



HAL
open science

Méthodes numériques appliquées à la conception par éléments finis

David Dureisseix

► **To cite this version:**

David Dureisseix. Méthodes numériques appliquées à la conception par éléments finis. DEA. ENS cachan - Université de Montpellier 2, 2008, pp.98. cel-00351713v1

HAL Id: cel-00351713

<https://cel.hal.science/cel-00351713v1>

Submitted on 10 Jan 2009 (v1), last revised 3 Jan 2014 (v2)

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

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

Méthodes numériques appliquées à la conception par éléments finis

David Dureisseix



28/09/2008

TABLE DES MATIÈRES

INTRODUCTION	1
1 RÉSOLUTION DES SYSTÈMES LINÉAIRES	5
1.1 Conditionnement	6
1.1.1 Définition	6
1.1.2 Influence des perturbations	7
1.2 Techniques de stockage et coût	8
1.2.1 Paramètres influents	8
1.2.2 Coût	9
1.3 Méthodes directes	10
1.3.1 Méthode d'élimination de Gauss	10
1.3.2 Factorisation de Cholesky	12
1.3.3 Factorisation de Crout	13
1.3.4 Lien avec la condensation de Schur	13
1.3.5 Autres techniques	14
1.4 Quelques particularités en éléments finis	14
1.4.1 Conditions aux limites en déplacement	14
1.4.2 Relations linéaires entre degrés de liberté	16
1.4.3 Système singulier : semi défini positif	16
1.5 Méthodes itératives stationnaires	17
1.5.1 Méthode de Jacobi	18
1.5.2 Méthode de Gauss-Seidel	19
1.5.3 Méthode de relaxation	19
1.5.4 Méthode d'Uzawa	20
1.6 Méthodes itératives non-stationnaires	21
1.6.1 Méthode générique de projection	21
1.6.2 Méthode du gradient	21
1.6.3 Méthode du gradient conjugué	22
1.6.4 Préconditionnement	23
1.6.5 Autres méthodes itératives	23
1.7 Critères d'arrêt	23
2 SYSTÈMES AUX VALEURS PROPRES	25
2.1 Préliminaires	25
2.1.1 Rappels	25
2.1.2 Domaine d'emploi des différentes méthodes	26
2.2 Méthode de Jacobi	27
2.3 Méthode des puissances	28
2.4 Méthode des itérations inverses	30
2.5 Méthode des itérations inverses avec décalage spectral	30
2.5.1 Principe et convergence	30
2.5.2 Exemples d'utilisation dans un code éléments finis	31
2.6 Transformation de Householder	31
2.6.1 Factorisation QR	31
2.6.2 Détermination des valeurs propres par la méthode QR	32
2.6.3 Forme de Hessenberg	33

2.7	Exemple d'utilisation dans un code éléments finis	34
2.7.1	Réduction dans un sous-espace	34
2.7.2	Technique itérative	35
2.8	Méthode de Lanczos	35
3	RÉSOLUTION DE SYSTÈMES NON-LINÉAIRES PARTICULIERS	37
3.1	Non-linéarités matérielles de type plasticité	37
3.1.1	Formulation du problème	37
3.1.2	Résolution pour un incrément de temps	38
3.1.3	Méthode de la plus grande pente	39
3.1.4	Méthodes de gradient conjugué	39
3.1.5	Méthode de Newton-Raphson	39
3.1.6	Méthode de Newton-Raphson modifiée	39
3.1.7	Mises à jour de type Quasi-Newton	40
4	RÉSOLUTION DE SYSTÈMES DIFFÉRENTIELS	43
4.1	Équation scalaire du premier ordre	43
4.1.1	Problème de Cauchy	43
4.1.2	Ré-écriture du problème	43
4.1.3	Schémas explicites à un pas	44
4.1.4	Quelques problèmes pratiques	45
4.1.5	A-stabilité	46
4.1.6	Méthode des trapèzes généralisée	47
4.1.7	Exemple de méthode à pas multiples : Adams	48
4.1.8	Exemple de méthode de prédiction correction	50
4.1.9	Méthode de Runge-Kutta	51
4.1.10	Autres méthodes	52
4.2	Systèmes différentiels du premier ordre	53
4.2.1	Emploi des méthodes précédentes	53
4.2.2	Choix d'un schéma : influence de l'ordre d'intégration	54
4.3	Systèmes différentiels du deuxième ordre	55
4.3.1	Méthode de Newmark	55
4.3.2	Méthode de Gear implicite à deux pas	59
4.3.3	θ -méthode	60
4.3.4	α -HHT	61
4.3.5	Autres méthodes	61
4.4	Intégration de modèles de comportement	61
4.4.1	Exemple de modèle de plasticité	61
4.4.2	Méthodes de type retour radial	62
4.4.3	Exemple de modèle de viscoplasticité	64
4.5	Cas de la dynamique non régulière	66
	RÉFÉRENCES	69
A	FACTORISATION DE CHOLESKY - ORTHOGONALISATION D'UN SOUS-ESPACE	73
B	BORNAGE PAR LES VALEURS PROPRES ÉLÉMENTAIRES	75
B.1	Système aux valeurs propres généralisé	75
B.2	Système aux valeurs propres	76
C	TECHNIQUE D'ACCÉLÉRATION D'AITKEN	77
D	REPÈRES BIOGRAPHIQUES	79
	INDEX	85

INTRODUCTION

Ce document de travail (avec renvoi à des ressources de références s'il y a besoin d'approfondir) est utile à ceux qui sont intéressés ou amenés à utiliser des méthodes de simulation (essentiellement par éléments finis) dans le but de dimensionner des structures mécaniques.

Il ne s'agit pas pour autant d'un document sur les éléments finis, puisqu'on suppose que les problèmes ont déjà été formulés dans ce cadre, et que l'on se place en aval, lorsque la résolution des problèmes éléments finis est en jeu. On se ramènera ainsi fréquemment aux techniques effectivement utilisées dans des codes de calcul industriels. Pour cela, on sera amené à utiliser des méthodes classiques d'algèbre linéaire, de calcul matriciel, adaptées aux calculs sur ordinateurs (algorithmique).

Selon les problèmes qui se posent, on est amené à utiliser un grand nombre d'outils différents mais ce document ne se veut pas exhaustif. Il ne traite en particulier pas les cas suivants :

- la recherche de racine de $f(U) = 0$;
- les problèmes de minimum non quadratiques (bisection, point fixe, Newton, Newton-Raphson, quasi-Newton, sécante, BFGS) ;
- les problèmes non linéaires de type grands déplacement ou contact unilatéral.

Structuré en quatre chapitres, il aborde les techniques de résolution de grands systèmes linéaires (de façon directe ou itérative), de problèmes de recherche de valeurs et vecteurs propres, et enfin de systèmes différentiels.

Pour un aperçu moins détaillé mais plus large des méthodes numériques, le lecteur pourra se reporter à [1].

Il semble aussi utile de préciser que ce document est en perpétuelle évolution et que de ce fait, son organisation ainsi que sa table des matières peuvent afficher quelque incohérence.

En cas d'erreur

Les différents algorithmes et techniques présentés ne sont pas garantis contre les erreurs et autres coquilles... De la même façon, tout document est perfectible, et celui-ci ne fait pas exception à la règle. Aussi, si vous en trouvez, signalez-les : c'est aussi une méthode itérative... pour réduire le nombre de bugs. J'ai déjà eu des retours qui m'ont permis d'améliorer la chose. Merci donc en particulier à Martin Genet.

Tout lecteur ayant des remarques et propositions constructives est amicalement invité à envoyer un retour à l'auteur par courriel à David.Dureisseix@lmgc.univ-montp2.fr.

Notations

Un vecteur est représenté par une lettre romane minuscule (comme \mathbf{x}), une matrice par une lettre romane majuscule (comme \mathbf{A}), un scalaire par une lettre grecque minuscule (comme α). La composante i du vecteur \mathbf{x} est notée x_i (c'est un scalaire) et la composante i, j de la matrice \mathbf{A} , $A_{i,j}$: c'est un scalaire également.

La colonne j de \mathbf{A} , par exemple, sera classiquement notée $\mathbf{A}_{1:n,j}$. Les vecteurs de la base dans laquelle les différentes matrices sont écrites seront notés \mathbf{e}_j .

À chaque fois, l'itéré numéro i est désigné par une mise en exposant (comme $\mathbf{x}^{(i)}$) et en général, une même lettre désignera en fait un même emplacement en mémoire ; par exemple, $\mathbf{x}^{(i)}$ écrase $\mathbf{x}^{(i-1)}$ en mémoire.

Pour approfondir

Pour des renseignements généraux complémentaires sur les techniques employées, il est utile de ce référer par exemple à [2, 3]... et pour les algorithmes et leur implantation, aux bibliothèques de netlib...

De façon générale, la bibliographie présentée regroupe deux types de références : les anciennes, voire très anciennes, en tant que base historique, et les plus récentes, en tant qu'ouvrages contemporains. Les biographies utilisées se sont largement inspirées des bases de données disponibles aux adresses suivantes :

- <http://www-history.mcs.st-andrews.ac.uk/>
- <http://www.math.bme.hu/html/index.html>
- <http://www.wikipedia.org>

Cadre de travail

Actuellement, et pour ce qui nous intéresse, nous avons tout intérêt à utiliser des bibliothèques scientifiques existantes de façon intensive : on ne peut prétendre à reprogrammer de façon plus efficace et plus robuste des algorithmes standards que des spécialistes de ce domaine et qui travaillent sur ces bibliothèques depuis un grand nombre d'années. Il est cependant nécessaire de connaître le principe d'un certain nombre de ces méthodes, leur domaine d'emploi, leur limite et d'avoir une idée des cas à problème pour aider au choix d'une méthode. C'est dans cette direction que se place ce document.

Parmi ces bibliothèques, en grande partie du domaine public, certaines émergent ou commencent à émerger comme des standards de fait. Par exemple, la bibliothèque LAPACK (Linear Algebra PACKage), [4] :

<http://www.netlib.org/lapack/>

qui succède en grande partie aux anciennes bibliothèques EISPACK, pour la résolution de systèmes aux valeurs propres :

<http://www.netlib.org/eispack/>

et LINPACK pour la résolution de systèmes linéaires :

<http://www.netlib.org/linpack/>

Sa description est celle de la FAQ (*Frequently Asked Questions*) :

LAPACK provides routines for solving systems of simultaneous linear equations, least-squares solutions of linear systems of equations, eigenvalue problems, and singular value problems. The associated matrix factorizations (LU, Cholesky, QR, SVD, Schur, generalized Schur) are also provided, as are related computations such as reordering of the Schur factorizations and estimating condition numbers. Dense and banded matrices are handled, but not general sparse matrices. In all areas, similar functionality is provided for real and complex matrices, in both single and double precision.

Citons aussi sa version parallélisée, ScaLAPACK (Scalable LAPACK, [5]) :

The ScaLAPACK project is a collaborative effort involving several institutions (Oak Ridge National Laboratory, Rice University, University of California, Berkeley, University of California, Los Angeles, University of Illinois, University of Tennessee, Knoxville), and comprises four components : dense and band matrix software (ScaLAPACK), large sparse eigenvalue software (PARPACK and ARPACK), sparse direct systems software (CAPSS and MFACT), preconditioners for large sparse iterative solvers (ParPre). The ScaLAPACK (or Scalable LAPACK) library includes a subset of LAPACK routines redesigned for distributed memory MIMD parallel computers. It is currently written in a Single-Program-Multiple-Data style using explicit message passing for interprocessor communication. It assumes matrices are laid out in a two-dimensional block cyclic decomposition. Like LAPACK, the ScaLAPACK routines are based on block-partitioned algorithms in order to minimize the frequency of data movement between different levels of the memory hierarchy. (For such machines, the memory hierarchy includes the off-processor memory of other processors, in addition to the hierarchy of registers, cache, and local memory on each processor.) The fundamental building blocks of the ScaLAPACK library are distributed memory versions (PBLAS) of the Level 1, 2 and 3 BLAS, and a set of Basic Linear Algebra Communication Subprograms (BLACS) for communication tasks that arise frequently in parallel linear algebra computations. In the ScaLAPACK routines, all interprocessor communication occurs within the PBLAS and the BLACS. One of the design goals of ScaLAPACK was to have the ScaLAPACK routines resemble their LAPACK equivalents as much as possible.

D'un niveau d'utilisation plus bas, et utilisées dans les bibliothèques précédentes, on trouve aussi des utilitaires intéressants dans les bibliothèques ARPACK :

<http://www.caam.rice.edu/software/ARPACK/>

ARPACK is a collection of Fortran77 subroutines designed to solve large scale eigenvalue problems. The package is designed to compute a few eigenvalues and corresponding eigenvectors of a general $n \times n$ matrix \mathbf{A} . It is most appropriate for large sparse or structured matrices \mathbf{A} where structured means that a matrix-vector product $w \leftarrow Av$ requires order n rather than the usual order n^2 floating

point operations. This software is based upon an algorithmic variant of the Arnoldi process called the Implicitly Restarted Arnoldi Method (IRAM). When the matrix \mathbf{A} is symmetric it reduces to a variant of the Lanczos process called the Implicitly Restarted Lanczos Method (IRLM). These variants may be viewed as a synthesis of the Arnoldi/Lanczos process with the Implicitly Shifted QR technique that is suitable for large scale problems. For many standard problems, a matrix factorization is not required. Only the action of the matrix on a vector is needed.

et BLAS (Basic Linear Algebra Subprograms) :

<http://www.netlib.org/blas/>

The BLAS (Basic Linear Algebra Subprograms) are high quality “building block” routines for performing basic vector and matrix operations. Level 1 BLAS do vector-vector operations, Level 2 BLAS do matrix-vector operations, and Level 3 BLAS do matrix-matrix operations. Because the BLAS are efficient, portable, and widely available, they are commonly used in the development of high quality linear algebra software, LINPACK and LAPACK for example.

Les développements concernant ces bibliothèques sont toujours en cours, en particulier les versions creuses SPBLAS (SParse BLAS) et parallèles PBLAS (Parallel BLAS).

Enfin, des bibliothèques scientifiques généralistes sont aussi disponibles, comme la GSL (GNU Scientific Library) :

<http://www.gnu.org/software/gsl/>

1

RÉSOLUTION DES SYSTÈMES LINÉAIRES

On est souvent, voire tout le temps, amené à résoudre des systèmes linéaires de la forme $\mathbf{Ax} = \mathbf{b}$ pour la simulation sur ordinateur ; par exemple en éléments finis [6], une analyse statique conduit à résoudre un système que l'on notera $\mathbf{Ku} = \mathbf{f}$ où \mathbf{K} est la matrice de rigidité, \mathbf{u} le déplacement en certains points de la structure et \mathbf{f} , les forces généralisées.

Une analyse non-linéaire implicite, la recherche d'une base modale... sont souvent basées sur la *brique de base* qu'est la résolution de systèmes linéaires. Dans ces cas de figure, comme dans d'autres, les matrices \mathbf{A} et \mathbf{b} correspondantes ont souvent des particularités dont il sera judicieux de tirer parti pour les avantages, ou qu'il faudra avoir en tête pour les inconvénients. Par exemple :

1. une grande dimension $\dim \mathbf{A} = n$, classiquement de 10^3 à 10^6 , donc :
 - on n'inverse pas \mathbf{A} (ce serait résoudre n systèmes linéaires) ;
 - il faut prendre en compte des considérations de stockage adapté, à la fois pour réduire les ressources en mémoire et le nombre d'opérations à effectuer.
2. la nécessité de traiter plusieurs second membres :
 - simultanément ;
 - successivement.
3. des propriétés souvent rencontrées :
 - des coefficients réels, une matrice \mathbf{A} carrée,
 - \mathbf{A} inversible $\Leftrightarrow \mathbf{A}$ régulière $\Leftrightarrow \det \mathbf{A} \neq 0$
 - \mathbf{A} symétrique $\Leftrightarrow \mathbf{A} = \mathbf{A}^T$ avec $\forall \mathbf{x}, \mathbf{y}, \mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y}, \mathbf{x}^T \mathbf{A} \mathbf{y} = \mathbf{y}^T \mathbf{A}^T \mathbf{x}$
 - \mathbf{A} positive $\Leftrightarrow \forall \mathbf{x}, \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$
 - \mathbf{A} définie $\Leftrightarrow \mathbf{x}^T \mathbf{A} \mathbf{x} = 0 \Rightarrow \mathbf{x} = \mathbf{0}$
 - \mathbf{A} mal conditionnée, typiquement $\text{cond} \mathbf{K} = \mathcal{O}(\frac{1}{h^2})$ où h est la taille de la maille en éléments finis massifs.

Pour plus d'informations concernant ces particularités et les techniques mises en œuvre, on pourra se référer à [7, 8].

1.1 Conditionnement

1.1.1 Définition

Pour \mathbf{A} symétrique et inversible,

$$\text{cond}\mathbf{A} = \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\| \quad (1.1)$$

Avec une norme matricielle cohérente avec celle définie pour les vecteurs :

$$\|\mathbf{A}\| = \sup_{\|\mathbf{x}\| \neq 0} \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} \quad (1.2)$$

où $\|\mathbf{x}\|^2 = \mathbf{x}^T \mathbf{x} = \sum_i x_i^2$ est la norme euclidienne, on a :

$$\text{cond}\mathbf{A} = \frac{|\lambda|_{\max}}{|\lambda|_{\min}} \quad (1.3)$$

où les λ sont les valeurs propres de \mathbf{A} .

En général, il est très coûteux d'atteindre le conditionnement de \mathbf{A} , par contre un certains nombre de techniques peuvent permettre de l'estimer. On pourra en particulier s'inspirer des techniques développées dans le chapitre concernant la recherche de valeurs propres. Un conditionnement grand conduit à des problèmes de *précision* (cumul d'arrondis numériques). Pour les illustrer, on est amené à considérer un système perturbé (par des perturbations généralement supposées petites) $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$.

Exemple 1

Considérons le système $\mathbf{Ax} = \mathbf{b}$ suivant :

$$\begin{bmatrix} 1 & 1 \\ 1 & 1 + \varepsilon \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix} \Rightarrow \mathbf{x} = \begin{pmatrix} 2 \\ 0 \end{pmatrix} \quad (1.4)$$

Le polynôme caractéristique $\lambda^2 - (2 + \varepsilon)\lambda + \varepsilon = 0$ permet de déterminer les valeurs propres, à savoir :

$$\lambda_{1,2} = \frac{2 + \varepsilon \pm \sqrt{4 + \varepsilon^2}}{2} \quad (1.5)$$

Avec ε petit, le conditionnement est donc très grand puisque $\text{cond}\mathbf{A} \approx 4/\varepsilon$. Considérons maintenant le système $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$ perturbé tel que :

$$\begin{bmatrix} 1 & 1 \\ 1 & 1 + \varepsilon \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 + \alpha \end{pmatrix} \Rightarrow \mathbf{x} = \begin{pmatrix} 2 - \frac{\alpha}{\varepsilon} \\ \frac{\alpha}{\varepsilon} \end{pmatrix} \quad (1.6)$$

Ce résultat entraîne la discussion suivante :

- si $\alpha \ll \varepsilon \ll 1$ alors $\mathbf{x} \approx [2 \ 0]^T$;
- si $\alpha \approx \varepsilon \ll 1$, ce qui est tout à fait possible avec les arrondis numériques, la solution est très différente puisque $\mathbf{x} \approx [1 \ 1]^T$.

Exemple 2

Examinons le système suivant :

$$\mathbf{A} = \begin{bmatrix} 1,2969 & 0,8648 \\ 0,2161 & 0,1441 \end{bmatrix} \quad \mathbf{b} = \begin{pmatrix} 0,8642 \\ 0,1440 \end{pmatrix} \quad \mathbf{x} = \begin{pmatrix} 0,8642 \\ 0,1440 \end{pmatrix} \quad (1.7)$$

Le conditionnement de \mathbf{A} est $\text{cond}\mathbf{A} = 3,3 \cdot 10^8$; le résidu est petit puisque $\mathbf{b} - \mathbf{Ax} = [-10^{-8} \ 10^{-8}]^T$ alors que la solution exacte est très différente de \mathbf{x} puisque c'est $[2 \ -2]^T$. Il semble donc que le conditionnement du système soit très fortement lié à la qualité de la solution (en terme d'influence des arrondis numériques) que l'on puisse espérer avoir.

Dans le cas des éléments finis massifs (autre que plaques en flexion, poutres...), le conditionnement varie en $\mathcal{O}(\frac{1}{h^2})$. La figure 1.1 illustre cela dans un cas 2D très simple. Le conditionnement a été obtenu en calculant explicitement la plus grande et la plus petite valeur propre de \mathbf{K} avec les méthodes qui seront présentées ultérieurement.

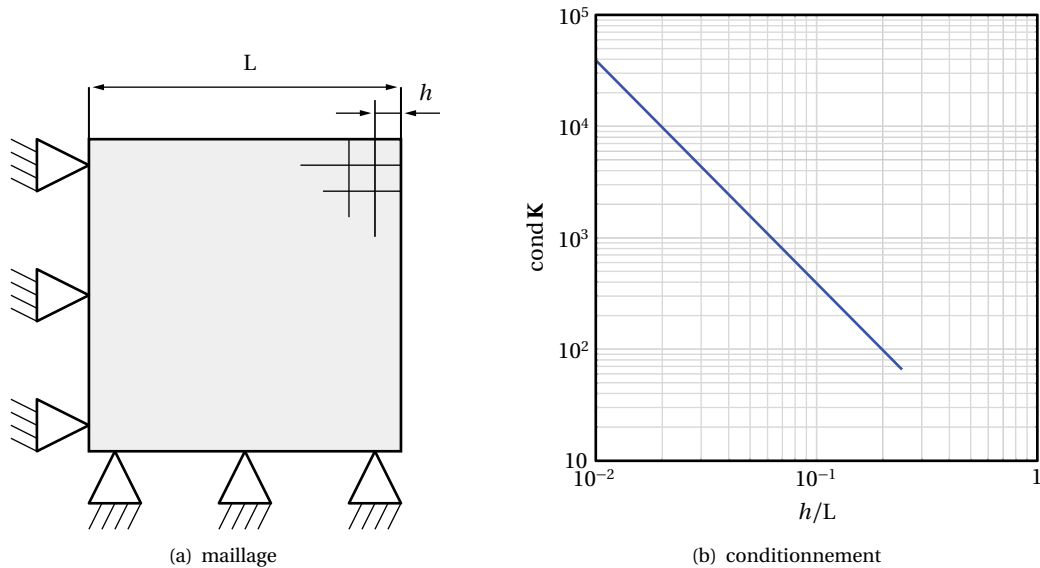


fig. 1.1 - Conditionnement en éléments finis

1.1.2 Influence des perturbations

Considérons le système $\mathbf{Ax} = \tilde{\mathbf{b}}$. Si $\tilde{\mathbf{b}} = \mathbf{b} + \delta\mathbf{b}$ avec $\|\delta\mathbf{b}\| \ll \|\mathbf{b}\|$, a-t-on $\|\delta\mathbf{x}\| \ll \|\mathbf{x}\|$ pour $\tilde{\mathbf{x}} = \mathbf{x} + \delta\mathbf{x}$? Calculons simplement, en utilisant l'inégalité de Milkowski :

$$\mathbf{Ax} = \mathbf{b} \Rightarrow \|\mathbf{Ax} = \mathbf{b}\| \Rightarrow \|\mathbf{A}\| \cdot \|\mathbf{x}\| \geq \|\mathbf{b}\| \quad (1.8)$$

autrement dit :

$$\frac{1}{\|\mathbf{x}\|} \leq \frac{\|\mathbf{A}\|}{\|\mathbf{b}\|} \quad (1.9)$$

et maintenant :

$$\mathbf{A}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b} + \delta\mathbf{b} \Rightarrow \mathbf{A}\delta\mathbf{x} = \delta\mathbf{b} \Rightarrow \delta\mathbf{x} = \mathbf{A}^{-1}\delta\mathbf{b} \Rightarrow \|\delta\mathbf{x}\| \leq \|\mathbf{A}^{-1}\| \cdot \|\delta\mathbf{b}\| \quad (1.10)$$

soit :

$$\|\delta \mathbf{x}\| \leq \|\mathbf{A}^{-1}\| \cdot \|\delta \mathbf{b}\| \quad (1.11)$$

par conséquent :

$$\frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} \leq \text{cond} \mathbf{A} \frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|} \quad (1.12)$$

Considérons maintenant le système perturbé $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \mathbf{b}$. Si $\tilde{\mathbf{A}} = \mathbf{A} + \delta \mathbf{A}$ et $\tilde{\mathbf{x}} = \mathbf{x} + \delta \mathbf{x}$, alors :

$$\begin{aligned} (\mathbf{A} + \delta \mathbf{A})(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} &\Rightarrow \mathbf{A}\delta \mathbf{x} + \delta \mathbf{A}(\mathbf{x} + \delta \mathbf{x}) = \mathbf{0} \Rightarrow \delta \mathbf{x} = -\mathbf{A}^{-1}\delta \mathbf{A}(\mathbf{x} + \delta \mathbf{x}) \\ \Rightarrow \|\delta \mathbf{x}\| &\leq \|\mathbf{A}^{-1}\| \cdot \|\delta \mathbf{A}\| \cdot \|\mathbf{x} + \delta \mathbf{x}\| \end{aligned} \quad (1.13)$$

autrement dit :

$$\frac{\|\delta \mathbf{x}\|}{\|\mathbf{x} + \delta \mathbf{x}\|} \leq \text{cond} \mathbf{A} \frac{\|\delta \mathbf{A}\|}{\|\mathbf{A}\|} \quad (1.14)$$

Il n'est donc pas étonnant qu'on se trouve face à une difficulté lorsque le système à résoudre est très mal conditionné, autant avec les méthodes directes qui fournissent une solution erronée qu'avec les méthodes itératives qui ne convergeront pas.

1.2 Techniques de stockage et coût

1.2.1 Paramètres influents

Dans les applications envisagées, en particulier pour les systèmes de grande taille, il est nécessaire d'envisager des techniques de stockage adaptées. Par exemple, si la taille du système est $n = 10^3$, le stockage de \mathbf{A} pleine nécessite 8 Mo de mémoire ; avec $n = 10^6$, 8 000 Go sont nécessaires.

Pire encore, le nombre d'opérations pour résoudre le système $\mathbf{Ax} = \mathbf{b}$ est lié au nombre de termes non nuls (en fait au nombre de termes stockés, donc supposés non nuls) de la matrice, et aussi de la numérotation des inconnues, comme on le verra. Par exemple, une méthode de Cramer sur une matrice pleine demande un nombre d'opérations en virgule flottante de $n_{\text{op}} \approx (n+1)!$, ce qui devient rapidement inextricable, même avec les ordinateurs les plus puissants [9, 10].

Dans le cas des systèmes issus d'une discrétisation par éléments finis, le système est dit creux¹. Il faut donc en tenir compte au niveau du stockage, en utilisant par exemple un stockage :

- bande (largeur de bande, *bandwidth*, l_b) ;
- ligne de ciel, *skyline* ;
- morse ;
- semi-morse.

Quant à l'influence de la renumérotation des degrés de liberté (ddl), des méthodes heuristiques sont employées car la détermination de la numérotation optimale est un problème extrêmement coûteux (en fait, plus coûteux que le système que l'on cherche à résoudre). Les techniques les plus employées se nomment :

1. en anglais, un système creux est dit *sparse*

- Cuthill-MacKee inverse [11] ;
- *nested dissection* [12] ;
- degré minimum [13, 14]. . .

1.2.2 Coût

Le coût d'un algorithme s'estime de différentes manières. Une première indication est la taille de la mémoire requise pour traiter un problème donné. Une autre est le temps mis pour effectuer la résolution du problème. Deux types de prise de temps sont généralement employés : le temps CPU (Central Processing Unit), lié aux opérations effectuées, et le temps horloge, lié au temps extérieur écoulé pendant cette résolution. Il est certain que ce dernier est une mesure très globale puisqu'elle fait intervenir la machine utilisée², les compétences du programmeur, et peuvent dépendre fortement du problème traité. Elle ne caractérise donc pas uniquement l'algorithme mais plutôt globalement le logiciel réalisé pour un problème donné sur une machine donnée.

D'autres indicateurs sont également disponibles, comme le nombre d'opérations en virgule flottante nécessaire à la résolution. C'est ce que l'on appelle la *complexité* de l'algorithme. Elle devient plus liée à l'algorithme même en faisant des estimations sur les types de problèmes à traiter, mais elle ne fait pas intervenir les éventuels transferts en mémoire, opérations sur les entiers. . .

La vraie qualification d'une technique doit donc faire appel à tous ces indicateurs. Ici, seule la complexité sera présentée par la suite. L'estimation du coût effectif en temps CPU des différentes opérations en virgule flottante peut ensuite être proposée par des essais numériques et des prises de temps sur des architectures de machines cibles.

À titre d'indication, un petit test a été réalisé dans ce sens, qui fait intervenir essentiellement les performances du processeur. Les coûts indiqués dans le tableau 1.1 ont été obtenus en janvier 2000. Entre parenthèses, on trouvera les coûts obtenus l'année ou les deux années précédentes. Le code est écrit en Fortran 77 ; les différentes machines testées ainsi que les options de compilation sont les suivantes :

- fronsac est une station de travail HP du LMT-Cachan (HP 9000/879, système B.10.20), la compilation est faite avec f77 +O2 ;
- bacchus est une machine parallèle SGI 02000 du Pôle Parallélisme Île de France Sud (SGI IP27 système 6.5), et la compilation, f90 -O2 ;
- artemis est une machine IBM SP de l'Antenne de Bretagne de l'ENS de Cachan (IBM AIX système 3), et la compilation, f77 -O2.

Les opérations dénommées daxpy (série de 2 opérations en fait : addition et multiplication), dscal (série de multiplications), dgemv (calcul d'un produit matrice \mathbf{A} par vecteur \mathbf{x} , pondéré par un coefficient β et ajouté à un précédent vecteur \mathbf{y} , lui même multiplié par un coefficient α) sont appelées avec la librairie BLAS, déjà mentionnée. L'opération dénommée $\alpha\mathbf{y} + \beta\mathbf{Ax}$ est donc l'équivalent de dgemv, programmé directement en Fortran. À partir de

2. les performances du système, du compilateur, les accès aux disques

type d'opération (moyennée sur 10^8 termes)	coût (10^{-8} s CPU)		
	fronsac	bacchus	artemis
addition	5,2 (5,7 ; 7)	5,54 (6,2)	2,67
multiplication	5,2 (5,7 ; 7)	5,54 (6,2)	2,67
division	9	9 (11,7)	2,73
racine carrée	9	14,5 (18,3)	18,2
daxpy	8,9 (9,6 ; 10)	4,66 (5,5)	2,02
dscal	4,97 (5,3 ; 6)	2,56 (3,9)	1,35
type d'opération (moyennée sur 100 vecteurs de 1000 termes)	coût (10^{-2} s CPU)		
	fronsac	bacchus	artemis
$\alpha\mathbf{y} + \beta\mathbf{Ax}$	4,95 (5,2)	3,4 (3,9)	0,84
dgemv	-	2 (2,3)	1,1

tab. 1.1 - estimation du coût des opérations élémentaires

ces quelques résultats, on peut tirer quelques conclusions : outre le fait qu'utiliser les options d'optimisation des compilateurs est intéressante, l'utilisation de bibliothèques (ici la BLAS) est plus efficace que reprogrammer des routines élémentaires soi-même. En fait, l'utilisation de routines standards augmente la lisibilité du programme, permet de programmer à un niveau plus élevé et donc réduit les temps de « debuggage », et permet souvent d'atteindre de meilleures performances, puisque ces bibliothèques sont programmées de façon plus efficace (souvent directement en assembleur) et sont optimisées pour la machine sur laquelle elles sont utilisées.

D'autre part, pour l'utilisateur normal, le coût effectif (en secondes CPU) dépend d'un grand nombre de paramètres comme l'efficacité du compilateur, de l'environnement de développement, des temps de latence de la mémoire, des caches, de l'accès aux disques, etc. En particulier, le changement de version du système d'exploitation d'une année sur l'autre influe sur les résultats du test présenté précédemment.

1.3 Méthodes directes

Encore majoritairement employée dans les codes actuels pour leur robustesse, et d'une précision fine du coût (nombre d'opérations, stockage nécessaire), elles deviennent coûteuses en temps et en encombrement lorsque le système à résoudre devient très grand. Pour plus d'informations sur ces méthodes, on pourra avantageusement consulter [15].

