan étant formé des points  $(r, z)$  satisfaisant  $r + \langle z, z' \rangle = k$  où  $k$  est une constante. Le th. 10 dit que pour  $k = \varphi(z')$ , cet hyperplan supporte l'ensemble  $[w, \Gamma]$ , la région située au dessus du graphe de  $w$ . De plus, en  $z = 0$ , on obtient  $r = \varphi(z')$ ; donc  $\varphi(z')$  est égal à l'intersection de cet hyperplan avec l'axe vertical, ce qui nous ramène à la discussion géométrique faite plus haut.

Nous sommes à même maintenant, de prouver le théorème de dualité de Lagrange.

### Théorème I.11

Soit  $f : \Omega \subset X \rightarrow \mathbb{R}$  convexe, avec  $\Omega$  convexe

Soit  $g : X \rightarrow \mathbb{R}$  normé, avec  $g$  convexe

Si  $\exists x_0$ ,  $g(x_0) < 0$  et  $\mu_0 = \inf\{f(x) + g(x) \mid x \in \Omega\}$  est fini.

Alors

$$\inf_{\substack{g(x) < 0 \\ x \in \Omega}} f(x) = \max_{z' \geq 0} \varphi(z')$$
(5)

et le maximum à droite est atteint en  $z'_0 \geq 0$ .

Si, de plus, l'infimum à gauche est atteint en  $x_0 \in \Omega$ , alors

$$\langle g(x_0), z'_0 \rangle = 0$$

et  $x_0$  minimise  $f(x) + \langle g(x), z'_0 \rangle$  sur  $\Omega$

Démonstration: Pour tout  $z' \geq 0$ , on a

$$\inf_{\substack{x \in \Omega \\ g(x) < 0}} \{f(x) + \langle g(x), z' \rangle\} \leq \inf_{\substack{x \in \Omega \\ g(x) \leq 0}} \{f(x) + \langle g(x), z' \rangle\} \leq \inf_{x \in \Omega} f(x) = \mu_0$$

Donc, le membre de droite de (5) est  $\leq \mu_0$ .

Cependant, le théorème I.7 établit l'existence d'un  $z'_0$  qui donne l'égalité.

Le même théorème I.7 fournit d'ailleurs le reste de la preuve.  $\blacksquare$

## I Optimisation avec contraintes : version locale.

Examinons maintenant le point de vue local de l'optimisation avec contraintes.

Introduisons d'abord la différentielle de Fréchet.

Soient  $X$  et  $Y$  deux espaces linéaires normés.

Soit  $D$  ouvert de  $X$

Soit  $T : D \rightarrow Y$

Si  $\forall x \in D \quad \forall h \in X \quad \exists ST(x, h) \in Y$  linéaire et continue en  $h$  tel que

$$\lim_{\|h\| \rightarrow 0} \frac{\|T(x+h) - T(x) - ST(x, h)\|}{\|h\|} = 0$$

alors,  $T$  est Fréchet différentiable en  $x$  et  $ST(x, h)$  est la différentielle de Fréchet de  $T$  en  $x$  avec  $h$  comme incrément.

Cette différentielle peut encore s'écrire  $ST(x, h) = T'(x)h$  où  $T'(x)$  est la dérivée de  $T$  au sens de Fréchet.

Si  $T$  est une transformation continûment différentiable au sens de Fréchet de l'ouvert  $D \subseteq X$  Banach vers  $Y$  Banach, si  $x_0 \in D$  est tel que  $T'(x_0)$  est surjectif sur  $Y$ , alors le point  $x_0$  est appelé point régulier de  $T$ .

Nous énonçons ensuite le théorème de l'inverse généralisé, sans en donner de preuve :

## Théorème

Soit  $x_0$  un point régulier de la transformation  $T$  du Banach  $X$  dans le Banach  $Y$ . Alors, il existe un voisinage  $N(y_0)$  du point  $y_0 = T(x_0)$  et une constante  $K$  telle que l'équation  $T(x) = y$  admette une solution pour tout  $y \in N(y_0)$ , satisfaisant

$$\|x - x_0\| \leq K \|y - y_0\|$$

Une démonstration détaillée se trouve dans Ref (L-1) p240 sq.

### §1

## Contraintes sous forme d'égalités.

Notre but, dans ce paragraphe, est de développer des conditions nécessaires pour un extremum de  $f(x)$  sous les contraintes  $H(x) = \Theta$  où  $f: X \rightarrow \mathbb{R}$  avec  $X$  Banach et  $H: X \rightarrow Z$  avec  $Z$  Banach.

Commençons par le

Théorème I.12

Si  $f$  admet un minimum local (sous la contrainte  $H(x) = \Theta$ ) au point  $x_0$  et si  $f$  et  $H$  sont continûment différentiables au sens de Fréchet dans un ouvert contenant  $x_0$  et si, enfin,  $x_0$  est un point régulier de  $H$ , alors  $f'(x_0)h = 0$  pour tout  $h$  tel que  $H'(x_0)h = 0$ .

Démonstration : Considérons la transformation  $T: X \rightarrow \mathbb{R} \times Z$  définie comme suit:  $T(x) = (f(x), H(x))$ . Si il existe un  $h$  tel que  $H'(x_0)h = 0$  et  $f'(x_0)h \neq 0$ , alors  $T'(x_0) = (f'(x_0), H'(x_0)): \mathbb{R}^n \rightarrow \mathbb{R} \times Z$

est surjective car  $H'(x_0)$  l'est sur  $Z$  par hypothèse. Par le théorème de l'inverse quelconque, on pourrait déduire que  $\forall \varepsilon > 0 \exists x \exists \delta > 0$  tel que  $\|x - x_0\| < \delta$  et  $T(x) = (f(x) - \varepsilon, \Theta)$ ; ce qui contredit l'hypothèse que  $x_0$  minimise  $f(x)$  localement. ■

Nous pourrons maintenant parler à nouveau des multiplicateurs de Lagrange.

### Théorème I.1.3

Si la fonctionnelle  $f$  est continûment différentiable au sens de Fréchet admet un extremum local, sous la contrainte  $H(x)=\Theta$ , au point  $x_0$ , alors il existe  $z'_0 \in Z'$  tel que

$$L(x) = f(x) + \langle H(x), z'_0 \rangle$$

est stationnaire en  $x_0$ ; i.e.  $f'(x_0) + z'_0[H'(x_0)] = \Theta$

Démonstration: D'après le thI.12 il est clair que  $f'(x_0)$  est orthogonale au noyau de  $H'(x_0)$ . Comme l'image de  $H'(x_0)$  est fermée, on peut déduire  $f'(x_0) \in R[(H'(x_0))^\circ]$ .

Donc, il existe  $z'_0 \in Z'$  tel que  $f'(x_0) = -(H'(x_0))^\circ z'_0$ , ce qui est équivalent à  $f'(x_0) + \langle H'(x_0), z'_0 \rangle = 0$ . ■

Supposons maintenant que  $x_0$  n'est plus un point régulier.

### Théorème I.1.4

Dans les hypothèses du thI.13., sauf que  $R(H'(x_0))$  est seulement fermé dans  $Z$ , il existe un élément normal de  $R \times Z'$ , soit  $(r_0, z'_0)$  tel que  $r_0 f'(x_0) + z'_0[H'(x_0)]$  est stationnaire en  $x_0$ .

Démonstration: Si  $x_0$  est régulier, on choisit  $r_0 = 1$  et on applique le th 13. Si  $x_0$  n'est pas régulier, soit  $H = R(H'(x_0))$ . Il existe alors un point  $z \in Z$  tel que  $\inf_{m \in H} \|z - m\| > 0$  et, donc, on peut déduire

qu'il existe  $z'_0 \in H^{\perp}$ ,  $z'_0 \neq 0$ . Or  $H^{\perp} = R(H'(x_0))^{\perp} = cl^{\circ}(H'(x_0))$  (Rappel:  $H^{\perp} = \{z' \mid \forall x \in H \quad \langle x, z' \rangle = 0\}$ ). Il suffit alors de choisir  $(0, z'_0)$  ■

## §2 Contraintes sous forme d'inégalités.

Nous nous intéressons à présent au problème suivant

$$\begin{cases} \text{minimiser } f(x) \\ \text{sous la contrainte } G(x) \leq 0 \end{cases}$$

où  $f$  est définie sur  $X$  espace linéaire et  $G: X \rightarrow Z$ , espace linéaire normé de cône positif  $P$ .

Définissons aussi les points réguliers. Dans ce but, nous introduisons la différentielle de Gâteaux

Soit  $x \in \text{dom}(T) \subset X$ , et  $h \in X$  qgn.

Si  $\lim_{t \rightarrow 0} \frac{T(x+th) - T(x)}{t}$  existe, elle est appelée la différentielle de Gâteaux de  $T$  en  $x$  direction  $h$ .

Remarque:  $T$  est Fréchet-différentiable  $\Leftrightarrow T$  est Gâteaux-différentiable

Soit  $P \neq \emptyset$ ,  $G$  Gâteaux-différentiable telle que  $\delta_G G(x, h)$  est linéaire en  $h$ .  $x_0 \in X$  est un point régulier de  $G(x) \leq 0$  si  $G(x_0) \leq 0$  et  $\exists h \mid G(x_0) + \delta_G G(x_0, h) < 0$

Il est à noter que traduire  $H(x) = 0$  par  $P = \{0\}$  et  $-H(x) \leq 0$  devient impossible car il faut  $\overset{\circ}{P} \neq \emptyset$ . Les contraintes exprimées sous forme d'égalité ne peuvent donc être reprises dans ce cadre.

Prenons maintenant le théorème important

### Théorème 15 (Kuhn-Tucker généralisé)

Dans les hypothèses précédentes, supposons que  $f$  soit une fonction sur  $\mathbb{R}$  6-différentiable, et que  $G: X \rightarrow \mathbb{Z}$  soit 6-différentiable, et que  $\delta_{\mathbb{R}} f$  et  $\delta_G G$  soient linéaires en  $h$ . Supposons de plus que  $x_0$  minimise  $f$  sous la contrainte  $G(x) \leq 0$  et que  $x_0$  est un point régulier de  $G(x) \leq 0$ . Alors,  $\exists z'_0 \in \mathbb{Z}', z'_0 \geq 0$  tel que

$$L = f(x) + \langle G(x), z'_0 \rangle$$

soit stationnaire en  $x_0$ . De plus,  $\langle G(x_0), z'_0 \rangle = 0$ .

Démonstration Dans  $W = \mathbb{R} \times \mathbb{Z}$ , définissons

$$A = \{(r, z) \mid r \geq \delta_G G(x_0, h), z \geq G(x_0) + \delta_G G(x_0, h) \text{ pour un } h \in X\}$$

$$B = \{(r, z) \mid r \leq 0, z \leq 0\}$$

A et B sont évidemment convexes. Puisque  $\overset{\circ}{P} \neq \emptyset$ , alors  $\overset{\circ}{B} \neq \emptyset$ . A ne contient pas de point intérieur à B. En effet, si  $(r, z) \in A$  avec  $r < 0$  et  $z < 0$ , alors  $\exists h \in X$  tel que

$$\delta_{\mathbb{R}} f(x_0) < 0 \quad G(x_0) + \delta_G G(x_0, h) < 0$$

Le point  $G(x_0) + \delta_G G(x_0, h)$  est alors le centre d'un sphère ouverte contenue dans le cône négatif  $N$  de  $Z$ . Soit  $\rho$  le rayon de cette sphère. Alors, pour  $\alpha \in ]0, 1[$ , le point  $\alpha[G(x_0) + \delta_G G(x_0, h)]$  est le centre d'une sphère ouverte de rayon  $\alpha\rho$  contenue dans  $N$ ; donc le point  $(1-\alpha)G(x_0) + \alpha[G(x_0) + \delta_G G(x_0, h)] = G(x_0) + \alpha\delta_G G(x_0, h)$  s'y trouve aussi. Puisque pour  $h$  fixé

$$\|G(x_0 + ah) - G(x_0) - a\delta_G G(x_0, h)\| = o(a)$$

Il est donc clair que  $b(x_0 + ah) < 0$  pour  $a$  suffisamment petit.

On peut reproduire un argument similaire pour prouver que  $f(x_0 + ah) < f(x_0)$  pour  $a$  petit, ce qui contredit l'optimalité de  $x_0$ . Donc  $A \cap B = \emptyset$ .

Employons à nouveau le théorème de l'hyperplan séparant A et B.

On déduit  $\exists r_0, z_0', \delta$  tels que

$$\forall (r, z) \in A \quad r_0 r + \langle z, z_0' \rangle \geq \delta$$

$$\forall (r, z) \in B \quad r_0 r + \langle z, z_0' \rangle \leq \delta$$

Comme  $(0, 0) \in A \cap B$ , on a  $\delta = 0$ .

De la nature de B, on déduit que  $r_0 z_0 \leq 0$  et  $z_0' \geq 0$ . De plus l'hyperplan ne peut être vertical puisque  $\exists h \quad b(x_0) + \delta_b b(x_0, h) < 0$ . Donc, posons  $r_0 = 1$ .

On dérive aussi de la propriété de séparation

$$\forall h \in X \quad \delta_b f(x_0, h) + \langle b(x_0) + \delta_b b(x_0, h), z_0' \rangle \geq 0$$

En choisissant  $h = 0$ , on obtient  $\langle b(x_0), z_0' \rangle \geq 0$ . Mais comme  $b(x_0) \leq 0$  et  $z_0' \geq 0$ , on a aussi  $\langle b(x_0), z_0' \rangle \leq 0$ . Donc  $\langle b(x_0), z_0' \rangle = 0$ .

Reste donc

$$\forall h \in X \quad \delta_b f(x_0, h) + \langle \delta_b b(x_0, h), z_0' \rangle \geq 0$$

Comme les b-différentielles sont linéaires en h, il est nécessaire que  $\delta_b f(x_0, h) + \langle \delta_b b(x_0, h), z_0' \rangle = 0$

### III Le principe du maximum de Pontryagin.

Dans ce paragraphe, nous allons aborder un théorème qui est à la base de bien des applications, spécialement en contrôle optimal, comme nous le verrons plus tard.

Nous développerons la théorie d'un point de vue abstrait, en expliquant les opérateurs duals.

Posons  $X$  et  $U$  des espaces linéaires normés et  $g[x, u]$  une fonction (de coût) sur  $X \times U \rightarrow \mathbb{R}$ . Considérons de plus une contrainte de la forme  $A[x, u] = \Theta$  où  $A : X \times U \rightarrow X$ . Cet opérateur  $A$  peut représenter un ensemble d'équations différentielles, intégrales, aux différences, aux dérivées partielles, etc.

Supposons que  $A[x, u] = \Theta$  définit une fonction  $x(u)$  implicite unique. Supposons, de plus, que  $A$  et  $g$  sont différentiables au sens de Fréchet par rapport à  $x$  et que  $A_x[x, u]$  et  $g_x[x, u]$  sont continues sur  $X \times U$ . Supposons enfin que

$$\|x(u) - x(v)\| \leq K \|u - v\|$$

Notre problème est de minimiser  $J = g[x, u]$  sous la contrainte  $A[x, u] = \Theta$  avec  $u \in \Omega \subseteq U$ . Il est clair que  $J = g[x, u] = g[x(u), u] = J[u]$ .

Introduisons maintenant le Lagrangien

$$\begin{aligned} L[x, u, \lambda'] &= \lambda' A[x, u] + g[x, u] \\ &= \langle A[x, u], \lambda' \rangle + g[x, u] \end{aligned}$$

pour  $x \in X, u \in U, \lambda' \in X'$ . Nous pouvons alors démontrer le

### Théorème I.6

Pour tout  $u \in \Omega$ , soit  $\lambda'$  une solution de l'équation

$$\lambda' A_x[x, u] + g_x[x, u] = \Theta$$

Alors, pour  $v \in \Omega$ ,

$$J(u) - J(v) = L[x(u), u, \lambda'] - L[x(v), v, \lambda'] + o(\|u - v\|)$$

Démonstration: Par définition

$$\begin{aligned} J(u) - J(v) &= g[x(u), u] - g[x(v), v] \\ &= g[x(u), u] - g[x(u), v] + g[x(u), v] - g[x(v), v] \\ &= g[x(u), u] - g[x(u), v] + g_x[x(u), u] \cdot [x(u) - x(v)] \\ &\quad + (g_x[x(v), v] - g_x[x(u), u]) [x(u) - x(v)] + o(\|x(u) - x(v)\|) \end{aligned}$$

$$J[u] - J[v] = g[x(u), u] - g[x(u), v] + g_x[x(u), u][x(u) - x(v)] \\ + o(\|u - v\|)$$

par la continuité de  $g_x$  et la condition  $\|x(u) - x(v)\| \leq k \|u - v\|$ .

De même

$$\|A[x(u), u] - A[x(u), v] - A_x[x(u), u][x(u) - x(v)]\| = o(\|v - u\|)$$

Donc

$$\begin{aligned} L[x(u), u, \lambda] - L[x(v), v, \lambda] + o(\|u - v\|) &= \lambda' A[x(u), u] + g[x(u), u] - \lambda' A[x(v), v] \\ &\quad - g[x(v), v] + o(\|v - u\|) \\ &= \lambda' A_x[x(u), u][x(u) - x(v)] + g[x(u), u] - g[x(v), v] + o(\|u - v\|) \\ &= g[x(u), u] - g[x(u), v] + g_x[x(u), u][x(u) - x(v)] + o(\|v - u\|) \\ &= J[u] - J[v] \end{aligned}$$

Appliquons maintenant ce résultat au système d'équations différentielles ordinaires de la forme

$$\dot{x}(t) = f(x(t), u(t)) \quad x(t_0) = x_0$$

et à la fonction de coût

$$J = \int_{t_0}^t l(x(t), u(t)) dt$$

où nous supposons que les fonctions  $f$  et  $l$  sont continûment différentiables par rapport à  $x$ , que  $f$  satisfait

$$\|f(x, u) - f(y, v)\| \leq M(\|x - y\| + \|u - v\|)$$

sachant que  $f, g$  et  $u$  prennent leurs valeurs dans l'espace  $\mathbb{R}^n$  et  $\mathbb{R}^m$  respectivement. En fait nous restreignons  $u(t) \in \Omega \subseteq \mathbb{R}^m$  et  $u(t)$  est différentiable par morceaux sur  $[t_0, t_1]$ .

Dans ces conditions, nous pouvons aborder le principe de Pontryagin.

### Théorème I.17

Soyons  $x_0$  et  $u_0$  optimaux pour le problème

$$\min J = \min \int_{t_0}^t l(x, u) dt$$

sous la contrainte  $\dot{x}(t) = f(x(t), u(t))$  avec  $x(t_0)$  donné,  $u(t) \in \mathcal{L}$ .

Soit  $\lambda$  la solution de l'équation

$$-\dot{\lambda}(t) = f_x^\dagger \lambda(t) + l_x^\dagger \quad \lambda(t_1) = 0$$

où les dérivées partielles sont évaluées le long de la trajectoire optimale, et définissons le Hamiltonien

$$H[x, u, \lambda, t] = \lambda^\dagger(t) f(x, u) + l(x, u).$$

Alors, pour tout  $t \in [t_0, t_1]$

$$H(x_0(t), u_0(t), \lambda(t), t) \leq H(x_0(t), u(t), \lambda(t), t)$$

pour tout  $u$  à valeurs dans  $\mathcal{L}$ .

Démonstration : Pour une question de facilité, nous supposons  $m = 1$ , i.e. les contrôles sont des fonctions scalaires.

Prenons  $X = C^1[t_0, t_1]$  et  $U = \{f \in C^1 \text{ continues par morceaux sur } [t_0, t_1]\}$  muni de la norme  $L^1$

$$\text{Prenons } A[x, u] = x(t) - x(t_0) - \int_{t_0}^t f(x(\tau), u(\tau)) d\tau$$

$$\text{et } g[x, u] = \int_{t_0}^t l(x, u) d\tau$$

On peut voir que  $A$  et  $g$  sont continûment différentiables au sens de Fréchet par rapport à  $x$ .

Si  $x, x + \delta x$  correspondent à  $u, u + \delta u$  dans  $\{(x, u) \mid A[x, u] = 0\}$ , alors

$$\|\delta g(t)\|_{L^1} \leq \int_{t_0}^t |H| \|\delta x(\tau)\|_{L^1} + |\delta u(\tau)| d\tau$$

$\alpha$  qui fournit [ par le théorème de Gronwall ] ,

$$\|\delta u(t)\|_{L^2} \leq M e^{M(t-t_0)} \int_{t_0}^t |\delta u(\tau)| d\tau$$

et donc  $\|\delta u\| \leq K \|f\|$  et la transformation A satisfait la condition du théorème précédent.

De plus on a bien que  $X \subseteq L^2$  qui est Hilbert . Donc  $\lambda' \mapsto \lambda$  l'élément de  $X$  associé à  $\lambda'$  dans  $X'$  par le théorème de Riesz . La forme de dualité se réduit alors au produit scalaire .

On a alors que  $-\dot{\lambda}(t) = f_x^+ \lambda + \ell_x^+$  est bien équivalent , après intégration , à  $\lambda'' A_x[x(u), u] + g_x[x(u), u] = \Theta$  , pour  $t \in [t_0, t]$  . De même  $\int_{t_0}^t H[x, u, \lambda] dt$  est identique au Lagrangien  $L[x, u, \lambda']$  à part un terme  $\int_{t_0}^t x(t) \dot{\lambda}(t) dt$  qui n'est pas important puisqu'il ne dépend pas de  $u$  explicitement .

Le théorème précédent nous donne

$$J[u_0] - J[u] = \int_{t_0}^t [H[x_0, u_0, \lambda] - H[x_0, u, \lambda]] dt + o(\|u - u_0\|)$$

Voyons maintenant que cette équation implique la minimisation du Hamiltonien . Supposons l'inverse : il existe un  $\tilde{t} \in [t_0, t]$  et  $\tilde{u}(\tilde{t}) \in \Omega$  tel que

$$H[x(\tilde{t}), u_0(\tilde{t}), \lambda(\tilde{t})] > H[x(\tilde{t}), \tilde{u}(\tilde{t}), \lambda(\tilde{t})]$$

Comme  $u$  est continue par morceaux et que  $x$  est continue , ainsi que  $\lambda$  ,  $f$  et  $\ell$  , il existe un intervalle  $[t', t'']$  qui contient  $\tilde{t}$  et il existe un  $\epsilon > 0$  tel que

$$H[x(t'), u_0(t'), \lambda(t')] - H[x(t), \tilde{u}(t), \lambda(t)] > \epsilon$$

pour tout  $t \in [t', t'']$  .

Choisissons  $u(t)$  la fonction continue par morceaux égale à  $u_0(t)$  en dehors de  $[t', t'']$  et à  $\tilde{u}(t)$  sur  $[t', t'']$  . Nous obtenons alors

$$J[u_0] - J[u] > \epsilon (t'' - t') + o(\|u - u_0\|).$$

Or  $t'' - t'$  est tel que  $\|u - u_0\| \approx 0(t'' - t')$ . Donc, en choisissant  $[t', t'']$  suffisamment petit, on peut rendre  $J[u_0] - J[u]$  positif, ce qui contredit l'optimalité de  $u_0$ . ■

### Remarque

Il est clair que l'équation  $-\dot{\lambda}(t) = \int_a^t \lambda + \dot{\lambda} dt$  équivaut à

$$-\frac{\partial H}{\partial x} = \dot{\lambda}$$

ce qui fournit parfois un moyen de calculer  $\lambda$  effectivement.

## Chapitre II

### Les systèmes dynamiques contrôlés.

#### I Introduction

Comme annoncé, ce chapitre sera consacré principalement à la description des systèmes dynamiques contrôlés, ou systèmes dynamique à contrôle. D'abord d'un point de vue intuitif et pratique, ensuite d'un point de vue plus formel, nous tenterons de définir cette classe de systèmes si fréquents dans les applications, en même temps que théoriquement intéressants pour eux-mêmes.

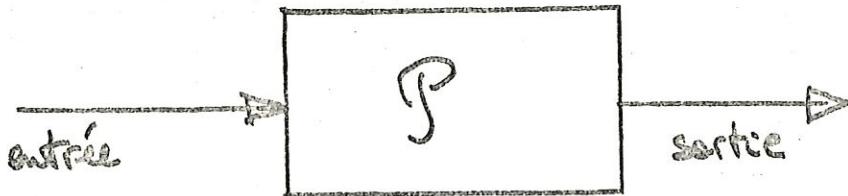
Les outils développés dans les pages précédentes ne nous serviront pas directement dans cette description, mais plutôt dans le paragraphe consacré aux méthodes de résolution et surtout dans le chapitre plus technique, traitant de la méthode de Ritz Galerkin.

Mais abordons d'abord la description.

## II Les systèmes dynamiques contrôlés : Vision intuitive

La connaissance du monde physique est basée sur l'expérience et l'abstraction. Tous les scientifiques connaissent la signification pratique des mots: induction et déduction. D'autre part, autant que la synthèse, l'analyse est un élément indispensable de cette connaissance. Le scientifique, en face d'un phénomène, tente d'abord de le cerner pour mieux le comprendre. C'est ce que nous allons faire.

Commençons par nous donner un système physique  $P$ , une de ces "boîtes noires" proverbiales, auquel nous pouvons fournir des "signaux" d'entrée et duquel nous pouvons observer le comportement par le biais de "signaux" de sortie.



Notre objectif final est la détermination d'un signal d'entrée, ou "input", qui produira, à travers le système  $P$ , un signal de sortie, ou "output", possédant certaines caractéristiques que nous nous donnons d'avance, tentant ainsi de minimiser "le coût" de notre opération, dans la mesure où le signal de sortie est conforme à nos désirs.

Dans les pages qui vont suivre, nous nous intéresserons à une classe de systèmes pour lesquels on peut construire des modèles mathématiques adéquats, appelés systèmes dynamiques.

Effectuons d'abord une petite expérience mentale. Partant d'un "moment initial", nous appliquons un signal à l'entrée de notre système  $P$ , et cela jusqu'à un instant futur. Pendant ce laps de temps, que nous appellerons intervalle d'observation, nous observons la sortie du système. Une question naturelle se pose alors : "Partant du même instant initial, nous appliquons à  $P'$ , système en tous points identique à  $P$ , un signal d'entrée semblable à celui appliqué à  $P$ . Observons nous, pour  $P'$ , les mêmes signaux de sortie que pour  $P$  ?"

L'expérience nous enseigne qu'il faut, en plus, tenir compte de l'état du système ( $P$  ou  $P'$ ) au moment initial. L'output du système dépend donc de son état initial.

Nous venons d'introduire, de façon naturelle, le concept supplémentaire d'"état du système". En fait, dans la description du système dynamique, nous utiliserons les concepts d'input, d'output et d'état. Formalisons un peu : appelons  $u(t)$  le signal d'input en fonction du temps; de même appelons  $y(t)$  le signal d'output en fonction du temps et  $\alpha(t)$  l'état du système en fonction du temps.

Au temps initial  $t_0$ , nous commençons donc d'introduire  $u(t_0, t) = u(t)$  pour tout  $t$  dans  $[t_0, t_f]$  où  $t_f$  est l'instant final, et nous observons  $y(t_0, t)$  sur le même intervalle.

De plus, nous allons faire quelques hypothèses sur le système lui-même.

1. Nous supposons que le système est entièrement déterministe.
2. Nous supposons qu'il n'est pas anticipatif, i.e. que l'input, l'output et l'état à l'instant  $t$  ne dépendent pas de leurs valeurs respectives aux instants postérieurs.

Precisons un peu plus encore notre modèle : donnons nous des relations qui unissent l'output à l'input, via l'état et l'état initial.

$$y(t) = g[x(t_0), u(t_0, t), (t_0, t)]$$

$$x(t) = f[x(t_0), u(t_0, t), (t_0, t)]$$

où la première équation est appelée équation d'output et la seconde équation d'état.

A présent, passons à une définition plus formelle de ces concepts.

### III Les systèmes dynamiques contrôlés : définition formelle.

Soit  $T \subseteq \mathbb{R}$ ,  $(\Sigma, d)$  un espace métrique, de même de  $(\Omega, \hat{d})$ , soit  $U$  un ensemble de fonctions continues par morceaux  $T \rightarrow \Omega$ .

Soit  $x(t) : T \rightarrow \Sigma$ , soit  $g : \Sigma \times \Omega \times T \rightarrow \mathbb{R}^P$ .

Si  $t_0 \in T$  et  $t_1 \in T$  et  $t_0 \leq t_1$ , soit  $[t_0, t_1] = \{t \in T \mid t_0 < t \leq t_1\}$

Nous utiliserons la notation  $u_{[t_0, t_1]}$  pour la restriction de  $u \in U$  à l'intervalle  $[t_0, t_1]$ .

Si  $t \in T$ ,  $u \in U$ ,  $x(t) \in \Sigma$  alors  $g[x(t), u(t), t]$  est un élément de  $\mathbb{R}^P$  bien défini que nous notons  $y(t)$ , i.e.

$$y(t) = g[x(t), u(t), t]$$

Nous posons aussi

$$y_{[t_0, t_1]} = \hat{g}[x(t_0); u_{[t_0, t_1]}]$$

Après ces quelques précisions, nous sommes maintenant en mesure d'envisager les axiomes qui définissent notre système avec exactitude.

### Axiome 1

Pour  $x(t_0) \in \Sigma$ , pour tout  $t_0 \in T$ , pour tout  $t \in T$  avec  $t \geq t_0$ , et pour tout  $u \in J_{[t_0,t]}$  avec  $u \in U$ , la connaissance de  $x(t_0)$  et  $u \in J_{[t_0,t]}$  détermine  $y \in J_{[t_0,t]}$  de manière unique.

En particulier, si  $u$  et  $v$  appartiennent à  $U$  et satisfont

$$\text{alors } \begin{aligned} u \in J_{[t_0,t]} &= v \in J_{[t_0,t]} \\ \hat{g}[x(t_0); u \in J_{[t_0,t]}] &= \hat{g}[x(t_0); v \in J_{[t_0,t]}] \end{aligned}$$

### Axiome 2

Si  $t_0 < \hat{t} < t$  sont des éléments de  $T$ , si  $x(t_0) \in \Sigma$ , si  $\sum[x(t_0), u, \hat{t}]$  représente l'ensemble des  $x(\hat{t})$  de  $\Sigma$  qui satisfont l'équation

$y \in J_{[\hat{t},t]} = \hat{g}[x(t_0); u \in J_{[t_0,\hat{t}]}] = \hat{g}[x(\hat{t}), u \in J_{[\hat{t},t]}]$   
et si  $u^*$  est un élément élément de  $U$ , alors l'intersection des ensembles  $\sum[x(t_0), u, \hat{t}]$  où  $u \in U$  avec  $u \in J_{[t_0,\hat{t}]} = u^* \in J_{[t_0,\hat{t}]}$  n'est pas vide, i.e.

$$\bigcap_{u \in U} \sum[x(t_0), u, \hat{t}] \neq \emptyset$$

$u \in J_{[t_0,\hat{t}]} = u^* \in J_{[t_0,\hat{t}]}$

En particulier, aucun  $\sum[x(t_0), u, \hat{t}]$  n'est vide. L'axiome assure l'existence d'au moins un élément de  $\Sigma$  pour toute paire  $(u \in J_{[\hat{t},t]}, y \in J_{[\hat{t},t]})$ .

On peut alors montrer que ces deux axiomes impliquent l'existence d'une fonction  $\phi[t, u \circ_{t_0, t}], x(t_0)]$  telle que

$$x(t) = \phi[t, u \circ_{t_0, t}], x(t_0)]$$

(cf Référence (Z-1))

### Axiome 3

Les fonctions  $g$ ,  $\hat{g}$  et  $\phi$  sont continues en toutes leurs variables.  
Pour ce qui est de la dépendance de  $g$  et de  $\phi$  en fonction du  $u \circ_{t_0, t}$ , cela signifie que si  $u, v \in U$  et vérifient

$$d_u(u \circ_{t_0, t}, v \circ_{t_0, t}) = \sup_{\tau \in [t_0, t] \cap T} \{d(u(\tau), v(\tau))\} \text{ petit}$$

$$\text{alors } d\{\phi[t; u \circ_{t_0, t}], x(t_0)], \phi[t; v \circ_{t_0, t}], x(t_0)]\}$$

$$\text{et } \sup_{\tau \in [t_0, t] \cap T} \{ \| \hat{g}[x(t_0); u \circ_{t_0, t}] - \hat{g}[x(t_0); v \circ_{t_0, t}] \| \}$$

sont aussi petits, et où  $\hat{g}[x(t_0); u \circ_{t_0, t}] = g[x(t_0), u(t+), t_0]$   
et  $\hat{g}[x(t_0); v \circ_{t_0, t}] = g[x(t_0), v(t+), t_0]$  par définition.

### Axiome 4

La fonction  $\phi$  satisfait les conditions suivantes :

- 1.) Pour tout  $t, t_0 \in T, u \in U$  et  $x(t_0) \in \Sigma$ ,  $\phi[t_0, u \circ_{t_0, t}], x(t_0)] = x(t_0)$   
au sens  $\varprojlim_{t \rightarrow t_0} \phi[t, u \circ_{t_0, t}], x(t_0)] = x(t_0)$

- 2.) Pour tout  $t_0 < \hat{t} < t \in T, u \in U$  et  $x(t_0) \in \Sigma$ ,

$\phi[t; u]_{t_0, t} x(t_0) = \phi[t, u]_{t_0, t}^1, \phi[\hat{t}; u]_{t_0, \hat{t}} x(t_0)]]$   
 où  $t_0 < \hat{t} \leq t$ . Cette propriété est appellée la propriété de transition  
 ou de semi-groupe.

- 3) Pour tout  $\tau, t_0, t$  avec  $\tau \in [t_0, t] \cap T$ , et pour tout  $x(t_0) \in \Sigma$ ,  
 si  $u, v \in U$  avec  $u|_{[t_0, t]} = v|_{[t_0, t]}$ , alors
- $$\phi[\tau, u]_{t_0, t} x(t_0) = \phi[\tau, v]_{t_0, t} x(t_0)]$$

On peut maintenant poser la définition :

Un système dynamique contrôlé, ou système dynamique S est  
 la donnée de  $T, \Sigma, \Omega$  et  $U$ , d'une variable  $x(t)$  et d'une  
 fonction  $g$  tels que les axiomes 1 à 4 soient vérifiés. Dans  
 ce cas,  $T$  est appellé le domaine du système ;  $\Sigma$  l'espace d'état  
 du système,  $U$  l'espace d'input du système et  $x(t)$  la variable  
 d'état ;  $u|_{[t_0, t]}$  est l'input pendant l'intervalle d'observa-  
 tion  $[t_0, t]$  et  $y|_{[t_0, t]} = g[x(t_0); u|_{[t_0, t]}]$  est l'output du système.  
 La fonction  $\phi[t, u|_{[t_0, t]}, x(t_0)]$  est la fonction de transition du  
 système. L'ensemble  $\{x(\tau) | x(\tau) = \phi[\tau, u|_{[t_0, \tau]}, x(t_0)]\}$  pour  $\tau \in [t_0, t] \cap T$   
 est la trajectoire du système pendant  $[t_0, t] \cap T$ , au départ  
 de l'état initial  $x(t_0)$  et générée par le contrôle (input)  $u|_{[t_0, t]}$ .  
 Enfin les équations

$$y(t) = g[x(t), u(t), t]$$

$$x(t) = \phi[t, u|_{[t_0, t]}, x(t_0)]$$

sont respectivement les équations de sortie et d'état du système.

Donnons maintenant une brève interprétation des axiomes.

Axiome 1. Cet axiome dit principalement que si on connaît  $x(t_0)$  et  $t_0$ , et si l'on applique un contrôle donné pendant  $[t_0, t]$  avec  $t > t_0$ , l'output du système est déterminé de manière unique. Pour prévoir la sortie pendant  $[t_0, t]$ , il n'est pas nécessaire de connaître l'input avant  $t_0$ .

$x(t_0)$  et  $u[t_0, t]$  suffisent. On note aussi que les valeurs futures du contrôle n'affectent pas  $y[t_0, t]$ , ce qui implique que le système est non anticipatif.

Axiome 2. L'axiome 2 assure l'existence d'"assez" d'états possibles du système pour pourvoir tout couple input, output. Il implique aussi que la connaissance de  $x(t_0)$  et de  $u[t_0, t]$  suffit à déterminer, non seulement  $y(t)$  pour  $t \in [t_0, t]$ , mais aussi  $x(\hat{t})$ ,  $t_0 < \hat{t} < t$ . L'état concerne, d'une certaine façon, toute l'information du passé requise pour prédire les états futurs de l'output et de  $x(t)$ .

Axiome 3 C'est un axiome de continuité qui concerne la "stabilité" du système pour des modifications de  $x(t_0)$  ou  $u[t_0, t]$  au plus de la continuité de la trajectoire.

Axiome 4 Il faut que 1)  $x(t_0)$  soit le point de départ de la trajectoire, 2) si un contrôle mène  $x(t)$  de  $x(t_0)$  à  $x(t)$  en passant par  $x(\hat{t})$ , ce même contrôle mène  $x(\hat{t})$  de  $x(t_0)$  en  $x(\hat{t})$ , 3) Le système soit non anticipatif.

Nous voyons donc qu'il s'agit bien de la traduction formelle des idées exprimées dans le chapitre précédent, qui correspondent, en gros, à une description "réaliste" de systèmes physiques.

Nous pouvons, après ces généralités, aborder maintenant la description des systèmes qui seront étudiés plus en détail dans le reste de ce travail.

## IV Les systèmes dynamiques contrôlés: cas particuliers.

### 1§ La fonctionnelle de coût

Maintenant que nous disposons d'un système dont nous pouvons influencer l'output par un contrôle, revenons à notre expérience du début. Comment produire, à partir d'un contrôle bien choisi, un output possédant cette qualité supplémentaire d'avoir des caractéristiques que nous déterminons d'avance ?

En d'autres termes, si nous assignons un coût déterminé à une paire input, output, quel input choisir pour minimiser ce coût. Représentons le coût sous la forme suivante :

$$J[u, y] = f(y(t)) + \int_{t_0}^T g(y, u, t) dt$$

Le terme  $f(y(t))$  s'introduit de manière naturelle. Il est en effet possible que le coût associé à la trajectoire de l'output soit influencé par la valeur de cet output à la fin du processus.

Désormais, notre problème sera de minimiser  $J[u, y]$  en travaillant sur un système dynamique contrôlé  $S$  donné.

$$y(t) = g[x(t), u(t), t]$$

$$x(t) = \phi[t, u, t_0, x(t_0)]$$

$$\min_{u \in U} J[u, y] = f[y(T)] + \int_{t_0}^T g_0[y(t), u(t), t] dt$$

28 Système différentiel non linéaire.

Nous considérons ici le problème suivant:

$$[t_0, t] \cap T = [0, T] \subset \mathbb{R}$$

$$G = \{W_2^\alpha[0, T], \mathbb{R}^n\} \quad \Sigma = \mathbb{R}^k$$

$$U = \{W_2^\alpha[0, T], \mathbb{R}^r\} \quad \Omega = \mathbb{R}^r$$

où  $G$  représente l'ensemble de  $x(t)$  et où  $\{W_2^\alpha[0, T], \mathbb{R}^k\}$  est l'espace de Sobolev de toutes les fonctions définies sur  $[0, T]$  et à valeurs dans  $\mathbb{R}^k$  telle qu'elles soient du carré intégrable sur  $[0, T]$  ainsi que toutes leurs dérivées au sens des distributions, jusqu'à l'ordre  $\alpha$ . Ce espace est muni de la norme

$$\|z\|_{2,\alpha}^2 = \sum_{i=1}^k \int_0^T \sum_{j=0}^{\alpha} |z_i^{(j)}(t)|^2 dt$$

les équations de ce système sont

$$\dot{x}(t) = f(x, u, t)$$

$$y(t) = x(t)$$

$$\text{avec } x(0) = x_0$$

et le coût associé  $J[u] = \int_0^T g(x, u, t) dt$

où  $f: G \times U \times [0, T] \rightarrow \mathbb{R}^n$

$g: G \times U \times [0, T] \rightarrow \mathbb{R}$

Il est clair que notre système est maintenant bien défini.

Nous étudierons son comportement dans les pages qui suivent, lorsqu'il sera question de la méthode de Ritz-Galerkin.

### 35 Le régulateur d'état différentiel linéaire

Choisissons cette fois

$$[t_0, t] \cap T = [0, T] \subset \mathbb{R} \quad \Sigma = \mathbb{R}^n$$

$G = \Phi^n = \{ \text{fonctions } [0, T] \rightarrow \mathbb{R}^n, \text{ continues par morceaux et à dérivée bornée} \}$

$$U = \Phi^r \quad L = \mathbb{R}^r$$

Les équations seront

$$\dot{x}(t) = A(t)x(t) + B(t)u(t)$$

$$y(t) = x(t)$$

$$\text{avec } x(0) = x_0$$

où  $A(t)$  et  $B(t)$  sont deux matrices respectivement  $n \times n$  et  $n \times r$ , continues par morceaux sur  $[0, T]$ .

Le coût est donné par la fonction suivante, où  $F$  est  $n \times n$  définie positive.

$$J[u, x] = \frac{1}{2} \int_0^T \{ \langle Q(t)x(t), Q(t)x(t) \rangle_n + \langle u(t), R(t)u(t) \rangle_r \} dt + \frac{1}{2} \langle x(T), Fx(T) \rangle_n$$

où  $Q(t)$  est une matrice  $n \times n$  symétrique, définie positive et  $R(t)$   $r \times r$ , symétrique et définie positive, toutes deux continues sur  $[0, T]$ .

La notation  $\langle \cdot, \cdot \rangle_p$  représente le produit scalaire habituel de  $\mathbb{R}^P$ .

Le système sera aussi étudié dans les pages suivantes.

§ 4 Le régulateur d'état parabolique linéaire

Introduisons un dernier problème :

$$[t_0, t] \cap T = [0, T] \subset \mathbb{R} \quad \Sigma = \mathbb{R}^n \quad \Omega = \mathbb{R}^r$$

$$\mathcal{G} = \left\{ V_2^1[(0,1) \times (0,T)], \mathbb{R}^n \right\}$$

$$\mathcal{U} = \left\{ V_2^0[(0,1) \times (0,T)], \mathbb{R}^r \right\}$$

où  $\{V_p^q[(0,1) \times (0,T)], \mathbb{R}^k\}$  est l'ensemble des fonctions à valeurs dans  $\mathbb{R}^k$  telles que

$$\frac{\partial^{i+j} f}{\partial x^i \partial t^j} \in L^{\frac{p}{p+1}}[(0,1) \times (0,T)]$$

pour  $i = 1-q$  et  $j = 1-q-1$ , muni de la norme

$$\|f(x,t)\|_{V_q^p}^p = \int_0^T \left[ \sum_{j=0}^q \sum_{k=0}^{q-1} \left| \frac{\partial^{j+k} f}{\partial x^j \partial t^k}(x,t) \right|^2 \right]^{p/2} dt$$

Les équations, condition initiale et conditions frontières sont

$$\frac{\partial v}{\partial t} = A(x,t) \frac{\partial^2 v}{\partial x^2} + B(x,t) u(x,t)$$

$$\underset{t \rightarrow \infty}{\lim} v(x,t) = v_0(x)$$

$$\alpha v(0,t) + \frac{\partial v(0,t)}{\partial x} = c f_1(t)$$

$$\beta v(1,t) + \frac{\partial v(1,t)}{\partial x} = d f_2(t)$$

$$g(x,t) = v(x,t)$$

où  $f_1$  et  $f_2 \in \{W_2^0[0,T], \mathbb{R}^n\}$ , A et B sont deux matrices respectivement  $n \times n$  et  $n \times r$ , toutes deux bornées sur  $C^\gamma(0,1) \times C^{\gamma-1}(0,T)$  pour  $\gamma \geq 1$ , et où  $v_0(x) \in C^\delta(0,1)$

La fonctionnelle de coût s'écrit

$$\begin{aligned} J[u, f_1, f_2] = & \frac{1}{2} \int_0^T \left\{ \langle v, Q(x,t)v \rangle + \langle u, R(x,t)u \rangle \right\} dx dt \\ & + \frac{1}{2} \int_0^T \left\{ \langle f_1, r(t) f_1 \rangle + \langle f_2, s(t) f_2 \rangle \right\} dt \end{aligned}$$

où  $Q, R, s$  et  $r$  sont des matrices  $n \times n$  définies positives et par morceaux  $C^\gamma(0,1) \times C^{\gamma-1}(0,T)$  pour  $\gamma \geq 1$  et  $\langle \cdot, \cdot \rangle$  est le produit scalaire.

Le système un peu plus complexe, où l'espace intervient en plus du temps, sera, lui aussi, étudié dans ce travail.

## § 5 Remarques sur les régulateurs d'état.

D'abord, les deux systèmes introduits sans ce nom, doivent leur appellation au fait que le but de l'optimisation est de ramener l'état du système au repos (au zéro). On désire de plus que le contrôle ne prenne pas des valeurs trop élevées (dans plusieurs cas pratiques, l'énergie dépensée pour le contrôle est proportionnelle à  $\int_0^T u(t)^2 dt$ ). Ces considérations expliquent la forme de la fonctionnelle de coût où les variables  $x$  et  $u$  apparaissent dans des formes quadratiques toujours positives.

Il est à remarquer qu'aucune contrainte explicite n'est là pour limiter la taille du contrôle. Seule la pénalisation dans le coût limite cette taille.

## II Les méthodes de résolution

Dans ce paragraphe, nous aborderons d'abord une méthode générale de résolution du problème du régulateur d'état différentiel linéaire. Ensuite nous parlerons de la méthode de Ritz-Galerkin, pour introduire le chapitre qui l'étudiera plus en détail.

### 16 Le problème du régulateur d'état.

Rappelons les données générales du problème

$$\begin{cases} \dot{x}(t) = Ax(t) + Bu(t) \\ x(0) = x_0 \end{cases}$$

$$J[u] = \frac{1}{2} \langle x(T), Fx(T) \rangle + \frac{1}{2} \int_0^T \{ \langle x(t), Q(t)x(t) \rangle + \langle u(t), R(t)u(t) \rangle \} dt$$

comme plus haut.

Nous allons montrer que  $u(t) = G(t)x(t)$ , où  $G(t)$  est une matrice  $r \times n$ .

Supposons d'abord que le contrôle optimal existe pour tout état initial. Nous allons nous servir du principe de Pontryagine pour obtenir des conditions nécessaires sur le contrôle optimal et tenter de le calculer à partir de là.

Ecrivons l'Hamiltonien du système

$$H = \frac{1}{2} \langle x(t), Qx(t) \rangle + \frac{1}{2} \langle u(t), Ru(t) \rangle + \langle Ax(t), \lambda(t) \rangle + \langle B(t)u(t), \lambda(t) \rangle$$

Le multiplicateur  $\lambda(t)$  est solution de  $\dot{\lambda}(t) = -\frac{\partial H}{\partial x(t)}$

$$\text{soit } \dot{\lambda}(t) = -Q(t)x(t) - A^T(t)\lambda(t)$$

Le long de la trajectoire optimale, on doit avoir  $\frac{\partial H}{\partial u(t)} = 0$   
ce qui équivaut à

$$\frac{\partial H}{\partial u(t)} = R(t)u(t) + B^t(t)\lambda(t) = 0$$

d'où on déduit  $u(t) = -R^{-1}(t)B^t(t)\lambda(t)$ .

Le fait que  $R(t)$  est définie positive pour tout  $t \in [0, T]$  implique l'existence de  $R^{-1}(t)$  pour toutes ces valeurs.

Nous savons que l'Hamiltonien doit être minimum. La condition  $\frac{\partial H}{\partial u(t)} = 0$  assure uniquement un extrémum. Pour un minimum, il faut que  $\frac{\partial^2 H}{\partial u^2(t)}$  soit une matrice définie positive.

Or

$$\frac{\partial^2 H}{\partial u^2(t)} = R(t)$$

qui est définie positive. Ce qui prouve donc que le contrôle  $u(t)$  trouvé plus haut minimise bien  $H$ .

Remplaçons maintenant  $u(t)$  par  $-R^{-1}(t)B^t(t)\lambda(t)$

$$\dot{x}(t) = A(t)x(t) - B(t)R^{-1}(t)B^t(t)\lambda(t)$$

Si nous définissons  $S(t) = B(t)R^{-1}(t)B^t(t)$ , qui est évidemment une matrice  $n \times n$  symétrique, nous pouvons écrire

$$\begin{bmatrix} \dot{x}(t) \\ \dot{\lambda}(t) \end{bmatrix} = \begin{bmatrix} A(t) & -S(t) \\ -g(t) & -A^t(t) \end{bmatrix} \begin{bmatrix} x(t) \\ \lambda(t) \end{bmatrix}$$

qui est un système de  $2n$  équations différentielles homogène linéaire. Nous savons obtenir une solution unique de ce système pourvu qu'il nous ayons  $2n$  conditions frontières. La condition initiale  $x(0) = x_0$  nous en fournit  $n$ . D'autre part, la valeur  $\lambda(T)$  n'étant pas encore déterminée

nous pouvons la choisir de manière à minimiser  $\frac{1}{2} \langle x(T), Fx(T) \rangle$ , soit

$$\lambda(T) = \frac{\partial}{\partial x(T)} \frac{1}{2} \langle x(T), Fx(T) \rangle = Fx(T)$$

ce qui nous fournit les  $n$  conditions nécessaires. (C'est la condition de transversalité, cf. (A-1) pp 270-271)

Soit  $\Omega(t, 0)$  la matrice fondamentale ( $2n \times 2n$ ) du système différentiel, on obtient donc

$$\begin{bmatrix} x(t) \\ \lambda(t) \end{bmatrix} = \Omega(t, 0) \begin{bmatrix} x(0) \\ \lambda(0) \end{bmatrix}$$

où  $\lambda(0)$  est inconnue. En  $t=T$ , on obtient

$$\begin{bmatrix} x(T) \\ \lambda(T) \end{bmatrix} = \Omega(T, t) \begin{bmatrix} x(t) \\ \lambda(t) \end{bmatrix}$$

Partageons maintenant  $\Omega(T, t)$  en 4 sous-matrices  $n \times n$

$$\Omega(T, t) = \begin{bmatrix} \Omega_{11}(T, t) & \Omega_{12}(T, t) \\ \Omega_{21}(T, t) & \Omega_{22}(T, t) \end{bmatrix}$$

Réécrivons l'équation précédente

$$x(T) = \Omega_{11}(T, t)x(t) + \Omega_{12}(T, t)\lambda(t)$$

$$\lambda(T) = \Omega_{21}(T, t)x(t) + \Omega_{22}(T, t)\lambda(t) = Fx(T)$$

ce qui fournit, après un bref calcul

$$\lambda(t) = [\Omega_{22}(T, t) - F\Omega_{12}(T, t)]^{-1} [F\Omega_{11}(T, t) - \Omega_{21}(T, t)]x(t)$$

pourvu que l'inverse considéré existe.

Pasons maintenant

$$K(t) = [\Omega_{22}(T, t) - F\Omega_{12}(T, t)]^{-1} [F\Omega_{11}(T, t) - \Omega_{21}(T, t)]$$

une matrice  $n \times n$  dépendant de  $T$  et  $F$ , mais pas de l'état initial

Notons maintenant que  $[\Omega_{22}(T, t) - F\Omega_{12}(T, t)]^{-1}$  existe. En  $t=T$  nous avons  $\Omega(T, T) = I$ , ce qui implique

$$\Omega_{11}(T,T) = \Omega_{22}(T,T) = I$$

$$\Omega_{12}(T,T) = \Omega_{21}(T,T) = 0$$

Donc  $\Omega_{22}(T,T) - F\Omega_{12}(T,T) = I$  est inversible. De plus, on voit immédiatement que

$$K(T) = [\Omega_{22}(T,T) - F\Omega_{12}(T,T)]^{-1} [F\Omega_{11}(T,T) - \Omega_{21}(T,T)] = F$$

Comme nous venons de le voir, tout marche bien entre T. On peut voir (Ref (K-1)) que l'inverse existe pour tout  $t \in [t_0, T]$ ; nous pouvons alors écrire

$$\lambda(t) = K(t)x(t)$$

Quelques commentaires s'imposent ici. Quand les matrices A, S et Q dépendent du temps, il est généralement impossible d'obtenir une expression analytique de K(t). Même lorsque ces matrices sont constantes, l'évaluation de K(t) est un calcul très laborieux, spécialement lorsque l'ordre du système est élevé.