1.3.1 Méthode d'élimination de Gauss

Principe

On simplifie le système pour le rendre triangulaire supérieur par utilisation de combinaison linéaires de lignes (c'est-à-dire d'équations).

Préliminaire : élimination de Gauss par bloc

Écrivons le système à résoudre sous une forme de blocs :

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{pmatrix} \quad (1.15)$$

Si \mathbf{A}_{11} est inversible, la première équation conduit à exprimer \mathbf{x}_1 en fonction de \mathbf{x}_2 :

$$\mathbf{x}_1 = \mathbf{A}_{11}^{-1}(\mathbf{b}_1 - \mathbf{A}_{12}\mathbf{x}_2) \quad (1.16)$$

En utilisant cette expression dans la deuxième équation, on se ramène à un système avec \mathbf{x}_2 comme seule inconnue :

$$\underbrace{\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}}_{\mathbf{S}}\mathbf{x}_2 = \underbrace{\mathbf{b}_2 - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{b}_1}_{\mathbf{c}} \quad (1.17)$$

\mathbf{S} est appelé complément de Schur ou matrice condensée sur les degrés de liberté \mathbf{x}_2 . Connaissant \mathbf{x}_2 , on remonte à \mathbf{x}_1 avec (1.16). Le système de départ (1.15) est équivalent à :

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{c} \end{pmatrix} \quad (1.18)$$

et on peut le réécrire :

$$\underbrace{\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{0} & \mathbf{S} \end{bmatrix}}_{\mathbf{A}} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{I} \end{bmatrix} \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{c} \end{pmatrix} \quad (1.19)$$

Utilisation pour la factorisation LU

Si on prend comme bloc 11 la première équation scalaire du système, il faut avoir $\mathbf{A}_{11} \neq 0$ (c'est le premier *pivot*). En continuant sous la forme de l'équation (1.19), récursivement sur les sous-matrices qui restent, on obtient la factorisation $\mathbf{A} = \mathbf{LU}$ (en $n - 1$ étapes) avec \mathbf{U} (pour *upper*), matrice triangulaire supérieure et \mathbf{L} (pour *lower*), matrice triangulaire inférieure avec diagonale unitaire.

Pour des raisons d'efficacité de stockage, \mathbf{L} et \mathbf{U} écrasent en mémoire \mathbf{A} si on ne stocke pas la diagonale de \mathbf{L} qu'on sait unitaire. À chaque itération k , la matrice de départ \mathbf{A} est donc modifiée en une matrice $\mathbf{A}^{(k)}$.

Conditions d'emploi

Si à chaque étape de la factorisation, le pivot $A_{kk}^{(k-1)}$ est non nul, on peut progresser sans problème. La seule condition est que \mathbf{A} soit *inversible*, mais quand on trouve un pivot nul, il faut procéder à des échanges d'équations. En pratique, on a bien entendu aussi des problèmes si le pivot est petit en valeur absolue ; il est parfois conseillé de procéder systématiquement à un pivotage. Cette opération est cependant coûteuse en terme d'accès et de transfert de données en mémoire, surtout si elle nécessite des accès aux disques. Si la matrice assemblée est initialement bande, les permutations détruisent généralement cette propriété, ce qui peut alors devenir très pénalisant.

Stratégies de choix du pivot

Il peut s'agir du pivotage partiel ou du pivotage total :

pivotage partiel

on choisit dans la même colonne k que celle du pivot à changer, la ligne l telle que :

$$|A_{l,k}| = \max_{m=k,\dots,n} |A_{m,k}| \quad (1.20)$$

pour remplacer le pivot ;

pivotage total

on peut aussi permuter les colonnes ; on cherche alors le nouveau pivot dans toute la sous-matrice restant à traiter, on choisit donc la ligne l et la colonne c telles que :

$$|A_{l,c}| = \max_{m,p=k,\dots,n} |A_{m,p}| \quad (1.21)$$

Particularités

Sans recours au pivotage, il y a unicité de la factorisation. Dans le cas où il est nécessaire de traiter plusieurs second membres, on opère simultanément sur les \mathbf{b} . Si les résolutions doivent être réalisées successivement, et pour chaque autre second membre \mathbf{b}' , on résout :

$$\mathbf{LU}\mathbf{x}' = \mathbf{b}' \quad (1.22)$$

- par résolution de $\mathbf{L}\mathbf{y}' = \mathbf{b}'$ appelée descente ;
- puis par résolution de $\mathbf{U}\mathbf{x}' = \mathbf{y}'$ appelée montée.

À chaque fois, il ne s'agit plus que de résoudre un système triangulaire (supérieur ou inférieur). On peut d'autre part calculer à faible coût le déterminant de \mathbf{A} une fois que l'on connaît sa décomposition \mathbf{LU} :

$$\det \mathbf{A} = \det(\mathbf{LU}) = \det \mathbf{L} \det \mathbf{U} = \det \mathbf{U} = \prod_{k=1,n} U_{k,k} \quad (1.23)$$

Complexité

Dans le cas d'une matrice stockée pleine, la factorisation \mathbf{LU} nécessite $\approx \frac{n^3}{3}$ opérations ; une montée (ou 1 descente) nécessite quant à elle $\approx \frac{n^2}{2}$ opérations.

1.3.2 Factorisation de Cholesky

Cette technique a été décrite pour la première fois dans [16] en 1924.

Caractéristiques de la méthode

Même idée que la factorisation \mathbf{LU} quand \mathbf{A} est *symétrique et inversible* : $\mathbf{A} = \mathbf{B}^T \mathbf{B}$ avec \mathbf{B} à diagonale positive. Dans le cas, le stockage et le calcul ne portent que sur un peu plus de la moitié de la matrice, par exemple sur sa partie inférieure.

Propriété

Si \mathbf{A} est symétrique, définie, positive (donc inversible), il n'est pas utile de pivoter.

Complexité

Pour \mathbf{A} pleine, la factorisation $\mathbf{B}^T \mathbf{B}$ requière environ $n^3/6$ opérations. Pour \mathbf{A} bande, la factorisation $\mathbf{B}^T \mathbf{B}$ requière approximativement $n l_b^2/2$ opérations, une montée (ou une descente) coûte de l'ordre de $n l_b$ opérations. Le stockage nécessite aux alentours de $n l_b$ réels (généralement 8 octets). De façon courante, on a $l_b \ll n^3$.

Un exemple d'utilisation de cette factorisation pour l'orthogonalisation des vecteurs d'un sous-espace est présentée dans l'annexe A.

3. $\frac{l_b}{n} = \frac{1}{10}$ à $\frac{1}{100}$

1.3.3 Factorisation de Crout

C'est la méthode la plus employée. La factorisation se met sous la forme $\mathbf{A} = \mathbf{LDL}^T$ avec \mathbf{L} triangulaire inférieure à diagonale unitaire et \mathbf{D} diagonale. Cette factorisation s'applique pour \mathbf{A} symétrique, définie positive, mais traite aussi les cas où les valeurs propres sont négatives ou nulles.

Complexité

Elle est identique à celle de la factorisation de Cholesky, mais on n'a plus besoin d'extraire n racines carrées (pour la diagonale).

L'algorithme 1.3.3 présente cette factorisation dans le cas d'une matrice symétrique \mathbf{A} stockée pleine. Il requière l'utilisation d'un vecteur \mathbf{v} supplémentaire. Sans pivotage, comme pour Cholesky et pour Gauss, la factorisation conserve le profil de la matrice écrasée en mémoire, comme on le verra dans la suite.

```

boucle sur  $j = 1, 2, \dots, n$ 
  boucle sur  $i = 1, 2, \dots, j - 1$ 
     $v_i \leftarrow A_{j,i} A_{i,i}$ 
  fin
   $v_j \leftarrow A_{j,j} - \mathbf{A}_{j,1:j-1} \mathbf{v}_{1:j-1}$ 
   $A_{j,j} \leftarrow v_j$ 
   $\mathbf{A}_{j+1:n,j} \leftarrow \frac{1}{v_j} (\mathbf{A}_{j+1:n,j} - \mathbf{A}_{j+1:n,1:j-1} \mathbf{v}_{1:j-1})$ 
fin

```

alg. 1.2 - Factorisation LDL^T stockée pleine

1.3.4 Lien avec la condensation de Schur

Super-éléments

Comme c'était le cas pour la factorisation LU, la factorisation de Cholesky ou Crout est très liée à la condensation de Schur. Dans le cas des éléments finis, on peut bâtir ainsi la notion de super-élément. Considérons une partie d'une structure qui représente une zone d'intérêt particulier (par exemple, une zone qui reste en élasticité alors que le reste de la structure plastifie, où une zone dans laquelle des fonctions de bases particulières ont été ajoutées comme dans la méthode de partition de l'unité...) et qui a été maillée par éléments finis, figure 1.3. L'interface avec le reste de la structure se fait par le bord de cette zone. Notons avec un indice i les degrés de libertés des nœuds strictement internes à cette zone et \mathbf{b} ceux des nœuds de l'interface, alors :

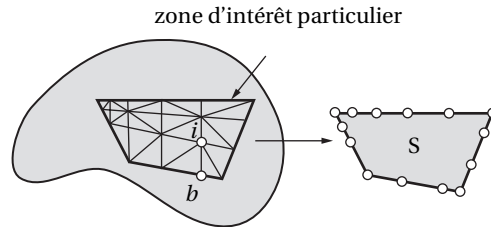


fig. 1.3 - Utilisation d'un super-élément

$$\mathbf{K} = \begin{bmatrix} \mathbf{K}_{ii} & \mathbf{K}_{ib} \\ \mathbf{K}_{bi} & \mathbf{K}_{bb} \end{bmatrix} \quad (1.24)$$

Si on condense sur les nœuds *bord* on obtient :

$$\mathbf{S} = \mathbf{K}_{bb} - \mathbf{K}_{bi} \mathbf{K}_{ii}^{-1} \mathbf{K}_{ib} \quad (1.25)$$

c'est la matrice de rigidité du super-élément considéré. Elle est similaire à une matrice de rigidité élémentaire, pour un *élément* qui peut être lui-même une petite structure.

Conservation de la bande

La figure 1.4 montre, à l'itération k , c'est-à-dire lors de la condensation de l'inconnue k sur les suivantes, les termes de la matrice qui vont être touchés, autrement dit, ceux qui *a priori* sont ou deviennent non-nuls. Clairement, la bande est conservée lors de la factorisation. On peut de la même façon montrer que seuls les termes sous le profil sont modifiés, et donc que la factorisation conserve le profil de la matrice.

1.3.5 Autres techniques

L'adaptation de la factorisation de Crout au cas où A est non-symétrique est classique, c'est la factorisation LDM^T .

D'autres variantes des techniques précédentes permettent de tirer mieux parti des conditions de stockage, en particulier pour optimiser les échanges de données avec la partie des informations stockées sur disque, comme la méthode frontale [17].

D'autres versions encore permettent de réaliser des factorisations en parallèle de façon assez performante, comme les méthodes multi-frontales [18].

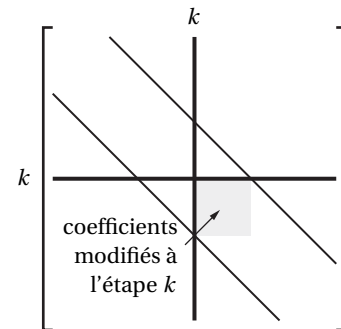


fig. 1.4 - Étape k de la factorisation

1.4 Quelques particularités en éléments finis

1.4.1 Conditions aux limites en déplacement

La prise en compte de déplacements imposés non-nuls a une influence sur l'utilisation des solveurs dans le sens où la matrice de rigidité est modifiée par ces contraintes.

Séparation des degrés de liberté

Le problème correspondant de minimisation du potentiel sous contrainte est le suivant :

$$\min_{\mathbf{C}\mathbf{u}=\mathbf{u}_d} J_1(\mathbf{u}) = \frac{1}{2} \mathbf{u}^T \mathbf{K} \mathbf{u} - \mathbf{f}^T \mathbf{u} \quad (1.26)$$

\mathbf{C} est une matrice booléenne dont le but est d'extraire une partie des degrés de libertés, ceux sur lesquels la condition d'égalité aux déplacements imposés \mathbf{u}_d est imposée.

Reprenons l'écriture par blocs déjà présentés dans l'équation (1.15) en séparant les degrés de liberté imposés $\mathbf{u}_2 = \mathbf{C}\mathbf{u} = \mathbf{u}_d$. Le problème (1.26) est ainsi équivalent à :

$$\min_{\mathbf{u}_1} \frac{1}{2} \mathbf{u}_1^T \mathbf{K}_{11} \mathbf{u}_1 - \mathbf{f}_1^T \mathbf{u}_1 + \mathbf{u}_1^T \mathbf{K}_{12} \mathbf{u}_d + \text{constante} \quad (1.27)$$

dont l'équation d'Euler correspondante est $\mathbf{K}_{11} \mathbf{u}_1 = \mathbf{f}_1 - \mathbf{K}_{12} \mathbf{u}_d$. On a ainsi un problème de taille réduite où les charges extérieures sont modifiées par le terme $-\mathbf{K}_{12} \mathbf{u}_d$ ⁴. Cette méthode est effectivement employée dans des codes comme abaqusTM. L'inconvénient de cette méthode réside dans le tri explicite nécessaire dans les degrés de liberté qui remplit donc le

4. ce sont les réactions à l'appui qui permettent d'avoir un déplacement bien égal à \mathbf{u}_d

terme \mathbf{K}_{12} , et qui entraîne un surcoût important lorsque le nombre de degrés de liberté bloqués est grand.

Multiplicateurs de Lagrange

Une autre technique consiste à reformuler le problème (1.26) en relaxant la contrainte à l'aide de multiplicateur de Lagrange \mathbf{l} et donc de transformer le problème de minimum en celui de la recherche du point selle du Lagrangien :

$$\text{extr}_{\mathbf{u}, \mathbf{l}} J_1(\mathbf{u}) - \mathbf{l}^T (\mathbf{C}\mathbf{u} - \mathbf{u}_d) \quad (1.28)$$

Il n'est plus nécessaire de séparer les degrés de liberté, mais la taille du problème a augmenté ; en effet les équations d'Euler sont cette fois :

$$\underbrace{\begin{bmatrix} \mathbf{K} & -\mathbf{C}^T \\ -\mathbf{C} & \mathbf{0} \end{bmatrix}}_{\mathbf{A}} \begin{pmatrix} \mathbf{u} \\ \mathbf{l} \end{pmatrix} = \begin{pmatrix} \mathbf{f} \\ -\mathbf{u}_d \end{pmatrix} \quad (1.29)$$

On peut remarquer que la première équation montre que les multiplicateurs de Lagrange sont bien les réactions à l'appui.

L'inconvénient vient du fait que la taille du problème augmente, et du fait que \mathbf{A} est une matrice symétrique mais plus définie positive ; elle possède des valeurs propres négatives. On a cependant le résultat suivant si \mathbf{K} est symétrique, définie positive :

$$\mathbf{A} \text{ inversible} \Leftrightarrow \mathbf{C} \text{ injective} \Leftrightarrow \text{kerr } \mathbf{C} = \emptyset \quad (1.30)$$

c'est-à-dire que l'on n'écrit que des liaisons indépendantes. Une factorisation de Gauss risque de détruire la structure bande de la matrice \mathbf{K} et la factorisation de Crout va demander des pivotages.

Une modification de cette technique permet, avec une numérotation ad hoc qui conserve la bande, de ne pas avoir besoin de pivoter. Il s'agit d'une technique de double multiplicateurs utilisée dans le code Cast3MTM. Elle peut être illustrée de la manière suivante :

$$\text{extr}_{\mathbf{u}, \mathbf{l}_1, \mathbf{l}_2} J_1(\mathbf{u}) + \frac{1}{2} \underbrace{\begin{pmatrix} \mathbf{l}_1 \\ \mathbf{l}_2 \\ \mathbf{u} \end{pmatrix}^T \begin{bmatrix} -1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{pmatrix} \mathbf{l}_1 \\ \mathbf{l}_2 \\ \mathbf{u}_2 \end{pmatrix} - \begin{pmatrix} \mathbf{u}_d \\ \mathbf{u}_d \end{pmatrix}^T \begin{pmatrix} \mathbf{l}_1 \\ \mathbf{l}_2 \end{pmatrix}}_{-\frac{1}{2}(l_1 - l_2)^2 + (l_1 + l_2)^T (\mathbf{C}\mathbf{u} - \mathbf{u}_d)} \quad (1.31)$$

dont les équations d'Euler sont :

$$\begin{aligned} \mathbf{K}\mathbf{u} &= \mathbf{f} - \mathbf{C}^T(l_1 + l_2) \\ l_1 &= l_2 \\ \mathbf{C}\mathbf{u} &= \mathbf{u}_d \end{aligned} \quad (1.32)$$

Bien entendu, en cas de nombreux blocages, la taille du problème croît en conséquence.

Valeurs très grandes

On peut aussi considérer une condition aux limites en déplacement imposé comme une modélisation de l'action de l'extérieur. Celui-ci pourrait aussi réagir comme un ressort ayant une grande raideur k , avec un déplacement imposé à la base. Le potentiel élastique de l'ensemble structure + ressort extérieur est alors :

$$J_1(\mathbf{u}) + \frac{1}{2}(\mathbf{C}\mathbf{u} - \mathbf{u}_d)^T k(\mathbf{C}\mathbf{u} - \mathbf{u}_d) \quad (1.33)$$

Les nouvelles équations d'Euler de cette formulation sont alors :

$$(\mathbf{K} + \mathbf{C}^T k \mathbf{C})\mathbf{u} = \mathbf{f} + \mathbf{C}^T k \mathbf{C}\mathbf{u}_d \quad (1.34)$$

On ne fait donc que rajouter des termes dans la matrice de raideur sur les degrés de libertés incriminés et dans le vecteur des forces généralisées sur les efforts duaux.

Le problème est le choix de k : pour des valeurs trop petites, la condition est mal imposée, pour des valeurs trop grandes, le système devient extrêmement mal conditionné, voire numériquement singulier.

1.4.2 Relations linéaires entre degrés de liberté

Cette situation se présente par exemple dans le cas d'utilisation de symétrie par rapport à un plan non parallèle à un des plans des axes globaux de l'analyse, dans le cas de conditions aux limites périodiques sur une cellule de base... À chaque fois, il s'agit de contraintes de la forme $\mathbf{C}\mathbf{u} = \mathbf{u}_d$ où \mathbf{C} n'est cette fois-ci pas forcément booléenne. Les techniques précédentes s'appliquent bien entendu sur ces cas de figure tout aussi bien... et avec les mêmes limitations.

1.4.3 Système singulier : semi défini positif

Exemple

Les exemples sont nombreux : il s'agit de structure avec des modes rigides, de problèmes de sous-intégration... il y a alors des déplacements non nuls à énergie de déformation nulle, autrement dit des modes rigides :

$$\frac{1}{2}\mathbf{u}^T \mathbf{K}\mathbf{u} = 0 \quad (1.35)$$

Dans la suite, on s'intéressera au cas où \mathbf{A} est symétrique, sachant que la généralisation est tout à fait possible.

Noyau

L'ensemble des solutions indépendantes, stockées dans les colonnes de \mathbf{R} , telles que $\mathbf{A}\mathbf{R} = \mathbf{0}$, forme le noyau de \mathbf{A} . S'il y en a p , le rang de \mathbf{A} est $n - p$.

Existence de solutions

Pour qu'il existe des solutions, une condition nécessaire et suffisante est que \mathbf{b} soit orthogonal au noyau de \mathbf{A} , c'est-à-dire $\mathbf{R}^T \mathbf{b} = 0$. Exemple : les efforts ne sollicitent pas les modes de solide rigide.

- Si \mathbf{A} est non symétrique, la condition est que \mathbf{b} appartienne au noyau adjoint de \mathbf{A} .
- Si \mathbf{x}_0 est une solution de ce système, l'ensemble des solutions est de la forme $\mathbf{x}_0 + \mathbf{R}\alpha$ où α est quelconque.

Inverses généralisées (ou pseudo-inverses)

\mathbf{A}^+ telle que $\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$. Si \mathbf{A} est inversible, il n'y a qu'une pseudo-inverse $\mathbf{A}^+ = \mathbf{A}^{-1}$. Pour \mathbf{A} symétrique, toutes les pseudo-inverses sont de la forme $\mathbf{A}^+ + \mathbf{B}^T\mathbf{R}^T + \mathbf{R}\mathbf{B}$, avec \mathbf{B} quelconque [19].

Factorisation et construction d'une inverse généralisée

Supposons que le système puisse s'écrire sous la forme de blocs à l'image de (1.15), où \mathbf{A}_{11} est la plus grande sous-matrice régulière. On a alors $\mathbf{S} = \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12} = \mathbf{0}$. Une pseudo-inverse de \mathbf{A} est :

$$\begin{bmatrix} \mathbf{A}_{11}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (1.36)$$

et la solution est :

$$\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \underbrace{\begin{pmatrix} \mathbf{A}_{11}^{-1}\mathbf{b}_1 \\ \mathbf{0} \end{pmatrix}}_{\mathbf{x}_0} + \underbrace{\begin{bmatrix} -\mathbf{A}_{12} \\ \mathbf{I}_d \end{bmatrix}}_{\mathbf{R}} \underbrace{\mathbf{x}_2}_{\alpha} \quad (1.37)$$

La factorisation fait progresser l'élimination tant qu'un pivot non nul existe. Dans notre cas, avec pivotage, on tombe sur une sous-matrice nulle. Pour les matrices semi-définies, il n'est pas nécessaire de pivoter car un pivot nul apparaît lorsque la ligne du pivot est une combinaison linéaire des précédentes. L'inconnue correspondante est repérée et fixée à $\mathbf{0}$, voir \mathbf{x}_0 dans (1.37), c'est-à-dire qu'on ajoute des liaisons complémentaires pour fixer les modes de solide rigide — à condition que le terme correspondant au second membre condensé soit nul aussi. La colonne correspondante de la matrice stocke le vecteur associé, voir \mathbf{R} dans (1.37).

1.5 Méthodes itératives stationnaires

De manière générale avec les techniques itératives, il faut se poser des questions telles que :

- quelles sont les conditions d'emploi de ces techniques afin d'assurer la convergence vers la solution cherchée ?
- quel critère d'arrêt des itérations employer ?
- quels sont les coûts de ces méthodes ?

Dans le cas des méthodes itératives stationnaires, on rencontre les particularités suivantes :

- on accède à \mathbf{A} uniquement par le produit $\mathbf{A}\mathbf{w}$ (phase la plus coûteuse). Éventuellement, comme en éléments finis, on peut le faire sous forme élémentaire :

$$\mathbf{A}\mathbf{w} = \sum_e \mathbf{A}^e \mathbf{w}^e \quad (1.38)$$

– on construit une suite de solutions approchées :

$$\mathbf{x}^{(k)} \xrightarrow[k \rightarrow \infty]{} \mathbf{A}^{-1} \mathbf{b} \quad (1.39)$$

à partir d'une récurrence linéaire $\mathbf{x}^{(k)} = \mathbf{B}\mathbf{x}^{(k-1)} + \mathbf{c}$ où \mathbf{B} et \mathbf{c} sont indépendantes de k ⁵. Ce sont des méthodes de type *point fixe* : si on converge, c'est vers :

$$\mathbf{x} = \mathbf{B}\mathbf{x} + \mathbf{c} \Rightarrow (\mathbf{x}^{(k)} - \mathbf{x}) = \mathbf{B}(\mathbf{x}^{(k-1)} - \mathbf{x}) \quad (1.40)$$

Pour avoir $\|\mathbf{x}^{(k)} - \mathbf{x}\| \rightarrow_{k \rightarrow \infty} 0$, il faut et il suffit que $\rho(\mathbf{B}) < 1$ où $\rho(\mathbf{B})$ est la *rayon spectral* de \mathbf{B} : c'est le sup des valeurs absolues des valeurs propres de \mathbf{B} . Le taux de convergence asymptotique est défini par :

$$\tau = -\ln \rho(\mathbf{B}) \quad (1.41)$$

Il peut servir à estimer le nombre d'itérations nécessaire pour atteindre le niveau de convergence souhaité ; en effet on a :

$$(\mathbf{x}^{(k)} - \mathbf{x}) = \mathbf{B}^k (\mathbf{x}^{(0)} - \mathbf{x}) \Rightarrow \frac{\|\mathbf{x}^{(k)} - \mathbf{x}\|}{\|\mathbf{x}^{(0)} - \mathbf{x}\|} \leq \|\mathbf{B}^k\| \Rightarrow -\ln \frac{\|\mathbf{x}^{(k)} - \mathbf{x}\|}{\|\mathbf{x}^{(0)} - \mathbf{x}\|} \leq k\tau \quad (1.42)$$

L'intérêt est de pouvoir estimer le nombre d'itérations nécessaires pour diviser le niveau d'erreur par un coefficient donné (utilisation d'un critère d'arrêt), connaissant le taux de convergence asymptotique de la méthode. Pour plus d'informations sur les méthodes itératives stationnaires, le lecteur pourra consulter [20].

1.5.1 Méthode de Jacobi

Pour les systèmes linéaires, cette méthode, d'algorithme 1.5, est décrite pour mémoire seulement, car elle possède un faible taux de convergence et dans nos applications, est remplacée par d'autres techniques plus efficaces. Elle illustre cependant bien les techniques de point fixe : si on suppose connues toutes les inconnues sauf l'inconnue i , cette dernière peut être déterminée à l'aide de l'équation i . Procédant ainsi sur toutes les inconnues, on obtient un nouvel itéré \mathbf{x}^k à partir de l'ancien \mathbf{x}^{k-1} . Cependant, comme la méthode de Gauss-Seidel, elle peut présenter un intérêt dans le cas de systèmes non-linéaires. Sous forme matricielle, cela se traduit par :

$$\mathbf{x}^k = (\mathbf{I} - \mathbf{D}^{-1}\mathbf{A})\mathbf{x}^{k-1} + \mathbf{D}^{-1}\mathbf{b} \quad (1.43)$$

où \mathbf{D} est la matrice diagonale, contenant la diagonale de \mathbf{A} . La matrice d'itération est donc ici $\mathbf{B} = \mathbf{I} - \mathbf{D}^{-1}\mathbf{A}$. On peut noter l'utilisation d'un vecteur de stockage $\bar{\mathbf{x}}$ supplémentaire.

```

Vecteur de départ  $\mathbf{x}^{(0)}$ 
tant que  $R > \epsilon$ 
  itération sur  $k$ 
  boucle sur  $i = 1, n$ 
     $\bar{x}_i \leftarrow \mathbf{A}_{i,1:i-1} \mathbf{x}_{1:i-1}^{(k-1)} + \mathbf{A}_{i,i+1:n} \mathbf{x}_{i+1:n}^{(k-1)}$ 
     $\bar{x}_i \leftarrow (\mathbf{b}_i - \bar{x}_i) / \mathbf{A}_{i,i}$ 
  fin
   $\mathbf{x}^{(k)} \leftarrow \bar{\mathbf{x}}$ 
fin

```

alg. 1.5 - Méthode de Jacobi

5. Cette caractéristique n'existe pas pour les méthodes non stationnaires.

Convergence

Si \mathbf{A} est symétrique, définie positive, alors $\rho_J = \rho(\mathbf{B}) < 1$, mais en éléments finis, on a une évaluation de ce rayon spectral qui varie comme :

$$\rho_J = 1 - \frac{h^2}{2} + \mathcal{O}(h^4) \quad (1.44)$$

d'où un taux de convergence faible lorsque le maillage est raffiné (donc lorsque h est faible) :

$$\tau \approx \frac{h^2}{2} \quad (1.45)$$

1.5.2 Méthode de Gauss-Seidel

```

Vecteur de départ  $\mathbf{x}^{(0)}$ 
tant que  $R > \varepsilon$ 
  itération sur  $k$ 
  boucle sur  $i = 1, 2, \dots, n$ 
     $\mathbf{x}_i^{(k-1)} \leftarrow \mathbf{0}$ 
     $\boldsymbol{\beta} \leftarrow \mathbf{A}_{i,1:n} \mathbf{x}_{1:n}^{(k-1)}$ 
     $\mathbf{x}_i^{(k)} \leftarrow (\mathbf{b}_i - \boldsymbol{\beta}) / A_{i,i}$ 
  fin
fin

```

alg. 1.6 - Méthode de Gauss-Seidel

L'idée est similaire à celle de la méthode de Jacobi, à ceci près que l'on utilise les nouvelles valeurs estimées dès qu'elles sont disponibles, et non pas à la fin de l'itération. Sous forme matricielle, cela se traduit par :

$$\mathbf{x}^k = (\mathbf{D} - \mathbf{E})^{-1} (\mathbf{F} \mathbf{x}^{k-1} + \mathbf{b}) \quad (1.46)$$

où $-\mathbf{E}$ est la partie triangulaire inférieure de \mathbf{A} , et $-\mathbf{F}$ sa partie triangulaire supérieure. La matrice d'itération est ici :

$$\mathbf{B} = (\mathbf{D} - \mathbf{E})^{-1} \mathbf{F} = (\mathbf{I} - \mathbf{D}^{-1} \mathbf{E})^{-1} \mathbf{D}^{-1} \mathbf{F} \quad (1.47)$$

Quand \mathbf{A} est symétrique, définie positive, on est assuré de la convergence puisque $\rho_{GS} < 1$. Mais l'estimation de la valeur du rayon spectral est difficile car il dépend en particulier de la numérotation des degrés de liberté employée.

1.5.3 Méthode de relaxation

Cette technique, SOR pour *successive over relaxation*, est attribuée à Southwell [21] et date de 1946. Elle utilise comme base la méthode de Gauss-Seidel ainsi qu'une technique maintenant souvent utilisée avec d'autres approches : la technique de *relaxation*.

```

Vecteur de départ  $\mathbf{x}^{(0)}$ 
tant que  $R > \varepsilon$ 
  itération sur  $k$ 
  boucle sur  $i = 1, 2, \dots, n$ 
     $\boldsymbol{\alpha} \leftarrow \mathbf{x}_i^{(k-1)}$ 
     $\mathbf{x}_i^{(k-1)} \leftarrow \mathbf{0}$ 
     $\boldsymbol{\beta} \leftarrow \mathbf{A}_{i,1:n} \mathbf{x}_{1:n}^{(k-1)}$ 
     $\mathbf{x}_i^{(k)} \leftarrow (1 - \omega) \boldsymbol{\alpha} + \omega \boldsymbol{\beta}$ 
  fin
fin

```

alg. 1.7 - Méthode SOR

Le principe, expliqué par l'algorithme 1.7, consiste à modifier (par relaxation) l'itéré \mathbf{x}^k fourni par l'algorithme de Gauss-Seidel à chaque itération. Pour cela, il faut conserver l'itéré précédent \mathbf{x}^{k-1} et modifier l'itéré courant comme suit :

$$\mathbf{x}^k \leftarrow \omega \mathbf{x}^k + (1 - \omega) \mathbf{x}^{k-1} \quad (1.48)$$

où ω est appelé le paramètre de relaxation. Si $\omega < 1$, on parle de sous-relaxation, et si $\omega > 1$, de sur-relaxation. On peut montrer que la matrice d'itération correspondante est :

$$\mathbf{B} = (\mathbf{I} - \omega \mathbf{D}^{-1} \mathbf{E})^{-1} [(1 - \omega) \mathbf{I} + \omega \mathbf{D}^{-1} \mathbf{F}] \quad (1.49)$$

Concernant la convergence, le théorème de Kahan donne une condition nécessaire : $0 < \omega < 2$. Si $\omega = 0$, l'algorithme ne progresse plus. Le théorème d'Ostrowski-Reich permet de conclure à la convergence $\forall \omega \in]0, 2[$ quand \mathbf{A} est symétrique, définie positive.

Le choix de ω est fait à partir d'heuristiques, comme $\omega = 2 - \mathcal{O}(h)$ pour une discrétisation éléments finis dont la taille de maille est h . Il existe aussi des techniques d'optimisation de ω en cours de route, comme la technique d'accélération d'Aitken, voir annexe C.

On peut aussi s'inspirer de résultats dans des cas particuliers. Par exemple dans un cas tridiagonal bloc, on montre que la valeur optimale de ω est :

$$\omega_{\text{opt}} = \frac{2}{1 + \sqrt{1 - \rho_j^2}} \quad (1.50)$$

La valeur de ρ_j est ainsi à estimer sur le problème en question. On a alors $\rho_{\text{SOR}} = \omega_{\text{opt}} - 1$, et, dans le cas éléments finis, $\tau \approx 2h$ (un ordre est gagné par rapport à la méthode de Jacobi).

La technique de relaxation peut être appliquée à d'autres algorithmes. Par exemple, un algorithme de Jacobi relaxé s'écrit sous la forme donnée dans l'algorithme 1.8.

```

Vecteur de départ  $\mathbf{x}^{(0)}$ 
tant que  $R > \varepsilon$ 
  itération sur  $k$ 
    boucle sur  $i = 1, 2, \dots, n$ 
       $\tilde{\mathbf{x}}_i \leftarrow \mathbf{A}_{i,1:i-1} \mathbf{x}_{1:i-1}^{(k-1)} + \mathbf{A}_{i,i+1:n} \mathbf{x}_{i+1:n}^{(k-1)}$ 
       $\tilde{\mathbf{x}}_i \leftarrow (\mathbf{b}_i - \tilde{\mathbf{x}}_i) / A_{i,i}$ 
    fin
     $\mathbf{x}^{(k)} \leftarrow \mathbf{x}^{(k-1)} + \omega(\tilde{\mathbf{x}} - \mathbf{x}^{(k-1)})$ 
  fin

```

alg. 1.8 - Méthode de Jacobi relaxée

1.5.4 Méthode d'Uzawa

Dans certains cas de figure comme les problèmes de point selle ou les liaisons traitées par multiplicateur de Lagrange par exemple, on peut être amené à résoudre des problèmes de la forme :

$$\begin{bmatrix} \mathbf{K} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{0} \end{bmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{l} \end{pmatrix} = \begin{pmatrix} \mathbf{f} \\ \mathbf{g} \end{pmatrix} \quad (1.51)$$

Pour que le problème continue d'être inversible, il suffit que \mathbf{C} soit injective.

La méthode d'Uzawa développée en 1958 consiste à résoudre itérativement de façon découplée, comme le montre l'algorithme 1.9 où α est un paramètre de la méthode. Pour analyser cet algorithme, on peut condenser la variable $\mathbf{x}^{(k)}$ sur la variable $\mathbf{l}^{(k)}$ pour trouver :

$$\mathbf{l}^{(k)} = \underbrace{(\mathbf{I} - \alpha \mathbf{C}^T \mathbf{K}^{-1} \mathbf{C})}_{\mathbf{B}} \mathbf{l}^{(k-1)} + \alpha (\mathbf{C}^T \mathbf{K}^{-1} \mathbf{f} - \mathbf{g}) \quad (1.52)$$

Si on appelle $\lambda_1, \dots, \lambda_p$, les valeurs propres de la matrice $\mathbf{C}^T \mathbf{K}^{-1} \mathbf{C}$, on en déduit l'équivalence suivante :

$$\varrho(\mathbf{B}) < 1 \Leftrightarrow \forall i, \quad |1 - \alpha \lambda_i| < 1 \quad (1.53)$$

```

Vecteurs de départ  $\mathbf{x}^{(0)}, \mathbf{l}^{(0)}$ 
tant que  $R > \varepsilon$ , itérer sur  $k$ 
  résoudre  $\mathbf{K} \mathbf{x}^{(k)} = \mathbf{b} - \mathbf{C} \mathbf{l}^{(k-1)}$ 
  réactualiser  $\mathbf{l}^{(k)} \leftarrow \mathbf{l}^{(k-1)} + \alpha (\mathbf{C}^T \mathbf{x}^{(k)} - \mathbf{c})$ 
fin

```

alg. 1.9 - Méthode d'Uzawa

Avec $\mathbf{C}^T \mathbf{K}^{-1} \mathbf{C}$ symétrique, définie positive, $\lambda_i > 0$ donc il faut prendre :

$$0 < \alpha < \frac{2}{\max_i \lambda_i} \quad (1.54)$$

1.6 Méthodes itératives non-stationnaires

Ces méthodes se décrivent de la même façon que les précédentes, à ceci près que les opérateurs \mathbf{B} et \mathbf{c} dépendent des itérations.

1.6.1 Méthode générique de projection

Reformulons le problème de départ en définissant le résidu :

$$\mathbf{r} = \mathbf{b} - \mathbf{Ax} \quad (1.55)$$

Il s'agit maintenant d'avoir $\mathbf{r} = \mathbf{0}$. L'idée est ici de se trouver avec un problème de plus petite taille, de structure plus simple (par exemple, tridiagonal) par plusieurs étapes de *projection*.

Considérons un sous-espace \mathcal{U} de l'espace dans lequel la solution est cherchée. Le problème approché est défini par $\mathbf{x} = \mathbf{Q}\hat{\mathbf{x}} \in \mathcal{U}$ avec $\mathbf{b} - \mathbf{Ax} \perp \mathcal{U}$, comme décrit sur la figure 1.10. Il est clair que si \mathcal{U} est l'espace complet, ce problème est le problème de départ. Dans tous les cas, il conduit à :

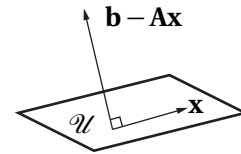


fig. 1.10 - Réduction dans un sous-espace

$$\mathbf{Q}^T(\mathbf{b} - \mathbf{AQ}\hat{\mathbf{x}}) = \mathbf{0} \Rightarrow \underbrace{(\mathbf{Q}^T \mathbf{A} \mathbf{Q})}_{\hat{\mathbf{A}}} \hat{\mathbf{x}} = \underbrace{\mathbf{Q}^T \mathbf{b}}_{\hat{\mathbf{b}}} \quad (1.56)$$

Une méthode de projection peut donc se schématiser de la façon suivante :

1. $\mathbf{r} \leftarrow \mathbf{b} - \mathbf{Ax}$
2. $\hat{\mathbf{x}} = \hat{\mathbf{A}}^{-1} \mathbf{Q}^T \mathbf{r}$
3. $\mathbf{x} \leftarrow \mathbf{x} + \mathbf{Q}\hat{\mathbf{x}}$
4. itération suivante

Le choix de \mathbf{Q} à chaque itération dépend bien entendu de la méthode employée.

1.6.2 Méthode du gradient

Pour \mathbf{A} symétrique, définie positive, la solution approchée est modifiée de la manière suivante :

$$\mathbf{x}^{(k+1)} \leftarrow \mathbf{x}^{(k)} + \mu^{(k)} \mathbf{d}^{(k)} \quad (1.57)$$

où $\mathbf{d}^{(k)}$ est la direction de recherche dans la direction du gradient :

$$\mathbf{d}^{(k)} = \mathbf{Ax}^{(k)} - \mathbf{b} \quad (1.58)$$

$\mu^{(k)}$ est le paramètre d'optimalité de la nouvelle solution dans la direction $\mathbf{d}^{(k)}$. Il se calcule comme :

$$\mu^{(k)} = \frac{\mathbf{d}^{(k)\top}(\mathbf{b} - \mathbf{A}\mathbf{x}^{(k)})}{\mathbf{d}^{(k)\top}\mathbf{A}\mathbf{d}^{(k)}} = -\frac{\mathbf{d}^{(k)\top}\mathbf{r}^{(k)}}{\mathbf{d}^{(k)\top}\mathbf{A}\mathbf{d}^{(k)}} \quad (1.59)$$

Il s'agit d'une méthode de projection pour laquelle la taille du sous-espace est 1 car on choisit pour \mathbf{Q} le sous-espace généré par le seul vecteur résidu $\mathbf{Q} = \mathbf{r}$.

La convergence de cette méthode est assez lente, d'où l'intervention d'une technique de conjugaison pour la méthode suivante.

1.6.3 Méthode du gradient conjugué

Cette méthode est généralement attribuée à Hestenes et Stiefel [22] en 1952. Elle est utilisée dans le cas où \mathbf{A} est symétrique, définie positive.

Description de la méthode

L'initialisation consiste à choisir un vecteur de départ $\mathbf{x}^{(0)}$, à calculer le résidu associé $\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}$ et de prendre comme première direction de recherche $\mathbf{d}^{(0)} = \mathbf{r}^{(0)}$.

Les itérations sont ensuite de la forme :

1. itéré $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mu^{(k)}\mathbf{d}^{(k)}$ avec le pas optimal

$$\mu^{(k)} = \frac{\mathbf{r}^{(k)\top}\mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top}\mathbf{A}\mathbf{d}^{(k)}} = \frac{\mathbf{r}^{(k)\top}\mathbf{r}^{(k)}}{\mathbf{d}^{(k)\top}\mathbf{A}\mathbf{d}^{(k)}} \quad (1.60)$$

2. résidu $\mathbf{r}^{(k+1)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k+1)}$

3. direction de recherche $\mathbf{d}^{(k+1)} = \mathbf{r}^{(k+1)} + \lambda^{(k)}\mathbf{d}^{(k)}$ qu'on orthogonalise par rapport à la précédente :

$$\mathbf{d}^{(k+1)\top}\mathbf{d}^{(k)} = 0 \Rightarrow \lambda^{(k)} = -\frac{\mathbf{r}^{(k+1)\top}\mathbf{A}\mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top}\mathbf{A}\mathbf{d}^{(k)}} = \frac{\mathbf{r}^{(k+1)\top}\mathbf{r}^{(k+1)}}{\mathbf{r}^{(k)\top}\mathbf{r}^{(k)}} \quad (1.61)$$

4. itération suivante

Propriétés

Si $i < j$, alors $\mathbf{r}^{(i)\top}\mathbf{r}^{(j)} = 0 = \mathbf{d}^{(i)\top}\mathbf{r}^{(j)} = \mathbf{d}^{(i)\top}\mathbf{A}\mathbf{d}^{(j)} = \mathbf{r}^{(i)\top}\mathbf{A}\mathbf{d}^{(j)}$. On construit donc des vecteurs orthogonaux au sens de \mathbf{A} , $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots$ qui forment donc une base. On minimise à chaque fois sur un sous-espace de taille de plus en plus grande.

En théorie, cette méthode est donc une méthode directe en n étapes. En pratique, on a progressivement perte d'orthogonalité due aux arrondis, c'est donc en ce sens une méthode itérative.

On a alors la propriété de convergence suivante :

$$\|\mathbf{x}^{(k)} - \mathbf{x}\|_{\mathbf{A}} \leq 2 \left[\frac{\sqrt{\text{cond}\mathbf{A}} - 1}{\sqrt{\text{cond}\mathbf{A}} + 1} \right]^{(k)} \|\mathbf{x}^{(0)} - \mathbf{x}\|_{\mathbf{A}} \quad (1.62)$$

d'où, en éléments finis :

$$\tau = \mathcal{O}\left(\frac{1}{\sqrt{\text{cond}\mathbf{A}}}\right) = \mathcal{O}(h) \quad (1.63)$$

1.6.4 Préconditionnement

```

Vecteur de départ  $\mathbf{x}^{(0)}$ 
Résidu d'équilibre  $\mathbf{r}^{(0)} \leftarrow \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}$ 
tant que  $R > \varepsilon$ , itérer sur  $k$ 
  résoudre  $\mathbf{M}\mathbf{z}^{(k-1)} = \mathbf{r}^{(k-1)}$ 
   $\rho_{k-1} \leftarrow \mathbf{r}^{(k-1)\top} \mathbf{z}^{(k-1)}$ 
  si  $k = 1$  alors
     $\mathbf{d}^{(1)} \leftarrow \mathbf{z}^{(0)}$ 
  sinon
     $\beta_{k-1} \leftarrow \frac{\rho_{k-1}}{\rho_{k-2}}$ 
     $\mathbf{d}^{(k)} \leftarrow \mathbf{z}^{(k-1)} + \beta_{k-1} \mathbf{d}^{(k-1)}$ 
  fin
   $\mathbf{q}^{(k)} \leftarrow \mathbf{A}\mathbf{d}^{(k)}$ 
   $\alpha_k \leftarrow \frac{\rho_{k-1}}{\mathbf{d}^{(k)\top} \mathbf{q}^{(k)}}$ 
   $\mathbf{x}^{(k)} \leftarrow \mathbf{x}^{(k-1)} + \alpha_k \mathbf{d}^{(k)}$ 
   $\mathbf{r}^{(k)} \leftarrow \mathbf{b}^{(k)} - \mathbf{A}\mathbf{x}^{(k)} = \mathbf{r}^{(k-1)} - \alpha_k \mathbf{q}^{(k)}$ 
fin

```

alg. 1.11 - Gradient conjugué
préconditionné

le problème de départ !

La convergence est cette fois-ci en :

$$\tau = \mathcal{O} \left(\frac{1}{\sqrt{\text{cond} \mathbf{M}^{-1} \mathbf{A}}} \right) \quad (1.66)$$

Comme exemples de préconditionneurs classiques, on a la factorisation de Cholesky incomplète, SSOR, etc. L'algorithme 1.11 présente cette méthode. Le choix et le développement de préconditionneurs adaptés font encore aujourd'hui l'objet de recherches.

1.6.5 Autres méthodes itératives

Pour information, on peut citer :

- les méthodes dites « rapides » dédiées à des applications particulières : Fourier ($\mathcal{O}(n \log n)$), multigrilles ($\mathcal{O}(n)$);
- pour \mathbf{A} symétrique et définie : Chebyshev [23];
- pour \mathbf{A} symétrique mais non définie : MinRES (résidu minimum, *minimal residual*), SymmLQ (*symmetric LQ*) [24];
- pour \mathbf{A} non symétrique connaissant \mathbf{A}^T : QMR (*quasi minimal residual*) [25];
- pour \mathbf{A} non symétrique : CGS (*conjugate gradient square*) [26], BiCGStab (*bi-conjugate gradient stabilized*) [27], GMRES (*generalized minimal residual*) [28], GMRES restarted...

1.7 Critères d'arrêt

Le critère d'arrêt recherché est celui de l'écart en solution tel que $\mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}$ soit suffisamment petit. Malheureusement, \mathbf{x} est la solution à convergence et on ne peut atteindre

Si on peut avoir facilement une approximation spectrale \mathbf{M} de \mathbf{A} , l'idée est de chercher à résoudre :

$$\mathbf{M}^{-1} \mathbf{A} \mathbf{x} = \mathbf{M}^{-1} \mathbf{b} \quad (1.64)$$

\mathbf{M} symétrique, définie positive. Ce problème devient non symétrique, mais on peut lui substituer :

$$(\mathbf{M}^{-1/2} \mathbf{A} \mathbf{M}^{-1/2}) \mathbf{M}^{1/2} \mathbf{x} = \mathbf{M}^{-1/2} \mathbf{b} \quad (1.65)$$

Le choix du préconditionneur dépend du problème à traiter. Dans le meilleur des cas, il est facile à inverser et représente bien la matrice \mathbf{A} . En fait, si $\mathbf{M} = \mathbf{I}$ il est facile à inverser, mais est très différent de \mathbf{A} , et si $\mathbf{M} = \mathbf{A}$ il y a convergence en une itération, mais dans laquelle il faut résoudre

l'erreur en solution au cours des itérations. La quantité que l'on peut espérer atteindre est liée au résidu :

$$\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{Ax}^{(k)} = -\mathbf{Ae}^{(k)} \quad (1.67)$$

Un critère souvent employé est le suivant :

$$\frac{\|\mathbf{r}^{(k)}\|}{\|\mathbf{Ax}^{(k)}\| + \|\mathbf{b}\|} < \varepsilon \quad (1.68)$$

où ε est un seuil fixé *a priori*, que l'on peut quantifier par rapport à la précision souhaitée et la précision machine. Parfois, on remplace au dénominateur $\|\mathbf{Ax}^{(k)}\|$ par $\|\mathbf{b}\|$. Le problème vient du fait, comme on a déjà pu le voir dans un exemple, que l'erreur en solution et le résidu sont liés de la façon suivante :

$$\|\mathbf{e}^{(k)}\| \leq \|\mathbf{A}^{-1}\| \quad (1.69)$$

donc :

$$\frac{\mathbf{r}^{(k)}}{\mathbf{b}} < \varepsilon \Rightarrow \frac{\mathbf{e}^{(k)}}{x} < \varepsilon \text{ cond}\mathbf{A} \quad (1.70)$$

et, en tout état de cause, le paramètre ε devrait être choisi aussi fonction du conditionnement de \mathbf{A} , qu'il faut donc estimer pour chaque problème traité.

Des critères d'arrêt en cas de non-convergence sont aussi mis en place. Par exemple, $\|\mathbf{r}^{(k)}\|$ décroît *trop lentement* ou ne décroît plus (mais il faut faire attention aux plateaux possibles avec certaines méthodes), ou la résolution est trop longue (critère en nombre maxi d'itérations). Il s'agit alors de sortie en cas d'erreur.

2

SYSTÈMES AUX VALEURS PROPRES

On considère ici les problèmes, dits *aux valeurs propres*, de la forme : trouver les caractéristiques propres, c'est-à-dire les couples (λ, \mathbf{x}) où $\mathbf{x} \neq \mathbf{0}$ tels que $\mathbf{Ax} = \lambda \mathbf{x}$ et l'on se restreindra aux cas où \mathbf{A} est symétrique.

Un problème aux valeurs propres généralisé est de la forme $\mathbf{Kx} = \omega^2 \mathbf{Mx}$, avec $\mathbf{x} \neq \mathbf{0}$, \mathbf{K} symétrique, définie positive et \mathbf{M} symétrique positive. Dans ce cas, on peut utiliser la factorisation de Cholesky de \mathbf{M} , soit $\mathbf{M} = \mathbf{LL}^T$ et avec $\mathbf{y} = \mathbf{L}^T \mathbf{x}$, on est ramené au problème précédent $\mathbf{L}^{-1} \mathbf{KL}^{-T} \mathbf{y} = \omega^2 \mathbf{y}$.

En général, dans les problèmes qui nous concernent, on ne veut qu'environ 30 valeurs et/ou vecteurs propres, de préférence à basse fréquence (dynamique lente, vibrations [29, 30, 31], flambage...). Dans le cas où il y a de l'amortissement, le système à résoudre est :

$$\mathbf{M}\ddot{\mathbf{x}} + \mathbf{C}\dot{\mathbf{x}} + \mathbf{K}\mathbf{x} = \mathbf{f}(t) \quad (2.1)$$

La matrice d'amortissement, \mathbf{C} , est souvent mal maîtrisée et si l'amortissement est peu influent à basse fréquence, ce qui n'est pas évident pour tous les cas, on est ramené à :

$$\mathbf{M}\ddot{\mathbf{x}} + \mathbf{K}\mathbf{x} = \mathbf{f}(t) \quad (2.2)$$

Avec une discrétisation éléments finis standard, les hautes fréquences obtenues numériquement sont éloignées des vraies.

Enfin, on prévoit parfois d'avoir des problèmes si les valeurs propres sont mal séparées, proches ou même confondues en cas de symétries, par exemple. Une référence de base concernant ces problèmes est [32, 33].

2.1 Préliminaires

2.1.1 Rappels

Matrice quelconque

λ est valeur propre de \mathbf{A} ssi $\det(\mathbf{A} - \lambda \mathbf{I}) = 0$. Dans \mathbb{C} , \mathbf{A} admet n valeurs propres distinctes ou non, réelles ou complexes. \mathbf{A} a n valeurs propres distinctes ssi \mathbf{A} est diagonalisable.

Matrice réelle

\mathbf{A} est symétrique ssi \mathbf{A} est diagonalisable dans \mathbb{R} . \mathbf{A} est symétrique, définie positive ssi :

- toutes les valeurs propres sont strictement positives ;
- $\exists \mathbf{B}$ symétrique, définie positive telle que $\mathbf{A} = \mathbf{B}^2$.

Dans \mathbb{C} , \mathbf{A} admet n valeurs propres distinctes ou non, réelles ou complexes. \mathbf{A} a n valeurs propres distinctes implique que \mathbf{A} est diagonalisable.

Quotient de Rayleigh pour une matrice réelle, symétrique

Le quotient de Rayleigh est défini comme étant, pour $\mathbf{x} \neq \mathbf{0}$:

$$R(\mathbf{x}) = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \quad \text{ou} \quad R(\mathbf{x}) = \frac{\mathbf{x}^T \mathbf{K} \mathbf{x}}{\mathbf{x}^T \mathbf{M} \mathbf{x}} \quad (2.3)$$

Il vient alors l'équivalence suivante :

$$(\lambda, \mathbf{x}) \text{ est caractéristique propre de } \mathbf{A} \Leftrightarrow \begin{cases} R \text{ est extremum en } \mathbf{x} \neq \mathbf{0} \\ R(\mathbf{x}) = \lambda \end{cases} \quad (2.4)$$

La i^{e} valeur propre λ_i est le minimum de $R(\mathbf{x})$, $\mathbf{x} \neq \mathbf{0}$ sur le sous espace orthogonal aux précédents vecteurs propres. De plus, en *dimension finie*, la valeur propre maxi est $\lambda_n = \max_{\mathbf{x} \neq \mathbf{0}} R(\mathbf{x})$ dans l'espace complet de dimension n .

2.1.2 Domaine d'emploi des différentes méthodes

En comparaison de la résolution d'un système linéaire, la résolution d'un système aux valeurs propres est beaucoup plus coûteuse. En fait, il s'agit d'un problème non-linéaire, et beaucoup de méthodes s'appuient sur la résolution d'une succession de systèmes linéaires, en interne. Elles souffriront alors des problèmes affectant ces solveurs, en particulier le mauvais conditionnement. L'ordre de grandeur de la taille des systèmes aux valeurs propres que l'on va considérer dans cette partie est de l'ordre de $n = 250$ à $n = 25000$.

La première méthode est celle de l'équation caractéristique. On peut en effet rechercher les zéros du polynôme caractéristique de la matrice, mais si les coefficients de ce polynôme peuvent être assez facilement atteints, une petite perturbation sur ceux-ci conduit à d'énormes variations sur ses racines. Cette méthode est donc réservée à des tailles de système réduites, jusqu'à $n \approx 10$. On peut montrer mathématiquement que, dès que $n \geq 5$, la méthode employée ne peut être qu'itérative : le problème ne peut plus être résolu numériquement en un nombre fini d'opérations.

Les méthodes de Jacobi, des puissances et de Householder vont ensuite pouvoir être employées pour des systèmes dont la taille va jusqu'à $n \approx 250$.

Au delà, elle doivent être adaptées à un stockage bande, et le remplissage des matrices au cours des itérations doit être étudié. Avec les méthodes de réduction, ces techniques permettent d'atteindre des systèmes dont la taille va jusqu'à $n \approx 2500$.

Jusqu'à $n \approx 25000$, les méthodes des puissances, plus ou moins couplées aux méthodes de sous-espaces, la méthode de Lanczos avec stratégie de réorthogonalisation permettent encore de traiter le problème.

Pour dépasser ces tailles de problème, il est ensuite nécessaire de faire appel à des techniques particulières. Par exemple, les techniques de parallélisation, vectorisation, optimisation des entrées-sorties, les méthodes de condensation et sous-structuration (Guyan-Irons, Craig-Bampton...) Ces dernières méthodes ne seront pas étudiées ici.

2.2 Méthode de Jacobi

Cette méthode, datant de 1846 [34], est décrite ici, dans le cas particulier où \mathbf{A} est symétrique, à des fins de rappel uniquement puisqu'elle est peu efficace pour les problèmes qui nous concernent.

```

tant que  $R > \varepsilon$ , itérer sur  $k$ 
  choix du terme à annuler  $A_{p,q}$ 
   $v \leftarrow \frac{A_{q,q}^{(k-1)} - A_{p,p}^{(k-1)}}{2A_{p,q}^{(k-1)}}$ 
  calcul de la tangente  $t$ 
  calcul du cosinus  $c \leftarrow 1/\sqrt{1+t^2}$ 
  calcul du sinus  $s \leftarrow ct$ 
   $\mathbf{A}^{(k)} \leftarrow \mathbf{A}^{(k-1)}$ 
   $A_{p,1:n}^{(k)} \leftarrow cA_{p,1:n}^{(k-1)} - sA_{q,1:n}^{(k-1)}$ 
   $A_{q,1:n}^{(k)} \leftarrow sA_{p,1:n}^{(k-1)} + cA_{q,1:n}^{(k-1)}$ 
   $A_{p,p}^{(k)} \leftarrow A_{p,p}^{(k-1)} - tA_{p,q}^{(k-1)}$ 
   $A_{q,q}^{(k)} \leftarrow A_{q,q}^{(k-1)} + tA_{p,q}^{(k-1)}$ 
   $A_{p,q}^{(k)} \leftarrow 0$ 
   $A_{q,p}^{(k)} \leftarrow 0$ 
fin

```

alg. 2.1 - Méthode de Jacobi

Dans cette méthode, on cherche à ramener la matrice \mathbf{A} à une forme diagonale pour trouver tous ses valeurs propres et vecteurs propres, et ce par une séquence de transformations orthogonales :

$$\mathbf{A}^{(k+1)} = \mathbf{R}^{(k)\text{T}} \mathbf{A}^{(k)} \mathbf{R}^{(k)} \quad (2.5)$$

avec les rotations $\mathbf{R}^{(k)}$ vérifiant $\mathbf{R}^{(k)\text{T}} \mathbf{R}^{(k)} = \mathbf{I}$, transformations qui sont des changements de base de type isométries ne modifiant pas les valeurs propres.

On commence par chercher l'élément hors diagonale de module maximal :

$$|A_{p,q}^{(k)}| = \max_{i \neq j} |A_{i,j}^{(k)}| \quad (2.6)$$

On construit alors la matrice de rotation d'un angle θ *ad hoc* dans le plan $(\mathbf{e}_p, \mathbf{e}_q)$ de façon à annuler dans l'itéré suivant les termes $A_{q,p}^{(k+1)}$ et $A_{p,q}^{(k+1)}$:

$$\mathbf{R}^{(k)} = 1 + (\cos \theta - 1)(\mathbf{e}_p \mathbf{e}_p^{\text{T}} + \mathbf{e}_q \mathbf{e}_q^{\text{T}}) + (\sin \theta - 1)\mathbf{e}_p \mathbf{e}_p^{\text{T}} - (\sin \theta + 1)\mathbf{e}_q \mathbf{e}_q^{\text{T}} \quad (2.7)$$

avec :

$$A_{p,q}^{(k+1)} = \frac{1}{2} \sin 2\theta [A_{p,p}^{(k)} - A_{q,q}^{(k)}] + \frac{1}{2} \cos 2\theta A_{p,q}^{(k)} = 0 \quad (2.8)$$

ce qui entraîne :

$$\cotan 2\theta = \frac{A_{p,p}^{(k)} - A_{q,q}^{(k)}}{A_{p,q}^{(k)}} \quad (2.9)$$

Seuls changent les termes sur les lignes ou colonnes p ou q . Dans l'algorithme 2.1, la tangente t est calculée comme suit :

$$t = \begin{cases} 1 & \text{si } v = 0 \\ -v + \text{sgn}(v)\sqrt{1+v^2} & \text{sinon} \end{cases} \quad (2.10)$$

Au cours des itérations, un terme nul peut redevenir non nul mais :

$$a^{(k)} = \sum_{i \neq j} A_{i,j}^{(k)2} \xrightarrow{k \rightarrow \infty} = 0 \quad (2.11)$$

En effet :

$$\begin{aligned} a^{(k+1)} &= \sum_{\substack{i \neq j, p \\ q, j \neq p, q}} \underbrace{[A_{i,j}^{(k+1)}]^2}_{[A_{i,j}^{(k)}]^2} + \sum_{i \neq p, q} \underbrace{\left([A_{i,p}^{(k+1)}]^2 + [A_{i,q}^{(k+1)}]^2 \right)}_{[A_{i,p}^{(k)}]^2 + [A_{i,q}^{(k)}]^2} + \sum_{j \neq p, q} \underbrace{\left([A_{p,j}^{(k+1)}]^2 + [A_{q,j}^{(k+1)}]^2 \right)}_{[A_{p,j}^{(k)}]^2 + [A_{q,j}^{(k)}]^2} \\ &= a^{(k)} - 2 [A_{p,q}^{(k)}]^2 \end{aligned} \quad (2.12)$$

donc $a^{(k+1)} < a^{(k)}$.

Cet algorithme est stable et simple, mais la convergence est lente et il s'avère coûteux en accès mémoire si n est grand. On a les estimations suivantes :

$$a^{(k+1)} \leq \left[1 - \frac{2}{n(n-1)} \right] a^{(k)} \quad \text{et} \quad \tau \approx \frac{2}{n^2} \quad (2.13)$$

2.3 Méthode des puissances

On se place ici dans le cas où \mathbf{A} est symétrique, définie positive et dont les valeurs propres sont toutes simples, rangées par ordre de module décroissant λ_i .

$x^{(0)}, \|x^{(0)}\| = 1$
tant que $R > \varepsilon$, itérer sur k
 $v \leftarrow \mathbf{K}x^{(k-1)}$
 résoudre $\mathbf{M}x^{(k)} = v$
 $\alpha \leftarrow \|x^{(k)}\|$
 $\lambda^{(k)} \leftarrow R(x^{(k)}) = \frac{x^{(k)T} \mathbf{K}x^{(k)}}{x^{(k)T} \mathbf{M}x^{(k)}}$
 normalisation $x^{(k)} \leftarrow \frac{x^{(k)}}{\alpha}$
fin

alg. 2.2 - Méthode des puissances

On construit un itéré de départ $\mathbf{x}^{(0)}$ puis les itérés successifs par récurrence de la manière suivante :

$$\mathbf{x}^{(k+1)} = \mathbf{A}\mathbf{x}^{(k)} \quad (2.14)$$

Si $(\mathbf{v}_i)_i$ est la base propre, on peut écrire les itérés dans cette base. En particulier :

$$\mathbf{x}^{(0)} = \sum_i \alpha_i \mathbf{v}_i \quad (2.15)$$

alors :

$$\mathbf{x}^{(k)} = \sum_i \lambda_i^k \alpha_i \mathbf{v}_i = \lambda_1^k \left[\alpha_1 \mathbf{v}_1 + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k \mathbf{v}_i \right] \quad (2.16)$$

et si $\alpha_1 \neq 0$ puisque $\frac{|\lambda_i|}{|\lambda_1|} < 1$ alors :

$$\mathbf{x}^{(k)} \xrightarrow{k \rightarrow \infty} \lambda_1^k \alpha_1 \mathbf{v}_1 \quad (2.17)$$

En pratique, il est utile de normer $\mathbf{x}^{(k)}$ à chaque itération pour éviter l'explosion et dans ce cas obtenir $\mathbf{x}^{(k)} \rightarrow \mathbf{v}_1$.

Pour comprendre le comportement de la convergence, prenons $\mathbf{x}^{(0)} = \sum_i \mathbf{v}_i$, alors :

$$\frac{\mathbf{x}^{(k)}}{\lambda_1^k} = \mathbf{v}_1 + \sum_{i=2}^n \left(\frac{\lambda_i}{\lambda_1} \right)^k \mathbf{v}_i \quad (2.18)$$

donc :

$$\frac{\left\| \frac{\mathbf{x}^{(k+1)}}{\lambda_1^{k+1}} - \mathbf{v}_1 \right\|}{\left\| \frac{\mathbf{x}^{(k)}}{\lambda_1^k} - \mathbf{v}_1 \right\|} = \frac{\left[\sum_{i=2}^n \left(\frac{\lambda_i}{\lambda_1} \right)^{2(k+1)} \right]^{1/2}}{\left[\sum_{i=2}^n \left(\frac{\lambda_i}{\lambda_1} \right)^{2k} \right]} \xrightarrow{k \rightarrow \infty} \frac{\lambda_2}{\lambda_1} \quad (2.19)$$

Le taux de convergence est alors estimé comme suit :

$$\tau \approx -\ln \frac{\lambda_2}{\lambda_1} = \ln \frac{\lambda_1}{\lambda_2} \quad (2.20)$$

Bien entendu, on peut avoir des ennuis si au départ $\mathbf{x}^{(0)}$ est orthogonal à \mathbf{v}_1 . Par rapport à la méthode de Jacobi, la convergence est maintenant indépendante de la taille du problème, n .

Cas de valeurs propres multiples ou proches

Dans le cas de valeurs propres multiples λ_1 , la convergence se fait vers λ_1 et vers une combinaison linéaire des vecteurs propres associés.

Dans le cas de valeurs propres proches, $\lambda_2 = \lambda_1 + \delta\lambda$, on a $\tau \approx \frac{\delta\lambda}{\lambda_1}$ et la convergence est donc très lente.