Nous allons maintenant mettre en évidence une propriété de K(t) qui permettra de la calculer, sans devoir inverser une matrice n × n.

Supposons que x(t) et λ(t) sont solutions du système différentiel avec  $\lambda(t) = K(t)x(t)$ . On obtient

$$\dot{\lambda}(t) = \dot{K}(t)x(t) + K(t)\dot{x}(t)$$

$$\text{Or nous avons } \dot{x}(t) = A(t)x(t) - S(t)\lambda(t)$$

$$\dot{\lambda}(t) = -Q(t)x(t) - A^*(t)\lambda(t)$$

En substituant, on obtient

$$\dot{x}(t) = [A(t) - S(t)K(t)]x(t)$$

$$\dot{\lambda}(t) = [K(t) + K(t)A(t) - K(t)S(t)K(t)]x(t)$$

$$\text{et encore } \dot{\lambda}(t) = [-Q(t) - A^*(t)K(t)]x(t)$$

ce qui fournit

$$[K(t) + K(t)A(t) - K(t)S(t)K(t) + A^*(t)K(t) + Q(t)]x(t) = 0$$

pour tout  $t \in [t_0, T]$

Comme cette équation doit être vraie pour tout choix des conditions initiales, car  $K(t)$  ne dépend pas de celles-ci, on doit avoir

$$\dot{K}(t) + K(t)A(t) + A^t(t)K(t) - K(t)S(t)K(t) + q(t) = 0$$

or  $S(t) = B(t)R^{-1}(t)B^t(t)$ , donc

$$\dot{K}(t) = -K(t)A(t) - A^t(t)K(t) + K(t)B(t)R^{-1}(t)B^t(t)K(t) - q(t)$$

Nous avons donc établi

Si  $x(t)$  et  $\lambda(t)$  sont les solutions du système différentiel et si  $\lambda(t) = K(t)x(t)$  pour tout  $t \in [t_0, T]$  et tout  $x(t)$ , alors  $K(t)$  doit satisfaire

$$\dot{K}(t) = -K(t)A(t) - A^t(t)K(t) + K(t)B(t)R^{-1}(t)B^t(t)K(t) - q(t)$$

De plus, si  $t = T$ , on doit avoir (cf supra)  $\lambda(T) = Fx(T)$ , mais aussi  $\lambda(T) = K(T)x(T)$ . Ce qui fournit

$$[K(T) - F]x(T) = 0$$

pour tout  $x(T)$ . Donc  $K(T) = F$

(analogue à ce qui a déjà été vu).

Nous avons donc une équation de Riccati pour le matrice  $K(t)$  avec une condition limite : nous avons donc une seule solution.

Pour réduire l'ordre de l'équation, voyons de plus que

$K(t)$  est symétrique, pour  $t \in [t_0, T]$

Démonstration : Transposons les deux membres de l'équation de Riccati :

$$\left[ \frac{d}{dt} K(t) \right]^+ = -K^T(t) A(t) - A^T(t) K^T(t) + K^T(t) B(t) R^{-1}(t) B^T(t) K^T(t) - Q(t)$$

puisque  $Q(t)$  et  $B(t) R^{-1}(t) B^T(t)$  sont symétriques.

D'autre part, pour toute matrice  $K(t)$ , il est vrai que

$$\left[ \frac{d}{dt} K(t) \right]^+ = \frac{d}{dt} [K^T(t)]$$

En comparant alors l'équation originale et l'équation transposée, on voit que  $K(t)$  et  $K^T(t)$  satisfont la même équation. Comme  $F$  est symétrique, on conclut  $K(T) = K^T(T) = F$ .

Comme la solution de l'équation est unique, on a bien  $K(t) = K^T(t)$  □

## Théorème II.1

Soit le système  $\dot{x}(t) = A(t)x(t) + B(t)u(t)$   
et la fonctionnelle de coût

$$J_1 = \frac{1}{2} \langle x(T), Fx(T) \rangle + \frac{1}{2} \int_0^T \{ \langle x(t), Q(t)x(t) \rangle + \langle u(t), R(t)u(t) \rangle \} dt$$

où  $u(t)$  n'est soumis à aucune contrainte,  $T$  est donné,  $F$  et  $Q(t)$  et  $R(t)$  sont définies positives et symétriques. Alors, l'unique contrôle optimal est donné par

$$u(t) = -R^{-1}(t) B^T(t) K(t) x(t)$$

où la matrice  $K(t)$  ( $n \times n$ ) est symétrique et seule solution de l'équation

$$\dot{K}(t) = -K(t)A(t) - A^T(t)K(t) + K(t)B(t)R^{-1}(t)B^T(t)K(t) - Q(t)$$

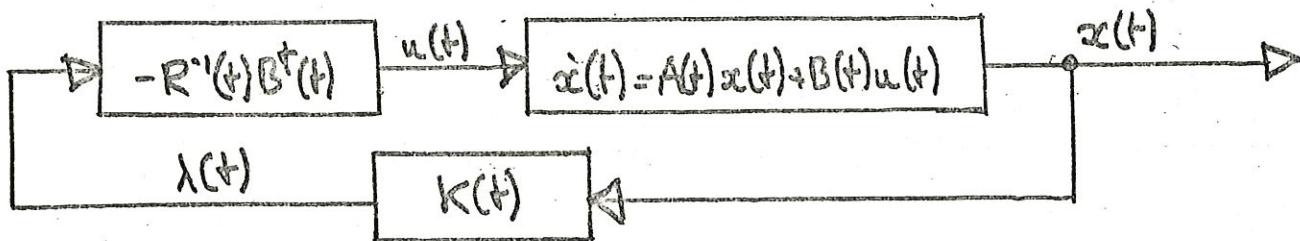
satisfaisant la condition limite  $K(T) = F$

L'état  $x(t)$  est donné par l'équation

$$\dot{x}(t) = [A(t) - B(t)R^{-1}(t)B^T(t)K(t)]x(t) \quad x(0) = x_0$$

Nous avons déjà établi la condition nécessaire. La condition suffisante, qui nécessite de longs développements, ne sera pas traitée ici. Une démonstration extensive peut être trouvée dans (A-1), pp 763 sq.

Pour terminer ce paragraphe, nous donnons ci-dessous un schéma montrant la structure du système étudié.



## §2 L'approximation de Ritz Galerkin

Dans des cas plus généraux que celui traité dans le §1, les calculs deviennent évidemment nettement plus complexes. Même dans le cas précédent, la masse de calcul peut devenir extrêmement lourde.

La difficulté réside principalement dans le fait que tous les espaces fonctionnels considérés (principalement  $U$  et  $G$ ) sont de dimension infinie : on ne peut donc exprimer les fonctions considérées comme des combinaisons linéaires finies de fonctions qu'on se donne à priori.

Pour remédier à cet inconvénient, nous introduirons l'approximation suivante. Choisissons dans  $U$  et dans  $G$  deux sous-espaces de dimension finie, soit  $U_m$  et  $G_m$  respectivement. Résolvons le problème sur ces espaces uniquement. Nous espérons obtenir ainsi une solution approchée, d'autant meilleure que  $U_m$  et  $G_m$  seront plus proches de  $U$  et  $G$ . Comme nous construisons  $U_m$  et  $G_m$  de manière à pouvoir

Si on choisit des bases finies, nous pouvons écrire

$$\left[ \begin{array}{l} x(t) = \sum_{j=1}^m \alpha_j G_j^m(t) \\ \lambda(t) = \sum_{j=1}^m \beta_j G_j^m(t) \\ u(t) = \sum_{j=1}^m \gamma_j U_j^m(t) \end{array} \right]$$

et, dans ces, nous écrivons la contrainte de l'équation d'état, lorsque cette écriture a un sens,

$$\int_0^t \langle G_j^m(\hat{t}), x(\hat{t}) - \phi[\hat{t}, u]_{t_0, \hat{t}}, x(t_0) \rangle d\hat{t} = 0$$

$$\text{avec } x(0) = x_0$$

L'idée qui consiste à choisir des sous-espaces de dimension finie est-elle bonne? Les solutions trouvées par cette méthode convergeront-elles, dans un sens encore à définir, vers les solutions du problème non approché?

Voilà quelques questions auxquelles nous tenterons de répondre dans les pages qui suivent, en considérant successivement les trois problèmes décrits dans le paragraphe IV.

## Chapitre III

### La méthode de Ritz Galerkin sur des modèles de contrôle optimal

#### I Introduction

Nous allons maintenant considérer successivement les différents problèmes abordés dans le paragraphe IV du Chapitre II. Nous tenterons de trouver, pour l'approximation de Ritz Galerkin, des bornes d'erreur à priori, qui permettront de prouver la convergence des solutions approchées vers "la" solution du problème non approché.

#### II Le régulateur d'état différentiel linéaire

Il s'agit du système présenté au Chapitre II, IV, §3. Rappelons brièvement les données du problème.

Soit  $Q(t)$  une matrice  $n \times n$ , symétrique, définie positive et  $R(t)$  une matrice  $r \times r$  symétrique, définie positive, toutes deux continues sur l'intervalle de temps  $[0, T]$  où  $T$  est fixé. Pour chaque  $k \geq 0$ , définissons

$$\Phi^k = \{ f : [0, T] \rightarrow \mathbb{R}^k \mid f \text{ est différentiable partout et à dérivée bornée} \}$$

Le problème est de trouver un contrôle optimal  $u^* \in \Phi^r$  et un état  $x^* \in \Phi^n$  qui minimisent

$$J[u, x] = \frac{1}{2} \int_0^T \{ \langle x(t), Q(t)x(t) \rangle_n + \langle u(t), R(t)u(t) \rangle_r \} dt \quad (\text{III.1})$$

permet tous les  $u \in \Phi^r$ , où  $x(t)$  est donné par

$$\begin{cases} \dot{x}(t) = A(t)x(t) + B(t)u(t) & t \geq 0 \\ x(0) = x_0 \end{cases} \quad (\text{III.2})$$

avec  $\|y\|_n^2 = \langle y, y \rangle_n = \sum_{i=1}^n y_i^2$  pour tout  $y \in \mathbb{R}^n$  et, de même

$\|z\|_r^2 = \langle z, z \rangle_r = \sum_{i=1}^r z_i^2$  pour tout  $z \in \mathbb{R}^r$ . De plus,  $A(t)$  est une matrice  $n \times n$  et  $B(t)$  une matrice  $n \times r$ , toutes deux continues par morceaux sur  $[0, T]$

### 1§ Procédure variationnelle

Comme vu précédemment, le problème ci-dessus est équivalent à celui de trouver  $\lambda^* \in \Phi^n$  qui extrémise le Lagrangien

$$\begin{aligned} L[u, x; \lambda, \gamma] = & J[u, x] + \int_0^T \langle \lambda(t), -\dot{x}(t) + A(t)x(t) + B(t)u(t) \rangle_n dt \\ & + \langle \gamma, x(0) - x_0 \rangle_n \end{aligned} \quad (\text{III.3})$$

sous la condition de transversalité  $\lambda(T) = 0$  (III.4), où  $\gamma, u(t)$  et  $x(t)$  sont donnés respectivement par

$$\gamma = -\lambda(0) \quad (\text{III.5})$$

$$u(t) = -R^{-1}(t) B^*(t) \lambda(t) \quad t \in [0, T] \quad (\text{III.6})$$

$$x(t) = -Q^{-1}(t) [\dot{\lambda} + A^*(t)\lambda(t)] \quad t \in [0, T] \quad (\text{III.7})$$

Effectuons à présent quelques changements d'écriture.

De (III.3) on tire :

$$\begin{aligned} L = & \frac{1}{2} \int_0^T \langle x, Qx \rangle + \frac{1}{2} \int_0^T \langle u, Ru \rangle + \int_0^T \langle \lambda, -\dot{x} \rangle + \int_0^T \langle \lambda, Ax \rangle + \int_0^T \langle \lambda, Bu \rangle \\ & + \langle -\lambda(0), x(0) - x_0 \rangle \end{aligned}$$

$$\begin{aligned}
 L &= \frac{1}{2} \int_0^T \langle \dot{x}, Qx \rangle + \frac{1}{2} \int_0^T \langle u, Ru \rangle - [\langle \lambda, x \rangle]_0^T + \int_0^T \langle \dot{\lambda}, x \rangle + \int_0^T \langle \lambda, Ax \rangle + \int_0^T \langle \lambda, Bu \rangle \\
 &\quad + \langle \lambda(0), x(0) - x_0 \rangle \\
 &= \frac{1}{2} \int_0^T \langle \dot{x}, Qx \rangle + \frac{1}{2} \int_0^T \langle u, Ru \rangle - \langle \lambda(T), x(T) \rangle + \langle \lambda(0), x(0) \rangle \\
 &\quad + \int_0^T \langle \dot{\lambda} + A^t \lambda, x \rangle + \int_0^T \langle B^t \lambda, u \rangle - \langle \lambda(0), x(0) \rangle + \langle \lambda(0), x_0 \rangle \\
 &= \frac{1}{2} \int_0^T \langle \dot{x}, Qx \rangle + \frac{1}{2} \int_0^T \langle u, Ru \rangle - \int_0^T \langle \dot{x}, Qx \rangle - \int_0^T \langle u, Ru \rangle + \langle \lambda(0), x_0 \rangle
 \end{aligned}$$

En effet,  $\lambda(T) = 0$  et (III.6) implique  $-Ru = B^t \lambda$ , tandis que de (III.7) on tire  $-Qx = \dot{\lambda} + A^t \lambda$

D'où

$$L = -\int_0^T \langle u, x \rangle + \langle \lambda(0), x_0 \rangle$$

Les conditions de continuité sur  $A, B, R, Q$  assurent l'existence des intégrales considérées.

Mais

$$\frac{1}{2} \int_0^T \langle u, Ru \rangle dt = \frac{1}{2} \int_0^T \langle R^{-1} B^t \lambda, RR^{-1} B^t \lambda \rangle dt = \frac{1}{2} \int_0^T \langle BR^{-1} B^t \lambda, \lambda \rangle dt$$

en utilisant (III.6).

De même, par (III.7), on obtient

$$\begin{aligned}
 \frac{1}{2} \int_0^T \langle \dot{x}, Qx \rangle dt &= \frac{1}{2} \int_0^T \langle Q^{-1} A^t \lambda + Q^{-1} \dot{\lambda}, A^t \lambda + \dot{\lambda} \rangle dt \\
 &= \frac{1}{2} \int_0^T \langle Q^{-1} \dot{\lambda}, \dot{\lambda} \rangle dt + \frac{1}{2} \int_0^T \langle Q^{-1} A^t \lambda, A^t \lambda \rangle dt \\
 &\quad + \frac{1}{2} \int_0^T \langle Q^{-1} A^t \lambda, \dot{\lambda} \rangle dt + \frac{1}{2} \int_0^T \langle Q^{-1} \dot{\lambda}, A^t \lambda \rangle dt \\
 &= \frac{1}{2} \int_0^T \langle Q^{-1} \dot{\lambda}, \dot{\lambda} \rangle dt + \frac{1}{2} \int_0^T \langle A Q^{-1} A^t \lambda, \lambda \rangle dt \\
 &\quad + \int_0^T \langle A Q^{-1} \dot{\lambda}, \lambda \rangle dt
 \end{aligned}$$

car  $Q$  est symétrique

Donc, si nous définissons, pour tout  $\lambda$  et  $\eta$  dans  $\mathbb{P}^n$

$$\begin{aligned} [\lambda, \eta] = & \int_0^T \langle Q^{-1}\dot{\lambda}, \dot{\eta} \rangle_n dt + \int_0^T \langle A Q^{-1} A^t \lambda, \eta \rangle_n dt \\ & + \int_0^T \langle B R^{-1} B^t \lambda, \eta \rangle_n dt + \int_0^T \{ \langle A Q^{-1} \dot{\lambda}, \eta \rangle_n + \langle A Q^{-1} \dot{\eta}, \lambda \rangle_n \} dt \end{aligned} \quad (\text{III.8})$$

nous pouvons bien écrire

$$\begin{aligned} F[\lambda] = -L[u, z; \lambda, \eta] = & \frac{1}{2} \int_0^T \langle Q^{-1} \dot{\lambda}, \lambda \rangle dt + \frac{1}{2} \int_0^T \langle A Q^{-1} \lambda, \lambda \rangle dt \\ & + \frac{1}{2} \int_0^T \langle B R^{-1} B^t \lambda, \lambda \rangle dt + \int_0^T \langle A Q^{-1} \dot{\lambda}, \lambda \rangle dt - \langle \lambda(0), z_0 \rangle_n \end{aligned}$$

$$F[\lambda] = +\frac{1}{2} [\lambda, \lambda] - \langle \lambda(0), z_0 \rangle_n \quad (\text{III.9})$$

qu'il faut minimiser par rapport à  $\lambda$ .

## 2§

### Quelques théorèmes

Définissons les notations suivantes

$$\|M\|_f = \max \{ |Mx|_f \mid x \in \mathbb{R}^n \text{ et } \|x\|_f = 1 \}$$

$$\lambda_g = \min_{t \in [0, T]} \{ \lambda(t) \mid \lambda(t) \text{ est une valeur propre de } Q(t) \}$$

$$\|Q\|_2^2 = \int_0^T \|Q(t)\|_n^2 dt$$

$$\|A^t\|_2^2 = \int_0^T \|A^t(t)\|_n^2 dt$$

$$\|A^t\|_\infty = \sup_{t \in [0, T]} \|A^t(t)\|_n$$

et enfin

$$\ell = \|Q\|_2 \left[ 2 \|A^T\|_\infty \right]^{-1/2} \exp \left[ \|A^T\|_\infty T \right]$$

On peut alors démontrer le théorème suivant

**Théorème III.1**

Le multiplicateur de Lagrange optimal  $\lambda$  existe et est l'unique solution dans  $\bar{\Phi}_0^n = \{ \phi \in \bar{\Phi}^n \mid \phi(T) = 0 \}$  du problème généralisé ( $[\lambda_m]$  est donné par (11))

$$[\lambda, \eta] = \langle \eta(0), x_0 \rangle_n \text{ pour tout } \eta \in \bar{\Phi}_0^n$$

De plus,  $[\lambda, \eta]$  est symétrique et

$$\|\dot{\lambda}\|_2^2 = \int_0^T \langle \dot{\lambda}, \dot{\lambda} \rangle_n dt \leq \tilde{\lambda}_Q^{-1} \left[ \|Q\|_2 + \ell \|A^T\|_2 \right]^2 [\lambda, \lambda]$$

$$\|\lambda\|_2^2 \leq \tilde{\lambda}_Q^{-1} \ell^2 [\lambda, \lambda]$$

(III.10)

(III.11)

(III.12)

Démonstration : La partie "existence" du théorème nous est déjà connue. Il s'agit d'ailleurs d'un résultat classique du calcul des variations (cf. Ref A-1 et B-1).

Si  $\eta \in \bar{\Phi}_0^n$  et  $\alpha \in \mathbb{R}$ , alors  $F[\lambda^* + \alpha\eta] \geq F[\lambda^*]$  avec l'égalité ssi  $\alpha = 0$ , où  $\lambda^*$  est le multiplicateur de Lagrange optimal.

Nous devons donc avoir

$$\frac{\partial F}{\partial \alpha} [\lambda^* + \alpha\eta](0) = 0$$

$$\begin{aligned} F[\lambda^* + \alpha\eta] &= \frac{1}{2} \left\{ \int_0^T \langle Q^{-1}(\lambda^* + \alpha\eta), \dot{\lambda} + \alpha\dot{\eta} \rangle + \int_0^T \langle A Q^{-1} A^T (\lambda^* + \alpha\eta), \lambda^* + \alpha\eta \rangle \right. \\ &\quad \left. + \int_0^T \langle B R^{-1} B^T (\lambda^* + \alpha\eta), \lambda^* + \alpha\eta \rangle \right\} + \int_0^T \langle A Q^{-1} (\lambda^* + \alpha\eta), \lambda^* + \alpha\eta \rangle \\ &\quad - \langle \lambda(0) + \alpha\eta(0), x_0 \rangle \end{aligned}$$

Dérivons terme à terme sous le signe "intégrale"

a) Soit  $T_1 = \langle Q^{-1}\lambda^*, \dot{\lambda}^* \rangle + \alpha \langle Q^{-1}\lambda^*, \dot{\gamma} \rangle + \alpha \langle Q^{-1}\dot{\gamma}, \lambda^* \rangle + \alpha^2 \langle Q^{-1}\dot{\gamma}, \dot{\gamma} \rangle$   
 or  $Q$  est symétrique, ce qui implique  $\langle Q^{-1}\lambda^*, \dot{\gamma} \rangle = \langle Q^{-1}\dot{\gamma}, \lambda^* \rangle$

On a donc

$$\frac{\partial T_1}{\partial \alpha} = 2 \langle Q^{-1}\lambda^*, \dot{\gamma} \rangle + 2\alpha \langle Q^{-1}\dot{\gamma}, \lambda^* \rangle$$

$$\left. \frac{\partial T_1}{\partial \alpha} \right|_{\alpha=0} = 2 \langle Q^{-1}\lambda^*, \dot{\gamma} \rangle$$

b) Soit  $T_2 = \langle A Q^{-1} A^t \lambda^*, \lambda^* \rangle + \alpha \langle A Q^{-1} A^t \lambda^*, \gamma \rangle + \alpha \langle A Q^{-1} A^t \gamma, \lambda^* \rangle + \alpha^2 \langle A Q^{-1} A^t \gamma, \gamma \rangle$

Comme supra, on a  $\langle A Q^{-1} A^t \lambda^*, \gamma \rangle = \langle A Q^{-1} A^t \gamma, \lambda^* \rangle$

ce qui donne

$$\frac{\partial T_2}{\partial \alpha} = 2 \langle A Q^{-1} A^t \lambda^*, \gamma \rangle + 2\alpha \langle A Q^{-1} A^t \gamma, \gamma \rangle$$

$$\left. \frac{\partial T_2}{\partial \alpha} \right|_{\alpha=0} = 2 \langle A Q^{-1} A^t \lambda^*, \gamma \rangle$$

c) Soit  $T_3 = \langle B R^{-1} B^t \lambda^*, \lambda^* \rangle + \alpha \langle B R^{-1} B^t \lambda^*, \gamma \rangle + \alpha \langle B R^{-1} B^t \gamma, \lambda^* \rangle + \alpha^2 \langle B R^{-1} B^t \gamma, \gamma \rangle$

La symétrie de  $R$  implique  $\langle B R^{-1} B^t \lambda^*, \gamma \rangle = \langle B R^{-1} B^t \gamma, \lambda^* \rangle$

Nous obtenons donc

$$\left. \frac{\partial T_3}{\partial \alpha} \right|_{\alpha=0} = 2 \langle B R^{-1} B^t \lambda^*, \gamma \rangle$$

d) Soit  $T_4 = \langle A Q^{-1} \lambda^*, \lambda^* \rangle + \alpha \langle A Q^{-1} \dot{\lambda}^*, \gamma \rangle + \alpha \langle A Q^{-1} \dot{\gamma}, \lambda^* \rangle + \alpha^2 \langle A Q^{-1} \dot{\gamma}, \gamma \rangle$

On obtient

$$\left. \frac{\partial T_4}{\partial \alpha} \right|_{\alpha=0} = \langle A Q^{-1} \dot{\lambda}^*, \gamma \rangle + \langle A Q^{-1} \dot{\gamma}, \lambda^* \rangle + 2\alpha \langle A Q^{-1} \dot{\gamma}, \gamma \rangle$$

$$\left. \frac{\partial T_4}{\partial \alpha} \right|_{\alpha=0} = \langle A Q^{-1} \dot{\lambda}^*, \gamma \rangle + \langle A Q^{-1} \dot{\gamma}, \lambda^* \rangle$$

e) Soit enfin  $T_5 = \langle \lambda(0), x_0 \rangle + \alpha \langle \gamma(0), x_0 \rangle$

Donc  $\left. \frac{\partial T_5}{\partial \alpha} \right|_{\alpha=0} = \langle \gamma(0), x_0 \rangle$

Rassemblant tous les termes on obtient

$$\begin{aligned}\frac{\partial F}{\partial \alpha}[\lambda^*, \alpha \eta](0) &= \int_0^T \langle Q^{-1} \dot{\lambda}^*, \dot{\eta} \rangle + \int_0^T \langle A Q^{-1} A^t \lambda^*, \eta \rangle + \int_0^T \langle B R^{-1} B^t \lambda^*, \eta \rangle \\ &+ \int_0^T \{ \langle A Q^{-1} \dot{\lambda}^*, \eta \rangle + \langle A Q^{-1} \dot{\eta}, \lambda^* \rangle \} - \langle \eta(0), x_0 \rangle = 0\end{aligned}$$

ce qui fournit, en utilisant (III.8)

$$[\lambda, \eta] - \langle \eta(0), x_0 \rangle = 0 \quad \text{pour tout } \eta \in \mathbb{D}^n.$$

ce qui est équivalent à (III.10)

D'autre part, on voit aisement dans (III.8) que  $[\lambda, \eta]$  est symétrique en  $\lambda$  et  $\eta$  car  $Q$  et  $R$  sont symétriques.

$$\text{De plus } \frac{1}{2} [\lambda, \lambda] = J[u, x] = \frac{1}{2} \int_0^T \langle u, Q u \rangle dt + \frac{1}{2} \int_0^T \langle u, R u \rangle dt$$

$$\geq \frac{1}{2} \int_0^T \langle u, Q u \rangle dt \geq \frac{1}{2} \lambda_Q \int_0^T \langle u, u \rangle dt = \frac{1}{2} \lambda_Q \|x\|_2^2 \quad (\text{III.9})$$

De (III.7) et (III.4) on peut déduire

$$\lambda(t) = - \int_t^T \frac{d\lambda}{ds} ds = \int_t^T [Q u + A^t \lambda] ds$$

Donc

$$|\lambda(t)| = |\langle \lambda(t), \lambda(t) \rangle|^{\frac{1}{2}} = \left| \int_t^T [Q u + A^t \lambda] ds \right|$$

car  $|\cdot|$  est la norme associée à  $\langle \cdot, \cdot \rangle$ .

ce qui donne

$$\begin{aligned}|\lambda(t)| &\leq \left| \int_t^T Q u \right| + \left| \int_t^T A^t \lambda \right| \leq \int_t^T |Q u| + \int_t^T |A^t \lambda| \\ &\leq \|Q\|_2 \|x\|_2 + \int_t^T |A^t(s)| \cdot |\lambda(s)| ds\end{aligned}$$

en utilisant l'inégalité de Cauchy-Schwarz.