Contrôle de la convergence

Comme critère de convergence sur les valeurs propres, on peut utiliser le fait que $\|\mathbf{x}^{(k+1)}\|/\|\mathbf{x}^{(k)}\|$ se stabilise, et pour la convergence sur les vecteurs propres, $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < \varepsilon$ ou :

$$\frac{\|\mathbf{A}\mathbf{x}^{(k+1)} - \lambda^{(k+1)}\mathbf{x}^{(k+1)}\|}{\|\mathbf{A}\mathbf{x}^{(k+1)} + \lambda^{(k+1)}\mathbf{x}^{(k+1)}\|} < \varepsilon' \quad (2.21)$$

On conseille aussi d'utiliser la norme infinie pour ces tests.

Obtention des modes suivants

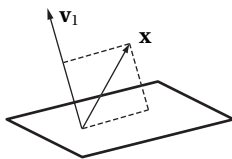


fig. 2.3 - Projection orthogonale à \mathbf{v}_1

La technique utilisée est la déflation orthogonale : on prend $\mathbf{x}^{(0)}$ orthogonal à \mathbf{v}_1 . Pour éviter les pertes d'orthogonalité par cumul d'arrondis, on réorthogonalise à chaque itération. Pour cela, on utilise un projecteur comme décrit sur la figure 2.3 :

$$\mathbf{P}_1 = \mathbf{I} - \frac{\mathbf{v}_1 \mathbf{v}_1^T}{\mathbf{v}_1^T \mathbf{v}_1} \quad (2.22)$$

d'où :

$$\mathbf{P}_1 \mathbf{x} = \mathbf{x} - \frac{\mathbf{v}_1^T \mathbf{x}}{\mathbf{v}_1^T \mathbf{v}_1} \mathbf{v}_1 \quad (2.23)$$

On a alors convergence vers $(\lambda_2, \mathbf{v}_2)$. On peut procéder de façon similaire pour accéder à la troisième caractéristique propre, et ainsi de suite.

2.4 Méthode des itérations inverses

```

 $\mathbf{x}^{(0)}, \|\mathbf{x}^{(0)}\| = 1$ 
tant que  $R > \varepsilon$ , itérer sur  $k$ 
   $\mathbf{v} \leftarrow \mathbf{M}\mathbf{x}^{(k-1)}$ 
  résoudre  $\mathbf{K}\mathbf{x}^{(k)} = \mathbf{v}$ 
   $\alpha \leftarrow \|\mathbf{x}^{(k)}\|$ 
   $\lambda^{(k)} \leftarrow R(\mathbf{x}^{(k)}) = \frac{\mathbf{x}^{(k)\top} \mathbf{K} \mathbf{x}^{(k)}}{\mathbf{x}^{(k)\top} \mathbf{M} \mathbf{x}^{(k)}}$ 
  normalisation  $\mathbf{x}^{(k)} \leftarrow \frac{\mathbf{x}^{(k)}}{\alpha}$ 
fin

```

alg. 2.4 - Itérations inverses

Comme décrit dans l'algorithme 2.4, on choisit encore un $\mathbf{x}^{(0)}$ arbitraire : les solutions classiques sont de prendre un vecteur de composantes aléatoires, normé, quitte à essayer avec plusieurs vecteurs de départ pour vérifier la convergence vers la même valeur.

On résoud ensuite $\mathbf{A}\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$ et on norme $\mathbf{x}^{(k+1)}$. Il s'agit bien de la méthode des puissances, appliquée à \mathbf{A}^{-1} , qui commence donc par converger vers la valeur propre de \mathbf{A} de module le plus élevé.

Bien évidemment, il est simple de démontrer que le taux de convergence est tel que :

$$\tau \approx \ln \frac{|\lambda_{n-1}|}{|\lambda_n|} \quad (2.24)$$

Comme pour la méthode des puissances, on converge d'autant mieux que les valeurs propres sont séparées... d'où la méthode du paragraphe suivant.

2.5 Méthode des itérations inverses avec décalage spectral

Appelée encore méthode de Wielandt, cette méthode suppose que l'on connaît une approximation μ de la valeur propre cherchée λ .

2.5.1 Principe et convergence

On utilise la matrice décalée $\mathbf{A}^* = \mathbf{A} - \mu\mathbf{I}$, on s'intéresse donc à $\mathbf{A}^*\mathbf{x} = \lambda^*\mathbf{x}$ et $\lambda^* = \lambda - \mu$. Comme on converge vers la valeur propre $|\lambda^*|$ la plus petite, on va donc trouver λ_i la plus proche de μ . Le taux de convergence est :

$$\tau = -\ln \left(\inf \left[\frac{|\lambda_i - \mu|}{|\lambda_{i-1} - \mu|}, \frac{|\lambda_i - \mu|}{|\lambda_{i+1} - \mu|} \right] \right) \quad (2.25)$$

et il est par conséquent très élevé si μ est très proche de λ_i . Voici quelques remarques :

- μ est proche de λ_i ; le système $\mathbf{A}^*\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$ est très mal conditionné. En fait, la méthode y est peu sensible : l'erreur en solution tend à s'aligner dans la direction du vecteur propre et est éliminée par la normalisation ;
- Cette méthode permet aussi d'accéder aux valeurs propres d'une région donnée du spectre ;
- On peut construire de cette manière la méthode des itérations inverses de Rayleigh : on change l'estimation μ par $\lambda^{(k)}$ à chaque itération. On obtient alors une convergence cubique :

$$\|\mathbf{x}^{(k+1)} - \mathbf{v}\| = \mathcal{O} \left(\|\mathbf{x}^{(k)} - \mathbf{v}\|^3 \right) \quad (2.26)$$

Cette méthode est coûteuse puisque l'on factorise à chaque itération.

2.5.2 Exemples d'utilisation dans un code éléments finis

On considère le système $\mathbf{Kx} = \omega^2 \mathbf{Mx}$.

Problème 1

On veut trouver le mode proche de la pulsation ω_0 . On procède alors de la façon suivante :

1. décalage : $\mathbf{K}^* \leftarrow \mathbf{K} - \omega_0^2 \mathbf{M}$
2. factorisation \mathbf{LDL}^T de \mathbf{K}^* ; soit N^* le nombre de termes strictement négatifs de \mathbf{D}
3. $\mathbf{x}^{(0)}$ aléatoire puis itérations inverses avec décalage de ω_0 sur $\mathbf{K}^* \mathbf{x} = \lambda^* \mathbf{Mx}$ pour obtenir (λ^*, \mathbf{x})
4. $\omega^2 \leftarrow \omega_0^2 + \lambda^*$
5. le numéro du mode trouvé est N^* si $\omega^2 < \omega_0^2$ ou $N^* + 1$ si $\omega^2 \geq \omega_0^2$.

Problème 2

On veut trouver tous les modes dans l'intervalle $[\omega_1, \omega_2]$. On procède alors de la façon suivante :

1. pour trouver le nombre de modes dans $[\omega_1, \omega_2]$, on compte le nombre de termes strictement négatifs de \mathbf{D} dans les factorisations de $\mathbf{K} - \omega_1^2 \mathbf{M}$ et $\mathbf{K} - \omega_2^2 \mathbf{M}$;
2. on procède par itérations inverses avec décalage et déflations orthogonales successives.

2.6 Transformation de Householder

Cette méthode date de 1958 [35].

2.6.1 Factorisation QR

On cherche à mettre \mathbf{A} sous la forme $\mathbf{A} = \mathbf{QR}$ où \mathbf{R} est triangulaire supérieure et \mathbf{Q} orthogonale, par applications successives de transformations de Householder à $\mathbf{A} \rightarrow \mathbf{HA}$ [36].

\mathbf{H} est de la forme $\mathbf{H} = \mathbf{I} - 2\mathbf{u}\mathbf{u}^T$ avec $\|\mathbf{u}\| = 1$. Il s'agit cette fois d'une symétrie par rapport à un plan de normale \mathbf{u} comme illustré sur la figure 2.5. En conséquence, \mathbf{H} est symétrique, orthogonale, et conserve les valeurs propres.

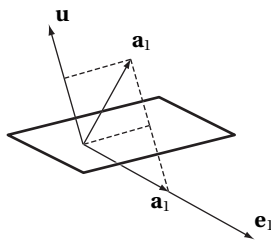


fig. 2.5 - Transformation de Householder de \mathbf{a}_1

Ces transformations de Householder sont appliquées sur les colonnes de $\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_n]$ successivement. Par exemple, sur la première, on veut $\mathbf{Ha}_1 = \alpha \mathbf{e}_1$ d'où le choix :

$$\alpha = \|\mathbf{a}_1\| = \|\mathbf{Ha}_1\| > 0 \quad (2.27)$$

alors :

$$\begin{aligned} \mathbf{a}_1 - 2\mathbf{u}(\mathbf{u}^T \mathbf{a}_1) &= \alpha \mathbf{e}_1 \\ \Rightarrow \mathbf{a}_1 - \alpha \mathbf{e}_1 &= 2\mathbf{u}(\mathbf{u}^T \mathbf{a}_1) \end{aligned} \quad (2.28)$$

En prenant le produit scalaire par \mathbf{a}_1 , on obtient :

$$\alpha^2 - \alpha \mathbf{a}_1^T \mathbf{e}_1 = 2(\mathbf{u}^T \mathbf{a}_1)^2 \quad (2.29)$$

d'où :

$$\mathbf{u}^T \mathbf{a}_1 = \sqrt{\frac{1}{2} \alpha (\alpha - \mathbf{a}_1^T \mathbf{e}_1)} \quad (2.30)$$

On peut noter que le terme $(\alpha - \mathbf{a}_1^T \mathbf{e}_1)$ est non nul si \mathbf{a}_1 n'est pas déjà parallèle à \mathbf{e}_1 . Finalement, on en déduit le vecteur cherché :

$$\mathbf{u} = \frac{\mathbf{a}_1 - \alpha \mathbf{e}_1}{2\mathbf{u}^T \mathbf{a}_1} \quad (2.31)$$

Le procédé se poursuit en considérant ensuite la deuxième colonne de \mathbf{A} et \mathbf{H} telle que :

$$\mathbf{H} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{H}_2 \end{bmatrix} \quad (2.32)$$

Le vecteur suivant \mathbf{u}_2 est déterminé à partir de \mathbf{e}_2 et \mathbf{a}_2 . C'est une méthode directe en $n - 1$ étapes.

On peut regarder les propriétés de remplissage de \mathbf{A} lors des transformations mais c'est assez délicat. En général, on ne stocke pas la matrice \mathbf{Q} ; l'application de \mathbf{Q} à un vecteur est coûteuse ($4n$ opérations), cette méthode n'est donc pas appliquée à la résolution de systèmes linéaires. Cet algorithme se trouve aussi sous le nom de Householder-Businger-Golub.

La méthode est présentée par l'algorithme 2.6. En ce qui concerne le stockage, \mathbf{v}_k écrase $\mathbf{A}_{1:k,k}$ en mémoire, $\mathbf{R}_{k+1,k+1:n}$ écrase $\mathbf{A}_{k+1,k+1:n}$. $\mathbf{R}_{k,k}$ doit être stocké ailleurs pour chaque itération, d'où un vecteur de dimension n supplémentaire. \mathbf{Q} n'est en général pas stockée, elle pourrait être reconstruite à partir des \mathbf{v}_k . Enfin, \mathbf{z} est un vecteur de stockage supplémentaire.

La complexité de cet algorithme pour une matrice pleine de taille n est de l'ordre de $\frac{2}{3}n^3$.

```

boucle sur  $k = 1, 2, \dots, n - 1$ 
 $\delta \leftarrow \mathbf{A}_{k:n,k}^{(k)T} \mathbf{A}_{k:n,k}^{(k)}$ 
 $\alpha \leftarrow \sqrt{\delta}$ 
 $\beta \leftarrow \frac{1}{\delta - \mathbf{A}_{k,k}^{(k)}}$ 
 $\mathbf{v}_{k:n}^{(k)} \leftarrow \mathbf{A}_{k:n,k}^{(k)}$ 
 $\mathbf{v}_k^{(k)} \leftarrow \mathbf{v}_k^{(k)} - \alpha$ 
 $\mathbf{z}_{k+1:n}^T \leftarrow \beta \mathbf{v}_{k:n}^{(k)T} \mathbf{A}_{k:n,k+1:n}^{(k)}$ 
 $\mathbf{A}_{k:n,k+1:n}^{(k+1)} \leftarrow \mathbf{A}_{k:n,k+1:n}^{(k)} - \mathbf{v}_{k:n}^{(k)} \mathbf{z}_{k+1:n}^T$ 
 $\mathbf{R}_{k,k} \leftarrow \alpha$ 
 $\mathbf{R}_{k,k+1:n} \leftarrow \mathbf{A}_{k,k+1:n}^{(k+1)}$ 
fin

```

alg. 2.6 - Transformation de Householder pour la factorisation QR

2.6.2 Détermination des valeurs propres par la méthode QR

Cette méthode est attribuée à Francis et Vera N. Kublanovskaya en 1961. On considère le problème $\mathbf{Ax} = \lambda \mathbf{x}$ avec des valeurs propres simples. On procède alors de la façon suivante :

1. $\mathbf{A}^{(1)} = \mathbf{A}$;
2. factorisation QR de $\mathbf{A}^{(k)} = \mathbf{Q}^{(k)} \mathbf{R}^{(k)}$;
3. calcul de $\mathbf{A}^{(k+1)} = \mathbf{R}^{(k)} \mathbf{Q}^{(k)}$;
4. itération suivante.

Propriétés

- si les valeurs propres sont toutes simples, $\mathbf{A}^{(k)}$ tend vers la matrice diagonale des valeurs propres. $\mathbf{Q}^{(k)}$ tend vers la matrice des vecteurs propres (propriété peu utilisée) ;
- il y a conservation des bandes de \mathbf{A} . C'est une méthode performante quand \mathbf{A} est tri-diagonale de petite taille... le coût d'une itération est alors d'environ $10n$. Quand \mathbf{A} est pleine, une factorisation coûte $\approx \frac{2}{3}n^3$ opérations et un produit \mathbf{QR} , $\approx \frac{1}{2}n^3$ opérations.

On peut aussi appliquer des translations (décalages) :

$$\begin{aligned}\mathbf{A}^{(k)} - \alpha^{(k)}\mathbf{I} &= \mathbf{Q}^{(k)}\mathbf{R}^{(k)} \\ \mathbf{A}^{(k+1)} &= \mathbf{R}^{(k)}\mathbf{Q}^{(k)} + \alpha^{(k)}\mathbf{I}\end{aligned}\tag{2.33}$$

Les matrices restent toujours semblables, donc conservent les valeurs propres.

Concernant la convergence de la méthode, avec des décalages, et en prenant $\alpha^{(k)} = A_{n,n}^{(k)}$:

$$|\lambda_n - A_{n,n}^{(k)}| \xrightarrow[k \rightarrow \infty]{} 0\tag{2.34}$$

alors :

$$\|(\mathbf{A} - \lambda_{k+1}\mathbf{I})\mathbf{q}^{(k)}\| = \mathcal{O}(\|(\mathbf{A} - \lambda_k\mathbf{I})\mathbf{q}^{(k-1)}\|^3)\tag{2.35}$$

Sans translation, la convergence reste linéaire.

2.6.3 Forme de Hessenberg

Comme son nom l'indique, cette méthode est due à Gerhard Hessenberg (1874-1925) pour qui on pourra consulter [37] et est dédiée aux matrices \mathbf{A} réelles symétriques.

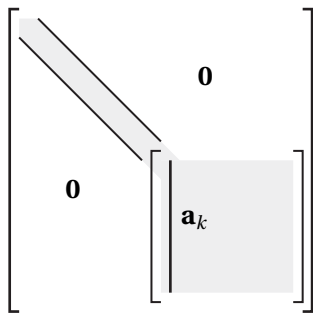


fig. 2.7 - Matrice en cours de tri-diagonalisation

On cherche cette fois à rendre \mathbf{A} tri-diagonale. Il s'agit de la même technique colonne par colonne que précédemment. On prend ici une transformation qui ne s'applique qu'à la deuxième colonne pour la première étape :

$$\mathbf{H} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{H}} \end{bmatrix}\tag{2.36}$$

mais comme \mathbf{A} est symétrique, il faut appliquer les transformations à gauche et à droite :

$$\mathbf{H}^{-1}\mathbf{A}\mathbf{H} = \mathbf{H}^T\mathbf{A}\mathbf{H} = \mathbf{H}\mathbf{A}\mathbf{H}\tag{2.37}$$

$\mathbf{A}^{(k)}$ prend donc la forme de la figure 2.7 et le principe consiste à travailler sur le vecteur $\begin{bmatrix} A_{k,k+1}^{(k)} & \dots & A_{k,n}^{(k)} \end{bmatrix}^T$ pour le rendre colinéaire à \mathbf{e}_{k+1} .

Cette méthode a la propriété intéressante de ne pas modifier les valeurs propres. C'est une méthode directe en $n - 2$ étapes et \mathbf{A} reste symétrique. Travaillons sur le premier itéré :

```

boucle sur  $k = 1, 2, \dots, n - 1$ 
 $\delta \leftarrow \mathbf{A}_{k+1:n,k}^{(k)T} \mathbf{A}_{k+1:n,k}^{(k)}$ 
 $\alpha \leftarrow \sqrt{\delta}$ 
 $\beta \leftarrow 1 / (\delta - \alpha \mathbf{A}_{k+1,k}^{(k)})$ 
 $\mathbf{v}_{k+1:n}^{(k)} \leftarrow \mathbf{A}_{k+1:n,k}^{(k)}$ 
 $\nu_{k+1}^{(k)} \leftarrow \nu_{k+1}^{(k)} - \alpha$ 
 $\mathbf{z}_{k+1:n} \leftarrow \beta \mathbf{A}_{k+1:n,k+1:n}^{(k)} \mathbf{v}_{k+1:n}^{(k)}$ 
 $\gamma \leftarrow \beta / 2 \mathbf{v}_{k+1:n}^{(k)T} \mathbf{z}_{k+1:n}$ 
 $\mathbf{z}_{k+1:n} \leftarrow \mathbf{z}_{k+1:n} - \gamma \mathbf{v}_{k+1:n}^{(k)}$ 
 $\mathbf{A}_{k+1:n,k+1:n}^{(k+1)} \leftarrow \mathbf{A}_{k+1:n,k+1:n}^{(k)}$ 
 $\quad - \mathbf{v}_{k+1:n}^{(k)} \mathbf{z}_{k+1:n}^T - \mathbf{z}_{k+1:n} \mathbf{v}_{k+1:n}^{(k)T}$ 
fin

```

alg. 2.8 - Transformation
de Householder

$$\mathbf{A}^{(1)} = \begin{bmatrix} \mathbf{A}_{1,1} & \tilde{\mathbf{a}}_1^T \\ \tilde{\mathbf{a}}_1 & \tilde{\mathbf{A}}^{(1)} \end{bmatrix} \quad (2.38)$$

on prend :

$$\tilde{\mathbf{H}}^{(1)} = \mathbf{I} - 2\tilde{\mathbf{u}}\tilde{\mathbf{u}}^T \quad (2.39)$$

en souhaitant avoir :

$$\tilde{\mathbf{H}}^{(1)} \tilde{\mathbf{a}}_1 = \alpha \tilde{\mathbf{e}}_1 \quad (2.40)$$

En fait, ici $\tilde{\mathbf{e}}_1 = \mathbf{e}_2$ et on procède comme précédemment pour avoir $\tilde{\mathbf{u}}$. On cherche donc $\tilde{\mathbf{H}}^{(1)} = \mathbf{I} - \beta \tilde{\mathbf{y}}\tilde{\mathbf{y}}^T$ avec $\alpha = \tilde{\mathbf{a}}_1^T \tilde{\mathbf{a}}_1$, $\beta = \frac{1}{\alpha(\alpha - A_{2,1})}$ et $\tilde{\mathbf{y}} = \tilde{\mathbf{a}}_1 - \alpha \tilde{\mathbf{e}}_1$. On calcule ensuite la transformée par :

$$\begin{aligned} \tilde{\mathbf{H}}^{(1)} \tilde{\mathbf{A}}^{(1)} \tilde{\mathbf{H}}^{(1)} &= (\mathbf{I} - \beta \tilde{\mathbf{y}}\tilde{\mathbf{y}}^T) \tilde{\mathbf{A}}^{(1)} (\mathbf{I} - \beta \tilde{\mathbf{y}}\tilde{\mathbf{y}}^T) \\ &= \tilde{\mathbf{A}}^{(1)} - \beta \tilde{\mathbf{y}}(\tilde{\mathbf{y}}^T \tilde{\mathbf{A}}^{(1)}) - \beta(\tilde{\mathbf{A}}^{(1)} \tilde{\mathbf{y}})\tilde{\mathbf{y}}^T + \beta^2 \tilde{\mathbf{y}}(\tilde{\mathbf{y}}^T \tilde{\mathbf{A}}^{(1)} \tilde{\mathbf{y}})\tilde{\mathbf{y}}^T \end{aligned} \quad (2.41)$$

et en utilisant les intermédiaires $\tilde{\mathbf{z}} = \beta \tilde{\mathbf{A}}^{(1)} \tilde{\mathbf{y}}$ et $\tilde{\mathbf{t}} = \tilde{\mathbf{z}} - \frac{\beta}{2}(\tilde{\mathbf{y}}^T \tilde{\mathbf{z}})\tilde{\mathbf{y}}$:

$$\begin{aligned} \tilde{\mathbf{H}}^{(1)} \tilde{\mathbf{A}}^{(1)} \tilde{\mathbf{H}}^{(1)} &= \tilde{\mathbf{A}}^{(1)} - \tilde{\mathbf{y}}\tilde{\mathbf{z}}^T - \tilde{\mathbf{z}}\tilde{\mathbf{y}}^T + \beta(\tilde{\mathbf{z}}^T \tilde{\mathbf{y}})\tilde{\mathbf{y}}\tilde{\mathbf{y}}^T \\ &= \tilde{\mathbf{A}}^{(1)} - \tilde{\mathbf{y}}\tilde{\mathbf{t}}^T - \tilde{\mathbf{t}}\tilde{\mathbf{y}}^T \end{aligned} \quad (2.42)$$

C'est la technique mise en œuvre dans l'algorithme 2.8, dans lequel \mathbf{v}_k n'est en général pas conservé et où \mathbf{z} est un vecteur de stockage supplémentaire. La réduction s'effectue en $n - 1$ étapes.

2.7 Exemple d'utilisation dans un code éléments finis

On considère $\mathbf{K}\mathbf{x} = \lambda\mathbf{M}\mathbf{x}$ avec \mathbf{K} symétrique, définie positive, \mathbf{M} symétrique positive pour l'obtention de m valeurs propres, en basse fréquence.

2.7.1 Réduction dans un sous-espace

On cherche à trouver toutes les m valeurs propres dans un sous-espace de dimension $p \geq m$, défini par p vecteurs rangés dans la matrice de base $\mathbf{V} = [\mathbf{v}_1 \ \dots \ \mathbf{v}_p]$. Éventuellement, les \mathbf{v}_i sont orthonormés au sens de \mathbf{K} , soit $\mathbf{V}^T \mathbf{K} \mathbf{V} = \tilde{\mathbf{K}} = \mathbf{I}$.

On cherche alors les vecteurs propres \mathbf{x} sous la forme $\mathbf{x} = \mathbf{V}\hat{\mathbf{x}}$, $\hat{\mathbf{x}}$ est de dimension p . Le problème est alors $\mathbf{K}\mathbf{V}\hat{\mathbf{x}} = \hat{\lambda}\mathbf{M}\mathbf{V}\hat{\mathbf{x}}$, d'où :

$$\underbrace{(\mathbf{V}^T \mathbf{K} \mathbf{V})}_{\tilde{\mathbf{K}}} \hat{\mathbf{x}} = \hat{\lambda} \underbrace{(\mathbf{V}^T \mathbf{M} \mathbf{V})}_{\tilde{\mathbf{M}}} \hat{\mathbf{x}} \quad (2.43)$$

Le problème $\tilde{\mathbf{K}}\hat{\mathbf{x}} = \hat{\lambda}\tilde{\mathbf{M}}\hat{\mathbf{x}}$ a été construit par une réduction de Ritz. Sa résolution complète permet de trouver p caractéristiques propres $(\hat{\lambda}, \hat{\mathbf{x}})$.

On retourne alors dans l'espace de départ où $\mathbf{x} = \mathbf{V}\hat{\mathbf{x}}$ est l'approximation d'un vecteur propre, $\hat{\lambda}$ est l'approximation de la valeur propre associée, approximations dont la qualité dépend du sous-espace choisi... Ce choix va être discuté dans les paragraphes suivants.

2.7.2 Technique itérative

Rapidement, il s'agit de :

1. choisir p vecteurs de départ ;
2. réduire à $\hat{\mathbf{K}}\hat{\mathbf{x}} = \hat{\lambda}\hat{\mathbf{M}}\hat{\mathbf{x}}$ dont on cherche les valeurs propres (si on a une approximation des valeurs propres, on fait un décalage) par une méthode QR ;
3. itérer en sens inverse avec décalage (souvent une seule itération suffit) pour avoir les approximations des vecteurs propres associés dans l'espace de départ ;
4. passer à l'itération suivante.

Comme la convergence, dépend du rapport des modules des valeurs propres, on a intérêt à conserver $p > m$ approximations. La règle conseillée est choisir une taille du sous-espace :

$$p = \min(2m, m + 8) \quad (2.44)$$

La convergence pour une valeur propre λ_i est proportionnelle à $\frac{|\lambda_i|}{|\lambda_{p+1}|}$.

2.8 Méthode de Lanczos

Cette méthode est due à Lanczos [38] en 1950. On travaille sur un sous-espace généré par la suite des itérés inverses d'un vecteur $\mathbf{v}^{(0)}$ de départ : $\left[\mathbf{v}^{(0)} \quad \mathbf{A}^{-1}\mathbf{v}^{(0)} \quad \dots \right]$ et sur une base orthonormée.

La convergence est très rapide : 3 à 4 itérations inverses par valeur propre convergée, et elle est indépendante de n ; cependant, la taille du sous-espace croît au cours des itérations.

Un inconvénient est la perte rapide d'orthogonalité sur les vecteurs successifs (elle se produit à chaque fois qu'une valeur propre converge), d'où des instabilités numériques. En pratique, pour pouvoir traiter de grands problèmes, il est nécessaire d'utiliser des stratégies de ré-orthogonalisation. Ces dernières ne seront pourtant pas décrites ici.

En termes simples, l'algorithme se résume à choisir un vecteur de départ $\mathbf{v}^{(0)}$ puis de résoudre le système $\mathbf{A}\mathbf{z}^{(k)} = \mathbf{v}^{(k)}$, d'orthogonaliser $\mathbf{z}^{(k)}$ par rapport aux $\mathbf{v}^{(j)}$, $j \leq k$ puis de normer $\mathbf{v}^{(k+1)} \leftarrow \frac{\mathbf{z}^{(k)}}{\|\mathbf{z}^{(k)}\|}$.

Pour démontrer le taux de convergence, il suffit de chercher $\mathbf{v}^{(k+1)}$ sous la forme :

$$\beta^{(k)}\mathbf{v}^{(k+1)} = \mathbf{A}^{-1}\mathbf{v}^{(k)} - \alpha^{(k)}\mathbf{v}^{(k)} - \beta^{(k-1)}\mathbf{v}^{(k-1)} \quad (2.45)$$

ce que l'on peut montrer par récurrence, avec $\mathbf{v}^{(k)\text{T}}\mathbf{v}^{(j)} = 0$ pour $j = k - 1$ et $j = k - 2$, et $\mathbf{v}^{(k)\text{T}}\mathbf{v}^{(k)} = 1$, pour avoir orthogonalité par rapport à tous les $\mathbf{v}^{(j)}$, $j < k$. Ces conditions

```

 $\mathbf{v}^{(0)} \leftarrow \mathbf{0}, \beta^{(0)} \leftarrow 0$ 
Vecteur de départ  $\mathbf{v}^{(1)}, \|\mathbf{v}^{(1)}\| = 1$ 
tant que  $R > \varepsilon$ , itérer sur  $k$ 
  résoudre  $\mathbf{A}\mathbf{z} = \mathbf{v}^{(k)}$ 
   $\alpha^{(k)} \leftarrow \mathbf{v}^{(k)\text{T}}\mathbf{z}$ 
   $\mathbf{d} \leftarrow \mathbf{z} - \alpha\mathbf{v}^{(k)} - \beta^{(k-1)}\mathbf{v}^{(k-1)}$ 
   $\beta^{(k)} \leftarrow \frac{1}{\sqrt{\mathbf{d}^T\mathbf{d}}}$ 
   $\mathbf{v}^{(k+1)} \leftarrow \beta^{(k)}\mathbf{d}$ 
  recherche de  $(\hat{\mathbf{x}}, \hat{\lambda})$  tel que  $\mathbf{T}\hat{\mathbf{x}} = \hat{\lambda}\hat{\mathbf{x}}$ 
fin

```

alg. 2.9 - Méthode de Lanczos

donnent les valeurs de $\alpha^{(k)} = \mathbf{v}^{(k)}\mathbf{A}^{-1}\mathbf{v}^{(k)}$ et $\beta^{(k-1)} = \mathbf{v}^{(k)}\mathbf{A}^{-1}\mathbf{v}^{(k-1)}$. Matriciellement, en notant \mathbf{V} de taille (n, k) la matrice qui contient tous les $\mathbf{v}^{(j)}$ disponibles à l'itération k , la relation (2.45) est équivalente à :

$$\mathbf{A}^{-1}\mathbf{V} - \mathbf{V} \underbrace{\begin{bmatrix} \alpha^{(0)} & \beta^{(0)} & & 0 \\ \beta^{(0)} & \ddots & \ddots & \\ & \ddots & \ddots & \beta^{(k-1)} \\ 0 & & \beta^{(k-1)} & \alpha^{(k)} \end{bmatrix}}_{\mathbf{T}} = \begin{bmatrix} 0 & \dots & 0 & \beta^{(k)} & \nu^{(k+1)} \end{bmatrix} \quad (2.46)$$

ce qui implique :

$$\mathbf{V}^T\mathbf{A}^{-1}\mathbf{V} - \mathbf{T} = \mathbf{0} \quad (2.47)$$

Par conséquent, \mathbf{T} a les mêmes valeurs propres que $\mathbf{V}^T\mathbf{A}^{-1}\mathbf{V}$ (voir le principe des sous-espaces). On cherche ensuite les caractéristiques propres de \mathbf{T} , $\mathbf{T}\hat{\mathbf{x}} = \hat{\lambda}\hat{\mathbf{x}}$, $\hat{\lambda}$ sont des approximations des valeurs propres de \mathbf{A}^{-1} et $\mathbf{V}\hat{\mathbf{x}}$, celles des vecteurs propres associés.

Afin d'illustrer le comportement de cette méthode, la figure 2.10 présente l'estimation des valeurs propres aux cours des itérations de la méthode de Lanczos sur le problème test de Wilkinson (de taille 21) comportant plusieurs valeurs propres multiples et plusieurs valeurs propres très proches (les cercles sont les valeurs propres exactes). On voit que la multiplicité et/ou proximité des valeurs propres échappe à l'algorithme, et que les valeurs propres se mettent à converger quasiment l'une après l'autre.

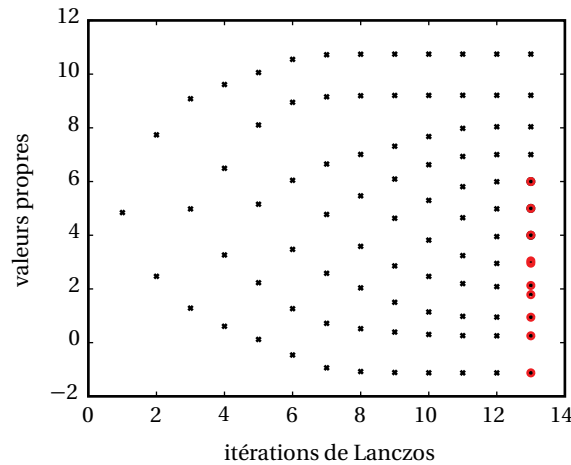


fig. 2.10 - Comportement de la méthode de Lanczos sur le test de Wilkinson

Critère d'arrêt

Classiquement, on borne la norme du résidu $\mathbf{r} = \mathbf{A}^{-1}\mathbf{V}\hat{\mathbf{x}} - \hat{\lambda}\mathbf{V}\hat{\mathbf{x}}$ avec :

$$\mathbf{A}^{-1}\mathbf{V}\hat{\mathbf{x}} = \mathbf{V}\mathbf{T}\hat{\mathbf{x}} + \begin{bmatrix} 0 & \dots & 0 & \beta^{(k)} & \nu^{(k+1)} \end{bmatrix} \hat{\mathbf{x}} \quad (2.48)$$

avec $\mathbf{T}\hat{\mathbf{x}} = \hat{\lambda}\hat{\mathbf{x}}$, d'où $\mathbf{r} = \beta^{(k)}\mathbf{v}^{(k+1)}\hat{\mathbf{x}}_k$.

L'algorithme 2.9 décrit cette méthode, pour lequel les $\alpha^{(k)}$ et les $\beta^{(k)}$ sont les composantes de \mathbf{T} .

RÉSOLUTION DE SYSTÈMES NON-LINÉAIRES PARTICULIERS

Dans cette partie, des préoccupations directement issues du calcul de structure par éléments finis sont abordées. Il s'agit ici d'appliquer les outils précédents à certains problèmes particuliers, et non pas de disposer d'un cours entier consacré aux éléments finis. Il en sera de même dans la partie 4.4 pour les relations de comportement du matériau. Pour une présentation dédiée aux éléments finis, le lecteur est invité à se reporter, par exemple, aux nombreux cours de qualité disponibles sur le réseau. Par exemple, ceux de C. Felippa : *Nonlinear Finite Element Methods*, essentiellement pour les non-linéarités d'origine géométriques, et *Advanced Finite Element Methods*, pour la technologie des éléments finis, donnés au *Aerospace Engineering Sciences Department* de l'*University of Colorado at Boulder*.

3.1 Non-linéarités matérielles de type plasticité

Le problème à résoudre est celui du calcul de structures dont le matériau possède un comportement non linéaire. Il peut par exemple s'agir d'un matériau élastique non linéaire (sans dissipation, donc) ou d'un problème d'évolution comme la plasticité ou la visco-plasticité. Dans ce dernier cas, une approche incrémentale transforme le problème d'évolution en une série de problèmes non-linéaires stationnaires à chaque incrément de temps. Ceux-ci sont alors résolus par exemple par des méthodes classiques de type Newton [39], décrites ici.

Pour plus d'informations, on pourra consulter [40].

3.1.1 Formulation du problème

Le problème de référence consiste à vérifier :

- les équations d'admissibilité cinématique : les déformations ε doivent être compatibles avec un champ de déplacement admissible U (régulier et satisfaisant les conditions aux limites en déplacement imposé)

$$\varepsilon = \varepsilon(U) \tag{3.1}$$

- les équations d'admissibilité statique : le champ de contraintes σ doit être en équilibre avec les charges extérieures ;

- les relations de comportement [41]. On considère un matériau dont l'état est décrit par (outre la déformation ou la contrainte) un jeu de variables internes \mathbb{X} , avec leurs variables thermodynamiques associées \mathbb{Y} .

La partition de la déformation en partie élastique et non élastique s'écrit :

$$\varepsilon = \varepsilon_e + \varepsilon_p \quad (3.2)$$

Les lois d'état reliant les variables cinématiques et statiques sont associées à l'énergie libre (qui dépend de l'état) $\psi(\varepsilon_e, \mathbb{X})$:

$$\begin{aligned} \sigma &= \frac{\partial \psi}{\partial \varepsilon_e} \\ \mathbb{Y} &= \frac{\partial \psi}{\partial \mathbb{X}} \end{aligned} \quad (3.3)$$

L'évolution de l'état, pour un modèle standard associé, se déduit du potentiel de dissipation $\Phi(\sigma, \mathbb{Y})$ par la loi de normalité :

$$\begin{aligned} \dot{\varepsilon}_p &= \dot{\lambda} \frac{\partial \Phi}{\partial \sigma} \\ -\dot{\mathbb{X}} &= \dot{\lambda} \frac{\partial \Phi}{\partial \mathbb{Y}} \end{aligned} \quad (3.4)$$

avec, dans le cas non visqueux, $\Phi \leq 0$, $\dot{\lambda} \geq 0$ et $\Phi \dot{\lambda} = 0$. λ est le multiplicateur plastique.

3.1.2 Résolution pour un incrément de temps

Pour les problèmes qui nous préoccupent, à chaque incrément de temps, on est amené à résoudre, après une discrétisation en déplacement, un problème de la forme :

$$\mathbf{r}(\mathbf{u}) = \mathbf{K}(\mathbf{u})\mathbf{u} - \mathbf{f} = \mathbf{0} \quad (3.5)$$

$\mathbf{r}(\mathbf{u})$ est le *gradient* à annuler. La matrice de rigidité dépend de la solution \mathbf{u} , par l'intermédiaire de la relation de comportement du matériau. Les méthodes de type Newton sont habituellement utilisées : si \mathbf{u}_0 est une approximation de la solution, cette dernière est $\mathbf{u} = \mathbf{u}_0 + \delta\mathbf{u}$ et un développement de Taylor de \mathbf{r} donne au voisinage de \mathbf{u} :

$$\frac{\partial \mathbf{r}}{\partial \mathbf{u}}(\mathbf{u}_0)\delta\mathbf{u} = -\mathbf{r}(\mathbf{u}_0) \quad (3.6)$$

ce qui suggère un algorithme itératif de la forme :

$$\mathbf{u}_{i+1} - \mathbf{u}_i = -s_i \mathbf{H}_i \mathbf{r}_i \quad (3.7)$$

où $\mathbf{r}_i = \mathbf{r}(\mathbf{u}_i)$, \mathbf{H}_i est une approximation de l'inverse de la matrice jacobienne $\mathbf{J}(\mathbf{u}_i) = \frac{\partial \mathbf{r}}{\partial \mathbf{u}}(\mathbf{u}_i)$. s_i est une amplitude destinée à réduire \mathbf{r}_{i+1} ¹.

En fait, les méthodes présentées ci-après sont issues d'une classe de problèmes plus vaste, celle de la minimisation de fonctionnelles non-linéaires $f(\mathbf{U})$; dans ce cadre, f est parfois appelée la fonction coût.

1. Dans un premier temps, on peut choisir $s_i = 1$ mais une procédure de type *line search* peut être utilisée comme une procédure itérative supplémentaire pour améliorer sa valeur

3.1.3 Méthode de la plus grande pente

Dans ce cas :

$$\mathbf{H} = \mathbf{I} \quad (3.8)$$

3.1.4 Méthodes de gradient conjugué

Ces méthodes n'imposent pas le stockage des matrices \mathbf{H}_i car seuls quelques vecteurs sont nécessaires au calcul de $\mathbf{H}_i \mathbf{r}_i$. Cela est réalisé avec une formule de récurrence sur une direction de recherche \mathbf{d} :

$$\mathbf{d}_i = \begin{cases} \mathbf{r}_i & \text{pour } i = 1 \\ \mathbf{r}_i + \zeta_i \mathbf{d}_{i-1} & \text{pour } i > 1 \end{cases} \quad (3.9)$$

et :

$$\mathbf{u}_{i+1} = \mathbf{u}_i - s_i \mathbf{d}_i \quad (3.10)$$

ζ_i est une scalaire utilisé pour conjuguer les directions de recherche, *i.e.* $\mathbf{d}_i^T \mathbf{d}_{i-1} = 0$, pour lequel deux expressions classiques sont :

1. Fletcher-Reeves [42] :

$$\zeta_i = \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{r}_{i-1}^T \mathbf{r}_{i-1}} \quad (3.11)$$

2. Polak-Ribière [43] :

$$\zeta_i = \frac{\mathbf{r}_i^T (\mathbf{r}_i - \mathbf{r}_{i-1})}{\mathbf{r}_{i-1}^T \mathbf{r}_{i-1}} \quad (3.12)$$

Si la fonctionnelle à minimiser est quadratique, donc si \mathbf{K} ne dépend pas de \mathbf{u} , alors toutes les directions de recherche successives sont conjuguées entre elles.

Pour mémoire, il existe d'autres méthodes de gradient conjugué, comme la méthode du *shortest residual* [44].

3.1.5 Méthode de Newton-Raphson

Ici, $s_i = 1$ et :

$$\mathbf{H}_i = \mathbf{J}(\mathbf{u}_i)^{-1} = \left[\frac{\partial \mathbf{r}}{\partial \mathbf{u}}(\mathbf{u}_i) \right]^{-1} = \left[\mathbf{K}(\mathbf{u}_i) + \frac{\partial \mathbf{K}}{\partial \mathbf{u}}(\mathbf{u}_i) \right]^{-1} \quad (3.13)$$

à chaque itération, il faut évaluer la matrice jacobienne et résoudre le précédent système linéaire.

3.1.6 Méthode de Newton-Raphson modifiée

Elle utilise un Jacobien figé, évalué à une précédente estimation de la solution \mathbf{u}^* :

$$\mathbf{H}_i = \mathbf{J}(\mathbf{u}^*)^{-1} \quad (3.14)$$

3.1.7 Mises à jour de type Quasi-Newton

L'idée est de modifier simplement \mathbf{H}_i à chaque itération plutôt que de le laisser fixé comme MNR, et plutôt que de le recalculer entièrement, comme pour NR.

$$\mathbf{H}_{i+1} = \mathbf{H}_i + \delta\mathbf{H}_i \quad (3.15)$$

en s'assurant que la condition sécante :

$$\mathbf{H}_{i+1}(\mathbf{u}_{i+1} - \mathbf{u}_i) = \mathbf{r}_{i+1} - \mathbf{r}_i \quad (3.16)$$

est satisfaite. Elle correspond à une approximation au premier ordre de la matrice jacobienne au voisinage de \mathbf{u}_i .

Si \mathbf{H}_i est mis à jour, on parle de mise à jour inverse, alors qu'une mise à jour directe modifie \mathbf{H}_i^{-1} .

Le rang de la mise à jour est le rang de la matrice de correction qui est utilisée. Pour une mise à jour de rang 1, $\delta\mathbf{H} = \mathbf{a}\mathbf{b}^T$ où \mathbf{a} et \mathbf{b} sont deux vecteurs, la formule de Sherman-Morrison-Woodbury [45, 46, 47] permet d'exprimer directement l'inverse, tant que \mathbf{H} n'est pas singulier et que $1 + \mathbf{b}^T\mathbf{H}^{-1}\mathbf{a} \neq 0$:

$$(\mathbf{H} + \mathbf{a}\mathbf{b}^T)^{-1} = \mathbf{H}^{-1} - \frac{\mathbf{H}^{-1}\mathbf{a}\mathbf{b}^T\mathbf{H}^{-1}}{1 + \mathbf{b}^T\mathbf{H}^{-1}\mathbf{a}} \quad (3.17)$$

L'algorithme générique consiste à itérer sur :

- direction de recherche $\mathbf{d}_i = -\mathbf{H}_i\mathbf{r}_i$;
- $\min_{s_i > 0} \|\mathbf{r}_{i+1}\|$ par exemple par une procédure de *line search* ;
- mise à jour de \mathbf{H}_i .

Le but est de construire une approximation de $\mathbf{J}^{-1}(\mathbf{u})$, qui s'améliore au cours des itérations, en utilisant l'information des gradients successifs. Par exemple, avec la précédente correction de rang 1, en vérifiant la symétrie et la condition sécante, on obtient :

$$\mathbf{H}_{i+1} = \mathbf{H}_i + \frac{(\delta_i - \mathbf{H}_i\boldsymbol{\gamma}_i)(\delta_i - \mathbf{H}_i\boldsymbol{\gamma}_i)^T}{\boldsymbol{\gamma}_i^T(\delta_i - \mathbf{H}_i\boldsymbol{\gamma}_i)} \quad (3.18)$$

où $\delta_i = \mathbf{u}_{i+1} - \mathbf{u}_i$ et $\boldsymbol{\gamma}_i = \mathbf{r}_{i+1} - \mathbf{r}_i$.

Une propriété intéressante pour l'efficacité est d'avoir \mathbf{H} symétrique, au moins quand le problème à résoudre l'est, et définie positive pour avoir effectivement une direction de descente telle que $\mathbf{d}_i^T\mathbf{r}_i < 0$. Avec la précédente mise à jour de rang 1, il n'est pas garanti d'avoir \mathbf{H}_i défini positif.

Méthode de Davidson-Fletcher-Powell

Cette méthode de rang 2, proposée dans [48] préserve le caractère défini positif. En démarant avec \mathbf{H}_0 symétrique, défini positif, la mise à jour est :

$$\mathbf{H}_{i+1}^{\text{DFP}} = \mathbf{H}_i + \frac{\delta_i\delta_i^T}{\delta_i^T\delta_i} - \frac{(\mathbf{H}_i\boldsymbol{\gamma}_i)(\mathbf{H}_i\boldsymbol{\gamma}_i)^T}{\boldsymbol{\gamma}_i^T\mathbf{H}_i\boldsymbol{\gamma}_i} \quad (3.19)$$

Dans le cas où la fonction coût est quadratique, les directions \mathbf{d}_i sont \mathbf{K} -orthogonales entre elles ; après n itérations, $\mathbf{H}_n = \mathbf{J}^{-1}$. Si de plus $\mathbf{H}_0 = \mathbf{I}$, il s'agit de la méthode du gradient conjugué.

Méthode de Broyden-Fletcher-Goldfarb-Shanno

Cette méthode de rang 2, proposée dans [49], préserve aussi le caractère défini positif. En démarrant avec \mathbf{H}_0 symétrique, défini positif, la mise à jour est :

$$\begin{aligned}\mathbf{H}_{i+1}^{\text{BFGS}} &= \mathbf{H}_{i+1}^{\text{DFP}} + (\boldsymbol{\gamma}_i^T \mathbf{H}_i \boldsymbol{\gamma}_i) \left(\frac{\boldsymbol{\delta}_i}{\boldsymbol{\delta}_i^T \boldsymbol{\gamma}_i} - \frac{\mathbf{H}_i \boldsymbol{\gamma}_i}{\boldsymbol{\gamma}_i^T \mathbf{H}_i \boldsymbol{\gamma}_i} \right) \left(\frac{\boldsymbol{\delta}_i}{\boldsymbol{\delta}_i^T \boldsymbol{\gamma}_i} - \frac{\mathbf{H}_i \boldsymbol{\gamma}_i}{\boldsymbol{\gamma}_i^T \mathbf{H}_i \boldsymbol{\gamma}_i} \right)^T \\ &= \mathbf{H}_i + \left(1 + \frac{\boldsymbol{\gamma}_i^T \mathbf{H}_i \boldsymbol{\gamma}_i}{\boldsymbol{\delta}_i^T \boldsymbol{\gamma}_i} \right) \frac{\boldsymbol{\delta}_i \boldsymbol{\delta}_i^T}{\boldsymbol{\delta}_i^T \boldsymbol{\gamma}_i} - \frac{\boldsymbol{\delta}_i \boldsymbol{\gamma}_i^T \mathbf{H}_i + \mathbf{H}_i \boldsymbol{\gamma}_i \boldsymbol{\delta}_i^T}{\boldsymbol{\delta}_i^T \boldsymbol{\gamma}_i}\end{aligned}\quad (3.20)$$

Cette méthode est réputée pour être l'une des plus robustes [50].

La famille de méthodes dite de Broyden est définie par :

$$\mathbf{H}_{i+1} = \alpha \mathbf{H}_{i+1}^{\text{BFGS}} + (1 - \alpha) \mathbf{H}_{i+1}^{\text{DFP}} \quad (3.21)$$

α étant tel que le caractère défini positif est assuré ($0 \leq \alpha \leq 1$ le fait).

Versions avec mémoire limitée

Même si elles sont symétriques, les corrections présentées précédemment souffrent d'une grande demande en mémoire, pour les problèmes éléments finis de grande taille (les matrices \mathbf{H}_i ne sont en effet pas creuses). Une solution consiste à ne stocker \mathbf{H}_i uniquement de façon implicite, c'est-à-dire à être capable de construire $\mathbf{H}_i \mathbf{r}_i$ avec \mathbf{H}_0 (qui lui est creux) et des vecteurs successifs. De plus, le calcul de la direction de recherche est alors souvent moins coûteux que dans le cas de l'utilisation explicite d'un produit matrice-vecteur.

Par exemple, le gradient conjugué préconditionné (PCG) par \mathbf{H}_0 , correspond à :

$$\mathbf{H}_{i+1} = \mathbf{V}_i \mathbf{H}_0 \quad (3.22)$$

avec :

$$\mathbf{V}_j = \mathbf{I} - \frac{\boldsymbol{\gamma}_j \boldsymbol{\delta}_j^T}{\boldsymbol{\delta}_j^T \boldsymbol{\gamma}_j} \quad (3.23)$$

et la méthode BFGS, à :

$$\mathbf{H}_{i+1} = \left(\prod_{j=0}^i \mathbf{V}_j \right)^T \mathbf{H}_0 \left(\prod_{j=0}^i \mathbf{V}_j \right) + \sum_{j=0}^i \left(\prod_{k=j+1}^i \mathbf{V}_k \right)^T \frac{\boldsymbol{\delta}_k \boldsymbol{\delta}_k^T}{\boldsymbol{\delta}_k^T \boldsymbol{\gamma}_k} \left(\prod_{k=j+1}^i \mathbf{V}_k \right) \quad (3.24)$$

Une troncature peut alors être effectuée pour limiter encore les ressources en mémoire, en éliminant les contributions des plus anciens termes. Elle correspond alors à la version *limited memory* de BFGS (L-BFGS) [51, 52].

4

RÉSOLUTION DE SYSTÈMES DIFFÉRENTIELS

Il est possible de consulter sur le sujet les références [53, 54]. Dans les applications qui nous intéressent, on est conduit à résoudre de tels systèmes, par exemple dans les cas suivants : intégration de relation de comportement, problème de thermique transitoire, problème de dynamique rapide...

4.1 Équation scalaire du premier ordre

4.1.1 Problème de Cauchy

Il s'agit de trouver la fonction $x(t)$, $t \in [0, T]$, telle que :

$$\begin{aligned} \frac{dx}{dt} &= f(x, t) \\ x(0) &= u_0 \end{aligned} \tag{4.1}$$

Théorème de Cauchy-Lipschitz

Il y a existence et unicité de la solution si f est assez régulière¹ en x :

$$\exists k > 0, \forall t \in [0, T], \forall (y, z) \in \mathbb{R}^2, |f(z, t) - f(y, t)| \leq k|z - y| \tag{4.2}$$

4.1.2 Ré-écriture du problème

Pour se ramener aux techniques d'intégration, il est plus intéressant d'écrire le problème sous forme intégrale :

$$x = u_0 + \int_0^t f(x, t) dt \tag{4.3}$$

Une *discrétisation* de l'intervalle de temps $[0, T]$ est la séquence d'instants $0 = t_0 < t_1 < \dots < t_m = T$. Le pas de temps est $h = t_{i+1} - t_i$. Le principe est de chercher à approcher la valeur de la fonction aux piquets de temps, $x(t_i)$, par la valeur x_i , ce qui permet de définir l'*erreur* $e_i = |x(t_i) - x_i|$. Dans la suite, on notera $f_i = f(x_i, t)$. Le problème incrémental est, ayant la solution approchée jusqu'à t_i , de la déterminer à t_{i+1} avec :

$$x = x_i + \int_{t_i}^t f(x, t) dt \tag{4.4}$$

1. Plus précisément en termes mathématiques, f est continue et lipschitzienne.

Exemple : schéma d'Euler explicite

La figure 4.3(a) illustre l'estimation de l'intégrale cherchée. Le schéma s'écrit :

$$\begin{aligned}x_{i+1} &= x_i + hf_i \\x_0 &= u_0\end{aligned}\tag{4.5}$$

4.1.3 Schémas explicites à un pas

Ce sont les schémas de la forme :

$$\begin{aligned}x_{i+1} &= x_i + h\Phi(x_i, t_i, h) \\x_0 &= u_0\end{aligned}\tag{4.6}$$

convergence

Si on a la propriété $\sup_i |e_i| \rightarrow_{h \rightarrow 0} 0$ (en l'absence d'erreur d'arrondi !)

stabilité

On considère le schéma perturbé :

$$\begin{aligned}\tilde{x}_{i+1} &= \tilde{x}_i + h\Phi(\tilde{x}_i, t_i, h) + \varepsilon_{i+1} \\ \tilde{x}_0 &= u_0 + \varepsilon_0\end{aligned}\tag{4.7}$$

Le schéma considéré est stable ssi :

$$\exists S > 0, \sup_i |\tilde{x}_i - x_i| \leq S \sum_{i=0}^m |\varepsilon_i|\tag{4.8}$$

consistance

Pour définir cette notion, on a besoin de l'erreur de consistance (*local truncation error*) :

$$c_{i+1} = x(t_{i+1}) - (x(t_i) + h\Phi(x(t_i), t_i, h))\tag{4.9}$$

Le schéma est consistant ssi :

$$\sum_i |c_i| \xrightarrow{h \rightarrow 0} 0\tag{4.10}$$

ordre

On dit d'un schéma qu'il est d'ordre p quand $\exists M > 0, \sup_i |c_i| \leq Mh^{p+1}$.

Les trois premières notions précédentes sont illustrées sur la figure 4.1.

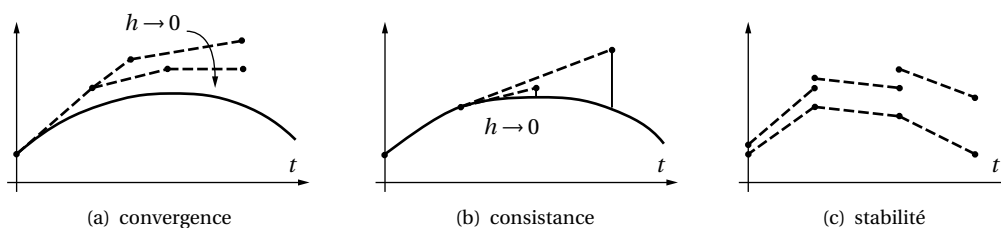


fig. 4.1 - Propriétés de base d'un schéma d'intégration

Propriétés

Il existe des liens entre les différentes caractéristiques précédentes, en particulier :

- stabilité et consistance \Leftrightarrow convergence
- Φ lipschitzienne en $x \Rightarrow$ stabilité
- consistance $\Leftrightarrow \Phi(x, t, 0) = f(x, t) \Leftrightarrow$ d'ordre $p \leq 1$
- d'ordre $p \Rightarrow |e_i| \leq Nh^p$

Exemple : schéma d'Euler explicite

Dans ce cas, $\Phi(x, t, h) = f(x, t)$ et par conséquent, le schéma est consistant et stable, donc convergent :

$$c_{i+1} = x(t_{i+1}) - (x(t_i) + hf(x(t_i), t_i)) = \mathcal{O}(h^2) \quad (4.11)$$

or :

$$\begin{aligned} x(t_{i+1}) &= x(t_i) + h\dot{x}(t_i) + \mathcal{O}(h^2) \\ \dot{x}(t_i) &= f(x(t_i), t_i) \end{aligned} \quad (4.12)$$

donc c'est un schéma d'ordre 1.

4.1.4 Quelques problèmes pratiques

Influence des erreurs d'arrondi

Perturbons le schéma :

$$\begin{aligned} \tilde{x}_{i+1} &= \tilde{x}_i + h\Phi(\tilde{x}_i, t_i, h) + \varepsilon_{i+1} \\ \tilde{x}_0 &= u_0 + \varepsilon_0 \end{aligned} \quad (4.13)$$

Si le schéma est convergent, $e_i = |x(t_i) - x_i| \xrightarrow{h \rightarrow 0} 0$. Qu'en est-il de $\tilde{e}_i = |x(t_i) - \tilde{x}_i|$?

$$\tilde{e}_i = |x(t_i) - x_i + x_i - \tilde{x}_i| \leq \underbrace{|x(t_i) - x_i|}_{\textcircled{1}} + \underbrace{|x_i - \tilde{x}_i|}_{\textcircled{2}} \quad (4.14)$$

avec, dans l'équation (4.14) :

- le schéma est d'ordre p donc le terme $\textcircled{1} \leq Nh^p$;
- le schéma est stable donc le terme $\textcircled{2} \leq S \left[\sum_{i=0}^m |\varepsilon| + \sum_{i=1}^m h|\varepsilon| \right]$.

Au bilan :

$$\tilde{e}_i \leq S[(m+1)|\varepsilon| + hm|\varepsilon|] + Nh^p \quad (4.15)$$

avec $hm = T$, d'où :

$$\tilde{e}_i \leq S|\varepsilon| \left(1 + T + \frac{T}{h} \right) + Nh^p \quad (4.16)$$

Il n'est donc pas forcément avantageux de réduire trop fortement le pas de temps pour la précision du schéma. En particulier, on peut dire que si le pas de temps n'est pas grand devant les erreurs numériques d'évaluation de la fonction Φ , la solution est entachée d'erreurs par cumul des erreurs numériques.

Problème numériquement mal posé

Considérons le problème suivant :

$$\begin{aligned}\frac{dx}{dt} &= 4x - 1 \\ x(0) &= \frac{1}{4}\end{aligned}\tag{4.17}$$

dont la solution est :

$$x(t) \equiv \frac{1}{4}\tag{4.18}$$

Perturbons maintenant les conditions initiales :

$$x(0) = \frac{1}{4} + \varepsilon\tag{4.19}$$

La nouvelle solution est alors très différente de la précédente :

$$x(t) = \frac{1}{4} + \varepsilon e^{4t}\tag{4.20}$$

Ce problème est numériquement mal posé. On risque d'avoir des difficultés à le résoudre car il est très sensible aux conditions initiales².

4.1.5 A-stabilité

De façon pratique, la propriété de stabilité d'un schéma n'est pas suffisante. Prenons par exemple le schéma d'Euler explicite déjà décrit comme stable, et résolvons le problème :

$$\begin{aligned}\frac{dx}{dt} &= -\omega^2 x \\ x(0) &= 1\end{aligned}\tag{4.21}$$

dont la solution exacte est :

$$x(t) = e^{-\omega^2 t}\tag{4.22}$$

La figure 4.2 présente les solutions obtenues numériquement en même temps que la solution exacte pour deux pas de temps. Ces solutions diffèrent fortement, et de plus en plus au cours du temps, si le pas de temps n'est pas assez petit. On parlera alors d'un schéma conditionnellement stable. Un schéma inconditionnellement stable³ [55] se définit pour les équations linéaires. On teste alors le schéma sur un problème du type :

$$\frac{dx}{dt} = cx\tag{4.23}$$

où c est une constante complexe et pour lequel le cas où le schéma est tel que $x_{i+1} = r(ch)x_i$, où r est un polynôme ou une fraction rationnelle. On définit alors le domaine de A-stabilité comme suit :

$$D_{\text{stab}} = \{z \in \mathbb{C}, |r(z)| \leq 1\}\tag{4.24}$$

Un schéma est alors dit A-stable (ou inconditionnellement stable) si et seulement si :

$$\{z \in \mathbb{C}, \text{Re}(z) \leq 0\} \subset D_{\text{stab}}\tag{4.25}$$

2. on peut d'ailleurs se poser des questions quant à la notion même de solution pour un tel problème
3. ou A-stable pour *absolute stability*

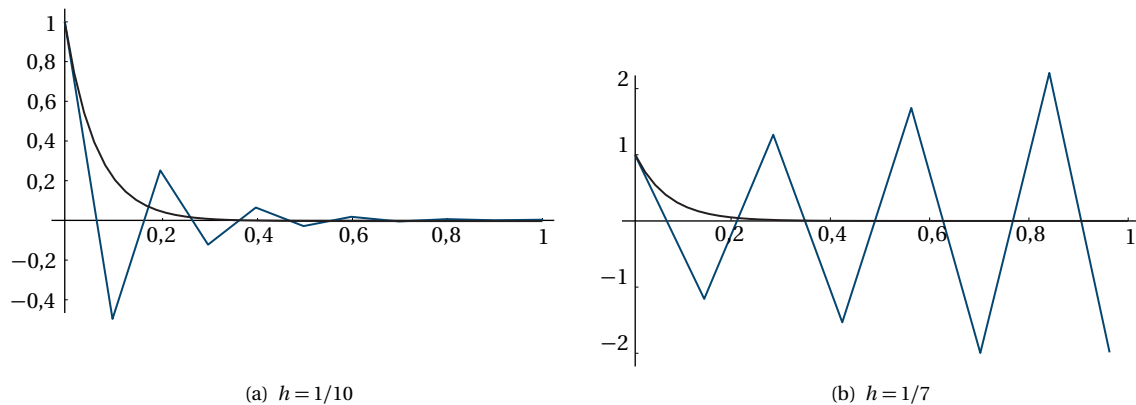


fig. 4.2 - Stabilité du schéma d'Euler explicite

Exemple : schéma d'Euler explicite

La relation est la suivante :

$$x_{i+1} = x_i + hc x_i = (1 + hc)x_i \Rightarrow r(z) = 1 + z \quad (4.26)$$

La condition de stabilité est donc :

$$|1 + hc| \leq 1 \quad (4.27)$$

Pour $c = -\omega^2 \in \mathbb{R}^-$, elle s'écrit $h \leq \frac{2}{\omega^2}$. Il s'agit donc d'un schéma conditionnellement stable. Pour l'exemple précédent, avec $c = -15$ (la solution tend vers 0), le critère de stabilité donne $0 \leq h \leq \frac{1}{7.5}$. Ceci est bien en accord avec les résultats de la figure 4.2.

La condition de stabilité dépend bien entendu du schéma, mais aussi du problème. Son interprétation est que, pour une solution bornée, il faut que le schéma donne une approximation bornée.

4.1.6 Méthode des trapèzes généralisée

Cette méthode est décrite avec un paramètre $\alpha \in [0, 1]$ par :

$$\begin{aligned} x_{i+1} &= x_i + h [(1 - \alpha)f_i + \alpha f_{i+1}] \\ x_0 &= u_0 \end{aligned} \quad (4.28)$$

Le calcul de x_{i+1} fait intervenir $\alpha f_{i+1} = \alpha f(x_{i+1}, t_{i+1})$. Il s'agit donc d'une méthode *implicite* si $\alpha \neq 0$.

Le cas $\alpha = 0$ conduit à la méthode d'Euler explicite ; dans le cas où $\alpha = 1$, il s'agit de la méthode d'Euler implicite ; enfin si $\alpha = \frac{1}{2}$, c'est la méthode des trapèzes, appelée aussi méthode de Crank-Nicolson et qui date de 1947. La figure 4.3 illustre l'approximation de l'intégration réalisée par chacun de ces cas particuliers.

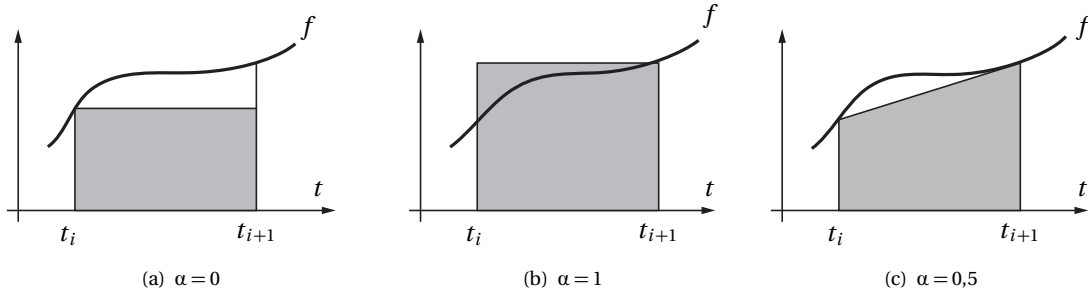


fig. 4.3 - Différents types d'intégration

Stabilité

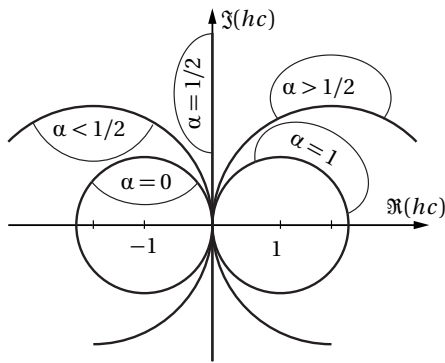


fig. 4.4 - Domaines de stabilité pour la méthode des trapèzes généralisés

On considère le problème :

$$\frac{dx}{dt} = cx \tag{4.29}$$

le schéma donne :

$$x_{i+1} = x_i + h[(1 - \alpha)cx_i + \alpha cx_{i+1}] \tag{4.30}$$

donc :

$$x_{i+1} = \underbrace{(1 - \alpha hc)^{-1}(1 + (1 - \alpha)hc)}_B x_i \tag{4.31}$$

Il faut avoir $\rho(B) < 1$, ce qui est le cas si toutes les valeurs propres de B sont simples,

c'est-à-dire si la condition suivante est vérifiée :

$$\frac{|1 + (1 - \alpha)hc|}{|1 - \alpha hc|} \leq 1 \tag{4.32}$$

- si $\alpha \geq \frac{1}{2}$ le schéma est inconditionnellement stable ;
- si $\alpha < \frac{1}{2}$ avec $c = -\omega^2$, la condition de stabilité s'écrit $h \leq \frac{2}{(1-2\alpha)\omega^2}$.