L'inégalité de Gronwall (Ref B-5) implique alors que

$$\|\lambda\|_2 \leq \|Q\|_2 [2 \|A^t\|_{\infty} T]^{\frac{1}{2}} \exp[\|A^t\|_{\infty} T] \|x\|_2 = \xi \|x\|_2$$

Donc  $[\lambda, \lambda] \geq \lambda_Q \xi^2 \|\lambda\|_2^2$ , ce qui prouve (III.12)

De plus

$$\|\lambda\|_2 \leq \|g\|_2 \|x\|_2 + \|A^t\|_2 \|\lambda\|_2 \leq [\|g\|_2 + \epsilon \|A^t\|_2] \|x\|_2$$

qui, avec (II.3), prouve (III.11).

Pour l'unicité, si  $\lambda$  et  $\mu$  vérifient (III.10), on a

$$0 = [\lambda - \mu, \lambda - \mu] \geq \lambda g^2 \|\lambda - \mu\|_2^2$$

ce qui implique que  $\lambda = \mu$

Comme nous allons travailler dans le cadre de l'approximation de Ritz-Galerkin, nous allons nous restreindre à  $S$  sous espace de dimension finie de  $\mathbb{P}_0^n$  pour chercher les solutions. Nous déterminons ainsi une approximation  $\lambda_S$  de  $\lambda^*$  en minimisant  $F$  sur  $S$  et nous déterminons une approximation  $u_S$  de  $u$  par (III.6). En appliquant le contrôle obtenu dans l'équation (III.2) on obtient l'état  $x_S$  qui est généralement différent de celui qu'on obtiendrait par (III.7)

Montrons d'abord que cette approximation fournit une solution unique

### Théorème III.2

Il existe un seul  $\lambda_S \in S$  qui minimise  $F[\lambda]$  sur  $S$

Démonstration : Choissons  $\{B_i(t)\}_{i=1, M}$  une base de  $S$  (de dimension  $M$ )

$$F\left[\sum_{i=1}^M \beta_i B_i\right] = \frac{1}{2} \left[ \sum_{i=1}^M \beta_i B_i, \sum_{i=1}^M \beta_i B_i \right] - \left< \sum_{i=1}^M \beta_i B_i(t), x_0 \right>$$

c'est une fonction de  $\beta \in \mathbb{R}^M$ , continûment différentiable, et donc F admet un minimum en  $\beta^*$  ssi

$$\left\{ \begin{array}{l} \frac{\partial F}{\partial \beta_i} [\beta^*] = 0 \quad 1 \leq i \leq M \\ H = \left[ \frac{\partial^2 F}{\partial \beta_i \partial \beta_j} \right] \end{array} \right. \quad (III.14)$$

$H = \left[ \frac{\partial^2 F}{\partial \beta_i \partial \beta_j} \right]$  est définie positive

En calculant (III.14) comme précédemment, on obtient

$$\frac{\partial F}{\partial \beta_i} [\beta^*] = \sum_{j=1}^M \beta_j [B_i, B_j] - \langle B_i(0), x_0 \rangle \quad 1 \leq i \leq M \quad (\text{III.15})$$

c'est à dire  $A\beta^* = k$  avec  $A = \{[B_i, B_j]\}$  et  $k_i = \langle B_i(0), x_0 \rangle$ .  
A est évidemment symétrique, car  $[B_i, B_j]$  l'est.

De plus, si  $\beta \neq 0$

$$\beta^T A \beta = \left[ \sum_{i=1}^M \beta_i B_i, \sum_{i=1}^M \beta_i B_i \right] \geq \lambda_g \gamma^2 \left\| \sum_{i=1}^M \beta_i B_i \right\|_2^2 > 0$$

en utilisant (III.13) et A est donc définie positive.

De plus, il résulte de (III.15) que  $H = A$ , et que  $\beta^*$  est, par conséquent,  
l'unique minimum de  $F$  sur  $\mathbb{R}^M$  ■

De plus, on peut obtenir une borne générale de l'erreur :

### Théorème III.3

Si  $\lambda_g$  est l'élément qui minimise  $F[\lambda]$  sur S  
alors  $|\lambda^* - \lambda_g| \equiv [\lambda^* - \lambda_g, \lambda^* - \lambda_g]^{\frac{1}{2}} = \inf_{w \in S} |\lambda^* - w|$

(III.16)

Démonstration : Soit  $w \in S$ , alors  $F[w] = \frac{1}{2} [w, w] - \langle w(0), x_0 \rangle$

$$\text{Donc } F[w] - F[\lambda^*] = \frac{1}{2} [w, w] - \frac{1}{2} [\lambda^*, \lambda^*] + \langle x_0, \lambda^*(0) - w(0) \rangle$$

$$\text{Posons } w = \eta \text{ dans (12). On obtient } [\lambda^* \langle w \rangle] = [\lambda^*, w] = \langle x_0, w(0) \rangle$$

$$\text{De même, en posant } \lambda^* = \eta, \text{ on obtient } [\lambda^*, \lambda^*] = \langle x_0, \lambda^*(0) \rangle$$

Donc

$$\begin{aligned} F[w] - F[\lambda^*] &= \frac{1}{2} [w, w] + \frac{1}{2} [\lambda^*, \lambda^*] + \langle x_0, -w(0) \rangle \\ &= \frac{1}{2} [w, w] + \frac{1}{2} [\lambda^*, \lambda^*] - [\lambda^*, w] \\ &= \frac{1}{2} [\lambda^* - w, \lambda^* - w] = \frac{1}{2} |\lambda^* - w|^2 \end{aligned}$$

en utilisant la linéarité de  $[\lambda, \eta]$

Par conséquent

$$|\lambda^* - \lambda_S|^2 = 2 [F[\lambda_S] - F[\lambda^*]] \leq 2 [F[w] - F[\lambda^*]] = |\lambda^* - w|^2$$

et on obtient

$$\inf_{w \in S} |\lambda^* - w| \leq |\lambda^* - \lambda_S| \leq \inf_{w \in S} |\lambda^* - w|$$

Le qui prouve le théorème ■

Des théorèmes précédents, on peut alors déduire le

### Théorème III.4

Si  $\lambda_S$  est l'élément qui minimise  $F[\lambda]$  sur  $S$   
alors  $\|\lambda^* - \lambda_S\|_2 \leq [\lambda_S]^{-1/2} \inf_{w \in S} |\lambda^* - w|$

(III.17)

$$\|\lambda^* - \lambda_S\|_2 \leq \lambda_S^{-1/2} [\|g\| + \|\Lambda^t\|_2] \inf_{w \in S} |\lambda^* - w|$$

(III.18)

Nous pouvons remarquer que (III.18) implique que  $\lambda_S$  est la projection de  $\lambda^*$  par rapport au produit scalaire  $[\cdot, \cdot]$ .

Du résultat précédent, on peut déduire les résultats suivants

### Théorème III.5

Si  $u_S(t) = -R^{-1}(t)B^t(t)\lambda_S(t)$ ,  $t \in [0, T]$ , est le contrôle calculé (optimal dans  $S$ ) : alors

$$\|u^* - u_S\|_2 \leq \|R^{-1}B^t\|_\infty \lambda_S^{-1/2} \inf_{w \in S} |\lambda^* - w|$$

(III.19)

$$\text{où } \|R^{-1}B^t\|_\infty = \sup_{t \in [0, T]} |R^{-1}B^t(t)|_r$$

Démonstration : Soit  $\delta_S(t) = u^*(t) - u_S(t)$

alors

$$\delta_s(t) = -R^{-1}(t)B^*(t)[\lambda^*(t) - \lambda_s(t)]$$

par linéarité de (III.6)

$$\text{Dès } \|\delta_s(t)\|_2 = \|R^{-1}B^*(\lambda^* - \lambda_s)\|_2 \leq \|R^{-1}B^*\|_\infty \|\lambda^* - \lambda_s\|_2$$

$$\text{et } \|\lambda^* - \lambda_s\|_2 \leq \lambda_g^{-1/2} \inf_{w \in S} |\lambda^* - w|$$

(III.17)

Ce qui prouve (III.19) ■

Définissons maintenant les quantités suivantes

$$\Gamma = T \|B\|_\infty \exp \left[ \int_0^T \|A(z)\|_n dz \right]$$

$$\|A\|_\infty = \sup_{t \in [0, T]} \|A(t)\|_n$$

$$\|B\|_\infty = \sup_{t \in [0, T]} \{ \|B(t)w\|_n \mid w \in \mathbb{R}^r \text{ et } \|w\|_r = 1 \}$$

Nous pouvons alors démontrer le

### Théorème III.6

Si  $\dot{x}_s(t) = A(t)x_s(t) + B(t)u_s(t) \quad t \in [0, T]$   
 et  $x_s(0) = x_0$ , alors

$$\|x^* - x_s\|_2 \leq \Gamma \|R^{-1}B^*\|_\infty \lambda_g^{-1/2} \inf_{w \in S} |\lambda^* - w| \quad (\text{III.20})$$

$$\|\dot{x}^* - \dot{x}_s\|_2 \leq [\Gamma \|A\|_\infty + \|B\|_\infty] \|R^{-1}B^*\|_\infty \lambda_g^{-1/2} \inf_{w \in S} |\lambda^* - w| \quad (\text{III.21})$$

Démonstration : Posons  $\varepsilon_s(t) = x^*(t) - x_s(t)$  pour  $t \in [0, T]$ .

Alors  $\dot{\varepsilon}_s(t) = A\varepsilon_s + B[u^* - u_s]$  pour  $t \in [0, T]$  par linéarité, et, de plus,  $\varepsilon_s(0) = 0$  (par hypothèse)

$$\text{Donc } \dot{\epsilon}_s(t) = \int_0^t A(z) \epsilon_s(z) dz + \int_0^t B(z) \delta_s(z) dz$$

En appliquant alors l'inégalité de Gronwall (Ref B-2), on obtient

$$\|\epsilon_s(t)\|_n \leq T^{1/2} \|B\|_\infty \|\delta_s\|_2 \exp \left[ \int_0^T \|A(z)\|_n dz \right]$$

ce qui donne

$$\|\epsilon_s\|_2^2 = \int_0^T \|\dot{\epsilon}_s(t)\|_n dt \leq T^2 \|B\|_\infty^2 \|\delta_s\|_2^2 \exp \left[ 2 \int_0^T \|A(z)\|_n dz \right]$$

$$\text{Soit } \|\epsilon_s\|_2^2 \leq \Gamma^2 \|u^* - u_s\|_2^2$$

qui, combiné avec (III.19), prouve (III.20).

De plus,

$$\|\dot{\epsilon}_s(t)\|_n \leq \|A\|_\infty \|\epsilon_s(t)\|_n + \|B\|_\infty \|u^*(t) - u_s(t)\|_r \text{ pour } t \in [0, T]$$

$$\begin{aligned} \text{Donc } \|\dot{\epsilon}_s\|_2^2 &= \int_0^T \|\dot{\epsilon}_s(t)\|_n dt \leq \int_0^T \{ \|A\|_\infty \|\epsilon_s\|_n + \|B\|_\infty \|u^* - u_s\|_r \}^2 dt \\ &\leq \{ \|A\|_\infty \|\epsilon_s\|_n + \|B\|_\infty \|u^* - u_s\|_r \} \|u^* - u_s\|_2^2 \end{aligned}$$

ce qui donne

$$\begin{aligned} \|\dot{\epsilon}_s\|_2 &\leq \sqrt{\|A\|_\infty \|\epsilon_s\|_n + \|B\|_\infty \|u^* - u_s\|_r} \|u^* - u_s\|_2 \leq \|A\|_\infty \|\epsilon_s\|_2 + \|B\|_\infty \|u^* - u_s\|_2 \\ &\leq [\Gamma \|A\|_\infty + \|B\|_\infty] \|u^* - u_s\|_2 \end{aligned}$$

qui, avec (III.19), prouve (III.21) ■

Ce théorème nous donne donc des bornes sur les erreurs de trajectoires et de vitesse du système approximé.

Nous pouvons étendre ce genre de résultat en donnant une borne d'erreur pour la fonctionnelle de coût  $J[u, x]$

Théorème III.7

Sous les hypothèses précédentes

$$0 \leq J[u_s, x_s] - J[u^*, x^*] \leq \frac{1}{2} \|R^{-1} B^*\|_\infty^2 \lambda_g^{-1} \gamma^2 \left[ \|g\|_\infty \Gamma^2 + \|R\|_\infty \right] \|u^* - u_s\|_2^2$$

(III.22)

Démonstration : Si  $\delta_s(t)$  et  $\epsilon_s(t)$  sont définis comme plus haut,  
alors :

$$\begin{aligned} J[u_s, x_s] &= \frac{1}{2} \int_0^T \langle x_s, Q x_s \rangle_n dt + \frac{1}{2} \int_0^T \langle u_s, R u_s \rangle_r dt \\ &= \frac{1}{2} \int_0^T \langle x^* + \epsilon_s, Q(x^* + \epsilon_s) \rangle_n dt + \frac{1}{2} \int_0^T \langle u^* + \delta_s, R(u^* + \delta_s) \rangle_r dt \\ &= J[u^*, x^*] + \int_0^T \langle \delta_s, R u^* \rangle_r dt + \int_0^T \langle \epsilon_s, Q x^* \rangle_n dt \\ &\quad + \frac{1}{2} \int_0^T \langle \delta_s, R \delta_s \rangle_r dt + \frac{1}{2} \int_0^T \langle \epsilon_s, Q \epsilon_s \rangle_n dt \end{aligned}$$

car  $Q$  est symétrique.

Mais, par (II.6), on trouve

$$\int_0^T \langle \delta_s, R u^* \rangle_r dt = - \int_0^T \langle \delta_s, Q \lambda^* \rangle_n dt = - \int_0^T \langle B \delta_s, \lambda^* \rangle_n dt.$$

Or  $\dot{\epsilon}_s(t) = A(t) \epsilon_s(t) + B(t) \delta_s(t)$ , on obtient

$$\begin{aligned} \int_0^T \langle \delta_s, R u^* \rangle_r dt &= - \int_0^T \langle \dot{\epsilon}_s - A \epsilon_s, \lambda^* \rangle_n dt \\ &= - [\langle \epsilon_s, \lambda^* \rangle]_0^T + \int_0^T \langle \epsilon_s, \dot{\lambda}^* \rangle_n dt + \int_0^T \langle A \epsilon_s, \lambda^* \rangle_n dt \\ &= 0 + \int_0^T \langle \epsilon_s, \dot{\lambda}^* + A^t \lambda^* \rangle_n dt = - \int_0^T \langle \epsilon_s, Q x^* \rangle_n dt \end{aligned}$$

car  $\epsilon_s(0) = 0$  et  $\dot{\lambda}^*(T) = 0$  car  $\lambda^* \in \overline{\Phi}_0$

Donc

$$\begin{aligned} J[u_s, x_s] &= J[u^*, x^*] + \frac{1}{2} \int_0^T \langle \delta_s, R \delta_s \rangle_r dt + \frac{1}{2} \int_0^T \langle \epsilon_s, Q \epsilon_s \rangle_n dt \\ &\leq J[u^*, x^*] + \frac{1}{2} \|R\|_{\infty} \|\delta_s\|_2^2 + \frac{1}{2} \|Q\|_{\infty} \|\epsilon_s\|_2^2 \\ &\leq J[u^*, x^*] + \frac{1}{2} \|R^t B^t\|_{\infty} \|\lambda^*\|_{\infty} \inf_{w \in S} \|u^* - w\| \left[ \|Q\|_{\infty} \Gamma^2 + \|R\|_{\infty} \right] \end{aligned}$$

en utilisant (II.18) et (II.20).  $\square$

§3

### Une particularisation de S

Il est clair que le choix d'une base finie pour engendrer le sous-espace  $S$  de  $\mathbb{P}_g^n$  est d'une grande importance pratique. Le calcul des solutions optimales en est, en effet, fort tributaire.

Nous nous contenterons de donner quelques renseignements supplémentaires pour ce choix : considérons l'espace (de dimension finie) des fonctions splines

$$S = S_d(\Delta) = \left\{ \sum_{i=1}^p \beta_i B_{d,i}(t) \mid \beta_i \in \mathbb{R}^n \right\}$$

En fait, si l'on choisit  $\Delta: 0 = x_0 < x_1 < \dots < x_{N+1} = T$  et  $h = \max_{0 \leq i \leq N} (x_{i+1} - x_i)$

$S_d(\Delta) = \left\{ S(t) \text{ polynômes par morceaux, de degré } d, \text{ tels que } S(x) \in C^d[0, T] \text{ et } S(T) = 0 \right\}$

Dans ce cas, on peut montrer (Cfr Ref S-1) que si chaque composante de  $\lambda^*$  est  $d+1$  fois continûment différentiable par morceaux par rapport à  $t$ , il existe une constante  $K_d$  positive et indépendante de la subdivision  $\Delta$  telle que

$$\inf_{w \in S} |\lambda^* - w| \leq K_d h^d \|D^{d+1} \lambda^*\|_2$$

où  $h$  est le "modèle" de la subdivision  $\Delta$ .  
On obtient alors le théorème suivant

Théorème II.8

Sous les hypothèses précédentes, et si chaque composante de  $\lambda^*$  est  $d+1$  fois continûment différentiable par morceaux par rapport à  $t$ , il existe une constante  $K_d \geq 0$  indépendante de la subdivision  $\Delta$  telle que

$$\|\lambda^* - \lambda_{S_d(\Delta)}\|_2 \leq \lambda_g^{-1/2} \cdot K_d h^d \|D^{d+1} \lambda^*\|_2$$

$$\|u^* - u_{S_d(\Delta)}\|_2 \leq \|R^{-1}B^T\|_\infty \lambda_g^{-1/2} \cdot K_d h^d \|D^{d+1} \lambda^*\|_2$$

$$\|x^* - x_{S_d(\Delta)}\|_2 \leq \Gamma \|R^{-1}B^*\|_\infty \lambda_g^{-1/2} \zeta K_d h^d \|D^{d+1}\lambda^*\|_2$$

$$\|\dot{x}^* - \dot{x}_{S_d(\Delta)}\|_2 \leq (\Gamma \|A\|_\infty + \|B\|_\infty) \|R^{-1}B^*\|_\infty \lambda_g^{-1/2} \zeta K_d h^d \|D^{d+1}\lambda^*\|_2$$

$$J[u_{S_d(\Delta)}, x_{S_d(\Delta)}] - J[u^*, x^*] \geq 0$$

$$\geq \frac{1}{2} \lambda_g^{-1} \zeta^2 \|R^{-1}B^*\|_\infty^2 (\|g\|_\infty \Gamma^2 + \|R\|_\infty) K_d h^d \|D^{d+1}\lambda^*\|_2$$

La démonstration est évidente en utilisant tous les théorèmes démontrés précédemment.

### III Système différentiel non linéaire

Nous allons maintenant considérer un système plus général : celui introduit au Chapitre II, II, § 2

Réposons brièvement le problème.  
Considérons la fonctionnelle de coût

$$J[u] = \int_0^T g(x, u, t) dt \quad (\text{III } 23)$$

et les équations du système

$$\dot{x}(t) = f(x, u, t) \quad x(0) = x_0, \quad t \in [0, T] \quad (\text{III } 24)$$

où  $x(t)$  est un vecteur à  $n$  dimensions,  $u(t)$  un vecteur à  $r$  dimensions,  $f(x, u, t)$  un vecteur à  $n$  dimensions et  $g(x, u, t)$  une fonction scalaire.

Nous supposons de plus que

$$U = \{u\} = \{W_2^\alpha[0, T], \mathbb{R}^r\} \quad (\text{III } 25)$$

$$C = \{x\} = \{W_2^\alpha[0, T], \mathbb{R}^n\} \quad (\text{III } 26)$$

où  $\{W_2^\alpha[0, T], \mathbb{R}^k\}$  est l'espace de Sobolev des fonctions vectorielles de dimension  $k$  définies sur  $[0, T]$  qui ont leurs dérivées au sens des distributions dans  $\{L^2[0, T], \mathbb{R}^k\}$  jusqu'à l'ordre  $\alpha$ .

On peut aussi caractériser  $\{W_2^\alpha[0, T], \mathbb{R}^k\}$  comme la collection des fonctions  $z$  définies sur  $[0, T]$  telles que  $z^{(j)}$ ,  $j \leq \alpha-1$ , est absolument continue dans  $[0, T]$  et  $z^{(\alpha)} \in \{L^2[0, T], \mathbb{R}^k\}$ . La norme de Sobolev sur cet espace est donnée par

$$\|z\|_{2,\alpha}^2 = \sum_{i=1}^k \int_0^{T,i} [z_i^{(\alpha)}(t)]^2 dt \quad (\text{III.27})$$

Nous utiliserons aussi la norme  $L^2$  ( $\|.\|_{2,0}$  ou  $\|.\|_2$ )

**15**

### Formulations du problème et hypothèses

Notre problème consiste donc à minimiser  $J[u]$  sur tous les  $u$  appartenant à  $\{W_2^\alpha[0, T], \mathbb{R}^r\}$

Nous pouvons l'envisager sous l'angle du calcul des variations et construire le Lagrangien :

$$L[u, x; \lambda, \gamma] = J[u] + \int_0^T \langle \lambda(t), -\dot{x}(t) + f(x, u, t) \rangle dt + \langle x(0) - x_0, \gamma \rangle \quad (\text{III.28})$$

où  $\lambda$  et  $\gamma$  appartiennent respectivement à  $\{W_2^\alpha[0, T], \mathbb{R}^n\}'$  et  $\mathbb{R}^n$ . Comme ces espaces sont des Hilbert,  $\gamma \in \mathbb{R}^n$  et  $\lambda \in \{W_2^\alpha[0, T], \mathbb{R}^n\}'$  par identification (th. de Riesz). Il s'agit donc maintenant de trouver  $x^*$ ,  $u^*$ ,  $\lambda^*$  et  $\gamma^*$  tels que

$$L[u^*, x^*; \lambda^*, \gamma^*] = \sup_{\substack{\lambda \in \mathcal{E} \\ \gamma \in \mathbb{R}^n}} \inf_{\substack{u \in \mathcal{U} \\ x \in \mathcal{X}}} L[u, x; \lambda, \gamma] \quad (\text{III.29})$$

Il faut donc que les dérivées partielles (au sens de Fréchet) vérifient

$$\frac{\partial L}{\partial u}[u^*, x^*; \lambda^*, \gamma^*] = 0 \quad (\text{III.30})$$

$$\frac{\partial L}{\partial x}[u^*, x^*; \lambda^*, \gamma^*] = 0 \quad (\text{III.31})$$

$$\frac{\partial L}{\partial \lambda}[u^*, x^*; \lambda^*, \gamma^*] = 0 \quad (\text{III.32})$$

$$\frac{\partial L}{\partial \gamma}[u^*, x^*; \lambda^*, \gamma^*] = 0 \quad (\text{III.33})$$

et l'hypothèse

Faisons maintenant des hypothèses sur la régularité du problème (sous la première forme). Posons la

### Définition

Le problème  $\in BC^{\alpha}[0,T]$  ( $\alpha$  est un entier positif) ssi

(i)  $f(x, u, t)$  et  $g(x, u, t)$  sont  $\alpha+1$  fois continûment différentiables en  $x$  et  $u$  et  $\alpha$  fois en  $t$  pour  $t \in [0, T]$ .

(ii) les opérateurs (non linéaires)  $\frac{\partial^{i+j} f}{\partial x^i \partial u^j}$  et  $\frac{\partial^{i+j} g}{\partial x^i \partial u^j}$  ( $i, j \in \mathbb{N}_0$ )

et ( $i+j \leq \alpha+1$ ) transforment des domaines bornés de  $G \times U$  en domaines bornés de  $\{L^2[0, T], \mathbb{R}^n\}$  et  $L^2[0, T]$  respectivement.

Nous supposons maintenant

Hypothèse 1 : le problème  $\in BC^{\alpha}[0, T]$

Hypothèse 2 : Le lagrangien est extrémisé en  $(u^*, x^*, \lambda^*, y^*)$  satisfaisant les conditions (III.30) à (III.33)

Le principe de Pontryagin nous fournit les conditions nécessaires (cf I, II)

$$\left\{ \begin{array}{l} \frac{\partial g^*}{\partial u} + \left[ \frac{\partial f^*}{\partial u} \right]^t \lambda^* = 0 \\ \lambda^* + \left[ \frac{\partial f^*}{\partial x} \right]^t \lambda^* + \frac{\partial g^*}{\partial x} = 0 \quad \text{avec } \lambda^*(T) = 0 \end{array} \right. \quad (\text{III.34})$$

$$\left\{ \begin{array}{l} -\dot{x}^* + f(x^*, u^*, t) = 0 \quad \text{avec } x^*(0) = x_0 \\ y^* = -\lambda^*(0) \end{array} \right. \quad (\text{III.35})$$

$$(\text{III.36})$$

$$(\text{III.37})$$

Hypothèse 3 : La seconde variation de  $L$  (par rapport à  $(u, x)$ ) est largement positive dans un voisinage convexe de  $(u^*, x^*, \lambda^*)$

i.e. il existe des voisinages  $N(u^*) \subseteq U$ ,  $N(x^*) \subseteq G$  et  $N(\lambda^*) \subseteq G$  tels que l'opérateur  $H$  défini par

$$H[u, \dot{x}; \lambda] = \begin{bmatrix} \frac{\partial^2 g}{\partial u^2} + \left[ \frac{\partial^2 f}{\partial u^2} \right]^+ \lambda & \frac{\partial^2 g}{\partial u \partial \dot{x}} + \left[ \frac{\partial^2 f}{\partial u \partial \dot{x}} \right]^+ \lambda \\ \frac{\partial^2 g}{\partial \dot{x} \partial u} + \left[ \frac{\partial^2 f}{\partial \dot{x} \partial u} \right]^+ \lambda & \frac{\partial^2 g}{\partial \dot{x}^2} + \left[ \frac{\partial^2 f}{\partial \dot{x}^2} \right]^+ \lambda \end{bmatrix} \quad (\text{III.38})$$

satisfasse

$$\int_0^T \left\langle \begin{bmatrix} \delta u \\ \delta \dot{x} \end{bmatrix}, H[\bar{u}, \bar{\dot{x}}; \bar{\lambda}] \begin{bmatrix} \delta u \\ \delta \dot{x} \end{bmatrix} \right\rangle dt \geq \tau \left[ \|\delta u\|_2^2 + \|\delta \dot{x}\|_2^2 \right] \quad (\text{III.39})$$

où  $\delta u = u - \bar{u}$  et  $\delta \dot{x} = \dot{x} - \bar{\dot{x}}$  avec  $(\bar{u}, \bar{\dot{x}}, \bar{\lambda}) \in N(u^*) \times N(\dot{x}^*) \times N(\lambda^*)$ , et où  $\tau$  est une constante strictement positive.