Les différents domaines de stabilité correspondants sont tracés sur la figure 4.4.

Ordre

En utilisant la même méthode que celle déjà employée pour le schéma d'Euler explicite, on peut facilement montrer que le schéma des trapèzes généralisés est d'ordre 1 si $\alpha \neq \frac{1}{2}$, et qu'il est d'ordre 2 pour $\alpha = \frac{1}{2}$, c'est-à-dire pour Crank-Nicolson.

4.1.7 Exemple de méthode à pas multiples : Adams

L'idée des méthodes à pas multiple est d'utiliser non seulement l'évaluation f_i mais aussi un certain nombre de valeurs passées $f_{i-1} \dots$

À titre d'illustration, la méthode d'Adams, datant de 1883, utilisant un développement de Taylor de la solution est maintenant décrite.

Adams explicite

On utilise pour cela un développement de Taylor (avant) de $x(t_{i+1})$:

$$x(t_{i+1}) = x(t_i) + h\dot{x}(t_i) + \frac{h^2}{2!}\ddot{x}(t_i) + \dots \quad (4.33)$$

avec $\dot{x}(t_i) = f(x(t_i), t_i)$ et $\ddot{x}(t_i) = \frac{df}{dt}(x(t_i), t_i)$. On garde ensuite les termes d'un ordre suffisant pour obtenir l'ordre voulu du schéma. Par exemple, pour Adams explicite d'ordre 1, on conserve les termes en $\mathcal{O}(h)$; on obtient $x_{i+1} = x_i + hf_i + \mathcal{O}(h^2)$, c'est Euler explicite !

Pour Adams explicite d'ordre 2, on conserve les termes en $\mathcal{O}(h^2)$ qui impliquent la dérivée première de f ; pour l'estimer on utilise encore le développement d'une fonction g :

$$g(t_{i+1}) = g(t_i) - h\dot{g}(t_i) + \mathcal{O}(h^2) \quad (4.34)$$

qui permet d'avoir l'expression de la dérivée \dot{g} :

$$\dot{g}(t_i) = \frac{1}{h}[g(t_i) - g(t_{i+1})] + \mathcal{O}(h) \quad (4.35)$$

Avec ce procédé appliqué à $\frac{df}{dt}$, on obtient la méthode cherchée :

$$x_{i+1} = x_i + \frac{h}{2}[3f_i - f_{i-1}] + \mathcal{O}(h^3) \quad (4.36)$$

La démarche se prolonge pour les ordres de schéma supérieurs. On obtient ainsi les formules d'Adams-Bashforth :

$$\begin{aligned} x_i &= x_{i-1} + hf_{i-1} \\ x_i &= x_{i-1} + \frac{h}{2}(3f_{i-1} - f_{i-2}) \\ x_i &= x_{i-1} + \frac{h}{12}(23f_{i-1} - 16f_{i-2} + 5f_{i-3}) \\ x_i &= x_{i-1} + \frac{h}{24}(55f_{i-1} - 59f_{i-2} + 37f_{i-3} - 9f_{i-4}) \\ x_i &= x_{i-1} + \frac{h}{720}(1901f_{i-1} - 2774f_{i-2} + 2616f_{i-3} - 1274f_{i-4} + 251f_{i-5}) \\ x_i &= x_{i-1} + \frac{h}{1440}(4277f_{i-1} - 7923f_{i-2} + 9982f_{i-3} - 7298f_{i-4} + 2877f_{i-5} - 475f_{i-6}) \end{aligned} \quad (4.37)$$

On ne peut cependant pas augmenter l'ordre du schéma impunément : l'inconvénient réside dans un rétrécissement du domaine de stabilité. Celui-ci est illustré sur la figure 4.5.

Adams implicite

L'idée est la même que précédemment, en utilisant un développement de Taylor arrière de :

$$x(t_i) = x(t_{i+1} - h) = x(t_{i+1}) - h\dot{x}(t_{i+1}) + \frac{h^2}{2!}\ddot{x}(t_{i+1}) + \dots \quad (4.38)$$

La méthode d'Adams implicite d'ordre 1 est ainsi :

$$x_i = x_{i+1} - hf_{i+1} + \mathcal{O}(h^2) \quad (4.39)$$

d'où :

$$x_{i+1} = x_i + hf_{i+1} + \mathcal{O}(h^2) \quad (4.40)$$

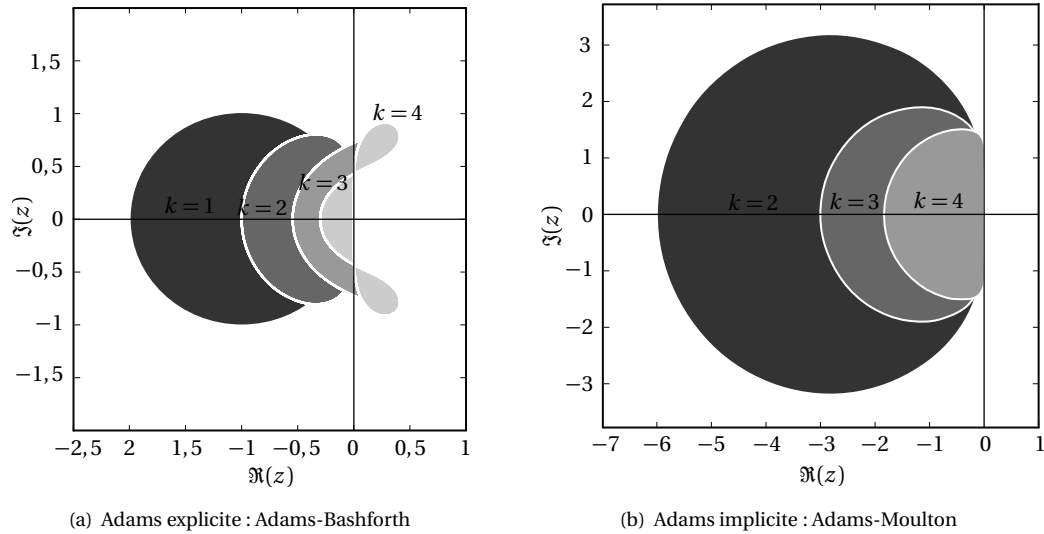


fig. 4.5 - Domaine de stabilité dans le plan complexe

c'est Euler implicite !

L'utilisation de la même technique que pour Adams explicite permet de construire les schémas implicites d'ordres plus élevés, et on obtient alors les formules d'Adams-Moulton :

$$\begin{aligned}
 x_i &= x_{i-1} + hf_i \\
 x_i &= x_{i-1} + \frac{h}{2}(f_i + f_{i-1}) \\
 x_i &= x_{i-1} + \frac{h}{12}(f_i + 8f_{i-1} - f_{i-2}) \\
 x_i &= x_{i-1} + \frac{h}{24}(9f_i + 19f_{i-1} - 5f_{i-2} + f_{i-3}) \\
 x_i &= x_{i-1} + \frac{h}{720}(251f_i + 646f_{i-1} - 264f_{i-2} + 106f_{i-3} - 19f_{i-4}) \\
 x_i &= x_{i-1} + \frac{h}{1440}(475f_i + 1427f_{i-1} - 798f_{i-2} + 482f_{i-3} - 173f_{i-4} + 27f_{i-5})
 \end{aligned} \tag{4.41}$$

Le fait d'utiliser des schémas implicites augmente la stabilité, mais la même tendance est observée lorsque l'ordre croît, voir figure 4.5.

Amorçage du schéma

Aux premiers pas de temps, on ne dispose pas forcément de tous les f_{i-q} . Dans ce cas, on peut utiliser une méthode à 1 pas pour démarrer l'intégration jusqu'à avoir suffisamment progressé pour pouvoir appliquer la méthode à pas multiples voulue.

4.1.8 Exemple de méthode de prédiction correction basée sur Adams

Le problème des méthodes implicites réside dans le fait, qu'en non-linéaire au moins, on ne peut exprimer explicitement la valeur au pas courant. On est alors obligé de mettre en œuvre une méthode itérative, généralement de type point fixe.

On part ainsi de la version implicite ; par exemple Adams d'ordre 2 : connaissant $x_{i+1}^{(0)}$, on évalue $f(x_{i+1}^{(0)}, t_{i+1})$ qui permet, en utilisant la formule d'intégration d'Adams implicite,

d'atteindre une valeur $x_{i+1}^{(1)}$. Cette dernière étape est appelée *correction*. On itère ainsi jusqu'à une convergence souhaitée. Pour avoir la première valeur $x_{i+1}^{(0)}$, la *prédiction*, on utilise la version implicite de la méthode, par exemple ici, Adams explicite d'ordre 2. En pratique, une à trois itérations suffisent.

On peut aussi apporter une amélioration à faible surcoût avec une correction du prédicteur avec l'erreur de troncature. Pour Adams d'ordre 2, on procède comme suit :

1. prédicteur (AE2) :

$$x_{i+1}^{(0)} \leftarrow x_i + \frac{h}{2}(3f_i - f_{i-1}) \quad (4.42)$$

2. modificateur :

$$\hat{x}_{i+1}^{(0)} \leftarrow x_{i+1}^{(0)} - \frac{1}{12} (x_i - x_i^{(0)}) \quad (4.43)$$

3. correcteur (AI2) :

$$x_{i+1} \leftarrow x_i + \frac{h}{2} (\hat{f}_{i+1} + f_{i-1}) \quad (4.44)$$

où $\hat{f}_{i+1} = f(\hat{x}_{i+1}^{(0)}, t_{i+1})$. Le coefficient $\frac{1}{12}$ est celui issu de la formule de (AI3).

4.1.9 Méthode de Runge-Kutta

Historiquement présentée dans [56] en 1901, cette méthode est de la forme :

$$x_{i+1} = x_i + h\Phi(x, t, h) \quad (4.45)$$

$$x_0 = u_0$$

avec :

$$\begin{aligned} \Phi(x, t, h) &= \sum_{l=1}^q a_l k_l(x, t, h) \\ k_1(x, t, h) &= f(x, t) \\ k_l(x, t, h) &= f\left(x + \sum_{j=1}^{l-1} \beta_j^{(l)} k_j(x, t, h), t + \alpha_l h\right) \end{aligned} \quad (4.46)$$

L'idée est d'intégrer en utilisant des valeurs en des points particuliers, avec des poids d'intégration particuliers. Les constantes a_l , α_l et $\beta_j^{(l)}$ sont ajustées pour obtenir l'ordre souhaité.

Exemple de schéma à l'ordre 2

On prend $\Phi = af(x, t) + bf(t + \alpha h, x + \beta hf(x, t))$ et on détermine a , b , α , β pour avoir :

$$\underbrace{\frac{1}{h}(x(t+h) - x(t))}_{\textcircled{1}} - \underbrace{\Phi(x(t), t, h)}_{\textcircled{2}} = \mathcal{O}(h^2) \quad (4.47)$$

avec, dans l'équation (4.47) :

$$\begin{aligned} \textcircled{1} &= f + \frac{h}{2} \left(\frac{\partial f}{\partial t} + \frac{\partial f}{\partial x} f \right) + \mathcal{O}(h^2) \\ \textcircled{2} &= af + b \left(f + \alpha h \frac{\partial f}{\partial t} + \beta hf \frac{\partial f}{\partial x} + \mathcal{O}(h^2) \right) \end{aligned}$$

Il faut donc avoir :

$$0 = 1 - a - b \quad 0 = \frac{h}{2} - b\alpha h \quad 0 = \frac{h}{2} - b\beta h \quad (4.48)$$

d'où les coefficients cherchés $a = 1 - b$, $\alpha = \beta = \frac{1}{2b}$ et le schéma :

$$\Phi(x, t, h) = (1 - b)f(x, t) + bf\left(x + \frac{h}{2b}f(x, t), t + \frac{h}{2b}\right) \quad (4.49)$$

où b est une constante arbitraire non nulle. Des cas particuliers existent :

- $b = 1$: il s'agit de la méthode d'Euler modifiée, d'ordre 2 ;
- $b = \frac{1}{2}$: il s'agit de la méthode de Heun, d'ordre 2.

Exemple de schéma à l'ordre 4

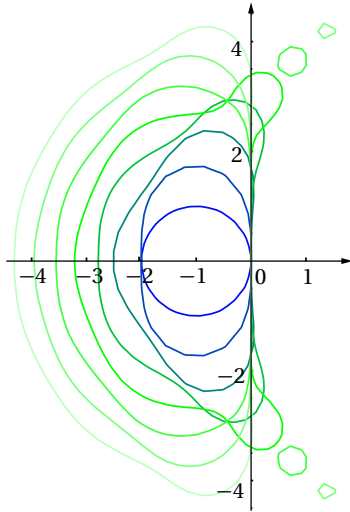


fig. 4.6 - Domaine de stabilité pour Runge-Kutta (ordre 1 à 8)

Bien entendu, cette technique se poursuit pour les ordres supérieurs. Un ordre souvent utilisé avec cette méthode est l'ordre 4. Le schéma se présente alors comme suit :

$$\begin{aligned} \Phi(x, t, h) &= \frac{1}{6}(k_1 + 2k_2 + 2k_3 + k_4) \\ k_1 &= f(x, t) \\ k_2 &= f\left(x + \frac{h}{2}k_1, t + \frac{h}{2}\right) \\ k_3 &= f\left(x + \frac{h}{2}k_2, t + \frac{h}{2}\right) \\ k_4 &= f(x + hk_3, t + h) \end{aligned} \quad (4.50)$$

Bien entendu, des méthodes de prédiction-correction basées sur la méthode de Runge-Kutta peuvent aussi être mises en place.

Cette méthode requière beaucoup d'évaluations de f , ce qui peut la pénaliser au niveau du coût. Concernant les domaines de stabilité, on peut observer la même tendance que pour les méthodes d'Adams, voir figure 4.6.

4.1.10 Autres méthodes

Citons pour mémoire les méthodes de Hanning, de Milne-Simpson, les méthodes de tir (*shooting*), les méthodes de collocation optimale, et une méthode bien adaptée à l'intégration des relations de comportements décrits par variables internes à 1 seul seuil, la méthode de retour radial [57].

Il existe aussi une famille de méthodes appelée Galerkin discontinu en temps, plus inspirée des formulations variationnelles en espace. Pour l'illustrer, considérons à nouveau le problème de Cauchy (4.1). Il peut être réécrit en utilisant une fonction test x^* à choisir dans un espace adéquat :

$$\int_0^T \left(\frac{dx}{dt} - f(x, t) \right) x^* dt = 0 \quad (4.51)$$

L'idée est de pouvoir utiliser des discrétisations en temps qui tolèrent *a priori* des discontinuités aux pas de temps : à $t = t_i$, on peut définir une valeur à gauche x_i^- et une valeur à droite x_i^+ différentes. Le saut est noté $\llbracket x \rrbracket_i = x_i^+ - x_i^-$. L'intégrale précédente est donc à prendre en un sens généralisé, qui fait « travailler » les sauts⁴ :

$$\int_0^T \left(\frac{dx}{dt} - f(x, t) \right) x^* dt = \sum_i \left[\int_{t_i}^{t_{i+1}} \left(\frac{dx}{dt} - f(x, t) \right) x^* dt + \llbracket x \rrbracket_i x_i^* \right] = 0 \quad (4.52)$$

La fonction test est choisie dans le même espace que x , et on définit sa valeur en t_i comme $x_i^* = \alpha x_i^{*+} + (1 - \alpha) x_i^{*-}$. Dans le cas standard, on prend $\alpha = 1$ ⁵. Le problème peut alors s'écrire sur chaque intervalle de temps :

$$\int_{t_i}^{t_{i+1}} \left(\frac{dx}{dt} - f(x, t) \right) x^* dt + \llbracket x \rrbracket_i x_i^* = 0 \quad (4.53)$$

Cette méthode est très ressemblante à une méthode de Runge-Kutta à un pas implicite (A-stable) d'ordre $2l + 1$ si on prend les espaces polynômiaux de degré inférieur ou égal à l par incrément de temps (Lesaint et Raviart 1974). Il faut aussi prévoir une méthode numérique de calcul d'intégrale en temps qui intègre exactement des polynômes de degré $2l$ (Cockburn et Shu 1989).

4.2 Systèmes différentiels du premier ordre

Cette situation est rencontrée en éléments finis, lors de l'intégration de relations d'un comportement décrit par variables internes et avec un système de taille bien plus grande, lors de calculs de diffusion, par exemple, un calcul de conduction thermique transitoire, qui amène à résoudre un problème de la forme :

$$\begin{aligned} \mathbf{M}\dot{\mathbf{x}} + \mathbf{K}\mathbf{x} &= \mathbf{g}(t) \\ \mathbf{x}(0) &= \mathbf{u}_0 \end{aligned} \quad (4.54)$$

où \mathbf{K} , généralement symétrique positive, est la matrice de conduction, et \mathbf{M} , généralement symétrique, définie positive, est la matrice de capacité. \mathbf{x} est le vecteur des températures inconnues aux nœuds de la discrétisation.

4.2.1 Emploi des méthodes précédentes

Les méthodes précédentes se généralisent sans grande difficulté à ce cas de figure.

Exemple

Prenons le schéma d'Euler explicite. Le problème précédent peut s'écrire :

$$\dot{\mathbf{x}} = -\mathbf{M}^{-1}\mathbf{K}\mathbf{x} + \mathbf{M}^{-1}\mathbf{g} \quad (4.55)$$

On applique alors le schéma voulu :

$$\mathbf{x}_{i+1} = \mathbf{x}_i + h(-\mathbf{M}^{-1}\mathbf{K}\mathbf{x}_i + \mathbf{M}^{-1}\mathbf{g}_i) \quad (4.56)$$

4. On suppose cependant $f(x, t)$ suffisamment régulière.

5. *one-sided upwinding*

Pour éviter de faire intervenir directement \mathbf{M}^{-1} , on préfère le mettre sous la forme :

$$\mathbf{M}\mathbf{x}_{i+1} = \mathbf{M}\mathbf{x}_i + h(\mathbf{K}\mathbf{x}_i + \mathbf{g}_i) \quad (4.57)$$

$$\mathbf{x}_0 = \mathbf{u}_0$$

Stabilité

On emploie la même démarche que dans le cas des équations scalaires, dans la base propre (pour les équations linéaires), avec la valeur propre la plus contraignante, c'est-à-dire la pulsation ω la plus grande.

Dans le cas d'Euler explicite, on a donc le critère de stabilité $h \leq \frac{2}{\omega_M^2}$ et on peut borner la pulsation propre maxi du système par la pulsation propre maxi élémentaire $\omega_M \leq \omega_{\max}^e$, voir annexe B.

Pour avoir un vrai schéma explicite, il faut ne pas avoir à résoudre de système linéaire couplé dans (4.57). Il faut donc avoir une matrice de masse diagonale. Une modification du problème (une approximation supplémentaire dans la modélisation) consiste à remplacer la matrice de masse par une matrice de masse « lumpée » diagonale. Pour cela, plusieurs techniques existent, comme une matrice proportionnelle à la masse diagonale, avec respect de la masse totale par élément, ou même la somme des lignes de la matrice de départ.

Outre une approximation supplémentaire [58], on peut se poser la question de la conservation des propriétés du schéma...

4.2.2 Choix d'un schéma : influence de l'ordre d'intégration

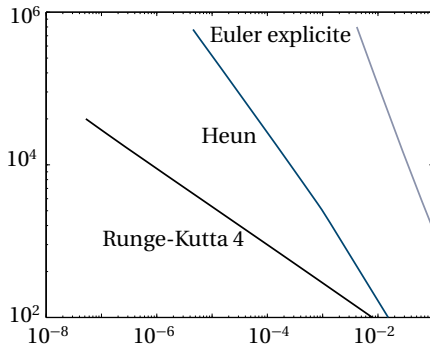


fig. 4.7 - Influence de l'ordre d'un schéma sur le coût

Si on considère que le coût est directement lié au nombre d'évaluations de la fonction f , on peut tracer le coût en fonction de l'erreur obtenue pour différents schémas.

Par exemple, si on choisit le cas test suivant :

$$\begin{pmatrix} \dot{x} \\ \dot{y} \end{pmatrix} = \begin{pmatrix} y \\ (1-x^2)y - x \end{pmatrix} \quad (4.58)$$

$$\begin{pmatrix} x(0) \\ y(0) \end{pmatrix} = \begin{pmatrix} u_0 \\ 0 \end{pmatrix}$$

Il s'agit d'un système différentiel du premier ordre, non linéaire, pour $t \in [0, T]$. Si $\alpha \approx$

2,008619861986087484313650940188 et avec $T \approx 6,6632868593231301896996820305$, la solution est périodique, de période T .

Une mesure de l'erreur peut donc être de prendre à la fin de l'intervalle d'étude :

$$e = \left\| \begin{pmatrix} x_T \\ y_T \end{pmatrix} - \begin{pmatrix} x_0 \\ y_0 \end{pmatrix} \right\| \quad (4.59)$$

La figure 4.7 reporte en échelles logarithmiques le nombre d'évaluations de f en fonction du niveau d'erreur e pour plusieurs schémas d'ordres différents : Euler explicite d'ordre 1, Heun

d'ordre 2 et Runge-Kutta d'ordre 4. On voit donc l'intérêt de mettre en œuvre des schémas d'ordre élevé, en gardant une stabilité raisonnable.

4.3 Systèmes différentiels du deuxième ordre

On trouve ce genre de problème par exemple en dynamique transitoire des structures :

$$\begin{aligned} \mathbf{M}\ddot{\mathbf{x}} + \mathbf{C}\dot{\mathbf{x}} + \mathbf{K}\mathbf{x} &= \mathbf{g}(t) \\ \dot{\mathbf{x}}(0) &= \mathbf{v}_0 \\ \mathbf{x}(0) &= \mathbf{u}_0 \end{aligned} \quad (4.60)$$

\mathbf{M} est la matrice de masse, \mathbf{C} est la matrice d'amortissement et \mathbf{K} la matrice de rigidité.

De façon générale, l'ordre d'un système linéaire peut être abaissé, ici à l'ordre 1, mais en doublant la taille du système. Ici, on s'intéressera au traitement direct du système au second ordre, dans le cas où on n'a pas d'amortissement, c'est-à-dire où $\mathbf{C} = \mathbf{0}$. L'équilibre est donc :

$$\mathbf{M}\ddot{\mathbf{x}} + \mathbf{K}\mathbf{x} = \mathbf{g}(t) \quad (4.61)$$

4.3.1 Méthode de Newmark

Il s'agit d'une famille de méthodes implicites à un pas [59] et définie en 1959. Elle se présente sous la forme suivante :

$$\begin{aligned} \mathbf{x}_{i+1} &= \mathbf{x}_i + h\dot{\mathbf{x}}_i + \frac{1}{2}h^2((1-2\beta)\ddot{\mathbf{x}}_i + 2\beta\ddot{\mathbf{x}}_{i+1}) \\ \dot{\mathbf{x}}_{i+1} &= \dot{\mathbf{x}}_i + h((1-\gamma)\ddot{\mathbf{x}}_i + \gamma\ddot{\mathbf{x}}_{i+1}) \end{aligned} \quad (4.62)$$

On retrouve bien le caractère implicite par l'apparition de quantités indicées $i+1$ au second membre. On peut aussi remarquer que si $\beta = \gamma = 0$, on retrouve des développements de Taylor puisque le schéma de Newmark ressemble à un développement au troisième ordre :

$$\begin{aligned} \mathbf{x}_{i+1} &= \mathbf{x}_i + h\dot{\mathbf{x}}_i + \frac{h^2}{2}\ddot{\mathbf{x}}_i + \frac{h^3}{6}6\beta\frac{\ddot{\mathbf{x}}_{i+1} - \ddot{\mathbf{x}}_i}{h} \\ \dot{\mathbf{x}}_{i+1} &= \dot{\mathbf{x}}_i + h\ddot{\mathbf{x}}_i + \frac{h^2}{2}2\gamma\frac{\ddot{\mathbf{x}}_{i+1} - \ddot{\mathbf{x}}_i}{h} \end{aligned} \quad (4.63)$$

Il reste à écrire l'équilibre (4.61) à l'instant t_{i+1} avec le schéma précédent. On peut alors prendre l'accélération comme inconnue :

$$(\beta h^2 \mathbf{K} + \mathbf{M})\ddot{\mathbf{x}}_{i+1} = \mathbf{g}_{i+1} - \mathbf{K} \left(\mathbf{x}_i + h\dot{\mathbf{x}}_i + h^2 \left(\frac{1}{2} - \beta \right) \ddot{\mathbf{x}}_i \right) \quad (4.64)$$

et on obtient le système à résoudre pour atteindre $\ddot{\mathbf{x}}_{i+1}$; mais on peut aussi l'écrire avec le déplacement comme inconnue :

$$(\beta h^2 \mathbf{K} + \mathbf{M})\mathbf{x}_{i+1} = \beta h^2 \mathbf{g}_{i+1} + \mathbf{M} \left(\mathbf{x}_i + h\dot{\mathbf{x}}_i + h^2 \left(\frac{1}{2} - \beta \right) \ddot{\mathbf{x}}_i \right) \quad (4.65)$$

qui est le système à résoudre pour atteindre \mathbf{x}_{i+1} . Dans ce dernier cas, cependant, un problème subsiste si $\beta = 0$. En effet, soit on utilise l'équilibre pour revenir à l'accélération :

$$\ddot{\mathbf{x}}_{i+1} = \mathbf{M}^{-1}(\mathbf{g}_{i+1} - \mathbf{K}\mathbf{x}_{i+1}) \quad (4.66)$$

ce qui peut être coûteux ; soit on utilise la relation inverse de (4.62) :

$$\ddot{\mathbf{x}}_{i+1} = \frac{1}{\beta h^2}(\mathbf{x}_{i+1} - \mathbf{x}_i) - \frac{1}{\beta h}\dot{\mathbf{x}}_i - \frac{1-2\beta}{2\beta}\ddot{\mathbf{x}}_i \quad (4.67)$$

mais on s'aperçoit bien du problème quand $\beta = 0$.

Cas particuliers

- $\gamma = \frac{1}{2}, \beta = \frac{1}{4}$: méthode des accélérations moyennes ;
- $\gamma = \frac{1}{2}, \beta = \frac{1}{6}$: méthode des accélérations linéaires ;
- $\gamma = \frac{1}{2}, \beta = 0$: méthode des différences centrées, qui peut être explicite⁶ ;
- $\gamma = \frac{1}{2}, \beta = \frac{1}{12}$: méthode de Fox-Goodwin.

Ordres

On peut montrer qu'il s'agit d'une méthode d'ordre 1 pour $\gamma \neq \frac{1}{2}$ et d'ordre 2 pour $\gamma = \frac{1}{2}$.

Stabilité

Dans le cas d'une seule équation scalaire, $\ddot{x} = -\omega^2 x$, le schéma donne :

$$\underbrace{\begin{bmatrix} 1 + (h\omega)^2\beta & 0 \\ h\omega^2\alpha & 1 \end{bmatrix}}_{\mathbf{A}_1} \underbrace{\begin{pmatrix} x_{i+1} \\ \dot{x}_{i+1} \end{pmatrix}}_{\mathbf{y}_{i+1}} = \underbrace{\begin{bmatrix} 1 - (h\omega)^2(\frac{1}{2} - \beta) & h \\ -h\omega^2(1 - \alpha) & 1 \end{bmatrix}}_{\mathbf{A}_2} \underbrace{\begin{pmatrix} x_i \\ \dot{x}_i \end{pmatrix}}_{\mathbf{y}_i} \quad (4.68)$$

d'où la relation $\mathbf{y}_{i+1} = \mathbf{A}_1^{-1}\mathbf{A}_2\mathbf{y}_i$ et la matrice d'itération $\mathbf{B} = \mathbf{A}_1^{-1}\mathbf{A}_2$. La condition de stabilité est donc $\rho(\mathbf{B}) < 1$ ou bien $\rho(\mathbf{B}) \leq 1$ avec des valeurs propres simples. L'équation caractéristique permettant de déterminer les valeurs propres de \mathbf{B} est $\lambda^2 - 2b\lambda + c = 0$ avec :

$$2b = 2 - \left(\alpha + \frac{1}{2}\right) \frac{(h\omega)^2}{1 + \beta(h\omega)^2} \quad c = 1 - \left(\alpha - \frac{1}{2}\right) \frac{(h\omega)^2}{1 + \beta(h\omega)^2} \quad (4.69)$$

les conditions de stabilité sont donc :

$$c \leq 1 \quad 1 - 2b + c \geq 0 \quad 1 + 2b + c \geq 0 \quad (4.70)$$

sauf en ce qui concerne les points isolés $(b, c) = (-1, 1)$ et $(1, 1)$. D'où les conditions en terme de coefficients de la méthode :

$$\gamma \geq \frac{1}{2} \quad 2 + \frac{h\omega(2\beta - \gamma)}{1 + \beta(h\omega)^2} \geq 0 \quad (4.71)$$

En conclusion, si $2\beta \geq \gamma \geq \frac{1}{2}$, le schéma est inconditionnellement stable (IS) ; si $\gamma \geq \frac{1}{2}$ et $\gamma > 2\beta$, le schéma est conditionnellement stable (CS) et la condition de stabilité est :

$$h \leq \frac{1}{\omega\sqrt{\gamma/2 - \beta}} \quad (4.72)$$

La figure 4.8 illustre ces résultats.

6. Cette méthode est aussi appelée *Velocity Verlet algorithm* ou *Leap Frog* en dynamique moléculaire !

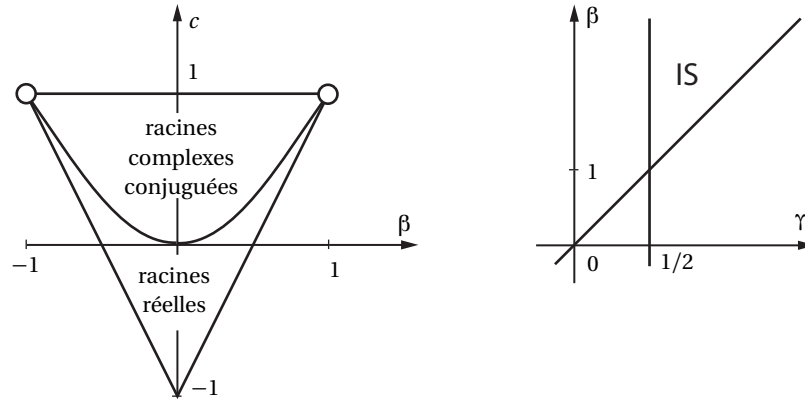


fig. 4.8 - Stabilité des schémas de Newmark

Influence de l'amortissement et comportement hautes fréquences

Considérons le cas où l'amortissement n'est pas négligé dans l'équation scalaire⁷ :

$$\ddot{x} + 2\zeta\omega\dot{x} + \omega^2x = 0 \quad (4.73)$$

On peut montrer dans ce cas que la condition de stabilité s'écrit :

$$\bar{h} \leq \frac{\zeta\left(\gamma - \frac{1}{2}\right) + \left(\frac{1}{2}\gamma - \beta + \zeta^2\left(\gamma - \frac{1}{2}\right)^2\right)^{1/2}}{\frac{1}{2}\gamma - \beta} \frac{1}{\omega} \quad (4.74)$$

Dans le cas où $\gamma > \frac{1}{2}$, l'amortissement a un effet stabilisant puisque sur les conditions de stabilité, on a $\bar{h} < h$.

Concernant les modes à haute fréquence, ils ne sont pas représentatifs du modèle continu mais sont néanmoins sollicités lors de l'intégration et conduisent à la présence de parasites non désirés. On utilise classiquement une *dissipation* numérique pour amortir ces hautes fréquences.