À ce stade, quelques remarques s'imposent. D'abord, les hypothèses 1, 2 et 3 constituent un ensemble de conditions suffisantes locales d'existence et d'unicité de la solution du problème (Cfr Réf A-1). Ensuite, on peut voir, à l'aide d'arguments classiques en équations différentielles, que  $x^*, \lambda^* \in C^1([0, T], \mathbb{R}^n)$  et  $u^* \in C^1([0, T], \mathbb{R}^m)$  (Cfr Réf C-1). Enfin, nous interchangerons sup et inf dans (II.29), pourvu qu'on reste dans un voisinage convexe convenable du point optimal (Cfr Réf L-1).

Désirons maintenant le comportement du Lagrangien au point optimal

### Théorème III.9

Si  $x \in N(x^*)$ ,  $u \in N(u^*)$ ,  $\lambda \in N(\lambda^*)$  et  $\gamma^* \in N(\gamma^*)$ , le Lagrangien a un point de selle en  $(u^*, x^*, \lambda^*, \gamma^*)$  avec

$$L[u^*, x^*, \lambda^*, \gamma^*] = L[u^*, x^*, \lambda^*, \gamma^*] \leq L[u, x, \lambda, \gamma] \quad (\text{III.40})$$

Démonstration : L'égalité de gauche est immédiate, si l'on se rappelle que le couple  $(x^*, u^*)$  doit satisfaire (III.24) pour tout couple  $(\lambda, \gamma)$   $\in G \times \mathbb{R}^n$  admissible. L'inégalité de droite s'obtient en développant  $L$  en série de Taylor autour de  $*$ .

$$L[u, x; \lambda^*, \gamma^*] = L[u^*, x^*; \lambda^*, \gamma^*]$$

$$+ \int_0^T \left\langle \frac{\partial g(x^*, u^*, t)}{\partial u} + \left[ \frac{\partial f(x^*, u^*, t)}{\partial u} \right]^t, \Delta u \right\rangle dt$$

$$+ \int_0^T \left\langle \frac{\partial g(x^*, u^*, t)}{\partial x} + \left[ \frac{\partial f(x^*, u^*, t)}{\partial x} \right]^t \lambda^* + \dot{\lambda}^*, \Delta x \right\rangle dt$$

$$+ \langle \gamma^* + \lambda^*(0), \Delta x(0) \rangle - \langle \lambda^*(T), \Delta x(T) \rangle$$

$$+ \int_0^1 (1-\tau) \int_0^T \left\langle \begin{bmatrix} \Delta u \\ \Delta x \end{bmatrix}, H[u, x, \lambda^*] \begin{bmatrix} \Delta u \\ \Delta x \end{bmatrix} \right\rangle dt dx \quad (\text{III.41})$$

où  $\tilde{u} = \tau u + (1-\tau)u^*$ ,  $\tilde{x} = \tau x + (1-\tau)x^*$  et  $\Delta u = u - u^*$ ,  $\Delta x = x - x^*$  avec  $\tau \in [0, 1]$ . Le développement provient de la Réf D-1. Notons que toutes les variations du premier ordre dans (III.41) sont nulles en vertu des conditions (III.34) à (III.37). De plus l'Hypothèse 3 implique

$$L[u, x; \lambda^*, \gamma^*] \geq L[u^*, x^*; \lambda^*, \gamma^*] + \frac{r}{2} \left[ \|\Delta u\|_2^2 + \|\Delta x\|_2^2 \right]$$

qui démontre le théorème ■

A partir de maintenant, nous éliminerons  $\gamma$  du Lagrangien en ne considérant plus que les variations qui satisfont  $\lambda(0) = -\gamma$ .

Nous allons maintenant choisir, comme dans le cas précédent, des sous-espaces de dimension finie

25

Sous espaces de dimension finie et approximation de Ritz-Galerkin

Définissons les ensembles admissibles pour  $u, x$  et  $\lambda$  par  $D_u, D_x$  et  $D_\lambda$ , ensembles de fonctions vectorielles de dimension  $r$  et  $n$ , qui sont continues et strictement différentiables par morceaux sur  $[0, T]$

Soient

$$C = \{u | u \in U \cap N(u^*) \cap D_u\}$$

$$S = \{x | x \in G \cap N(x^*) \cap D_x\}$$

$$H = \{\lambda | \lambda \in G \cap N(\lambda^*) \cap D_\lambda\}$$

où  $N(u^*)$  et  $N(x^*)$  et  $N(\lambda^*)$  sont déterminés par l'hypothèse de positivité forte.

Choisissons

$$C_m(w) = \{u | u(t) = \sum_{i=1}^m \alpha_i w_{i,m}(t), t \in [0, T]\} \subseteq U$$

$$S_m(w) = \{x | x(t) = \sum_{i=1}^m \beta_i w_{i,m}(t), t \in [0, T]\} \subseteq G$$

$$H_m(w) = \{\lambda | \lambda(t) = \sum_{i=1}^m \gamma_i w_{i,m}(t), t \in [0, T]\} \subseteq G$$

où  $\{w_{i,m}\}$  représente des bases qu'on peut toujours choisir orthonormées, et où les  $\alpha_i$  sont des vecteurs de  $\mathbb{R}^r$  et les  $\beta_i$  et  $\gamma_i$  de  $\mathbb{R}^n$ .

Choisissons  $\{w_{i,m}\}$  tels que

$$C_m = \{u | u \in C_m(w) \cap \bar{C}\}$$

$$S_m = \{x | x \in S_m(w) \cap \bar{S}\}$$

$$H_m = \{\lambda | \lambda \in H_m(w) \cap \bar{H}\}$$

soient non-vides. Comme  $C_m$ ,  $S_m$  et  $H_m$  sont des ensembles fermés de  $W_2^1$ , on peut voir que ces ensembles possèdent des meilleures approximations de  $u^*$ ,  $x^*$  et  $\lambda^*$  uniques. (Cf Rq L1, C-1)

On a donc

$$\| \hat{u} - u^* \|_{2,1} = \inf_{u \in C_m} \| u - u^* \|_{2,1}$$

$$\| \hat{x} - x^* \|_{2,1} = \inf_{x \in S_m} \| x - x^* \|_{2,1}$$

$$\| \hat{\lambda} - \lambda^* \|_{2,1} = \inf_{\lambda \in H_m} \| \lambda - \lambda^* \|_{2,1}$$

relations qui déterminent les "meilleures approximations"  $\bar{u}$ ,  $\bar{x}$  et  $\bar{\lambda}$ .

Posons maintenant

$$\epsilon(\bar{u}) = \bar{u} - u^*, \quad \epsilon(\bar{x}) = \bar{x} - x^*, \quad \epsilon(\bar{\lambda}) = \bar{\lambda} - \lambda^* \quad (\text{III}42)$$

Nous supposons évidemment que  $C_m$ ,  $S_m$  et  $H_m$  sont de "bons" sous-espaces approximatifs, c.-à-d que  $\|\epsilon(\bar{u})\|_{2,1}$ ,  $\|\epsilon(\bar{x})\|_{2,1}$  et  $\|\epsilon(\bar{\lambda})\|_{2,1}$  tendent vers zéro quand  $m$  tend vers l'infini.

Posons maintenant le problème dans les sous-espaces considérés : il faut minimiser  $J[\bar{u}]$  sur  $C_m$ . Soit

$$J[\bar{u}] = \inf_{u \in C_m} J[u]$$

sous la contrainte

$$\left\{ \begin{array}{l} \int_0^T \langle w_{j,m}(t), -\dot{\bar{x}}(t) + f(\bar{x}(t), \bar{u}(t), t) \rangle dt = 0 \\ \bar{x}(0) = x_0 \end{array} \right.$$

avec  $\bar{x} \in S_m$  et  $j=1, m$

En termes de sup inf, cela donne

$$L[\bar{u}, \bar{x}; \bar{\lambda}] = \sup_{\lambda \in H_m} \inf_{\substack{x \in S_m \\ u \in C_m}} L[u, x; \lambda]$$

Pour les dérivées, il faut trouver  $(\bar{u}, \bar{x}, \bar{\lambda}) \in C_m \times S_m \times H_m$  tel que

$$\frac{\partial L}{\partial u}(\bar{u}, \bar{x}, \bar{\lambda}) = 0$$

$$\frac{\partial L}{\partial x}(\bar{u}, \bar{x}, \bar{\lambda}) = 0$$

$$\frac{\partial L}{\partial \lambda}(\bar{u}, \bar{x}, \bar{\lambda}) = 0$$

Ces différentes formulations, bien sûr équivalentes, fournissent le système d'équations algébriques suivant

$$\int_0^T \left\langle w_{j,m}(t), \frac{\partial}{\partial u}(\bar{x}, \bar{u}, t) + \left[ \frac{\partial f}{\partial u}(\bar{x}, \bar{u}, t) \right]^t \bar{\lambda} \right\rangle dt = 0 \quad (\text{III-43})$$

$$\int_0^T \left\langle w_{j,m}(t), \dot{\bar{\lambda}} + \left[ \frac{\partial f}{\partial x}(\bar{x}, \bar{u}, t) \right]^t \bar{\lambda} + \frac{\partial g}{\partial x}(\bar{x}, \bar{u}, t) \right\rangle dt = 0 \quad (\text{III-44})$$

$$\bar{\lambda}(T) = 0 \quad (\text{III-45})$$

$$\int_0^T \left\langle w_{j,m}(t), -\dot{\bar{x}} + \dot{f}(\bar{x}, \bar{u}, t) \right\rangle dt = 0 \quad (\text{III-46})$$

$$\bar{x}(0) = x_0 \quad (\text{III-47})$$

où  $(\bar{u}, \bar{x}, \bar{\lambda}) \in C_m \times S_m \times H_m$  et  $j=1, m$ .

Exactement comme plus haut, on peut déduire le

### Théorème III-10

Soit  $(\bar{u}, \bar{x}, \bar{\lambda})$  la solution du système (III-43) à (III-47)

Si  $(u, x, \lambda) \in C_m \times S_m \times H_m$ , le Lagrangien a un point de selle dégénéré en  $(\bar{u}, \bar{x}, \bar{\lambda})$ , avec

$$L[\bar{u}, \bar{x}, \bar{\lambda}] = L[\bar{u}, \bar{x}, \bar{\lambda}] \leq L[u, x, \lambda] \quad (\text{III-48})$$

Nous allons maintenant démontrer la "qualité" de l'approximation de Ritz-Galerkin en dérivant des bornes d'erreurs données en termes des normes de  $\epsilon(\hat{u})$ ,  $\epsilon(\hat{x})$ ,  $\epsilon(\hat{\lambda})$ ,  $\epsilon(\hat{e})$  et  $\epsilon(\hat{s})$ , où les deux dernières sont définies par

$$\epsilon(\hat{e}) = \hat{e} - \bar{e}, \quad \epsilon(\hat{s}) = \hat{s} - \bar{s}$$

avec  $(\hat{e}, \hat{s})$  la meilleure approximation  $L^2$  de  $(\bar{e}, \bar{s})$ . La définition formelle de  $e$  et  $s$  sera faite ultérieurement.

Définissons d'abord les quantités suivantes :

$$\mathcal{J}^* = \mathcal{J}[u^*] = L[u^*, x^*, \lambda^*]$$

(III 49)

$$\mathcal{J}_m = \mathcal{J}[\bar{u}] = L[\bar{u}, \bar{x}; \bar{\lambda}]$$

(III 50)

où  $(\bar{u}, \bar{x}, \bar{\lambda})$  est la solution du système de Galerkin (III 43) à (III 47)

### 35 Convergence du couple $(x, u)$ .

Nous démontrons d'abord quelques résultats préparatoires

Théorème III 11

Soit  $(u, x, \lambda) \in P_m = C_m \times S_m \times M_m$ . Il existe alors des constantes  $k_1, k_2 > 0$  telles que

$$\left\| \frac{\partial g^*}{\partial u} + \left[ \frac{\partial f^*}{\partial u} \right]^t \lambda \right\|_2 \leq k_1 (\|\Delta u\|_2 + \|\Delta x\|_2) + \|\ell\|_2 \quad (\text{III 51})$$

et

$$\left\| \dot{\lambda} + \left[ \frac{\partial f^*}{\partial x} \right]^t \lambda + \frac{\partial \ell^*}{\partial x} \right\|_2 \leq k_2 (\|\Delta u\|_2 + \|\Delta x\|_2) + \|\dot{\ell}\|_2 \quad (\text{III 52})$$

$$\text{ou } \ell(t) = \ell(x, u, \lambda; t) = \frac{\partial g}{\partial u}(x, u, t) + \frac{\partial f}{\partial u}[(x, u, t)]^t \lambda \quad (\text{III 53})$$

$$\dot{\lambda}(t) = \dot{\lambda}(x, u, \lambda, t) = \dot{\lambda} + \left[ \frac{\partial f}{\partial x}(x, u, t) \right]^t \lambda + \frac{\partial g}{\partial x}(x, u, t) \quad (\text{III 54})$$

avec  $\Delta x = x - x^*$  et  $\Delta u = u - u^*$

Démonstration Nous avons évidemment

$$\frac{\partial g^*}{\partial u} + \left[ \frac{\partial f^*}{\partial u} \right]^t \lambda = \ell(t) + \left[ \frac{\partial g^*}{\partial u} - \frac{\partial g}{\partial u}(x, u, t) \right] + \left\{ \left[ \frac{\partial f^*}{\partial u} \right]^t - \left[ \frac{\partial f}{\partial u}(x, u, t) \right]^t \right\} \lambda$$

En appliquant le théorème de la valeur moyenne pour les opérateurs, dont les hypothèses font intervenir les hypothèses 1, (cf Réf D), on obtient

$$\left\| \frac{\partial g^*}{\partial u} + \left[ \frac{\partial f^*}{\partial u} \right]^t \lambda \right\|_2 \leq \|\ell\|_2 + k_g [\|\Delta u\|_2 + \|\Delta x\|_2] + k_f [\|\Delta u\|_2 + \|\Delta x\|_2]$$

qui se ramène à (III 51) en posant  $k_1 = k_g + k_f$ .

La preuve de (III 52) est identique.

### Théorème III 12

Si  $(\bar{u}, \bar{x}, \bar{\lambda})$  est la solution du système de Galerkin sur  $P_m$ . Soient  $\hat{e}$  et  $\hat{z}$  les meilleures approximations  $L^2$  de  $e(x, \bar{u}, \bar{\lambda}, t)$  et  $z(\bar{x}, \bar{u}, \bar{\lambda}, t)$  respectivement. Alors

$$\|\bar{e}\|_2 = \|e(\hat{e})\|_2$$

$$\|\bar{z}\|_2 = \|z(\hat{z})\|_2$$

(III 55)

(III 56)

Démonstration : Il suffit de voir que  $\hat{e}$  et  $\hat{z}$  sont nulles. Calculons les coefficients de Fourier pour  $\bar{e}$  et  $\bar{z}$

$$\langle \bar{e}, w_{i,m} \rangle_{L^2} = \int_0^T w_{i,m}(t) e(\bar{x}, \bar{u}, \bar{\lambda}, t) dt = 0$$

par (III 43) pour  $j=1, m$ .

$$\langle \bar{z}, w_{i,m} \rangle_{L^2} = \int_0^T w_{i,m}(t) z(\bar{x}, \bar{u}, \bar{\lambda}, t) dt = 0$$

par (III 44) pour  $j=1, m$ .

Le théorème a pour seul but d'uniformiser les notations.

Démontrons maintenant la convergence.

### Théorème II 13

Si  $(\bar{u}, \bar{x}, \bar{\lambda}) \in P_m$  est la solution du système de Galerkin, il existe des constantes  $k, k', k'' > 0$  indépendantes de  $m$ , telles que

$$\|\bar{u} - u^*\|_2 \leq \eta_m$$

$$\|\bar{x} - x^*\|_2 \leq \tau_m$$

(III 57)

(III 58)

où les constantes  $\eta_m$  et  $\tau_m$  sont données par les relations suivantes :

$$\eta_m = k'' [\|\epsilon(\hat{u})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_{2,1}] \\ + k [\|\epsilon(\hat{e})\|_2 \|\epsilon(\hat{u})\|_2 + \|\epsilon(\hat{z})\|_2 \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{e})\|_2^2 + \|\epsilon(\hat{u})\|_2^2 + \|\epsilon(\hat{\lambda})\|_{2,1}^2]^{1/2}$$
(III.5)

$$\Gamma_m = k' [\|\epsilon(\hat{u})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2] \\ + k [\|\epsilon(\hat{e})\|_2 \|\epsilon(\hat{u})\|_2 + \|\epsilon(\hat{z})\|_2 \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{u})\|_2^2 + \|\epsilon(\hat{e})\|_2^2 + \|\epsilon(\hat{\lambda})\|_{2,1}^2]^{1/2}$$
(III.6)

Démonstration : Utilisons le théorème III.10 en posant  $u = \hat{u}$  et  $u = \hat{u}$ . On obtient

$$J_m \leq L[\hat{u}, \hat{x}, \hat{\lambda}]$$

Développons maintenant le membre de droite en série de Taylor autour de  $(u^*, x^*, \lambda^*)$ . On obtient alors

$$J_m \leq L[u^*, x^*, \bar{\lambda}] + \int_0^T \left\langle \frac{\partial g^*}{\partial u} + \left[ \frac{\partial f^*}{\partial u} \right]^T \bar{\lambda}, \epsilon(\hat{u}) \right\rangle dt$$

$$+ \int_0^T \left\langle \bar{\lambda} + \left[ \frac{\partial f^*}{\partial x} \right]^T \bar{\lambda} + \frac{\partial g^*}{\partial x}, \epsilon(\hat{x}) \right\rangle dt$$

$$+ \int_0^1 (1-\tau) \int_0^T \left\langle \begin{bmatrix} \epsilon(\hat{e}) \\ \epsilon(\hat{z}) \end{bmatrix}, H(u^*, x^*, \lambda) \begin{bmatrix} \epsilon(\hat{u}) \\ \epsilon(\hat{x}) \end{bmatrix} \right\rangle dt d\tau$$

où  $\tilde{x} = \tau x^* + (1-\tau)\hat{x}$  et  $\tilde{u} = \tau u^* + (1-\tau)\hat{u}$ .

D'autre part, nous avons

$$L[u^*, x^*, \lambda^*] = J^* = L[u^*, x^*, \bar{\lambda}] = J[u^*, x^*, \hat{\lambda}] \quad (\text{III.6})$$

en vertu du théorème III.9

Prenons les normes

$$J_m \leq J^* + \left\| \frac{\partial g^*}{\partial u} + \left[ \frac{\partial f^*}{\partial u} \right]^T \bar{\lambda} \right\|_2 \|\epsilon(\hat{u})\|_2 \\ + \left\| \bar{\lambda} + \left[ \frac{\partial f^*}{\partial x} \right]^T \bar{\lambda} + \frac{\partial g^*}{\partial x} \right\|_2 \|\epsilon(\hat{x})\|_2 \\ + k [\|\epsilon(\hat{u})\|_2^2 + \|\epsilon(\hat{x})\|_2^2]$$

En utilisant le théorème III.11 on déduit immédiatement

$$\begin{aligned}
J_m &\leq J^* + c \left[ \|\bar{u} - u^*\|_2 (\|E(\bar{u})\|_2 + \|E(\hat{u})\|_2) \right. \\
&\quad \left. + \|\bar{x} - x^*\|_2 (\|E(\bar{x})\|_2 + \|E(\hat{x})\|_2) \right] \\
&\quad + \|E(\hat{u})\|_2 \|E(\hat{u})\|_2 + \|E(\hat{x})\|_2 \|E(\hat{x})\|_2 \\
&\quad + k (\|E(\bar{u})\|_2^2 + \|E(\hat{x})\|_2^2)
\end{aligned} \tag{III.62}$$

Utilisons maintenant l'égalité du théorème (III.10) où nous choisissons  $\lambda = \hat{\lambda}$ . Nous avons donc

$$J_m = L[\bar{u}, \bar{x}; \hat{\lambda}]$$

Développons le terme de droite en série de Taylor autour de  $(u^*, x^*, \hat{\lambda})$

$$\begin{aligned}
J_m &= L[u^*, x^*; \hat{\lambda}] + \int_0^T \left\langle \frac{\partial g^*}{\partial u} + \left[ \frac{\partial f^*}{\partial u} \right]^t \hat{\lambda}, \bar{u} - u^* \right\rangle dt \\
&\quad + \int_0^T \left\langle \dot{\hat{\lambda}} + \left[ \frac{\partial f^*}{\partial x} \right]^t \hat{\lambda} + \frac{\partial g^*}{\partial x}, \bar{x} - x^* \right\rangle dt \\
&\quad + \int_0^1 (1-\tau) \int_0^T \left\langle \begin{bmatrix} \bar{u} - u^* \\ \bar{x} - x^* \end{bmatrix}, H(u_\tau, x_\tau; \hat{\lambda}) \begin{bmatrix} \bar{u} - u^* \\ \bar{x} - x^* \end{bmatrix} \right\rangle dt d\tau
\end{aligned}$$

où  $u_\tau = \tau \bar{u} + (1-\tau)u^*$  et  $x_\tau = \tau \bar{x} + (1-\tau)x^*$ . En utilisant (III.61) et la positivité forte on obtient

$$\begin{aligned}
J_m &\geq J^* - \left| \int_0^T \left\langle \frac{\partial g^*}{\partial u} + \left[ \frac{\partial f^*}{\partial u} \right]^t \hat{\lambda}, \bar{u} - u^* \right\rangle dt \right| - \left| \int_0^T \left\langle \dot{\hat{\lambda}} + \left[ \frac{\partial f^*}{\partial x} \right]^t \hat{\lambda} + \frac{\partial g^*}{\partial x}, \bar{x} - x^* \right\rangle dt \right| \\
&\quad + \frac{1}{2} [\|\bar{u} - u^*\|_2^2 + \|\bar{x} - x^*\|_2^2] \\
&\geq J^* - \left\| \frac{\partial g^*}{\partial u} + \left[ \frac{\partial f^*}{\partial u} \right]^t \hat{\lambda} \right\|_2 \|\bar{u} - u^*\|_2 - \left\| \dot{\hat{\lambda}} + \left[ \frac{\partial f^*}{\partial x} \right]^t \hat{\lambda} + \frac{\partial g^*}{\partial x} \right\|_2 \|\bar{x} - x^*\|_2 \\
&\quad + \frac{1}{2} [\|\bar{u} - u^*\|_2^2 + \|\bar{x} - x^*\|_2^2]
\end{aligned}$$

Or, nous savons que

$$\begin{aligned}
\left\| \frac{\partial g^*}{\partial u} + \left[ \frac{\partial f^*}{\partial u} \right]^t \hat{\lambda} \right\|_2 &= \left\| -\frac{\partial f^*}{\partial u} - \left[ \frac{\partial f^*}{\partial u} \right]^t \lambda^* + \frac{\partial g^*}{\partial u} + \left[ \frac{\partial f^*}{\partial u} \right]^t \hat{\lambda} \right\|_2 \\
&= \left\| \left[ \frac{\partial f^*}{\partial u} \right]^t E(\hat{\lambda}) \right\|_2 \leq k_0 \|E(\hat{\lambda})\|_2
\end{aligned}$$

De même

$$\begin{aligned}\|\dot{\lambda} + \left[ \frac{\partial J^*}{\partial x} \right]^T \lambda + \frac{\partial g^*}{\partial x} \|_2 &= \left\| -\frac{\partial g^*}{\partial x} - \left[ \frac{\partial J^*}{\partial x} \right]^T \lambda^* - \lambda^* + \frac{\partial g^*}{\partial x} + \left[ \frac{\partial J^*}{\partial x} \right]^T \dot{\lambda} + \dot{\lambda} \right\|_2 \\ &= \|\varepsilon(\dot{\lambda}) + \left[ \frac{\partial J^*}{\partial x} \right]^T \varepsilon(\hat{\lambda})\|_2 \leq \|\varepsilon(\dot{\lambda})\|_2 + \left\| \left[ \frac{\partial J^*}{\partial x} \right]^T \varepsilon(\hat{\lambda}) \right\|_2 \\ &\leq k' [\|\varepsilon(\dot{\lambda})\|_2 + \|\varepsilon(\hat{\lambda})\|_2] \leq k'' \|\varepsilon(\hat{\lambda})\|_{2,1}\end{aligned}$$

en effet

$$\|\varepsilon(\hat{\lambda})\|_{2,1}^2 = \|\varepsilon(\hat{\lambda})\|_2^2 + \|\varepsilon(\dot{\lambda})\|_2^2$$

implique

$$\|\varepsilon(\hat{\lambda})\|_{2,1}^2 \geq \|\varepsilon(\hat{\lambda})\|_2^2 \text{ et donc } \|\varepsilon(\hat{\lambda})\|_{2,1} \geq \|\varepsilon(\hat{\lambda})\|_2$$

$$\|\varepsilon(\hat{\lambda})\|_{2,1}^2 \geq \|\varepsilon(\dot{\lambda})\|_2^2 \text{ et donc } \|\varepsilon(\hat{\lambda})\|_{2,1} \geq \|\varepsilon(\dot{\lambda})\|_2$$

$$\text{ce qui donne } \|\varepsilon(\dot{\lambda})\|_2 + \|\varepsilon(\hat{\lambda})\|_2 \leq 2 \|\varepsilon(\hat{\lambda})\|_{2,1}$$