En choisissant $\gamma > \frac{1}{2}$, on prend β pour maximiser l'amortissement (région des racines complexes conjuguées sur la figure 4.8), ce qui conduit à :

$$\beta \leq \frac{1}{4} \left(\gamma + \frac{1}{2} \right)^2 \quad (4.75)$$

Le point de dissipation Haute Fréquence maximal est déterminé par le couple ($\beta = 1$, $\gamma = \frac{3}{2}$), mais il présente le gros inconvénient d'être peu précis en Basse Fréquence. Un choix généralement recommandé pour les méthodes de Newmark est alors :

$$\gamma > \frac{1}{2} \quad \beta = \frac{1}{4} \left(\gamma + \frac{1}{2} \right)^2 \quad (4.76)$$

La figure 4.9 présente l'utilisation des schémas de Newmark sans amortissement pour le cas de réflexions d'ondes dans une poutre unidimensionnelle, ce qui met en évidence la présence de perturbations à hautes fréquences. La figure 4.10 présente quant à elle l'influence de l'amortissement numérique pour ces mêmes schémas.

7. Cette écriture dans la base propre d'un système entier correspond à un modèle d'amortissement de Rayleigh.

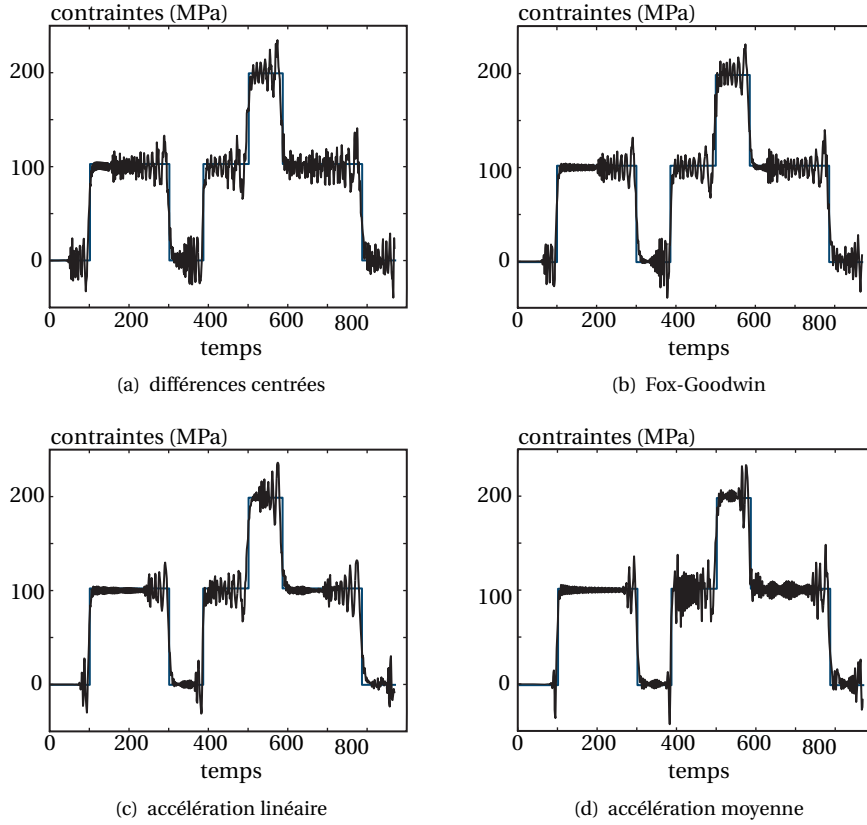


fig. 4.9 - Parasites hautes fréquences (d'après [60])

Toujours pour illustrer l'amortissement numérique, considérons un bilan d'énergie entre les instants t_i et t_{i+1} [61]. L'écriture du schéma de Newmark (4.62) en différence entre ces deux instants, en utilisant la notation suivante :

$$\begin{aligned} \langle \mathbf{x}_i \rangle &= \frac{1}{2}(\mathbf{x}_{i+1} + \mathbf{x}_i) \\ [\mathbf{x}_i] &= \mathbf{x}_{i+1} - \mathbf{x}_i \end{aligned} \quad (4.77)$$

conduit à :

$$\begin{aligned} [\mathbf{x}_i] &= h \langle \dot{\mathbf{x}}_i \rangle + \frac{h^2}{2}(2\beta - \gamma)[\ddot{\mathbf{x}}_i] \\ [\dot{\mathbf{x}}_i] &= h \langle \ddot{\mathbf{x}}_i \rangle + h \left(\gamma - \frac{1}{2} \right) [\ddot{\mathbf{x}}_i] \end{aligned} \quad (4.78)$$

De même, la différence des deux équilibres (4.61) devient :

$$\mathbf{M}[\ddot{\mathbf{x}}_i] + \mathbf{K}[\dot{\mathbf{x}}_i] - [\mathbf{g}_i] = 0 \quad (4.79)$$

En multipliant cette dernière relation par $[\dot{\mathbf{x}}_i]^T$ et en utilisant la propriété $\langle \mathbf{x}_i \rangle^T [\mathbf{x}_i] = [\frac{1}{2}\mathbf{x}_i^2]$, avec des matrices \mathbf{M} et \mathbf{K} positives, on obtient le bilan :

$$\left[\frac{1}{2} \dot{\mathbf{x}}_i^T \mathbf{M}^* \dot{\mathbf{x}}_i \right] + \left[\frac{1}{2} \dot{\mathbf{x}}_i^T \mathbf{K} \dot{\mathbf{x}}_i \right] = \frac{1}{h} [\dot{\mathbf{x}}_i]^T [\mathbf{g}_i] - \left(\gamma - \frac{1}{2} \right) [\ddot{\mathbf{x}}_i]^T \mathbf{M}^* [\ddot{\mathbf{x}}_i] \quad (4.80)$$

où $\mathbf{M}^* = \mathbf{M} + \frac{1}{2}h^2(2\beta - \gamma)\mathbf{K}$. On trouve donc, dans l'ordre, un terme lié à la variation de l'énergie cinétique, un terme lié à la variation de l'énergie interne, un terme lié au travail des

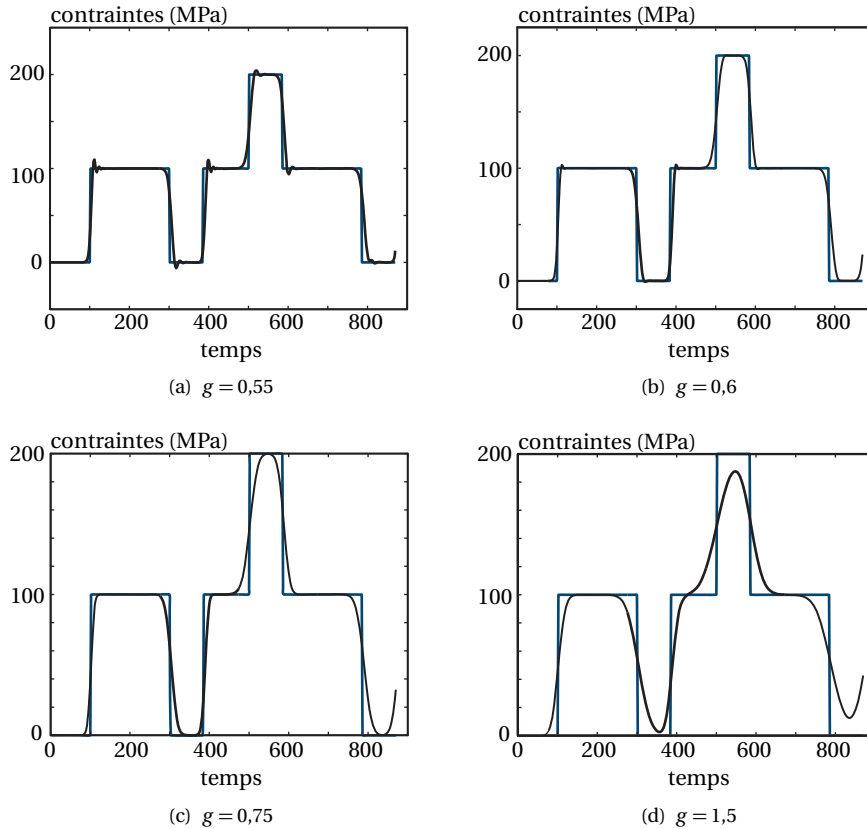


fig. 4.10 - Influence de l'amortissement numérique (d'après [60])

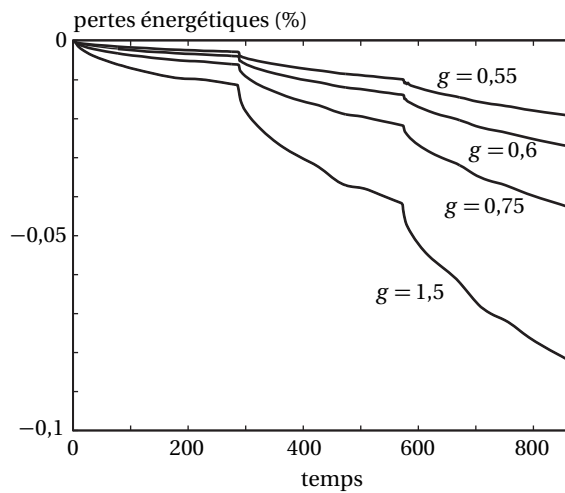


fig. 4.11 - Influence de l'amortissement numérique (suite) (d'après [60])

efforts extérieurs et un dernier terme lié à l'amortissement numérique (dissipation). Cela permet de conclure que si $\gamma \geq 0,5$, alors la solution est d'amplitude bornée et que si $\gamma = 0,5$, le schéma numérique ne dissipe pas d'énergie.

4.3.2 Méthode de Gear implicite à deux pas

Cette méthode [62] proposé en 1969, appelée aussi méthode BDF pour *Backward Differentiation Formulas*, est réputée être très stable. Pour information, elle utilise le développement

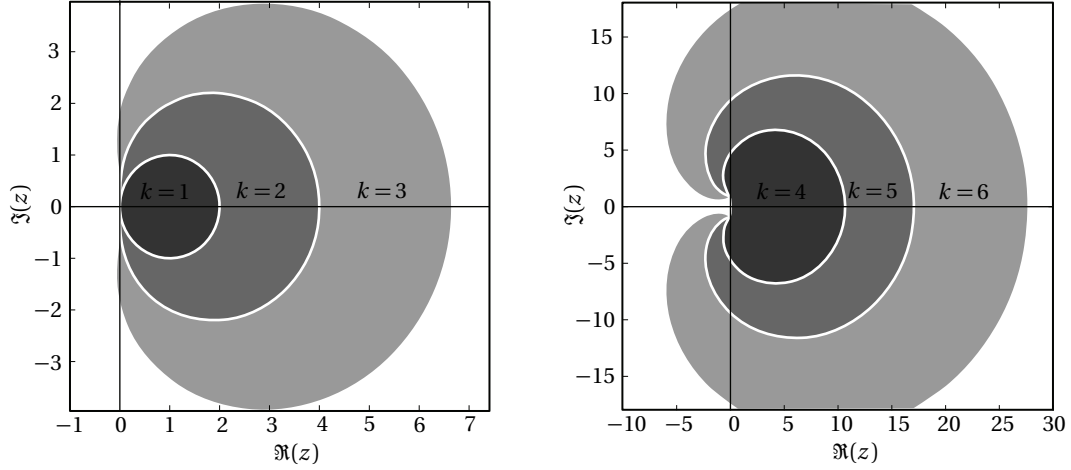


fig. 4.12 - A-stabilité dans le plan complexe pour la méthode de Gear : k est le nombre de pas

suisant :

$$\begin{aligned}\dot{\mathbf{x}}_{i+1} &= \frac{1}{2h}(3\mathbf{x}_{i+1} - 4\mathbf{x}_i + \mathbf{x}_{i-1}) \\ \ddot{\mathbf{x}}_{i+1} &= \frac{1}{2h}(3\dot{\mathbf{x}}_{i+1} - 4\dot{\mathbf{x}}_i + \dot{\mathbf{x}}_{i-1})\end{aligned}\quad (4.81)$$

De la même façon, l'équilibre (4.61) à l'instant t_{i+1} donne le système à résoudre pour atteindre $\ddot{\mathbf{x}}_{i+1}$:

$$\left(\mathbf{M} + \frac{4}{9}h^2\mathbf{K}\right)\ddot{\mathbf{x}}_{i+1} = \mathbf{g}_{i+1} - \mathbf{K}\left(\frac{4}{3}\mathbf{x}_i - \frac{1}{3}\mathbf{x}_{i-1} + \frac{8}{9}h\dot{\mathbf{x}}_i - \frac{2}{9}h\dot{\mathbf{x}}_{i-1}\right)\quad (4.82)$$

La figure 4.12 présente le domaine de stabilité de cette méthode.

4.3.3 θ -méthode

Historiquement attribuée à Wilson en 1968 [63], cette méthode estime une accélération linéaire entre les instants t_i et $t_{i+\theta} = t_i + \theta h$, $\theta \geq 0$ étant un paramètre de la méthode. On a donc :

$$\ddot{\mathbf{x}}(t_i + \tau) = \frac{\ddot{\mathbf{x}}_{i+\theta} - \ddot{\mathbf{x}}_i}{\theta h}\tau + \ddot{\mathbf{x}}_i\quad (4.83)$$

et par intégration, on a les expressions :

$$\begin{aligned}\dot{\mathbf{x}}_{i+\theta} &= \dot{\mathbf{x}}_i + \frac{1}{2}\theta h(\ddot{\mathbf{x}}_i + \ddot{\mathbf{x}}_{i+\theta}) \\ \mathbf{x}_{i+\theta} &= \mathbf{x}_i + \theta h\dot{\mathbf{x}}_i + \theta^2 h^2 \left(\frac{1}{3}\ddot{\mathbf{x}}_i + \frac{1}{6}\ddot{\mathbf{x}}_{i+\theta}\right)\end{aligned}\quad (4.84)$$

Cette fois-ci, l'équilibre (4.61) à l'instant $t_{i+\theta}$ conduit au système en $\mathbf{x}_{i+\theta}$:

$$\left(\mathbf{K} + \frac{6}{\theta^2 h^2}\mathbf{M}\right)\mathbf{x}_{i+\theta} = \mathbf{g}_{i+\theta} + \mathbf{M}\left(\frac{6}{\theta^2 h^2}\mathbf{x}_i + \frac{6}{\theta h}\dot{\mathbf{x}}_i + 2\ddot{\mathbf{x}}_i\right)\quad (4.85)$$

- $\theta \geq 1,37$: la méthode est inconditionnellement stable ;
- $\theta = 1$: méthode des accélérations linéaires.

4.3.4 α -HHT

Cette méthode [64] proposée en 1977, ressemble *a priori* à la méthode de Newmark, à ceci près qu'elle introduit une modification de l'équilibre (et donc en particulier du second membre) en utilisant un paramètre additionnel α .

L'écriture du déplacement, comme celle de la vitesse sont identique à la méthode de Newmark, et l'équilibre est écrit sous la forme :

$$\mathbf{M}\ddot{\mathbf{x}}_{i+1} + ((1 + \alpha)\mathbf{C}\dot{\mathbf{x}}_{i+1} - \alpha\mathbf{C}\dot{\mathbf{x}}_i) + ((1 + \alpha)\mathbf{K}\mathbf{x}_{i+1} - \alpha\mathbf{K}\mathbf{x}_i) = ((1 + \alpha)\mathbf{f}_{i+1} - \alpha\mathbf{f}_i) \quad (4.86)$$

Le démarrage est réalisé avec \mathbf{x}_0 et $\dot{\mathbf{x}}_0$ données ainsi que $\ddot{\mathbf{x}}_0 = \mathbf{M}^{-1}(\mathbf{f}_0 - \mathbf{C}\dot{\mathbf{x}}_0 - \mathbf{K}\mathbf{x}_0)$.

Pour le choix de paramètres de la forme $\beta = \frac{1}{4}(1 - \alpha)^2$ et $\gamma = \frac{1}{2} - \alpha$, le schéma est du second ordre pour les basses fréquences, il est inconditionnellement stable pour $-\frac{1}{3} \leq \alpha \leq 0$, la dissipation en hautes fréquences est maximum pour $\alpha = -\frac{1}{3}$, enfin, si on impose de plus $\alpha = 0$ on retrouve la méthode des trapèzes.

4.3.5 Autres méthodes

Il existe un grand nombre de méthodes d'intégration plus ou moins adaptées à un type de problème. En dynamique, citons encore la méthode de Park à trois pas [65].

4.4 Intégration de modèles de comportement

Si on considère le problème du calcul de structure non-linéaire de la section 3.1.1, on était conduit dans les méthodes de résolution à estimer le résidu $\mathbf{r}(\mathbf{u}) = \mathbf{K}(\mathbf{u})\mathbf{u} - \mathbf{f}$ qui dépend de l'état courant au travers du modèle de comportement du matériau :

$$\mathbf{K}(\mathbf{u})\mathbf{u} = \int_{\Omega} \text{Tr} \boldsymbol{\sigma} \boldsymbol{\varepsilon}(\mathbf{u}) d\Omega \quad (4.87)$$

puisque la contrainte $\boldsymbol{\sigma}$ dépend de l'état courant.

4.4.1 Exemple de modèle de plasticité

Utilisons les notations de la section 3.1.1. Les lois d'état sont de la forme :

$$\begin{aligned} \boldsymbol{\sigma} &= \mathbf{D}(\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}_p) \\ \mathbb{Y} &= \mathbf{G}(\mathbb{X}) \end{aligned} \quad (4.88)$$

où \mathbf{D} est l'opérateur de Hooke et \mathbf{G} la fonction d'écrouissage, caractéristiques du matériau. Dans le cadre d'un modèle associé, le potentiel de dissipation est pris égal au seuil de plasticité, sur lequel s'annule la fonction caractéristique f dépendant de l'état $f(\boldsymbol{\sigma}, \mathbb{Y})$.

Les lois d'évolution s'écrivent avec l'utilisation d'un multiplicateur plastique $\dot{\lambda}$ à cause du caractère non différentiable du potentiel de dissipation :

$$\dot{\boldsymbol{\varepsilon}}_p = \dot{\lambda} \frac{\partial f}{\partial \boldsymbol{\sigma}} \quad (4.89)$$

$$-\dot{\mathbb{X}} = \dot{\lambda} \frac{\partial f}{\partial \mathbb{Y}} \quad (4.90)$$

On a alors les conditions de complémentarité de Khun-Tucker :

$$f \leq 0, \quad \dot{\lambda} \geq 0, \quad f\dot{\lambda} = 0 \quad (4.91)$$

L'expression du multiplicateur est donnée dans le cas de l'écoulement plastique par la condition de cohérence :

$$\dot{f} = 0 = \frac{\partial f}{\partial \sigma} \dot{\sigma} + \frac{\partial f}{\partial \mathbb{Y}} \dot{\mathbb{Y}} \quad (4.92)$$

(dans le cas contraire $\dot{\lambda} = 0$ et le comportement est localement élastique). En remplaçant $\dot{\sigma}$ et $\dot{\mathbb{Y}}$ par leurs expressions issues de la dérivation des lois d'état, et en utilisant les lois d'évolution, on obtient l'expression du multiplicateur :

$$\dot{\lambda} = \frac{\frac{\partial f}{\partial \sigma} D\dot{\varepsilon}}{\frac{\partial f}{\partial \sigma} D\frac{\partial f}{\partial \sigma} + \frac{\partial f}{\partial \mathbb{Y}} \frac{\partial G}{\partial \mathbb{X}} \frac{\partial f}{\partial \mathbb{Y}}} \quad (4.93)$$

On peut donc formellement écrire (à partir des lois d'évolution, de l'expression précédente du multiplicateur et en remplaçant chaque fois que nécessaire σ et \mathbb{Y} par leurs expressions en ε , ε_p et \mathbb{X} issues des lois d'état) l'évolution du comportement sous la forme d'un système d'équations non-linéaires différentielles ordinaires (ODE) :

$$\begin{pmatrix} \dot{\varepsilon}_p \\ -\dot{\mathbb{X}} \end{pmatrix} = g \left(\begin{pmatrix} \varepsilon_p \\ \mathbb{X} \end{pmatrix}, \varepsilon, \dot{\varepsilon} \right) \quad (4.94)$$

Il est à noter que le temps n'intervient pas directement dans l'évolution, contrairement au cas des comportements visqueux. Remarquons aussi que ce système ne peut être résolu : il comporte trop d'inconnues par rapport au nombre d'équations disponibles. En effet, il faut une sollicitation pour que le comportement évolue. En utilisation dans un code aux éléments finis en déplacement, il est agréable de conserver une déformation admissible tout le temps. Le modèle de comportement est donc généralement piloté en intégration par le fait que l'évolution de la déformation totale est imposée. Dans le système précédent, ε et $\dot{\varepsilon}$ sont alors des données.

La résolution de ce système différentiel peut être effectuée avec les techniques décrites dans le chapitre 4. On peut ainsi utiliser des méthodes explicites (avec des sous-pas plus nombreux) ou implicites (avec parfois un seul sous-pas par pas de chargement de la structure, c'est-à-dire par résolution de système différentiel). Classiquement, Runge-Kutta d'ordre 4 avec contrôle de la taille du sous-pas et Euler implicite sont utilisés. En fait, ces méthodes sont utilisées sur une partie du problème à résoudre seulement, car souvent, des méthodes de type retour radial sont employées.

4.4.2 Méthodes de type retour radial

Ces méthodes sont employées car elles sont mieux adaptées au type de problème à traiter. On souhaite en effet retourner assez précisément sur le seuil après intégration du système différentiel. De plus, il faut pouvoir prendre en compte le cas de la décharge élastique. Pour

cela, on prend en compte la propriété de séparation en partie élastique et plastique des modèles de comportement, pour séparer en deux temps l'intégration. Considérons un incrément de charge sur la structure, et un pilotage en incrément de déformation totale, noté $\Delta\varepsilon$. Cette méthode est aussi considérée comme une méthode de partage d'opérateur (*operator split*, cf. prédiction élastique et correction plastique par un retour sur le seuil ou *return map*). La notation δ représente une variation sur un éventuel sous-pas d'intégration. La technique employée consiste en :

prédiction élastique : elle est telle que $\Delta\varepsilon_e + \Delta\varepsilon_p = \Delta\varepsilon$, avec $\Delta\varepsilon_p = 0$ et $\Delta\mathbb{X} = 0$, les lois d'état (4.88) donnent sur cet incrément :

$$\begin{pmatrix} \Delta\sigma \\ \Delta\mathbb{Y} \end{pmatrix} = \begin{bmatrix} \mathbb{D} & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} \Delta\varepsilon \\ 0 \end{pmatrix} \quad (4.95)$$

Si, au bout de cet incrément, le nouvel état $(\sigma + \Delta\sigma, \varepsilon_p, \mathbb{X}, \mathbb{Y} + \Delta\mathbb{Y})$ vérifie $f \leq 0$, on est à l'intérieur du seuil et le comportement est resté élastique. Sinon, il faut corriger la prédiction de la façon suivante.

correction : à partir de l'état déterminé précédemment, on procède à une nouvelle correction (toujours notée avec Δ) $\Delta\varepsilon_e + \Delta\varepsilon_p = 0$ et on cherche à vérifier les équations d'état, les équations d'évolution et, généralement, directement le retour sur le seuil $f = 0$ à la fin de la correction.

Cette phase de correction peut être réalisée par utilisation des méthodes d'intégration déjà citées. On peut aussi utiliser des méthodes plus dédiées comme celle du retour radial d'Ortiz et Simo [66, 57].

Retour radial d'Ortiz et Simo

Pour résoudre la phase de correction, l'idée est d'utiliser une méthode de Newton sur le seuil. On cherche à avoir :

$$f(\sigma + \Delta\sigma, \mathbb{Y} + \Delta\mathbb{Y}) = 0 \quad (4.96)$$

on procède par sous-pas (notés avec un δ) en utilisant la matrice jacobienne des dérivées partielles de f lors des développements successifs :

$$f(\sigma + \delta\sigma, \mathbb{Y} + \delta\mathbb{Y}) \approx f(\sigma, \mathbb{Y}) + \frac{\partial f}{\partial \sigma} \delta\sigma + \frac{\partial f}{\partial \mathbb{Y}} \delta\mathbb{Y} \quad (4.97)$$

pour lesquels on cherche à rester sur le seuil :

$$f + \frac{\partial f}{\partial \sigma} \delta\sigma + \frac{\partial f}{\partial \mathbb{Y}} \delta\mathbb{Y} = 0 \quad (4.98)$$

Comme $\delta\varepsilon = 0$ lors de la correction, les équations (4.89) et (4.90) donnent alors la valeur d'un « multiplicateur plastique » pour le sous-pas :

$$\delta\lambda = \frac{f}{\frac{\partial f}{\partial \sigma} \mathbb{D} \frac{\partial f}{\partial \sigma} + \frac{\partial f}{\partial \mathbb{Y}} \frac{\partial \mathbb{G}}{\partial \mathbb{X}} \frac{\partial f}{\partial \mathbb{Y}}} \quad (4.99)$$

Le potentiel de dissipation étant convexe, le dénominateur de l'expression précédente est toujours positif. Après avoir sommé ce sous-incrément à l'état courant, on itère sur la correction jusqu'à annuler le résidu $f(\sigma, \mathbb{Y})$.

Rigidité tangente cohérente

Il a été constaté que si l'on utilise la matrice de rigidité tangente continue précédente, pour itérer avec une méthode de Newton-Raphson (voir le chapitre 3.1.5), on perd la convergence quadratique de la méthode, et que pour la retrouver, il faut utiliser la matrice de rigidité tangente cohérente avec le schéma d'intégration employé.

Ce dernier peut être vu comme un outil pour produire, à partir d'un état donné :

$$s = (\sigma, \mathbb{Y}, \varepsilon_p, \mathbb{X}) \quad (4.100)$$

un nouvel état satisfaisant au comportement du matériau, piloté par exemple par un incrément de déformation totale $\Delta\varepsilon$. Ce nouvel état est noté :

$$s + \Delta s = (\sigma + \Delta\sigma, \mathbb{Y} + \Delta\mathbb{Y}, \varepsilon_p + \Delta\varepsilon_p, \mathbb{X} + \Delta\mathbb{X}) \quad (4.101)$$

Si la rigidité tangente continue à l'état s est D_T telle que $\dot{\sigma} = D_T \dot{\varepsilon}$, la matrice tangente cohérente D_C est définie par $\delta\Delta\sigma = D_C \delta\Delta\varepsilon$. Elle relie la variation de correction de contrainte obtenue, à la variation de pilotage en incrément de déformation totale : piloter avec $\Delta\varepsilon + \delta\Delta\varepsilon$ à partir du même état s conduira à un nouvel état suivant $s + \delta s + \delta\Delta s$.

On peut obtenir D_C par une méthode de perturbation (on procède à plusieurs intégrations avec des perturbations différentes $\delta\Delta\varepsilon$ sur le pilotage) mais cela est assez coûteux numériquement, et nécessite le choix de la valeur de la perturbation (trop petite, elle conduit à des problèmes de précision numériques, trop grande, à une mauvaise approximation de la dérivée). Une autre technique consiste à dériver analytiquement le schéma d'intégration, ce qui peut s'avérer assez pénible suivant le schéma et le modèle de comportement.

4.4.3 Exemple de modèle de viscoplasticité

Par rapport au modèle précédente, un grand nombre de modèles de viscoplasticité s'écrivent en n'imposant plus de signe sur la valeur de la fonction seuil $f(\sigma, \mathbb{Y})$, mais avec un terme de rappel sur le multiplicateur viscoplastique $\dot{\lambda}$: les conditions de complémentarité de Kuhn-Tucker (4.91) sont alors remplacées par :

$$\dot{\lambda} = \left(\left\langle \frac{f(\sigma, \mathbb{Y})}{m} \right\rangle_+ \right)^n \quad (4.102)$$

où les parties positive et négative sont respectivement notées :

$$\langle x \rangle_+ = \max(0, x) = \frac{1}{2}(x + |x|) \quad \text{et} \quad \langle x \rangle_- = \min(0, x) = \frac{1}{2}(x - |x|) = -\langle -x \rangle_+ \quad (4.103)$$

m est un module de viscoplasticité et n un exposant de viscosité.

Il est possible de définir l'équivalent du seuil de plasticité de la façon suivante : si on considère $\bar{f} = f - \langle f \rangle_+$, alors $\bar{f} = 0$ correspond à $f \geq 0$ c'est-à-dire à un écoulement viscoplastique, et $\bar{f} < 0$ correspond à $f \leq 0$ c'est-à-dire à une évolution élastique, et dans ce cas $\bar{f} = f$.

On peut aussi se ramener à une relation de complémentarité : si on note $\mu = \lambda$, on doit en effet avoir $\mu \geq 0$, $\tilde{f} = f - m\mu^{1/n} \leq 0$ et $\mu\tilde{f} = 0$. La principale différence par rapport au cas de la plasticité est la dépendance avec le temps au travers des vitesses des quantités cherchées. L'intégration de la loi de comportement ne sera plus pilotée par le seul incrément de déformation $\Delta\varepsilon$, mais aussi par l'incrément de temps Δt . Le principe consiste encore à prédire une évolution élastique :

$$\begin{aligned}\Delta\sigma &= D\Delta\varepsilon, \Delta\mathbb{Y} = 0 \\ \Delta\mathbb{X} &= 0, \Delta\varepsilon_p = 0\end{aligned}\tag{4.104}$$

et à estimer la valeur du seuil $f(\sigma, \mathbb{Y})$. Si celle-ci est négative, la solution est effectivement élastique. Sinon, il faut procéder à une correction viscoplastique ; cette étape est encore itérative, mais il faut bien réévaluer toutes les quantités (en particulier les dérivées par l'intermédiaire du schéma d'intégration en temps choisi, ici la θ -méthode) par rapport au précédent pas de temps t_i (elle sont ici indicées par i). La correction d'effectue toujours à déformation totale nulle, et l'estimation du seuil corrigé est :

$$\tilde{f}(\sigma + \delta\sigma, \mathbb{Y} + \delta\mathbb{Y}, \mu + \delta\mu) \approx \tilde{f}(\sigma, \mathbb{Y}, \mu) + \frac{\partial \tilde{f}}{\partial \sigma} \delta\sigma + \frac{\partial \tilde{f}}{\partial \mathbb{Y}} \delta\mathbb{Y} + \frac{\partial \tilde{f}}{\partial \mu} \delta\mu\tag{4.105}$$

pour lesquels on cherche à rester sur le seuil corrigé :

$$f(\sigma, \mathbb{Y}) - m\mu^{1/n} + \frac{\partial f}{\partial \sigma} \delta\sigma + \frac{\partial f}{\partial \mathbb{Y}} \delta\mathbb{Y} - \frac{m}{n} \mu^{1/n-1} \delta\mu = 0\tag{4.106}$$

avec $\delta\sigma = -D\delta\varepsilon_p$, $\delta\mathbb{Y} = \frac{\partial G}{\partial \mathbb{X}} \delta\mathbb{X}$ et les incréments :

$$\begin{aligned}\frac{\delta\varepsilon_p}{\theta\Delta t} &= (\mu + \delta\mu) \frac{\partial f}{\partial \sigma} - \dot{\varepsilon}_p \\ \frac{\delta\mathbb{X}}{\theta\Delta t} &= -(\mu + \delta\mu) \frac{\partial f}{\partial \mathbb{Y}} - \dot{\mathbb{X}}\end{aligned}\tag{4.107}$$

où les dérivées sont estimées par rapport au précédent pas de temps (et doivent être réévaluée à chaque itéré de la correction) :

$$\begin{aligned}\dot{\varepsilon}_p &= \frac{\varepsilon_p - \varepsilon_{pi}}{\theta\Delta t} - \frac{1-\theta}{\theta} \dot{\varepsilon}_{pi} \\ \dot{\mathbb{X}} &= \frac{\mathbb{X} - \mathbb{X}_i}{\theta\Delta t} - \frac{1-\theta}{\theta} \dot{\mathbb{X}}_i\end{aligned}\tag{4.108}$$

L'équation (4.105) permet alors d'obtenir l'incrément de multiplicateur viscoplastique :

$$\delta\mu = \frac{\frac{f - m\mu^{1/n}}{\theta\Delta t} + \frac{\partial f}{\partial \sigma} D\dot{\varepsilon}_p - \frac{\partial f}{\partial \mathbb{Y}} \frac{\partial G}{\partial \mathbb{X}} \dot{\mathbb{X}}}{\frac{\partial f}{\partial \sigma} D \frac{\partial f}{\partial \sigma} + \frac{\partial f}{\partial \mathbb{Y}} \frac{\partial G}{\partial \mathbb{X}} \frac{\partial f}{\partial \mathbb{Y}} + \frac{m}{n} \mu^{1/n-1} \frac{1}{\theta\Delta t}}\tag{4.109}$$

Ayant les incréments donnés $\Delta\varepsilon$ et Δt ,
prédiction élastique avec (4.104)
prédiction du seuil $f(\sigma, \mathbb{Y})$
si $f > 0$ alors
correction viscoplastique
 $\mu = 0, \tilde{f} = f$
tant que $\tilde{f} > 0$
estimer $\dot{\varepsilon}_p$ et $\dot{\mathbb{X}}$ avec (4.108)
calculer la correction $\delta\mu$ avec (4.109)
mise à jour de $\sigma, \mathbb{Y}, \varepsilon_p, \mathbb{X}$ avec (4.107)
mise à jour du multiplicateur $\mu \leftarrow \mu + \delta\mu$
mise à jour du seuil corrigé \tilde{f}
fin
fin

alg. 4.13 - Intégration du comportement viscoplastique

La valeur courante du multiplicateur doit alors aussi être mise à jour par $\mu + \delta\mu$, à chaque itéré de la correction, qui est poursuivie jusqu'à retourner sur le seuil corrigé \tilde{f} .