On a choisi les constantes comme suit :  $k' = \sup [1, k]$ ;  $k'' = 2k'$

Si on choisit  $c_0 = \sup [k_0, k'']$ , on obtient alors

$$\begin{aligned}J_m &\geq J^* - c_0 [\|\bar{u} - u^*\|_2 \|\varepsilon(\hat{\lambda})\|_2 + \|\bar{x} - x^*\|_2 \|\varepsilon(\hat{\lambda})\|_{2,1}] \quad (\text{III} \cdot 63) \\ &\quad + \frac{\sigma}{2} [\|\bar{u} - u^*\|_2^2 + \|\bar{x} - x^*\|_2^2]\end{aligned}$$

En combinant cette inégalité avec (III.62), on obtient

$$\begin{aligned}&\|\varepsilon(\hat{c})\|_2 \|\varepsilon(\hat{u})\|_2 + \|\varepsilon(\hat{s})\|_2 \|\varepsilon(\hat{x})\|_2 + k [\|\varepsilon(\hat{u})\|_2^2 + \|\varepsilon(\hat{x})\|_2^2] \\ &\geq - \|\bar{u} - u^*\|_2 [c_0 \|\varepsilon(\hat{\lambda})\|_2 + C (\|\varepsilon(\hat{u})\|_2 + \|\varepsilon(\hat{x})\|_2)] \\ &\quad - \|\bar{x} - x^*\|_2 [c_0 \|\varepsilon(\hat{\lambda})\|_{2,1} + C (\|\varepsilon(\hat{u})\|_2 + \|\varepsilon(\hat{x})\|_2)] \\ &\quad + \sigma [\|\bar{u} - u^*\|_2^2 + \|\bar{x} - x^*\|_2^2]\end{aligned}$$

Or

$$\begin{aligned}&k^2 (\|\varepsilon(\hat{c})\|_2 \|\varepsilon(\hat{u})\|_2 + \|\varepsilon(\hat{s})\|_2 \|\varepsilon(\hat{x})\|_2 + \|\varepsilon(\hat{u})\|_2^2 + \|\varepsilon(\hat{x})\|_2^2 + \|\varepsilon(\hat{\lambda})\|_{2,1}^2) \\ &\geq [\|\bar{u} - u^*\|_2 - k'' (\|\varepsilon(\hat{u})\|_2 + \|\varepsilon(\hat{x})\|_2 + \|\varepsilon(\hat{\lambda})\|_2)]^2 \quad (\text{III} \cdot 64) \\ &\quad + [\|\bar{x} - x^*\|_2 - k'' (\|\varepsilon(\hat{u})\|_2 + \|\varepsilon(\hat{x})\|_2 + \|\varepsilon(\hat{\lambda})\|_{2,1})]^2\end{aligned}$$

Le second membre est plus grand que

$$\begin{aligned}
 & [\|\bar{u} - u^*\|_2^2 - 2k' \|\bar{u} - u^*\|_2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2) \\
 & + k'^2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2)^2 + \|\bar{x} - x^*\|_2^2 \\
 & - 2k'' \|\bar{x} - x^*\|_2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2)_{2,1}^2 \\
 & + k''^2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2)_{2,1}^2] \\
 \geq & \|\bar{u} - u^*\|_2^2 - 2k' \|\bar{u} - u^*\|_2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2) \\
 & + \|\bar{x} - x^*\|_2^2 - 2k'' \|\bar{x} - x^*\|_2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2)_{2,1} \\
 & + k''^2 (\|\epsilon(\hat{\lambda})\|_{2,1}^2 + \|\epsilon(\hat{x})\|_2^2 + \|\epsilon(\hat{w})\|_2^2)
 \end{aligned}$$

Donc

$$\begin{aligned}
 & \|\epsilon(\hat{e})\|_2 \|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{z})\|_2 \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{w})\|_2^2 + \|\epsilon(\hat{x})\|_2^2 \\
 \geq & \left\{ \|\bar{u} - u^*\|_2^2 + \|\bar{x} - x^*\|_2^2 + (k''^2 - 1) \|\epsilon(\hat{\lambda})\|_{2,1}^2 + k''^2 (\|\epsilon(\hat{x})\|_2^2 + \|\epsilon(\hat{w})\|_2^2) \right. \\
 & - 2k' \|\bar{u} - u^*\|_2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2) \\
 & \left. - 2k'' \|\bar{x} - x^*\|_2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2)_{2,1} \right\} \frac{1}{k_0^2}
 \end{aligned}$$

C'est à dire

$$\begin{aligned}
 & \|\epsilon(\hat{e})\|_2 \|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{z})\|_2 \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{x})\|_2^2 + k (\|\epsilon(\hat{w})\|_2^2 + \|\epsilon(\hat{x})\|_2^2) \\
 \geq & (k + k''^2 - 1) (\|\epsilon(\hat{w})\|_2^2 + \|\epsilon(\hat{x})\|_2^2) + \frac{1}{k_0^2} \|\bar{u} - u^*\|_2^2 + \|\bar{x} - x^*\|_2^2 \\
 & - 2k' \|\bar{u} - u^*\|_2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2) \\
 & - 2k'' \|\bar{x} - x^*\|_2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2)_{2,1} \\
 \geq & \frac{1}{k_0^2} \left\{ \|\bar{x} - x^*\|_2^2 + \|\bar{u} - u^*\|_2^2 - 2k'' \|\bar{x} - x^*\|_2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2)_{2,1} \right. \\
 & \left. - 2k' \|\bar{u} - u^*\|_2 (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2 + \|\epsilon(\hat{\lambda})\|_2) \right\} \\
 \geq & \sigma \left( \|\bar{x} - x^*\|_2^2 + \|\bar{u} - u^*\|_2^2 \right) - \|\bar{u} - u^*\|_2 [C_0 \|\epsilon(\hat{\lambda})\|_2 + C (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2)] \\
 & - \|\bar{x} - x^*\|_2 [C_0 \|\epsilon(\hat{\lambda})\|_{2,1} + C (\|\epsilon(\hat{w})\|_2 + \|\epsilon(\hat{x})\|_2)]
 \end{aligned}$$

en choisissant  $k_0 \geq \sigma^{-1/2}$ ,  $k' = \sup \{C_0, C\} \frac{k_0^2}{2}$  et  $k'' = \sup \{1, k'\}$ , ce qui prouve (III.64).

Résolvons (III 14) pour  $\bar{u}-u^*$  et  $\bar{x}-x^*$  séparément. Nous obtenons alors rapidement la thèse. ■

45

### Convergence de la fonctionnelle de coût $J_m$

Nous allons maintenant chercher des bornes d'erreur pour la fonctionnelle de coût. Pour ce faire, nous allons introduire une approximation intermédiaire : nous cherchons dans  $N(u^*) \times N(x^*) \times M_m$  un triplet  $(\bar{u}, \bar{x}, \bar{\lambda})$  tel que

$$J_m = L[\bar{u}, \bar{x}; \bar{\lambda}] = \sup_{\lambda \in M_m} \inf_{\substack{u \in N(u^*) \\ x \in N(x^*)}} L[u, x; \lambda]$$

Nous pouvons alors obtenir les résultats suivants

Théorème III 14

Sous les hypothèses précédentes, il existe une constante  $k_1 > 0$  telle que

$$J^* - k_1 \| \epsilon(\hat{\lambda}) \|_{2,1}^2 \leq J_m \leq J^*$$

(III 15)

Démonstration : L'inégalité de droite est immédiate.

D'autre part

$$J_m \geq \inf_{\substack{x \in N(x^*) \\ u \in N(u^*)}} L[u, x; \bar{\lambda}]$$

où  $\bar{\lambda}$  est défini comme plus haut. Développons le membre du droit en série de Taylor autour de  $(u^*, x^*, \hat{\lambda})$ :

$$L[u, x; \bar{\lambda}] = L[u^*, x^*; \hat{\lambda}] + \int_0^T \left\langle \frac{\partial g^k}{\partial u} + \left[ \frac{\partial f^k}{\partial u} \right]^{t_k}, \Delta u \right\rangle dt$$

$$+ \int_0^T \left\langle \hat{\lambda} + \left[ \frac{\partial f^k}{\partial x} \right]^{t_k} \hat{\lambda} + \frac{\partial g^k}{\partial x}, \Delta x \right\rangle dt$$

(cont.)

$$+ \int_0^1 (1-\tau) \int_0^\tau \left\langle \begin{bmatrix} \Delta u \\ \Delta x \end{bmatrix}, H(\bar{x}, \bar{z}; \hat{\lambda}) \begin{bmatrix} \Delta u \\ \Delta x \end{bmatrix} \right\rangle dt d\tau$$

où  $\Delta x = x - x^*$ ,  $\Delta u = u - u^*$ ,  $\bar{x} = \tau x + (1-\tau)x^*$  et  $\bar{u} = \tau u + (1-\tau)u^*$ .

Exactement comme dans le théorème III.13 nous obtenons

$$\begin{aligned} L[u, x; \hat{\lambda}] &\geq L^* - k'_1 \|e(\hat{\lambda})\|_2 \|\Delta u\|_2 \\ &\quad - k'_2 \|e(\hat{\lambda})\|_2 \|\Delta x\|_2 + \frac{\epsilon}{2} (\|\Delta u\|_2^2 + \|\Delta x\|_2^2) \end{aligned}$$

Or

$$\begin{aligned} L^* - k'_1 \|e(\hat{\lambda})\|_{2,1}^2 + \frac{k'_1}{2} [(\|\Delta u\|_2 - \|e(\hat{\lambda})\|_{2,1})^2 + (\|\Delta x\|_2 - \|e(\hat{\lambda})\|_2)^2] \\ = L^* - k'_1 \|e(\hat{\lambda})\|_{2,1}^2 + \frac{k'_1}{2} [\|\Delta x\|_2^2 - 2\|\Delta x\|_2 \|e(\hat{\lambda})\|_{2,1} + \|e(\hat{\lambda})\|_{2,1}^2 \\ \quad + \|\Delta u\|_2^2 - 2\|\Delta u\|_2 \|e(\hat{\lambda})\|_2 + \|e(\hat{\lambda})\|_2^2] \\ = L^* + \frac{k'_1}{2} [\|\Delta x\|_2^2 - 2\|\Delta x\|_2 \|e(\hat{\lambda})\|_{2,1} - \|e(\hat{\lambda})\|_{2,1}^2 + \|\Delta u\|_2^2 \\ \quad - 2\|\Delta u\|_2 \|e(\hat{\lambda})\|_2 + \|e(\hat{\lambda})\|_2^2] \\ \leq L^* + \frac{k'_1}{2} [\|\Delta x\|_2^2 - 2\|\Delta x\|_2 \|e(\hat{\lambda})\|_{2,1} + \|\Delta u\|_2^2 - 2\|\Delta u\|_2 \|e(\hat{\lambda})\|_2] \\ \leq L^* - k'_1 \|e(\hat{\lambda})\|_2 \|\Delta u\|_2 - k'_2 \|e(\hat{\lambda})\|_{2,1} \|\Delta x\|_2 + \frac{\epsilon}{2} (\|\Delta u\|_2^2 + \|\Delta x\|_2^2) \end{aligned}$$

si l'on choisit  $k_1 = \inf[\sigma, \sup(k'_1, k'_2)]$

Nous avons donc

$$\begin{aligned} L[u, x; \hat{\lambda}] &\geq L^* - k_1 \|e(\hat{\lambda})\|_{2,1}^2 + \frac{k_1}{2} \{(\|\Delta x\|_2 - \|e(\hat{\lambda})\|_{2,1})^2 \\ &\quad + (\|\Delta u\|_2 - \|e(\hat{\lambda})\|_2)^2\} \end{aligned}$$

Comme les termes entre crochets sont toujours positifs, on peut les supprimer et l'on obtient bien la théorie.

### Théorème III.15

Si  $J_m = L[\bar{u}, \bar{x}, \bar{\lambda}]$  avec  $(\bar{u}, \bar{x}, \bar{\lambda}) \in P_m$ , il existe des constantes  $k_1, k_2 > 0$  telles que

$$-J_m + J^* \leq J_m \leq J^* + J_m$$

(III.16)

où  $\mu_m = k_1 \|\varepsilon(\hat{\lambda})\|_{2,1}^2$ ,  
et

$$\begin{aligned} J_m &= k_2 [(\gamma_m + \tau_m + \|\varepsilon(\hat{e})\|_2) \|\varepsilon(\hat{u})\|_2 \\ &\quad + (\gamma_m + \tau_m + \|\varepsilon(\hat{z})\|_2) \|\varepsilon(\hat{x})\|_2 + \|\varepsilon(\hat{u})\|_2^2 + \|\varepsilon(\hat{x})\|_2^2] \\ \text{avec } \gamma_m \text{ et } \tau_m \text{ définis par (III59) et (III60)} \end{aligned}$$
(III67)
(III68)

Démonstration. Commençons par remarquer que  
 $I_m \leq J_m$ .

Le théorème précédent nous permet alors d'affirmer

$$J^* \leq k_1 \|\varepsilon(\hat{\lambda})\|_{2,1}^2 \leq I_m \leq J_m$$

Ce qui prouve la borne inférieure.

D'autre part, en substituant (III57) et (III58) dans (III68) on obtient

$$\begin{aligned} J_m &\leq J^* + c [\gamma_m (\|\varepsilon(\hat{u})\|_2 + \|\varepsilon(\hat{x})\|_2) \\ &\quad + \tau_m (\|\varepsilon(\hat{u})\|_2 + \|\varepsilon(\hat{x})\|_2) + \|\varepsilon(\hat{e})\|_2 \|\varepsilon(\hat{u})\|_2 \\ &\quad + \|\varepsilon(\hat{z})\|_2 \|\varepsilon(\hat{x})\|_2 + k (\|\varepsilon(\hat{u})\|_2^2 + \|\varepsilon(\hat{x})\|_2^2)] \end{aligned}$$

ce qui donne, en rassemblant les termes, la borne supérieure ■

58

### Une particularisation de $P_m$

Nous allons maintenant indiquer quelques résultats pour un choix particulier du sous espace  $P_m$ : un espace de functions splines.

#### Définition

Supposons  $H_m^a$  ( $a \geq 1, m \geq 1$ ) un espace de dimension  $m$  de polynômes permanents d'ordre  $a-1$ , tel que

- 1)  $\exists L_m$  opérateur linéaire :  $PC^a[0, T] \rightarrow H_m^a$
- 2)  $\forall f \in PC^a[0, T] \quad \|L_m f - f\|_2 = O(m^{-a})$

- 3)  $\forall f \in PC^\alpha[0, T] \quad \| \frac{d^\beta}{dt^\beta} (\mathcal{L}_m f - f) \|_2 = O(h^{\alpha-\beta}), \quad 0 \leq \beta \leq \alpha$
- 4)  $(\mathcal{L}_m f)^{(i)}(0) = f^{(i)}(0)$  et  $(\mathcal{L}_m f)^{(i)}(T) = f^{(i)}(T), \quad i = 0, 1$
- 5)  $H_m^\alpha \subseteq \mathcal{C}^\alpha$  où  $\mathcal{C}^\alpha$  est un "ensemble admissible" (ex:  $G, U$ )  
avec  $PC^\alpha[0, T] = \{f \mid f^{(i)} \text{ est continue par morceaux sur } [0, T]\}$

Il est classique qu'il existe des sous-espaces existent, à condition d'avoir des "bonnes conditions" sur  $\mathcal{C}^\alpha$ . Comme, dans notre cas,  $G$  et  $U$  sont inclus dans  $W_2^\alpha$ , toutes ces conditions sont satisfaites.

Si nous définissons un découpage de  $[0, T]$  tel que

$$0 = t_0 < t_1 < \dots < t_q = T$$

avec  $h = \Delta t = \max_{i=1, q} |t_i - t_{i-1}|$

nous pouvons alors choisir  $H_m^\alpha(\Pi)$  l'espace des splines d'ordre  $\alpha$  pour lesquelles

$$\left\| \frac{d^\beta}{dt^\beta} (\mathcal{L}_m f - f) \right\|_2 = O(h^{\alpha-\beta}) \quad 0 \leq \beta \leq \alpha$$

On peut alors démontrer les résultats suivants

### Théorème III.6

Si  $u_S, x_S, \lambda_S \in H_m^\alpha(\Pi)$  ou  $(u_S, x_S, \lambda_S)$  est l'approximation spline de  $(u^*, x^*, \lambda^*)$ ,  $\alpha \geq 1$ , pour  $\alpha \geq 1$

$\  \epsilon(u_S) \ _{2, \beta} \leq O(h^{\alpha-\beta})$	$0 \leq \beta \leq \alpha$
$\  \epsilon(x_S) \ _{2, \beta} \leq O(h^{\alpha-\beta})$	$0 \leq \beta \leq \alpha$
$\  \epsilon(\lambda_S) \ _{2, \beta} \leq O(h^{\alpha-\beta})$	$0 \leq \beta \leq \alpha$

### Théorème III 17

Sous les hypothèses précédentes

$$\|\varepsilon(\hat{u})\|_{2,1} \leq \|\varepsilon(u_s)\|_{2,1}$$

$$\|\varepsilon(\hat{z})\|_{2,1} \leq \|\varepsilon(z_s)\|_{2,1}$$

$$\|\varepsilon(\hat{\lambda})\|_{2,1} \leq \|\varepsilon(\lambda_s)\|_{2,1}$$

### Théorème III 18

Si  $\bar{u}, \bar{z}$  et  $\bar{\lambda} \in H_m^k(\Pi)$  sont solutions du système de Galerkin, alors

$$\|\varepsilon(\bar{e})\|_2 = \|\bar{e} - \hat{e}\|_2 \leq O(h^\alpha)$$

$$\|\varepsilon(\bar{z})\|_2 = \|\bar{z} - \hat{z}\|_2 \leq O(h^{\alpha-1})$$

Ce dernier théorème fait intervenir, entre autres, l'hypothèse 1.

En utilisant ces résultats, nous obtenons les bornes suivantes

### Théorème III 19

Si  $\bar{u}, \bar{z}$  et  $\bar{\lambda} \in H_m^k(\Pi)$  sont solutions du système de Galerkin,

alors

$$\|\bar{z} - z^*\|_2 \leq O(h^{\alpha-1})$$

$$\|\bar{u} - u^*\|_2 \leq O(h^{\alpha-1})$$

$$J^* - O(h^{(k-1)}) \leq J_m \leq J^* + O(h^{2(k-1)})$$

## 65 Application au cas linéaire quadratique

Nous envisageons de nouveau la minimisation de

$$\mathcal{J}[u, x] = \frac{1}{2} \int_0^T \{ \langle x, g_x \rangle + \langle u, R u \rangle \} dt$$

sous la contrainte

$$\dot{x}(t) = A(t)x(t) + B(t)u(t) \quad x(0) = x_0$$

avec  $x(t) \in \mathbb{R}^n$ ,  $u(t) \in \mathbb{R}^r$  pour tout  $t \in [0, T]$

Nous supposons de plus

- 1)  $R(t)$  et  $Q(t)$  sont des matrices réelles symétriques et définies positives
- 2)  $R(t), Q(t), A(t)$  et  $B(t) \in PC^1[0, T]$

Les hypothèses 1 à 3 sont alors évidemment vérifiées. La théorie générale fournit, en implémentant le système de Galerkin à l'aide des expressions

$$\bar{x}(t) = \sum_{i=1}^m a_i w_i(t)$$

$$\bar{u}(t) = \sum_{i=1}^m b_i w_i(t)$$

$$\bar{\lambda}(t) = \sum_{i=1}^m c_i w_i(t)$$

on obtient, après quelques calculs simples, le système

$$G_a + A_a + B_b = -g$$

$$G_c + A_c + Q_a = 0$$

$$R_b + B^T c = 0$$

$$\text{où } G = (G_{ij}) \text{ avec } G_{ij} = \int_0^T \langle \dot{w}_j(t), w_i(t) \rangle dt$$

$$A = (A_{ij}) \text{ avec } A_{ij} = \int_0^T \langle w_j(t), A(t) w_i(t) \rangle dt$$

et  $\underline{B}$ ,  $\underline{Q}$  et  $\underline{R}$  définis de manière analogue à  $A$ .

Nous avons aussi choisi

$$\underline{a}^t = (a_1^t, a_2^t, \dots, a_m^t)$$

$$\underline{b}^t = (b_1^t, b_2^t, \dots, b_n^t)$$

$$\underline{c}^t = (c_1^t, c_2^t, \dots, c_m^t)$$

et

$$\underline{d}^t = (x_0^t w_1(0), x_0^t w_2(0), \dots, x_0^t w_m(0))^t$$

Comme  $R$  et  $Q$  sont définies positives,  $\underline{R}$  et  $\underline{Q}$  le sont aussi. Nous pouvons alors réduire le système en

$$[(G+A)\underline{Q}^{-1}(G+A)^t + \underline{B}\underline{R}^{-1}\underline{B}^t] \underline{c} = \underline{d}$$

ou plus simplement

$$F \underline{c} = \underline{d}$$

où nous avons posé

$$F = [(G+A)\underline{Q}^{-1}(G+A)^t + \underline{B}\underline{R}^{-1}\underline{B}^t]$$

une matrice  $(m,n) \times (m,n)$  symétrique et définie positive qui possède, quand on choisit les espaces "approximations" de manière adéquate (et possible!), des propriétés numériques intéressantes.

La résolution du système fournit alors tous les renseignements désirés.

## IV Le régulateur d'état parabolique linéaire

Rappelons brièvement les données du problème.

Nous considérons le système parabolique linéaire suivant

$$\frac{\partial v}{\partial t} = A(x,t) \frac{\partial^2 v}{\partial t^2} + B(x,t) u(x,t) \quad (\text{III cg})$$

avec la condition initiale  $\lim_{t \rightarrow 0^+} v(x,t) = v_0(x)$   
et les conditions-limites

$$\alpha v(0,t) + \frac{\partial v(0,t)}{\partial x} = c f_1(t) \quad (\text{III z0})$$

$$\beta v(1,t) + \frac{\partial v(1,t)}{\partial x} = d f_2(t) \quad (\text{III z1})$$

où  $x \in [0,1]$  et  $t \in [0,T]$ . La fonctionnelle du coût est

$$\begin{aligned} J[u, f_1, f_2] = & \frac{1}{2} \int_0^T \left[ \int_0^1 [ \langle v, Q(x,t) v \rangle + \langle u, R(x,t) u \rangle ] dx dt \right. \\ & \left. + \frac{1}{2} \int_0^T [ \langle f_1, r(t) f_1 \rangle + \langle f_2, s(t) f_2 \rangle ] dt \right] \quad (\text{III z2}) \end{aligned}$$

où  $v(\cdot)$  est un vecteur d'état à  $n$  dimensions,  $u(\cdot)$  un vecteur de contrôle à  $r$  dimensions et  $f_1(\cdot)$  et  $f_2(\cdot)$  sont des vecteurs de contrôle aux frontières à  $n$  dimensions. Nous supposons de plus que  $A(\cdot)$  est une matrice  $n \times n$ ,  $B(\cdot)$  une matrice  $n \times r$ ,  $Q(\cdot)$ ,  $R(\cdot)$ ,  $s(\cdot)$  des matrices  $n \times n$  définies positives, et que  $R(\cdot)$  est une matrice  $r \times r$  définie positive.  $\alpha, \beta, c$  et  $d$  sont des scalaires. De plus, nous exigeons que  $A, B, Q$  et  $R \in PC^\gamma[0,1] \times PC^{\gamma-1}[0,T]$  pour  $\gamma \geq 1$  et que  $v_0(x) \in C^0[0,1]$ .  $T$  est fixé,  $0 < r \leq n$ . Nous n'imposons pas de contrainte supplémentaire sur  $u(\cdot)$ ,  $f_1(\cdot)$  et  $f_2(\cdot)$ . De plus,  $Q$  et  $R$  sont symétriques.

Nous imposons encore que  $Au$  (notre ensemble de contrôles  $u$ ) soit une partie de  $\{V_2^0[0,1] \times (0,T)], Rf\}$  et que  $Af$  (l'ensemble des  $f_i$ ) soit une partie de  $W_2^0[0,T]$ . Nous désirons trouver  $(u^*, f_1^*, f_2^*, v^*)$  tel que

$$J[u^*, f_1^*, f_2^*] = \min_{\substack{u \in Au \\ f_i \in Af, f_i \in Ap}} J[u, f_1, f_2] \quad (\text{III z3})$$

**[16] Formulations du problème.**

Sous l'angle du calcul des variations, on peut attaquer le problème en introduisant des multiplicateurs de Lagrange vectoriels  $\lambda_1, \lambda_2, \lambda_3$  et  $\lambda_4$  et définir le lagrangien comme suit

$$\begin{aligned} L[\underline{u}, \underline{f}_1, \underline{f}_2, v; \lambda_1, \lambda_2, \lambda_3, \lambda_4] &= L[\underline{u}, v; \underline{\lambda}] \\ &= J[\underline{u}] + \int_0^T \left\{ \left\langle \lambda_1(x,t), -\frac{\partial v}{\partial t} + A \frac{\partial^2 v}{\partial x^2} + Bu \right\rangle dx dt \right. \\ &\quad + \left. \left\langle \lambda_3(t), Cf_1 - av(0,t) - \frac{\partial v}{\partial x}(0,t) \right\rangle dt \right. \\ &\quad + \left. \left\langle \lambda_4(t), df_2 - \beta v(1,t) - \frac{\partial v(1,t)}{\partial x} \right\rangle dt \right. \\ &\quad \left. + \int_0^1 \left\langle \lambda_2(t), v(x,0) - v_0(x) \right\rangle dx \right\} \quad (\text{III } 24) \end{aligned}$$

où  $\underline{u} = (u, f_1, f_2) \in V_2^0[(0,1) \times (0,T)] \times W_2^0[0,T]^2$  et  $\underline{\lambda} = (\lambda_1, \lambda_2, \lambda_3, \lambda_4)$   $\in V_2^1[(0,1) \times (0,T)] \times W_2^1[0,T]^2$ .

On peut aussi formuler le problème de manière équivalente comme suit : trouver  $v^*, \underline{u}^*$  et  $\underline{\lambda}^*$  tels que

$$L[\underline{u}^*, v^*; \underline{\lambda}^*] = \max_{\underline{\lambda} \in \underline{\Lambda}} \min_{\substack{u \in A\underline{u} \\ v \in Av}} L[\underline{u}, v, \underline{\lambda}] \quad (\text{III } 25)$$

où  $\underline{\Lambda} = \Lambda_{\lambda_1} \times \Lambda_{\lambda_2} \times \Lambda_{\lambda_3} \times \Lambda_{\lambda_4}$  avec  $\Lambda_{\lambda_i} \subseteq V_2^0[(0,1) \times (0,T)] \times W_2^0[0,T], \mathbb{R}^4$  et  $A\underline{u} = A_u \times A_f^2$

Nous pouvons aussi demander de trouver  $v^*, \underline{u}^*$  et  $\underline{\lambda}^*$  tels que

$$L[\underline{u}^*, \underline{v}^*, \lambda^*] = \max_{\lambda \in A\Lambda} L[\underline{u}, \underline{v}, \lambda] \quad (\text{III}7c)$$

sous les conditions

$$\frac{\partial L[\underline{u}, \underline{v}, \lambda]}{\partial v} = \frac{\partial L[\underline{u}, \underline{v}, \lambda]}{\partial u} = \frac{\partial L[\underline{u}, \underline{v}, \lambda]}{\partial \lambda_1} = \frac{\partial L[\underline{u}, \underline{v}, \lambda]}{\partial \lambda_2} = 0 \quad (\text{III}7f)$$

(dérivées au sens de Fréchet).

L'équivalence de ces formulations est un résultat du calcul des variations (cf. supra).

## 25 Conditions nécessaires d'optimalité.

Effectuons d'abord quelques transformations sur le Lagrangien

$$\begin{aligned}
 & \int_0^T \left\langle \lambda_1, -\frac{\partial v}{\partial t} + A \frac{\partial^2 v}{\partial x^2} + Bu \right\rangle dx dt = \int_0^T \left\langle \lambda_1, -\frac{\partial v}{\partial t} \right\rangle dx dt + \int_0^T \left\langle A^T \lambda_1, \frac{\partial^2 v}{\partial x^2} \right\rangle dx dt \\
 & + \int_0^T \int_0^1 \left\langle \lambda_1, Bu \right\rangle dx dt \\
 & = - \int_0^T \left[ \left\langle \lambda_1, v \right\rangle \right]_0^T dx + \int_0^T \int_0^1 \left\langle v, \frac{\partial \lambda_1}{\partial t} \right\rangle dx dt + \left[ \left\langle A^T \lambda_1, \frac{\partial v}{\partial x} \right\rangle \right]_0^T dt \\
 & - \int_0^T \int_0^1 \left\langle \frac{\partial A^T \lambda_1}{\partial x}, \frac{\partial v}{\partial x} \right\rangle dx dt + \int_0^T \int_0^1 \left\langle \lambda_1, Bu \right\rangle dx dt \\
 & = - \int_0^T \left\langle \lambda_1, v \right\rangle \int_0^T dx + \int_0^T \int_0^1 \left\langle v, \frac{\partial \lambda_1}{\partial t} \right\rangle dx dt + \int_0^T \int_0^1 \left\langle A^T \lambda_1, \frac{\partial v}{\partial x} \right\rangle dx dt \\
 & - \int_0^T \int_0^1 \left\langle v, \frac{\partial A^T \lambda_1}{\partial x} \right\rangle dx dt + \int_0^T \int_0^1 \left\langle v, \frac{\partial^2 A^T \lambda_1}{\partial x^2} \right\rangle dx dt + \int_0^T \int_0^1 \left\langle \lambda_1, Bu \right\rangle dx dt \\
 & = - \int_0^T \left\langle \lambda_1, v \right\rangle \int_0^T dx + \int_0^T \left\{ \left\langle A^T \lambda_1, \frac{\partial v}{\partial x} \right\rangle - \left\langle v, \frac{\partial A^T \lambda_1}{\partial x} \right\rangle \right\} \int_0^T dt \\
 & + \int_0^T \int_0^1 \left\{ \left\langle v, \frac{\partial \lambda_1}{\partial t} + \frac{\partial^2 A^T \lambda_1}{\partial x^2} \right\rangle + \left\langle \lambda_1, Bu \right\rangle \right\} dx dt
 \end{aligned}$$

ce qui donne évidemment

$$\begin{aligned}
 J[u, v, \lambda] &= J[u] - \int_0^T \langle \lambda_1, v \rangle |_0^T dx \\
 &\quad + \int_0^T \left\{ \left\langle A^t \lambda_1, \frac{\partial v}{\partial x} \right\rangle - \left\langle v, \frac{\partial A^t \lambda_1}{\partial x} \right\rangle \right\} |_0^T dt \\
 &\quad + \int_0^T \int_0^1 \left\{ \left\langle v, \frac{\partial \lambda_1}{\partial t} + \frac{\partial^2 A^t \lambda_1}{\partial x^2} \right\rangle + \langle \lambda_1, \partial_x v \rangle \right\} dx dt \\
 &\quad + \int_0^T \left\{ \left\langle \lambda_3, Cf_1 - \alpha v(0, t) - \frac{\partial v(0, t)}{\partial x} \right\rangle \right. \\
 &\quad \quad \left. + \left\langle \lambda_4, df_2 - \beta v(1, t) - \frac{\partial v(1, t)}{\partial x} \right\rangle \right\} dt \\
 &\quad + \int_0^1 \langle \lambda_2, v(x, 0) - v_0(x) \rangle dx \tag{III 78}
 \end{aligned}$$

De manière analogue aux problèmes précédents, on obtient les conditions nécessaires suivantes

$$\frac{\partial \lambda_1}{\partial t} + \frac{\partial^2 A^t \lambda_1}{\partial x^2} + Q(x, t) v(x, t) = 0 \tag{III 79}$$

$$R(x, t) u(x, t) + B^t \lambda_1(x, t) = 0 \tag{III 80}$$

$$\lambda_1(x, T) = 0 \quad \lambda_2(x) + \lambda_1(x, 0) = 0 \tag{III 81}$$

$$A^t(0, t) \lambda_1(0, t) + \lambda_3(t) = 0 \tag{III 82}$$

$$A^t(1, t) \lambda_1(1, t) - \lambda_4(t) = 0 \tag{III 83}$$

$$C \lambda_3(t) + c(t) \lambda_4(t) = 0 \tag{III 84}$$

$$d\lambda_4(t) + s(t) f_2(t) = 0 \quad (\text{III.85})$$

$$\alpha(A^t \lambda_1)(0,t) + \frac{\partial(A^t \lambda_1)(0,t)}{\partial x} = 0 \quad (\text{III.86})$$

$$\beta(A^t \lambda_1)(1,t) + \frac{\partial(A^t \lambda_1)(1,t)}{\partial x} = 0 \quad (\text{III.87})$$

à l'optimum cherché.

Modifions le Lagrangien. Remarquons que

$$\int_0^1 \langle -\lambda_1(x,0), v(x,0) \rangle + \int_0^1 \langle \lambda_1(x,0), v_0(x) \rangle = \int_0^1 \langle \lambda_2, v(x,0) \cdot v_0(x) \rangle$$

en utilisant (III.81); et donc

$$\int_0^1 \langle \lambda_2, v(x,0) \cdot v_0(x) \rangle + \int_0^1 \langle \lambda_1(x,0), v(x,0) \rangle = \int_0^1 \langle \lambda_1(x,0), v_0(x) \rangle$$

De plus (III.81) implique

$$\int_0^1 \langle \lambda_1(x,T), v(x,T) \rangle = 0$$

Nous avons aussi que

$$\int_0^T \left\{ \left\langle v, \frac{\partial \lambda_1}{\partial t} + \frac{\partial^2 A^t \lambda_1}{\partial x^2} \right\rangle + \langle B^t \lambda_1, u \rangle \right\} = - \int_0^T \left\{ \langle Kv, Qv \rangle + \langle u, Rv \rangle \right\}$$

en utilisant (III.79) et (III.80). De même, avec (III.84) et (III.85),

$$\int_0^T \left\{ \langle \lambda_3, cf_1 \rangle + \langle \lambda_4, df_2 \rangle \right\} = - \int_0^T \left\{ \langle f_1, cf_1 \rangle + \langle f_2, sf_2 \rangle \right\}$$

Nous pouvons aussi déduire

$$\int_0^T \left\{ \langle \lambda_3, -\alpha v(0,t) - \frac{\partial v(0,t)}{\partial x} \rangle - \langle (A^t \lambda_1)(0,t), \frac{\partial v(0,t)}{\partial x} \rangle + \langle v(0,t), \frac{\partial (A^t \lambda_1)(0,t)}{\partial x} \rangle \right\}$$

$$= \int_0^T \left\langle \alpha (A^t \lambda_1)(0,t) + \frac{\partial (A^t \lambda_1)(0,t)}{\partial x}, v(0,t) \right\rangle = 0$$

en utilisant (III 82) et (III 86).

De même

$$\int_0^T \left\{ \langle \lambda_4, -\beta v(1,t) - \frac{\partial v(1,t)}{\partial x} \rangle + \langle (A^t \lambda_1)(1,t), \frac{\partial v(1,t)}{\partial x} \rangle - \langle v(1,t), \frac{\partial (A^t \lambda_1)(1,t)}{\partial x} \rangle \right\}$$

$$= - \int_0^T \left\langle \beta (A^t \lambda_1)(1,t) + \frac{\partial (A^t \lambda_1)(1,t)}{\partial x}, v(1,t) \right\rangle = 0$$

en utilisant (III.83) et (III.87).

En rassemblant toutes ces simplifications, on obtient

$$\begin{aligned} L[u, v, \lambda] &= J[u] - \int_0^T \left\{ \langle v, g_v \rangle + \langle u, R_v \rangle \right\} dx dt \\ &\quad - \int_0^T \left\{ \langle f_1, r f_2 \rangle + \langle f_2, s f_2 \rangle \right\} dt + \int_0^1 \langle \lambda_1(x, 0) N_0(x) \rangle dx \\ &= - J[u] + \int_0^1 \langle \lambda_1(x, 0) N_0(x) \rangle dx \end{aligned} \tag{III 88}$$

**35**

Sous espace de dimension finie et approximation de Ritz Galerkin.

Comme dans les deux problèmes précédents, nous allons introduire un sous espace de dimension finie et tenter de résoudre le problème

initial sur ce sous-espace.

Soit  $S_h^{\gamma}$ ,  $\gamma \geq 1$  un espace de polynômes par morceaux d'ordre  $\gamma-1$  et d'une seule variable.  $h$  est ici le pas de la discréteisation. Alors,  $S_h^{\gamma}$  possède les propriétés suivantes

- 1.  $\exists L_h$  linéaire  $L_h: C^{\gamma}[a,b] \rightarrow S_h^{\gamma}$
- 2.  $\forall f \in C^{\gamma}[a,b]$   $\|L_h f - f\|_2 = O(h^{\gamma})$
- 3.  $\forall f \in C^{\gamma}[a,b]$   $\left\| \frac{d}{dx} (L_h f - f) \right\|_2 = O(h^{\gamma-1})$

Il est classique que de tels espaces existent, les espaces de fonctions splines par exemple.

Etendons maintenant, par un procédé simple, ces caractérisations, pour obtenir un espace de polynômes à deux variables.

Considérons d'abord  $\Omega$  un rectangle fermé dans  $\mathbb{R}^2$  défini par

$$\Omega = \{(p_1, p_2) \mid a_i \leq p_i \leq b_i, i=1,2\}$$

où  $a_i$  et  $b_i \in \mathbb{R}$  pour chaque  $i$ .

Soit  $\Gamma$  la frontière de  $\Omega$ . Introduisons une partition  $\Pi$  de  $\Omega$  comme suit

$$\Pi_i \mid a_i = x_i^0 < x_i^1 < \dots < x_i^{m_i} = b_i \quad i=1,2$$

Soit, de plus

$$h_i = \Pi_i = \max_{j \in J_{m_i}} \{x_j^i - x_{j-1}^i\}$$

Posons maintenant  $\sum_M^{\Gamma} = S_{m_1} \times S_{m_2}^{\Gamma}$

où  $\Gamma = \{\gamma_1, \gamma_2\}$  et  $M = \{h_1, h_2\}$ .

On en déduit immédiatement que  $\sum_M^{\Gamma}$  possède les propriétés

- 1.  $\exists L_s$  bilinéaire  $L_s: C^{\Gamma}(\Omega) \rightarrow \sum_M^{\Gamma}$
- 2.  $\forall f \in C^{\Gamma}(\Omega)$

$$\|L_s f - f\|_2 = \sum_{j=1}^2 O(h_j^{\gamma_j}) = O(h^{\gamma})$$

$$3. \quad \dot{f} \in C^{\gamma}[\Omega]$$

$$\left\| \frac{\partial}{\partial x_j} (L_h f - f) \right\|_2 = \sum_{j=1}^2 O(h_j^{\gamma_j-1}) = O(h^{\gamma-1})$$

où  $\gamma = (\gamma_1, \gamma_2)$  et  $h = (h_1, h_2)$ .

Notre espace est donc construit.

Définissons maintenant :

Le problème initial  $\in C^{\gamma}[(0,1) \times (0,T)]$



- (i)  $A, B, Q$  et  $R \in C^{\gamma}$  par rapport à  $t$
- (ii)  $A, B, Q$  et  $R \in C^{\gamma+1}$  par rapport à  $x$
- (iii)  $V_0(x) \in C^{\gamma+1}$  par rapport à  $x$
- (iv)  $r(\cdot)$  et  $s(\cdot) \in PC^{\gamma}$  par rapport à  $t$

Nous pouvons alors voir le théorème suivant

### Théorème III.20

Si le problème  $\in C^{\gamma-1}[(0,1) \times (0,T)]$ , alors

- (i)  $u^*, v^*, \lambda^* \in C^{\gamma}[(0,1) \times (0,T)]$
- (ii)  $f_1^*, f_2^*, \lambda_3^*, \lambda_4^* \in C^{\gamma}[0,T]$
- (iii)  $\lambda_2^* \in C^{\gamma}[0,1]$

Démonstration : Par (III.69) et (III.79), on voit immédiatement que  $v^*(\cdot)$  et  $\lambda^*$  satisfont le système.

$$\begin{bmatrix} \frac{\partial v^*}{\partial t} \\ \frac{\partial \lambda_1^*}{\partial t} \end{bmatrix} = \begin{bmatrix} A \frac{\partial^2 (\cdot)}{\partial x^2} & -S(x,t) \\ -Q(x,t) & -\frac{\partial^2 A^t(\cdot)}{\partial x^2} \end{bmatrix} \begin{bmatrix} v^* \\ \lambda_1^* \end{bmatrix} \quad (\text{III.89})$$

avec les conditions initiales et aux limites données par (III.69)-(III.70)-(III.71) et (III.81 à 87). (III.80) implique  $u^*(x,t) = -R^{-1}B^t \lambda_1^*$  et nous définissons alors

$$S(x,t) = B(x,t) R^{-1}(x,t) B^t(x,t)$$

Comme  $B, R \in C^{0,1}([0,1] \times (0,T])$  et  $R$  est définie positive,  $R^{-1} \in C^0([0,1] \times (0,T])$  et donc  $S \in C^{0,1}([0,1] \times (0,T])$ . Du système (III.89), on obtient évidemment que  $v^*, \lambda_1^* \in C^0([0,1] \times (0,T])$

De plus  $u^* = -R^{-1}B^t \lambda_1^*$  implique  $u^* \in C^0([0,1] \times (0,T])$ . (III.82) et  $A \in C^0$  implique  $\lambda_3 \in C^0[0,T]$ , puisque (III.83), on a aussi  $\lambda_4 \in C^0[0,T]$ . (III.84) et (III.85) indiquent respectivement  $f_1 \in C^0[0,T]$  et  $f_2 \in C^0[0,T]$ .

Enfin (III.81) implique  $\lambda_2 \in C^0[0,T]$  ■

Comme précédemment, nous pouvons poser le problème: trouver  $(\bar{u}, \bar{f}_1, \bar{f}_2)$   
 $= \underline{u}$  tel que

$$J[\underline{u}] = \min_{u \in \mathcal{A}_{\underline{u}}} J[u]$$

sous les contraintes

$$\int_0^1 \langle w_j(x) w_k(0), (\bar{v}(x,0) - v_0(x)) \rangle_i dx = 0$$

$$\int_0^T \langle w_j(0), w_k(t), (c \bar{f}_1(t) - \alpha \bar{v}(0,t) - \frac{\partial \bar{v}}{\partial x}(0,t)) \rangle_i dt = 0$$

$$\int_0^T \langle w_j(t) w_k(t), (\partial_t \bar{f}_2(t) - \beta \bar{v}(t,t) - \frac{\partial \bar{v}}{\partial x}(t,t))_j \rangle dt = 0$$

$$\int_0^T \int_0^1 \langle w_j(x) w_k(t), \left( -\frac{\partial \bar{v}}{\partial t} + A \frac{\partial^2 \bar{v}}{\partial x^2} + B \bar{v} \right)_j \rangle dx dt = 0$$

le système de Galerkin pour chaque fonction de base  $w_j(x), w_k(t) \in \Sigma_H^r$ ,  
 $j \in J_{m_1}, k \in J_{m_2}$  avec  $t \in J_n$

On peut évidemment formuler cela comme suit :

trouver  $\bar{u}, \bar{v}, \bar{f}_1$  et  $\bar{f}_2$  tels que

$$L[\bar{u}, \bar{v}, \bar{\lambda}] = \max_{\lambda \in \Sigma_H^r} \min_{\substack{u \in A_u \\ v \in A_v}} L[u, v, \lambda] \quad (\text{III}g_0)$$

où  $\bar{\lambda} \in \Sigma_H^r$  s'interprète naturellement comme  $(\lambda_1)_i \in \Sigma_{H_1}^r$ ,  $(\lambda_2)_i \in S_{m_1}^{k_1}$ ,  
 $(\lambda_3)_i \in S_{m_2}^{k_2}$  et  $(\lambda_4)_i \in S_{m_2}^{k_2}$  ( $i = 1 \dots n$ ).

On peut aussi chercher  $\bar{\lambda}$  tel que

$$L[\bar{u}, \bar{v}, \bar{\lambda}] = \max_{\lambda \in \Sigma_H^r} L[\bar{u}, \bar{v}, \lambda] \quad (\text{III}g_1)$$

sous les contraintes

$$\frac{\partial L[\bar{u}, \bar{v}, \lambda]}{\partial u} = \frac{\partial L[\bar{u}, \bar{v}, \lambda]}{\partial f_1} = \frac{\partial L[\bar{u}, \bar{v}, \lambda]}{\partial f_2} = \frac{\partial L[\bar{u}, \bar{v}, \lambda]}{\partial v} = 0 \quad (\text{III}g_2)$$

Un algorithme effectif sera développé plus bas, en utilisant cette dernière formulation.

Appelons  $\bar{u}$  et  $\bar{v}$  les solutions du système de Galerkin  
et abordons le point suivant.

46

## Convergence du coût pour le quadruplet $(\bar{u}, \bar{v}, \bar{f}_1, \bar{f}_2)$

Notre but est de voir comment  $J[\bar{u}]$  approxime  $J[u^*]$  où  $u^*$  est la solution optimale non approximée.

**Théorème III 21**

Si le problème initial  $\in C^{0,1}[(0,1) \times (0,T)]$  et si  $\bar{u}, \bar{v}, \bar{f}_1, \bar{f}_2$  et  $\bar{\lambda}$  sont les solutions du système de Galerkin pour  $\Sigma_H^T$ .

Si  $\lambda_s \in \Sigma_H^T$  où  $\lambda_s$  est une approximation de  $\lambda^*$  à  $O(h^\gamma)$ .

Si  $(u_s, v_s, f_{1s}, f_{2s})$  sont engendrés par  $\lambda_s$  d'après (III 79 à 87)

Alors, pour  $\gamma \geq 2$

$$\|\epsilon_u\|_2 = O(h^\gamma) \quad \|\epsilon_v\|_2 = O(h^{\gamma-2})$$

$$\|\epsilon_{f_1}\|_2 = O(h^\gamma) \quad \|\epsilon_{f_2}\|_2 = O(h^\gamma)$$

et

$$J[u^*] \geq J[\bar{u}, \bar{v}, \bar{\lambda}] \geq J[u^*] + L[\epsilon_u, \epsilon_v, \epsilon_\lambda] + \int_0^1 \langle \epsilon_\lambda, v_s(x) \rangle dx$$

où

$$\epsilon_u = u_s - u^*, \quad \epsilon_v = v_s - v^* \quad \text{et} \quad \epsilon_\lambda = \lambda_s - \lambda^*$$

Démonstration: Puisque le problème initial  $\in C^{0,1}$ , nous avons  $\lambda^* \in C^0[(0,1) \times (0,T)]$ . Or par hypothèse, il existe un  $\lambda_s \in \Sigma_H^T$  tel que  $\|\lambda_s - \lambda^*\|_2 = O(h^\gamma)$ . Donc, de (III 79) à (III 87) on tire

$$\|\lambda_{3s} - \lambda_3^*\|_2 = \|A^t(0,t)(\lambda_s(0,t) - \lambda_3^*(0,t))\|_2 \leq O(h^\gamma)$$

$$\|\lambda_{4s} - \lambda_4^*\|_2 = \|A^t(1,t)(\lambda_1^*(1,t) - \lambda_s(1,t))\|_2 \leq O(h^\gamma)$$

$$\|u_s - u^*\|_2 = \|R^{-1}B^t(\lambda^* - \lambda_s)\|_2 \leq O(h^\gamma)$$

$$\|v_s - v^*\|_2 = \|g^{-1}\left[\frac{2}{\delta t}(\lambda_s - \lambda_1^*) + \frac{\delta^2}{\partial t^2} A^t(\lambda_s - \lambda^*)\right]\|_2 \leq O(h^{\gamma-2})$$

$$\|f_{1s} - f_1^*\|_2 = c \|\tau^{-1}(\lambda_{3s} - \lambda_3^*)\|_2 \leq O(h^\gamma)$$

$$\|f_{2s} - f_2^*\|_2 = d \|s^{-1}(\lambda_{4s} - \lambda_4^*)\|_2 \leq O(h^\delta)$$

en se rappelant que  $r, s, A, B, R$  et  $Q$  sont continues sur des ensembles clos bornés. Nous avons donc prouvé (III.94)

Si nous nous souvenons qu'à l'extremum, le lagrangien est égal au coût, nous pouvons écrire, en utilisant (III.73)

$$J[u^*] = L[\underline{u}^*, v^*, \lambda^*] \geq L[\bar{u}, \bar{v}, \bar{\lambda}]$$

puisque  $\sum \underline{\lambda}_i \subseteq A \underline{\lambda}$ .

$$\text{De (III.91) on déduit que } L[\bar{u}, \bar{v}, \bar{\lambda}] \geq L[\underline{u}_s, v_s, \lambda_s]$$

Mais

$$\begin{aligned} L[\underline{u}_s, v_s, \lambda_s] &= L[\underline{u}^* + \epsilon_u, v^* + \epsilon_v, \lambda^* + \epsilon_\lambda] \\ &= \frac{1}{2} \int_0^T \left[ \left[ \langle v^* + \epsilon_v, Q(v^* + \epsilon_v) \rangle + \langle u^* + \epsilon_u, R(u^* + \epsilon_u) \rangle \right] \right. \\ &\quad \left. + \frac{1}{2} \int_0^T \left[ \langle \dot{v}_1^* + \epsilon \dot{v}_1, r(\dot{v}_1^* + \epsilon \dot{v}_1) \rangle + \langle \dot{f}_2^* + \epsilon \dot{f}_2, s(\dot{f}_2^* + \epsilon \dot{f}_2) \rangle \right] \right. \\ &\quad \left. + \int_0^T \left\langle \lambda_1^* + \epsilon \lambda_1, -\frac{\partial u^*}{\partial t} + A \frac{\partial^2 v^*}{\partial x^2} + B u^* - \frac{\partial \epsilon v}{\partial t} + A \frac{\partial^2 \epsilon v}{\partial x^2} + C \epsilon u \right\rangle \right. \\ &\quad \left. + \int_0^T \left\langle \lambda_3^* + \epsilon \lambda_3, c(\dot{v}_1^* + \epsilon \dot{v}_1) - \alpha(V^*(0,t) + \epsilon_v(0,t)) - \frac{\partial v^*}{\partial x}(0,t) - \frac{\partial \epsilon v}{\partial x}(0,t) \right\rangle \right. \\ &\quad \left. + \int_0^T \left\langle \lambda_4^* + \epsilon \lambda_4, d(\dot{f}_2^* + \epsilon \dot{f}_2) - \beta(V^*(t,t) + \epsilon_v(t,t)) - \frac{\partial f^*}{\partial x}(t,t) - \frac{\partial \epsilon f}{\partial x}(t,t) \right\rangle \right. \\ &\quad \left. + \int_0^T \left\langle \lambda_2^* + \epsilon \lambda_2, V^*(x,0) + \epsilon_v(x,0) - v_0(x) \right\rangle \right. \\ &= \frac{1}{2} \int_0^T \left[ \left[ \langle v^*, Q v^* \rangle + \langle u^*, R u^* \rangle \right] + \frac{1}{2} \int_0^T \left[ \langle \dot{v}_1^*, r \dot{v}_1^* \rangle + \langle \dot{f}_2^*, s \dot{f}_2^* \rangle \right] \right] \end{aligned}$$

cost.

$$\begin{aligned}
& + \int_0^T \int_0^1 \left[ \langle \lambda_1^*, -\frac{\partial v^*}{\partial t} + A \frac{\partial^2 v^*}{\partial x^2} + B u^* \rangle + \int_0^T \langle \lambda_3^*, c f_1^* - \alpha v^*(0,t) - \frac{\partial v^*}{\partial x}(0,t) \rangle \right. \\
& + \int_0^T \langle \lambda_4^*, d f_2^* - \beta v^*(1,t) - \frac{\partial v^*}{\partial x}(1,t) \rangle + \int_0^1 \langle \lambda_2^*, v^*(x,0) - v_0(x) \rangle \\
& + \frac{1}{2} \int_0^T \int_0^1 [\langle \epsilon_v, Q \epsilon_v \rangle + \langle \epsilon_u, R \epsilon_u \rangle] + \frac{1}{2} \int_0^T [\langle \epsilon_{f_1}, R \epsilon_{f_1} \rangle + \langle \epsilon_{f_2}, S \epsilon_{f_2} \rangle] \\
& + \int_0^T \int_0^1 \left[ \langle \epsilon_{\lambda_1}, -\frac{\partial \epsilon_v}{\partial t} + A \frac{\partial^2 \epsilon_v}{\partial x^2} + B \epsilon_u \rangle + \int_0^T \langle \epsilon_{\lambda_3}, c \epsilon f_1 - \alpha \epsilon_v(0,t) - \frac{\partial \epsilon_v}{\partial x}(0,t) \rangle \right. \\
& + \int_0^T \langle \epsilon_{\lambda_4}, d \epsilon f_2 - \beta \epsilon_v(1,t) - \frac{\partial \epsilon_v}{\partial x}(1,t) \rangle + \int_0^1 \langle \epsilon_{\lambda_2}, \epsilon_v(x,0) - v_0(x) \rangle \\
& + \int_0^1 \langle \epsilon_{\lambda_2}, v_0(x) \rangle + \frac{1}{2} \int_0^T \int_0^1 [\langle \epsilon_v, Q \epsilon_v \rangle + \langle \epsilon_u, R \epsilon_u \rangle] + \frac{1}{2} \int_0^T [\langle \epsilon_{f_1}, R \epsilon_{f_1} \rangle + \langle \epsilon_{f_2}, S \epsilon_{f_2} \rangle] \\
& + \int_0^T \int_0^1 \left[ \langle \epsilon_{\lambda_1}, -\frac{\partial \epsilon_v}{\partial t} + A \frac{\partial^2 \epsilon_v}{\partial x^2} + B \epsilon_u \rangle + \int_0^T \langle \epsilon_{\lambda_3}, c \epsilon f_1 - \alpha \epsilon_v(0,t) - \frac{\partial \epsilon_v}{\partial x}(0,t) \rangle \right. \\
& + \int_0^T \langle \epsilon_{\lambda_4}, d \epsilon f_2 - \beta \epsilon_v(1,t) - \frac{\partial \epsilon_v}{\partial x}(1,t) \rangle + \int_0^1 \langle \epsilon_{\lambda_2}, \epsilon_v(x,0) - v_0(x) \rangle \\
& + \frac{1}{2} \int_0^T \int_0^1 [\langle v^*, Q \epsilon_v \rangle + \langle u^*, R \epsilon_u \rangle] + \frac{1}{2} \int_0^T [\langle f_1^*, R \epsilon_{f_1} \rangle + \langle f_2^*, S \epsilon_{f_2} \rangle] \\
& + \int_0^T \langle \lambda_3^*, -\frac{\partial \epsilon_v}{\partial t} + A \frac{\partial^2 \epsilon_v}{\partial x^2} + B \epsilon_u \rangle + \int_0^T \langle \lambda_3^*, c \epsilon f_1 - \alpha \epsilon_v(0,t) - \frac{\partial \epsilon_v}{\partial x}(0,t) \rangle \\
& + \int_0^T \langle \lambda_4^*, d \epsilon f_2 - \beta \epsilon_v(1,t) - \frac{\partial \epsilon_v}{\partial x}(1,t) \rangle + \int_0^1 \langle \lambda_2^*, \epsilon_v(x,0) \rangle .
\end{aligned}$$

Les termes saillants sont nuls car  $v^*$  est optimal et satisfait les conditions massives (III 73) à (III 82).

D'autre part, puisque  $R$  et  $Q$  sont symétriques, on obtient bien que

$$L[u_s, v_s, \lambda_s] = L[u^*, v^*, \lambda^*] + L[\epsilon_u, \epsilon_v, \epsilon_\lambda]$$

$$\begin{aligned} &+ \int_0^T \left[ \left[ \langle v^*, g_{\epsilon_v} \rangle + \langle u^*, R_{\epsilon_u} \rangle \right] + \left[ \langle \ell_1^*, s \epsilon_{\ell_1} \rangle + \langle \ell_2^*, s \epsilon_{\ell_2} \rangle \right] \right. \\ &\quad \left. + \int_0^T \left[ \langle \lambda_2^*, \epsilon_v(x, t) \rangle + \langle \epsilon_{\lambda_2}, v_0(x) \rangle \right] + \int_0^T \left[ \langle \lambda_3^*, c \epsilon_{\ell_1} - \alpha \epsilon_v(0, t) - \frac{\partial \epsilon_v}{\partial x}(0, t) \rangle \right. \right. \\ &\quad \left. \left. + \int_0^T \left[ \langle \lambda_4^*, d \epsilon_{\ell_2} - \beta \epsilon_v(t, t) - \frac{\partial \epsilon_v}{\partial x}(t, t) \rangle + \int_0^T \left[ \langle \lambda^*, -\frac{\partial \epsilon_u}{\partial t} + A \frac{\partial^2 \epsilon_u}{\partial x^2} + B \epsilon_u \rangle \right] \right] \right] \right] \end{aligned}$$

Dans cette égalité, le terme souligné vaut  $2 J[\epsilon_u]$  par définition et les termes suivants  $-2 J[\epsilon_u]$  par un calcul exactement semblable à celui qui a permis de déduire (III 88) (On excepte  $\int_0^T \langle \epsilon_{\lambda_2}, v_0(x) \rangle$ )

Nous avons donc

$$L[u_s, v_s, \lambda_s] = J[u^*] + L[\epsilon_u, \epsilon_v, \epsilon_\lambda] + \int_0^T \langle \epsilon_{\lambda_2}, v_0(x) \rangle dx$$

ce qui achève la démonstration ■

Nous pouvons alors démontrer le corollaire suivant

### Théorème III 22

Sous les hypothèses du théorème III 21, on a

$$0 \leq J[u^*] - J[\bar{u}] \leq O(h^{2(\beta-2)})$$

(III 35)

Démonstration: Remarquons d'abord, comme précédemment, que  $J[\bar{u}] = L[\bar{u}, \bar{v}, \bar{\lambda}]$ .

De plus, le théorème précédent implique

$$0 \leq J[u^*] - L[\bar{u}, \bar{v}, \bar{\lambda}] \leq -L[\bar{e}_u, \bar{e}_v, \bar{e}_{\lambda}] - \int_0^1 \langle e_{\lambda_2}, v_0(x) \rangle dx$$

Mais, en utilisant (III 24),

$$\begin{aligned} L[\bar{e}_u, \bar{e}_v, \bar{e}_{\lambda}] &= \frac{1}{2} \int_0^T \left[ \langle \bar{e}_u, R\bar{e}_u \rangle + \langle \bar{e}_v, Q\bar{e}_v \rangle \right] + \frac{1}{2} \int_0^T \left[ \langle \bar{e}_{\lambda_1}, R\bar{e}_{\lambda_1} \rangle + \langle \bar{e}_{\lambda_2}, S\bar{e}_{\lambda_2} \rangle \right] \\ &\quad + \int_0^T \left\{ \langle \bar{e}_{\lambda_1}, -\frac{\partial \bar{e}_u}{\partial t} + A \frac{\partial^2 \bar{e}_u}{\partial x^2} + B \bar{e}_u \rangle + \int_0^T \left\{ \langle \bar{e}_{\lambda_3}, S\bar{e}_{\lambda_3} - \alpha \bar{e}_v(0,t) - \frac{\partial \bar{e}_v(0,t)}{\partial x} \rangle \right. \right. \\ &\quad \left. \left. + \langle \bar{e}_{\lambda_4}, d\bar{e}_{\lambda_2} - \beta \bar{e}_v(d,t) - \frac{\partial \bar{e}_v(d,t)}{\partial x} \rangle \right\} + \int_0^1 \langle \bar{e}_{\lambda_2}, \bar{e}_v(x,0) - v_0(x) \rangle \right\} \end{aligned}$$

En appliquant les conditions (III 79) à (III 85), car  $\bar{e}_v = v_s - v^*$ , on obtient

$$\begin{aligned} L[\bar{e}_u, \bar{e}_v, \bar{e}_{\lambda}] &= -\frac{1}{2} \int_0^T \left[ \langle \bar{e}_v, Q\bar{e}_v \rangle + \langle \bar{e}_u, R\bar{e}_u \rangle \right] - \frac{1}{2} \int_0^T \left[ \langle \bar{e}_{\lambda_1}, R\bar{e}_{\lambda_1} \rangle + \langle \bar{e}_{\lambda_2}, S\bar{e}_{\lambda_2} \rangle \right] \\ &\quad + \int_0^1 \left\{ \langle \bar{e}_{\lambda_1}(x,0), v_0(x) \rangle - \langle \bar{e}_{\lambda_1}(x,T), \bar{e}_v(x,T) \rangle \right\} \\ &\quad + \int_0^T \langle \bar{e}_v(0,t), \frac{\partial (A^t \bar{e}_{\lambda_1})(0,t)}{\partial x} + \alpha (A^t \bar{e}_{\lambda_1})(0,t) \rangle \\ &\quad - \int_0^T \langle \bar{e}_v(d,t), \frac{\partial (A^t \bar{e}_{\lambda_1})(d,t)}{\partial x} + \beta (A^t \bar{e}_{\lambda_1})(d,t) \rangle \end{aligned}$$

En appliquant maintenant l'inégalité de Cauchy-Schwartz aux deux dernières intégrales, on obtient

$$0 \leq J[u^*] - L[\bar{u}, \bar{v}, \bar{\lambda}] \leq \frac{1}{2} \| \bar{e}_v \|^2_S + \frac{1}{2} \| \bar{e}_u \|^2_R + \frac{1}{2} \| \bar{e}_{\lambda_1} \|^2_R + \frac{1}{2} \| \bar{e}_{\lambda_2} \|^2_S$$

$$+ \| \bar{e}_{\lambda_1}(x,0) \|_2 \| v_0(x) \|_2 + \| \bar{e}_{\lambda_1}(x,T) \|_2 \| \bar{e}_v(x,T) \|_2$$

$$+ \| \bar{e}_v(0,t) \|_2 \| \frac{\partial (A^t \bar{e}_{\lambda_1})(0,t)}{\partial x} + \alpha (A^t \bar{e}_{\lambda_1})(0,t) \|_2$$

cont.

$$\begin{aligned}
 & + \| \epsilon_0(1,t) \|_2 \| \frac{\partial(A^t \epsilon_\lambda)(1,t)}{\partial x} + \beta(A^t \epsilon_\lambda)(1,t) \|_2 \\
 & - \| \epsilon_\lambda \|_2 \| V_0(x) \| \\
 & \leq O(h^{2(\gamma-2)}) + O(h^{2\gamma}) + O(h^{2\gamma}) + O(h^{2\gamma}) + O(h^\gamma) \cdot O(h^{\gamma-2}) \\
 & + O(h^{\gamma-2}) O(h^{\gamma-2}) + O(h^{\gamma-2}) O(h^{\gamma-2}) = O(h^{2(\gamma-2)})
 \end{aligned}$$

Comme les normes  $\|\cdot\|_2$ ,  $\|\cdot\|_R$ ,  $\|\cdot\|_p$  et  $\|\cdot\|_S$  sont bien équivalentes aux normes  $\|\cdot\|_2$  respectives (cf th III 1), nous avons bien démontré la théorie ■

De la même manière, on peut voir le

### Théorème III 23

Supposons vérifiées les hypothèses du théorème III 21.

Imposons de plus à  $\sum_m^r$  les conditions suivantes

$$\forall \lambda_i \in \sum_m^r \lambda_i(x,T) = 0 \quad x \in (0,1)$$

$$\alpha(A^t \lambda_i)(0,t) + \frac{\partial(A^t \lambda_i)(0,t)}{\partial x} = 0 \quad t \in (0,T)$$

$$\beta(A^t \lambda_i)(1,t) + \frac{\partial(A^t \lambda_i)(1,t)}{\partial x} = 0 \quad t \in (0,T)$$

(III 9c)

$$\text{alors } 0 \leq J[\underline{u}^*] - J[\bar{u}] \leq O(h^{2(\gamma-2)})$$

où  $\bar{u}, \underline{u}, J$  est la solution de Galerkin sur le nouveau  $\sum_m^r$

La démonstration est identique à celle du théorème précédent. ■

56

### Convergence en norme du quadruplet $(\bar{u}, \bar{v}, \bar{\lambda}, \bar{f})$

Nous nous intéressons maintenant à la convergence en norme du quadruplet  $\bar{u}, \bar{v}, \bar{\lambda}, \bar{f}$  vers  $u^*, v^*, \lambda^*, f^*$ .

Prenons le théorème suivant

Théorème III 24

Si le problème initial  $\in C^1([0,1] \times (0,T))$ ,  
 si  $\bar{u}, \bar{v}$  et  $\bar{\lambda}$  sont les solutions sur le nouveau  $\Sigma_m^T$  (i.e.  
 $\Sigma_m^T$  muni des conditions du th. III 22), alors

$$\|\bar{u} - u^*\|_P \leq O(h^{j-2})$$

$$\|\bar{f}_1 - f_1^*\|_r \leq O(h^{j-2})$$

$$\|\bar{f}_2 - f_2^*\|_s \leq O(h^{j-2})$$

$$\|\bar{v} - v^*\|_Q \leq O(h^{j-2})$$

Démonstration Développons  $L[u^*, v, \lambda^*]$  en série de Taylor autour de  $(\bar{u}, \bar{v}, \bar{\lambda})$

$$L[u^*, v^*, \lambda^*] - L[\bar{u}, \bar{v}, \bar{\lambda}] = L[u^*, v^*, \bar{\lambda}] - L[\bar{u}, \bar{v}, \bar{\lambda}]$$

$$+ \frac{\partial L[\bar{u}, \bar{v}, \bar{\lambda}]}{\partial u}(u^* - \bar{u}) + \frac{\partial L[\bar{u}, \bar{v}, \bar{\lambda}]}{\partial f_1}(f_1^* - \bar{f}_1)$$

$$+ \frac{\partial L[\bar{u}, \bar{v}, \bar{\lambda}]}{\partial f_2}(f_2^* - \bar{f}_2) + \frac{\partial L[\bar{u}, \bar{v}, \bar{\lambda}]}{\partial v}(v^* - \bar{v})$$

$$+ \int_0^T \int_0^1 (u^* - \bar{u} - f_1^* + \bar{f}_1, f_2^* - \bar{f}_2, v^* - \bar{v}) H(x,t) \begin{bmatrix} u^* - \bar{u} \\ f_1^* - \bar{f}_1 \\ f_2^* - \bar{f}_2 \\ v^* - \bar{v} \end{bmatrix} dx dt$$

où  $M(x, t)$  est défini par

$$M(x, t) = \begin{bmatrix} R(x, t) & 0 & 0 & 0 \\ 0 & r(t) & 0 & 0 \\ 0 & 0 & s(t) & 0 \\ 0 & 0 & 0 & Q(x, t) \end{bmatrix}$$

(III.93) et (III.96) nous permettent d'assurer que les dérivées (au sens de Fréchet) du Lagrangien sont nulles en  $(\bar{x}, \bar{v}, \bar{\lambda})$ . (III.95) donne alors

$$0 \leq L[\bar{x}, \bar{v}, \bar{\lambda}^*] - L[\bar{x}, v, \bar{\lambda}] = \|u^* - \bar{u}\|_R^2 + \|f_1^* - f_1\|_r^2 + \|f_2^* - f_2\|_s^2$$

$$+ \|V^* - \bar{V}\|_g^2 \leq O(h^{2(8-2)})$$

La même remarque que supra s'impose à propos de l'équivalence des normes  $\|\cdot\|_R, \|\cdot\|_g, \|\cdot\|_r, \|\cdot\|_s$  avec la norme  $\|\cdot\|_2$ , puisque  $R, Q, S, r$  sont définies positives.

Nous avons donc prouvé tout ce que nous désirions.

**65**

### Algorithme numérique de résolution

Remarquons d'abord que (III.79) (III.80) et (III.88) nous permettent d'écrire  $L[\bar{x}, v, \bar{\lambda}] = L[\bar{\lambda}_i]$ .

Nous voulons

$$\frac{\partial L}{\partial \lambda_i} = 0$$

Comme tout  $(\lambda_i); i \in \sum_{\text{fin}}$ , on peut exprimer  $\lambda_i$  comme une combinaison linéaire

$$\lambda_i(x, t) = \sum_{j=1}^{m_1} \sum_{k=1}^{m_2} c_{jk} w_j(x) w_k(t) \quad x \in (0, 1), t \in (0, T)$$

où chaque  $\zeta_{jk}$  est un vecteur de dimension  $n$ .

On obtient alors successivement

$$\int_0^T \int_0^1 \langle u, R_u \rangle = \int_0^T \int_0^1 \langle R^{-1} B^t \lambda_1, B^t \lambda_1 \rangle = \int_0^T \int_0^1 \langle \lambda_1, B R^{-1} B^t \lambda_1 \rangle$$

$$= \sum_{j=1}^{m_1} \sum_{k=1}^{m_2} \sum_{p=1}^{m_1} \sum_{q=1}^{m_2} \zeta_{jk}^t \left[ \int_0^T \int_0^1 E_{jkpq} \, dx dt \right] c_{pq}$$

où  $E_{jkpq}(x,t) = w_j(x) w_k(t) w_p(x) w_q(t) B(x,t) R^{-1}(x,t) B^t(x,t)$   
 $(= \langle w_j(x) w_k(t), B(x,t) R^{-1}(x,t) B^t(x,t) w_p(x) w_q(t) \rangle)$

De même

$$\int_0^T \int_0^1 \langle v, Qv \rangle = \int_0^T \int_0^1 \langle Q^{-1} \left( \frac{\partial \lambda_1}{\partial t} + \frac{\partial^2 A^t \lambda_1}{\partial x^2} \right), \frac{\partial \lambda_1}{\partial t} + \frac{\partial^2 A^t \lambda_1}{\partial x^2} \rangle$$

$$= \sum_j \sum_k \sum_p \sum_q \zeta_{jk}^t \left[ \int_0^T \int_0^1 F_{jkpq} \, dx dt \right] c_{pq}$$

où  $F_{jkpq}(x,t) = w_j(x) \dot{w}_k(t) w_p(x) \dot{w}_q(t) Q^{-1}$

$$+ w_p(x) \dot{w}_k(t) w_q(t) Q^{-1} \frac{\partial^2 w_j(x) A^t(x,t)}{\partial x^2}$$

$$+ w_j(x) w_k(t) \dot{w}_q(t) Q^{-1} \frac{\partial^2 w_p(x) A^t(x,t)}{\partial x^2}$$

$$+ w_k(t) w_q(t) \frac{\partial^2 w_j(x) A^t(x,t)}{\partial x^2} Q^{-1} \frac{\partial^2 w_p(x) A^t(x,t)}{\partial x^2}$$

et aussi, de la même façon,

$$\int_0^T \langle f_1, r(t) f_1 \rangle = \sum_j \sum_k \sum_p \sum_q c_{jk}^+ \left[ \int_0^T G_{jklpq} dt \right] c_{pq}$$

$$\text{où } G_{jklpq} = c^2 w_j(0) w_k(t) w_p(0) w_q(t) A(0,t) r^{-1}(t) A^+(0,t)$$

et enfin

$$\int_0^T \langle f_2, s(t) f_2 \rangle = \sum_j \sum_k \sum_p \sum_q c_{jk}^+ \left[ \int_0^T G'_{jklpq} dt \right] c_{pq}$$

$$\text{où } G'_{jklpq} = d^2 w_j(1) w_k(t) w_p(1) w_q(t) A(1,t) s^{-1}(t) A^+(1,t)$$

et

$$\int_0^1 \langle \lambda_1(x,0), v_0(x) \rangle dx = \sum_{j=1}^{m_1} \sum_{k=1}^{m_2} c_{jk} \int_0^1 w_j(x) w_k(0) v_0(x) dx$$

Maximisons maintenant  $L[c_{11}, c_{12}, \dots, c_{21}, \dots, c_{m_1 m_2}]$

$$\begin{aligned} \frac{\partial L}{\partial c_{pq}} &= -\frac{1}{2} \sum_{j=1}^{m_1} \sum_{k=1}^{m_2} \left[ \int_0^T \left\{ E_{jklpq} + E_{jklpq}^+ + F_{jklpq} + F_{jklpq}^+ \right\} dt \right. \\ &\quad \left. + \int_0^T \left\{ G_{jklpq} + G_{jklpq}^+ + G'_{jklpq} + G'^{+}_{jklpq} \right\} dt \right] c_{jk}^+ \\ &\quad + \int_0^1 w_p(x) w_q(0) v_0(x) dx \end{aligned} \quad (\text{III 97})$$

pour  $p \in J_{m_1}$  et  $q \in J_{m_2}$ .

Définissons  $H_{i,l}$  par

$$H_{i,l} = \frac{1}{2} \int_0^T \int_0^1 [E_{jkpq} + E_{jkpq}^t + F_{jkpq} + F_{jkpq}^t] dx dt \\ + \int_0^T [G_{jkpq} + G_{jkpq}^t + G_{jkpq}^t + G_{jkpq}^t] dt$$

$$\text{et } b_i := \int_0^1 w_p(x) w_q(0) v_0(x) dx$$

où  $i = (p-1)m_2 + q$  et  $l = (j-1)m_2 + k$  pour  $j, p \in J_{m_1}$  et  $k, q \in J_{m_2}$ .  
Résoudre (III 97) revient alors à résoudre

$$Hy = b$$

(III 98)

où la matrice symétrique  $H ((m_1 m_2 n) \times (m_1 m_2 n))$  est définie par

$$H = \begin{bmatrix} H_{11} & & H_{1, m_1, m_2} \\ | & & | \\ H_{m_1, m_2, 1} & & H_{m_1, m_2, m_1, m_2} \end{bmatrix} \quad (\text{III 99})$$

avec  $H_{ij} = H_{ji}^t$ , et où  $y$  est donné par

$$y = \begin{bmatrix} c_{11} \\ | \\ c_{m_1, m_2} \end{bmatrix}$$

On peut, de plus, démontrer une propriété numérique supplémentaire du système algébrique (III 98)

### Théorème III 25

Si  $H$  est définie comme en (III 89), alors  $H$  est définie positive.

Démonstration: Supposons  $y \in \mathbb{R}^{m_1 m_2 n}$  et  $y \neq 0$ .

Sait  $\lambda$  correspondant à  $y$  par construction et soient  $u, f$  et  $v$  correspondants à  $\lambda$  par (III 79) à (II 87). Nous avons alors

$$\begin{aligned} g^t H y &= \frac{1}{2} y^t \left\{ \int_0^T \left\{ [E + E^t + F + F^t] d\alpha dt + \int_0^T \left\{ [G + G^t] d\alpha dt \right\} \right\} y \right. \\ &= \int_0^T \int_0^1 \langle u, R u \rangle d\alpha dt + \int_0^T \int_0^1 \langle v, Q v \rangle d\alpha dt + \int_0^T \langle f, S f \rangle dt \\ &= 2 J[u, v] \geq 0 \end{aligned}$$

par construction de  $J$ . ■

Ainsi se termine l'étude des trois problèmes annuels.

## Conclusion

Il est clair que les problèmes abordés dans ces quelques pages sont loin d'être épuisés. De nombreuses portes restent encore ouvertes pour une recherche plus approfondie.

Nous aurions personnellement beaucoup souhaité pouvoir traiter le côté plus numérique de la chose, mais les difficultés en matière d'informatique en 1973-1974 n'ont malheureusement pas permis l'implémentation directe des méthodes décrites.

Il nous reste à souhaiter que les circonstances futures permettant de revenir à ces problèmes qui, nous semble-t-il, peuvent avoir de fécondes répercussions non seulement dans l'application directe, mais aussi dans la théorie mathématique elle-même.

## Bibliographie

- A-1 ATHANS M. et FALB P.L.: "Optimal control: an introduction to the theory and its applications", Lincoln Laboratory Pub. (McGraw Hill), NY, 1966
- B-1 BLISS G.A.: "Lectures on the calculus of variations", Phoenix Science series, Chicago, 1968
- B-2 BOSARGE W.E. Jr, JOHNSON O.G., MCKNIGHT R.S. and THILAKES W.P.: "The Galerkin procedure for nonlinear control problems", SIAM J. Numer. Anal., Vol 10, No 1, March 1973, pp 94 sq
- B-3 BOSARGE W.E. Jr, JOHNSON O.G. and SMITH C.L.: "A direct method: Approximation to the linear parabolic regulator problem over multivariate spline bases", SIAM J. Numer. Anal., Vol 10, No 1, March 1973, pp 35 sq
- B-4 BRAUER F. and NOHEL J.: "Ordinary differential equations". W.A. Benjamin N.Y., 1966
- B-5 BRYSON A.E. and HO Y.C.: "Applied optimal control", Blaisdell Pub. Company, Waltham, Massachusetts, 1969
- B-6 BOURBAKI N. "Éléments de mathématique" (Livre II: Espaces vectoriels topologiques), Hermann, Paris, 1955

C-1 CHENEY E.W. "Introduction to Approximation theory", McGraw Hill,  
N.Y., 1966

C-2 CODDINGTON E.A. and LEVINSON N. "Theory of Ordinary Differential  
equations", McGraw Hill, N.Y., 1955

D-1 DIEUDONNÉ J.: "Les fondements de l'analyse moderne", Gauthier-  
Villars, Paris, 1965

E-1 EDWARDS R.E.: "Functional analysis: theory and applications.", Holt,  
Rinehart and Winston Inc., N.Y., 1965

K-1 KALMAN R.E. "Contributions to the theory of optimal control". Bol. Soc.  
Mat. Mex., vol 5, 1960, pp 102-119

L-1 LUENBERGER D.G.: "Optimization by vector space methods.", J. Wiley,  
N.Y., 1969

R-1 RAVIART P.A. "Méthodes des éléments finis", cours du III<sup>e</sup> cycle à  
Paris VI et CNRS, Paris, 1971

S-1 SCHULTZ M.H.: "Multivariate spline functions and elliptic problems",  
ed. I.J. SCHOENBERG, Academic Press, N.Y., 1969

S-2 SCHUL TZ M.H.: "A Ritz method for an optimal control problem", Journ.  
of Opt. theory and appl., vol 11, n° 3, 1973

T-1 YOSIDA K.: "Functional analysis", Springer Verlag, Berlin, 1971

Z-1 ZADEH L.A. and DESOER C.A.: "Linear System Theory: the state  
space approach", McGraw Hill, N.Y., 1963