Connaissant la solution à t_i , on procède à une succession de prédictions élastiques jusqu'à rééquilibrer les forces extérieures (c'est-à-dire à annuler le résidu des équations d'équilibre) et d'intégration du comportement selon l'algorithme décrit dans 4.13. La solution à t_{i+1} est celle qui doit vérifier à la fois l'équilibre et le comportement.

4.5 Cas de la dynamique non régulière

Dans les approches précédentes, on a souvent considéré que les évolutions temporelles étaient suffisamment régulières pour pouvoir les approcher avec un développement limité. Ce n'est pas forcément toujours le cas ; par exemple dans le cas de la fissuration dynamique, lorsque les structures sont soumises à des chocs ou impacts. Dans ce dernier cas, en particulier quand les corps sont des solides rigides, les vitesses ne sont plus dérivables (lors d'un choc, un corps rigide change peu en position, mais beaucoup en vitesse, par exemple dans le cas d'un rebond). Une technique (dite *event driven*) consiste à traiter de façon standard les phases de *vol libre* c'est-à-dire sans impact, à repérer les instants de choc, et à traiter ceux-ci de façon particulière, avec un saut de vitesse de part et d'autre de l'instant de choc.

Avec une modélisation en solides rigides, il manque une équation de comportement de l'impact pour fermer le problème. Un tel comportement peut être modélisé par la loi de Newton⁸. Elle postule que la vitesse normale du point de contact — d'où des difficultés pour modéliser un impact sur une surface entière — après impact (notée v_n^+) est opposée en direction à celle avant impact (notée v_n^-) :

$$v_n^+ = -e v_n^- \quad (4.110)$$

le saut de vitesse lors de l'impact étant $v_n^+ - v_n^-$. Le coefficient *matériau* e , appelé coefficient de restitution d'énergie, traduit le comportement de l'impact ; $0 \leq e \leq 1$. Le problème est qu'il dépend de beaucoup de paramètres du problème, puisqu'il représente de façon phénoménologique le comportement complexe des propagations / réflexions d'ondes lors de l'impact, ou des impacts multiples...

Sans présence de frottement, il n'est pas nécessaire d'écrire autre choses sur les glissements tangentiels lors de l'impact. On se contentera de décrire ce type de problème ici.

Une autre situation particulière est celle des chocs nombreux dans les matériaux granulaires ou soumis à la fragmentation... Dans ce cas, une technique *event driven* est inapplicable à cause des impacts qui peuvent être extrêmement nombreux. Une alternative consiste à employer une technique de *time-stepping*. Le pas de temps est fixé et on s'emploie à déterminer les vitesses immédiatement après les piquets de temps, en considérant l'impact *moyen* pendant le pas de temps courant. On ne s'intéresse alors pas vraiment à l'évolution du système continûment sur le pas de temps, mais plutôt à une configuration

8. Elle masque le fait que le solide est quand même légèrement déformable et que le rebond est dû à un échange entre l'énergie cinétique et l'énergie de déformation, lors de la propagation des ondes de chocs dans le solide

possible en chaque piquet de temps.

Considérons tout d'abord un système dynamique régulier (sans choc), rigide, évoluant sous l'action d'une force extérieure f . Son équation du mouvement est de la forme :

$$m\ddot{x} = f \quad (4.111)$$

où l'on suppose m constante. On se propose d'intégrer cette équation différentielle avec un schéma particulier en diminuant l'ordre du système⁹. Pour cela, commençons par intégrer l'équation du mouvement sur un pas de temps $[t_i, t_{i+1}]$; on obtient :

$$m(\dot{x}_{i+1} - \dot{x}_i) = \int_{t_i}^{t_{i+1}} f(t) dt \quad (4.112)$$

L'intégrale de la force extérieure (régulière) peut être approximée par un schéma d'intégration, par exemple la méthode des trapèzes généralisés, voir section 4.1.6. Le problème s'écrit alors :

$$m\dot{x}_{i+1} = m\dot{x}_i + (1 - \alpha)hf_i + \alpha hf_{i+1} \quad (4.113)$$

La position est récupérée par une intégration similaire de la vitesse :

$$x_{i+1} = x_i + \int_{t_i}^{t_{i+1}} \dot{x} dt \approx x_i + (1 - \alpha)h\dot{x}_i + \alpha h\dot{x}_{i+1} \quad (4.114)$$

Le lecteur pourra vérifier qu'il s'agit exactement de l'application de la méthode des trapèzes généralisés au système du premier ordre :

$$\begin{aligned} m\dot{v} &= f \\ \dot{x} &= v \end{aligned} \quad (4.115)$$

Cette façon de procéder va permettre de traiter les systèmes non réguliers comprenant des chocs. En effet, l'équation de la dynamique entre deux instants immédiatement suivant les piquets de temps t_i et t_{i+1} va s'écrire :

$$m(\dot{x}_{i+1} - \dot{x}_i) = \int_{t_i}^{t_{i+1}} f(t) dt + R \quad (4.116)$$

où R est l'impulsion totale, c'est-à-dire, la somme des impulsions lors des impacts éventuels sur le pas de temps traité.

Dans le cas d'une évolution statique, le contact unilatéral du système contre un mur rigide placé en $x = 0$ par exemple ferait intervenir un effort F et non pas une impulsion R , avec une loi de comportement du contact parfait du type $0 \leq x \perp F \geq 0$, qui est une version concise de $x \geq 0, F \geq 0, x \cdot F = 0$. Le pendant de ce comportement en dynamique est la condition de Signorini-vitesse issue du Lemme de Moreau :

- si $x > 0, R = 0$
- si $x = 0, 0 \leq \dot{x} \perp R \geq 0$

9. On travaillera alors à la fois sur les positions x_i et les vitesses \dot{x}_i aux piquets de temps t_i .

Pour dicrétiser en temps ce comportement, on utilise en général une version explicite sur la configuration — c'est-à-dire qu'on va utiliser un prédicteur pour la position x qui permettra de sélectionner le cas correspondant, et qu'on notera x^p :

- si $x^p > 0$, $R_{i+1} = 0$
- si $x^p \leq 0$, $0 \leq \dot{x}_{i+1} \perp R_{i+1} \geq 0$

Un prédicteur classique (saute-mouton) est $x^p = x_i + \frac{h}{2}\dot{x}_i$. L'inconvénient est la possibilité d'avoir une pénétration résiduelle, d'où l'inégalité utilisée sur les valeurs possiblement négatives de x^p , qui tend cependant vers 0 avec le pas de temps. Ce modèle de comportement correspond à un impact plastique (sans rebond), qui dissipe le plus l'énergie. Dans le cas de choc unique, on peut retrouver l'équivalent du modèle de Newton en substituant \dot{x}_{i+1} par une vitesse dite *formelle* :

$$\tilde{v} = \frac{\dot{x}_{i+1} - \dot{x}_i}{1 + e} \tag{4.117}$$

RÉFÉRENCES

- [1] B. Mohammadi and J.-H. Saiac. *Pratique de la simulation numérique*. Dunod / Industries et Technologies, 2003. ISBN 210006407X. 1
- [2] G.H. Golub and C.F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, Baltimore, MD, USA, 2^e edition, 1989. ISBN 0-8018-3772-3, 0-8018-3739-1. 2
- [3] K. J. Bathe and E. L. Wilson. *Numerical Analysis in Finite Element Analysis*. Prentice-Hall, 1982. 2
- [4] E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen. *LAPACK Users' Guide*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 3^e edition, 1999. ISBN 0-89871-447-8 (paperback). 2
- [5] L.S. Blackford, J. Choi, A. Cleary, E. d'Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R.C. Whaley. *ScaLAPACK Users' Guide*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 1997. ISBN 0-89871-397-8 (paperback). 3
- [6] H. Oudin. *Méthode des éléments finis*. École Centrale de Nantes, France, 2008. <http://cel.archives-ouvertes.fr/docs/00/34/17/72/PDF/MEF.pdf>. 5
- [7] Y. Saad. *Iterative Methods for Sparse Linear Systems*. 2000. 5
- [8] R. Barrett, M. Berry, T. F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, and H. Van der Vorst. *Templates for the Solution of Linear Systems : Building Blocks for Iterative Methods*. SIAM, 1994. http://www.netlib.org/linalg/html_templates/Templates.htm. 5
- [9] J. J. Dongarra, H. W. Meuer, and E. Strohmaier. TOP 500 supercomputer sites. In *Proceedings of Supercomputer '99 Conference*, number 17, Mannheim, 1999. Disponible sur <http://www.top500.org/> et <http://www.netlib.org/benchmark/top500.html>. 8
- [10] H. W. Meuer, E. Strohmaier, J. J. Dongarra, and H. Simon. Top500 supercomputer sites. In *16th International Supercomputer Conference : Applications — Architectures — Trends*, 2001. <http://www.netlib.org/benchmark/top500.html>, <http://www.top500.org/>. 8
- [11] E. Cuthill and J. McKee. Reducing the bandwidth of sparse symmetric matrices. In *Proc. of the 24th National Conference of the ACM*, volume P-69, pages 157–172, 1969. 9
- [12] A. George. Nested dissection of a regular finite-element mesh. *SIAM Journal on Numerical Analysis*, 10 :345–363, 1973. doi:10.1137/0710032. 9
- [13] W. F. Tinney and J. W. Walker. Direct solutions of sparse network equations by optimally ordered triangular factorization. In *Proc. of the IEEE*, volume 55, pages 1801–1809, 1967. 9
- [14] A. George and J. Liu. The evolution of the minimum degree ordering algorithm. *SIAM Review*, 31(1) :1–19, 1989. 9
- [15] I. S. Duff, A. M. Erisman, and J. K. Reid. *Direct Methods for Sparse Matrices*. Oxford University Press, Oxford, 1990. 10

- [16] Comm. Benoît. Note sur une méthode de résolution des équations normales provenant de l'application de la méthode des moindres carrés à un système d'équations linéaires en nombre inférieur à celui des inconnues — application de la méthode à la résolution d'un système d'équations linéaires (procédé du commandant cholesky). *Bulletin Géodésique*, 2 :5–77, 1924. 12
- [17] B. M. Irons. A frontal solution program for finite element analysis. *International Journal for Numerical Methods in Engineering*, 2(1) :5–32, 1970. doi:10.1002/nme.1620020104. 14
- [18] I. S. Duff. Parallel implementation of multifrontal schemes. *Parallel Computing*, 3 : 192–204, 1986. 14
- [19] C. Farhat and D. Rixen. *Encyclopedia of Vibration*, chapter Linear Algebra, pages 710–720. Academic Press, 2002. 17
- [20] L. Hagemana and D. Young. *Applied Iterative Methods*. Academic Press, New York, 1981. 18
- [21] R. Southwell. *Relaxation Methods in Theoretical Physics*. Clarendon Press, Oxford, 1946. 19
- [22] M. R. Hestenes and E. Stiefel. Methods of conjugate gradients for solving linear systems. *Journal of research of the National Bureau of Standards*, 49 :409–436, 1952. 22
- [23] G. Shortley. Use of Tschebyscheff-polynomial operators in the numerical solution of boundary-value problems. *Journal of Applied Physics*, 24 :392–396, 1953. doi:10.1063/1.1721292. 23
- [24] C. C. Paige and M. A. Saunders. Solution of sparse indefinite systems of linear equations. *SIAM Journal on Numerical Analysis*, 12 :617–629, 1975. doi:10.1137/0712047. 23
- [25] R. W. Freund and N. M. Nachtigal. QMR : A quasi-minimal residual method for nonhermitian matrices. *Numerische Mathematik*, 60 :315–339, 1991. 23
- [26] P. Sonneveld. CGS, a fast Lanczos-type solver for nonsymmetric linear systems. *SIAM Journal on Scientific and Statistical Computing*, 10 :36–52, 1989. 23
- [27] H. A. Van der Vorst. BiCGSTAB : A fast and smoothly converging variant of Bi-CG for the solution of nonsymmetric linear systems. *SIAM Journal on Scientific and Statistical Computing*, 13 :631–644, 1992. 23
- [28] Y. Saad and M. H. Schultz. GMRES : A generalized minimal residual algorithm for solving nonsymmetric linear systems. *SIAM Journal on Scientific and Statistical Computing*, 7 :856–869, 1986. 23
- [29] J. S. Rao. *Advanced Theory of Vibration*. John Wiley & Sons, 1989. 25
- [30] M. Géradin and D. Rixen. *Théorie des Vibrations. Application à la dynamique des structures*. Physique Fondamentale et Appliquée. Masson, 1993. 25
- [31] M. Géradin and D. Rixen. *Mechanical Vibrations. Theory and Application to Structural Dynamics*. Wiley, 1994. 25
- [32] J. H. Wilkinson. *The Algebraic Eigenvalue Problem*. Oxford University Press, 1965. ISBN 0-19-853403-5, 0-19-853418-3. 25
- [33] Z. Bai, J. Demmel, J. Dongarra, A. Ruhe, and H. van der Vorst, editors. *Templates for the Solution of Algebraic Eigenvalue Problems : A Practical Guide*. 1999. Draft. 25
- [34] C. G. J. Jacobi. Ueber ein leichtes verfahren, die in der theorie der säcularstörungen vorkommenden gleichungen numerish aufzulösen. *Journal für die reine und angewandle Mathematik*, pages 51–94, 1846. 27
- [35] R. A. Willoughby. Collection of articles honoring Alston S. Householder. *ACM*, 18 : 3–58, 1975. 31

-
- [36] A. S. Householder. Unitary triangularization of a nonsymmetric matrix. *Journal of the Association of Computing Machinery*, 5 :339–342, 1958. 31
- [37] G. Hessenberg. *Transzendenz von e und ; ein Beitrag zur höheren Mathematik vom elementaren Standpunkte aus. von 1912*. Johnson Reprint Corp., New York, 1965. 33
- [38] C. Lanczos. An iteration method dor the solution of the eigenvalue problem of linear differential and integral operators. *Journal of research of the National Bureau of Standards*, 45 :225–280, 1950. 35
- [39] C. T. Kelley. *Solving Nonlinear Equations with Newton's Method*, volume 1 of *Fundamentals of Algorithms*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2003. ISBN 0-89871-546-6. 37
- [40] M. A. Crisfield. *Non-Linear Finite Element Analysis of Solids and Structures*. John Wiley, 1991. 37
- [41] J. Lemaitre and J.-L. Chaboche. *Mechanics of Solid Materials*. Cambridge University Press, 1994. 38
- [42] R. Fletcher and C. Reeves. Function minimization by conjugate gradients. *Comput. J.*, 7 :149–154, 1964. 39
- [43] E. Polak and G. Ribière. Note sur la convergence de directions conjuguées. *Rev. Française Informat. Recherche Opérationnelle*, (16) :35–43, 1969. 39
- [44] Y. Dai and Y. Yuan. Global convergence of the method of shortest residuals. *Numerische Mathematik*, 83(4) :581–598, 1999. doi:10.1007/s002119900080. 39
- [45] W. J. Duncan. Some devices for the solution of large sets of simultaneous linear equations. *Philos. Mag. Ser.*, 35(7) :660–670, 1944. 40
- [46] W. W. Hager. Updating the inverse of a matrix. *SIAM Review*, 31(2) :221–239, 1989. 40
- [47] C. Huang and G. Verchery. An exact structural static reanalysis method. *Communications in Numerical Methods in Engineering*, 13 :103–112, 1997. 40
- [48] R. Fletcher and M. J. D. Powell. A rapidly convergent descent method for minimization. *Computer Journal*, 6 :163–168, 1963. 40
- [49] C. G. Broyden. The convergence of a class of double-rank minimization algorithms. *Journal of the Institute of Mathematics and its Applications*, 6 :76–90, 1970. 41
- [50] L. Nazareth. A relationship between BFGS and conjugate gradient algorithms and its implications for new algorithms. *SIAM Journal on Numerical Analysis*, 16 :794–800, 1979. doi:10.1137/0716059. 41
- [51] T. G. Kolda, D. P. O’Leary, and L. Nazareth. BFGS with update skipping and varying memory. *SIAM Journal on Optimization*, 8(4) :1060–1083, 1998. 41
- [52] D.C. Liu and J. Nocedal. On the limited memory BFGS method for large scale optimization. *Numerische Mathematik*, 45 :503–528, 1989. 41
- [53] T. J. R. Hughes. *The Finite Element Method, Linear Static and Dynamic Finite Element Analysis*. Prentice-Hall, Englewood Cliffs, 1987. 43
- [54] D. W. Jordan and P. Smith. *Nonlinear Ordinary Differential Equations*. Oxford applied mathematics and computing science series, 1987. 43
- [55] G. Dahlquist. A special stability problem for linear multistep methods. *BIT. Numerical Mathematics*, 3(1) :27–43, 1963. doi:10.1007/BF01963532. 46
- [56] W. Kutta. Beitrag zur naherungsweise integration totaler differentialgleichunge. *Zeit. Math. Physik*, 46 :435–453, 1901. 51
- [57] M. Ortiz and J. C. Simo. An analysis of a new class of integration algorithms for elasto-plastic constitutive relations. *International Journal for Numerical Methods in Engineering*, 23(3) :353–366, 1986. doi:10.1002/nme.1620230303. 52, 63

- [58] J.-P. Combe. *Sur le contrôle des calculs en dynamique rapide. Application aux problèmes d'impact*. PhD thesis, ENS de Cachan, LMT-Cachan, 2000. 54
- [59] N. M. Newmark. A method for computation of structural dynamics. In EMS, editor, *Proc. of the American Society of Civil Engineers*, volume 85, pages 67–94, 1959. 55
- [60] H. Lemoussu. *Approche non incrémentale pour le calcul de choc avec contact unilatéral avec frottement*. PhD thesis, ENS de Cachan, LMT-Cachan, 1999. 58, 59
- [61] T. J. R. Hughes and T. Belytschko. *Nonlinear Finite Element Analysis*. Zace sevice Ltd - ICE Division, 1995. 58
- [62] C. W. Gear. *Numerical Initial Value Problems in Ordinary Differential Equation*. Prentice-Hall, 1971. 59
- [63] K. J. Bathe and E. L. Wilson. Stability and accuracy analysis of direct integration method. *Earthquake Eng. Struct. Dyn.*, 1 :283–291, 1973. 60
- [64] H. M. Hilber, T. J. R. Hughes, and R. L. Taylor. Improved numerical dissipation for time integration algorithm in structural dynamics. *Earthquake Eng. Struct. Dyn.*, 5 : 283–292, 1977. 61
- [65] K. C. Park. An improved stiffly stable method for direct integration of non linear structural dynamics equations. *Journal of Applied Mechanics*, pages 164–470, 1975. 61
- [66] M. Ortiz and E. P. Popov. Accuracy and stability of integration algorithms for elastoplastic constitutive relations. *International Journal for Numerical Methods in Engineering*, 21(9) :1561–1576, 1985. doi:10.1002/nme.1620210902. 63
- [67] A. C. Aitken. On the iterative solution of a system of linear equations. *Proceedings of the Royal Society of Edinburgh*, 63 :52–60, 1950. 77



FACTORISATION DE CHOLESKY - ORTHOOGONALISATION D'UN SOUS-ESPACE

On considère un ensemble de m vecteurs de taille n \mathbf{v}_i indépendants, stockés dans les colonnes d'une matrice rectangulaire \mathbf{V} de taille (n, m) . Ils génèrent un sous-espace de l'espace de départ. Dans un certain nombre d'applications, il est utile de construire à partir de ces vecteurs, une nouvelle base orthogonale \mathbf{U} . Pour ce faire, on utilise généralement le procédé d'orthogonalisation de Gram-Schmidt. On cherche donc \mathbf{U} engendrant le même sous-espace que \mathbf{V} et tel que $\mathbf{U}^T \mathbf{U} = \mathbf{I}$.

Pour ce faire, on peut utiliser une technique basée sur la factorisation précédente : si on construit la matrice $\mathbf{M} = \mathbf{V}^T \mathbf{V}$, pleine et de taille réduite (n, n) et sa factorisation $\mathbf{M} = \mathbf{L} \mathbf{L}^T$, alors :

$$\mathbf{V}^T \mathbf{V} = \mathbf{L} \mathbf{L}^T \Rightarrow \mathbf{L}^{-1} \mathbf{V}^T \mathbf{V} \mathbf{L}^{-T} = \mathbf{I} \Rightarrow \mathbf{U}^T = \mathbf{L}^{-1} \mathbf{V}^T \quad (\text{A.1})$$

Il faut donc procéder à n résolutions du système triangulaire de taille n :

$$\mathbf{L} \mathbf{U}^T = \mathbf{V}^T \quad (\text{A.2})$$

On réalise ainsi une orthonormalisation des vecteurs de départ.

B

BORNAGE PAR LES VALEURS PROPRES ÉLÉMENTAIRES

On se place ici dans le cas de figure où les problèmes sont issus d'une analyse par éléments finis. L'élément fini courant est noté e et les quantités élémentaires associées porteront aussi l'indice e .

B.1 Système aux valeurs propres généralisé

Le problème est de la forme $\mathbf{K}\mathbf{x} = \lambda\mathbf{M}\mathbf{x}$, et on définit les matrices diagonales par bloc contenant les quantités élémentaires de la façon suivante :

$$\bar{\mathbf{K}} = \begin{bmatrix} \mathbf{K}_e & & \\ & \mathbf{K}_{e'} & \\ & & \ddots \end{bmatrix} \quad \text{et} \quad \bar{\mathbf{M}} = \begin{bmatrix} \mathbf{M}_e & & \\ & \mathbf{M}_{e'} & \\ & & \ddots \end{bmatrix} \quad (\text{B.1})$$

L'assemblage consiste à utiliser formellement une matrice d'assemblage \mathbf{P} booléenne telle que :

$$\mathbf{K} = \mathbf{P}\bar{\mathbf{K}}\mathbf{P}^T \quad (\text{B.2})$$

Comme \mathbf{M} est issue de la même discrétisation, la même matrice d'assemblage est utilisée :

$$\mathbf{M} = \mathbf{P}\bar{\mathbf{M}}\mathbf{P}^T \quad (\text{B.3})$$

Considérons le système, de taille supérieure à n :

$$\bar{\mathbf{K}}\bar{\mathbf{x}} = \tilde{\lambda}\bar{\mathbf{M}}\bar{\mathbf{x}} \quad (\text{B.4})$$

Il englobe le système initial car il suffit de l'écrire dans le sous-espace où $\bar{\mathbf{x}}$ est issu du *désassemblage* d'un vecteur \mathbf{x} de taille n : $\bar{\mathbf{x}} = \mathbf{P}^T\mathbf{x}$ ¹. Avec le quotient de Rayleigh :

$$\bar{R}(\bar{\mathbf{x}}) = \frac{\bar{\mathbf{x}}^T \bar{\mathbf{K}} \bar{\mathbf{x}}}{\bar{\mathbf{x}}^T \bar{\mathbf{M}} \bar{\mathbf{x}}} \quad (\text{B.5})$$

on sait que :

$$\tilde{\lambda}_{\max} = \max_{\bar{\mathbf{x}} \neq 0} \bar{R}(\bar{\mathbf{x}}) \quad (\text{B.6})$$

1. Les forces généralisées, elles, s'assembleraient de la façon suivante : $\mathbf{f} = \mathbf{P}\bar{\mathbf{f}}$.

donc :

$$\bar{\lambda}_{\max} \geq \frac{\bar{\mathbf{x}}^T \bar{\mathbf{K}} \bar{\mathbf{x}}}{\bar{\mathbf{x}}^T \bar{\mathbf{M}} \bar{\mathbf{x}}} \quad (\text{B.7})$$

Au vu du problème *désassemblé*, $\bar{\lambda}_{\max}$ est la valeur maximale des valeurs propres des éléments pris séparément et la relation (B.7) est en particulier vraie pour $\bar{\mathbf{x}} = \mathbf{P}^T \mathbf{x}$ donc :

$$\bar{\lambda}_{\max} \geq \frac{\bar{\mathbf{x}}^T \mathbf{P} \bar{\mathbf{K}} \mathbf{P}^T \bar{\mathbf{x}}}{\bar{\mathbf{x}}^T \mathbf{P} \bar{\mathbf{M}} \mathbf{P}^T \bar{\mathbf{x}}} = \frac{\mathbf{x}^T \mathbf{K} \mathbf{x}}{\mathbf{x}^T \mathbf{M} \mathbf{x}} = R(\mathbf{x}) \quad (\text{B.8})$$

donc :

$$\bar{\lambda}_{\max} \geq \max_{\mathbf{x} \neq 0} R(\mathbf{x}) = \lambda_{\max} \quad (\text{B.9})$$

B.2 Système aux valeurs propres

Cette fois-ci, le problème est de la forme $\mathbf{A} \mathbf{x} = \lambda \mathbf{x}$. Concernant les contributions élémentaire, ce problème n'est plus équivalent au problème précédent ; en particulier, si on écrit $\mathbf{A} = \mathbf{M}^{-1} \mathbf{K}$ et comme \mathbf{M}^{-1} est pleine, on n'a plus la connectivité éléments finis dans \mathbf{A} , et si on écrit $\mathbf{M} = \mathbf{I}$, \mathbf{M} n'est pas l'assemblage des matrices élémentaires identité. Le résultat précédent ne s'applique donc plus ici. Par contre, on a toujours :

$$\lambda_{\max} = \max_{\mathbf{x} \neq 0} R(\mathbf{x}) \quad (\text{B.10})$$

où :

$$R(\mathbf{x}) = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \quad (\text{B.11})$$

donc :

$$\lambda_{\max} \geq \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \quad (\text{B.12})$$

Prenons alors comme vecteur \mathbf{x} particulier un vecteur nul partout sauf sur les nœuds d'un élément quelconque e et notons \mathbf{x}_e sa restriction sur les degrés de liberté de cet élément (vecteur de petite taille, donc). On a $\mathbf{x}_e^T \mathbf{x}_e = \mathbf{x}^T \mathbf{x}$ et les contributions élémentaires $\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{P} \bar{\mathbf{A}} \mathbf{P}^T \mathbf{x} = \mathbf{x}_e^T \mathbf{A}_e \mathbf{x}_e + a$ où, si toutes les matrices élémentaires \mathbf{A}_e sont positives, $a \geq 0$, alors :

$$\lambda_{\max} \geq \frac{\mathbf{x}_e^T \mathbf{A}_e \mathbf{x}_e + a}{\mathbf{x}_e^T \mathbf{x}_e} \geq \frac{\mathbf{x}_e^T \mathbf{A}_e \mathbf{x}_e}{\mathbf{x}_e^T \mathbf{x}_e} \quad (\text{B.13})$$

donc :

$$\lambda_{\max} \geq \bar{\lambda}_{\max} \quad (\text{B.14})$$

C'est le résultat contraire au précédent, concernant les systèmes aux valeurs propres généralisés.

C

TECHNIQUE D'ACCÉLÉRATION D'AITKEN

Cette technique [67] se sert des itérés précédents pour construire une *meilleure* nouvelle approximation. Si on considère une série scalaire η_n qui tend vers $\bar{\eta}$, construite par un algorithme du premier ordre avec un taux de convergence κ ; en supposant les itérés $\eta_n, \eta_{n-1}, \eta_{n-2}$ calculés, l'approximation au premier ordre s'écrit :

$$\begin{aligned}\eta_n - \bar{\eta} &\approx \kappa(\eta_{n-1} - \bar{\eta}) \\ \eta_{n-1} - \bar{\eta} &\approx \kappa(\eta_{n-2} - \bar{\eta})\end{aligned}\tag{C.1}$$

Ces deux équations permettent de calculer une approximation de $\bar{\eta}$, qui va servir pour construire le nouvel itéré η'_n , et κ , qui permet donc d'estimer le taux de convergence en cours de route. On détermine ainsi l'expression du nouvel itéré :

$$\eta'_n = \eta_n - \frac{(\eta_{n-1} - \eta_n)^2}{\eta_n - 2\eta_{n-1} + \eta_{n-2}}\tag{C.2}$$

D

REPÈRES BIOGRAPHIQUES



Alexander Craig Aitken

Aitken left the Otago Boys' High School in Dunedin in 1913 having won a scholarship to Otago University. He began to study languages and mathematics with the intention of becoming a school teacher but his university career was interrupted by World War I. Back to New Zealand in 1917, Aitken followed his original intention and became a school teacher. At his old school Otago Boys' High School. Encouraged by the new professor of mathematics at Otago University, Aitken came to Scotland in 1923 and studied for a Ph.D. at Edinburgh under Whittaker. In 1925 he was appointed to Edinburgh where he spent the rest of his life.

- Born : 1 April 1895 in Dunedin, New Zealand.
- Died : 3 Nov 1967 in Edinburgh, Scotland.



André-Louis Cholesky

Cholesky was a French military officer involved in geodesy and surveying in Crete and North Africa just before World War I. He entered l'École Polytechnique at the age of 20 and was attached to the Geodesic Section of the Geographic Service, in June 1905. That was the period when the revision of the French triangulation had just been decided to be used as the base of a new cadastral triangulation. Cholesky solved the problem of the adjustment of the grid with the method now named after him to compute solutions to the normal equations for the least squares data fitting problem.

- Born : 15 Oct 1875 in Montguyon (Charente-Inférieure).
- Died : 31 Aug 1918.



Johann Carl Friedrich Gauss

Gaussian elimination, which first appeared in the text *Nine Chapters of the Mathematical Art* written in 200 BC, was used by Gauss in his work which studied the orbit of the asteroid Pallas. Using observations of Pallas taken between 1803 and 1809, Gauss obtained a system of six linear equations in six unknowns. Gauss gave a systematic method for solving such equations which is precisely Gaussian elimination on the coefficient matrix.

- Born : 30 April 1777 in Brunswick, Duchy of Brunswick (now Germany).
- Died : 23 Feb 1855 in Göttingen, Hanover (now Germany).



Issai Schur

In 1894, Schur entered the University of Berlin to read mathematics and physics. Schur made major steps forward in representation theory of groups, in collaboration with Frobenius, one of his teacher.

In 1916, in Berlin, he built his famous school and spent most of the rest of his life there. Schur's own impressive contributions were extended by his students in a number of different directions. They worked on topics such as soluble groups, combinatorics, and matrix theory.

- Born : 10 Jan 1875 in Mogilyov, Mogilyov province, Belarus.
- Died : 10 Jan 1941 in Jerusalem, Palestine.



Joseph-Louis Lagrange

Joseph-Louis Lagrange is usually considered to be a French mathematician, but the Italian Encyclopedia refers to him as an Italian mathematician. They certainly have some justification in this claim since Lagrange was born in Turin and baptised in the name of Giuseppe Lodovico Lagrangia.

The papers by Lagrange cover a variety of topics : beautiful results on the calculus of variations, work on the calculus of probabilities... In a work on the foundations of dynamics, Lagrange based his development on the principle of least action and on kinetic energy. The 'Mécanique analytique' which Lagrange had written in Berlin, was published in 1788. It summarized all the work done in the field of mechanics since the time of Newton and is notable for its use of the theory of differential equations. With this work Lagrange transformed mechanics into a branch of mathematical analysis.

« Lagrange, in one of the later years of his life, imagined that he had overcome the difficulty (of the parallel axiom). He went so far as to write a paper, which he took with him to the Institute, and began to read it. But in the first paragraph something struck him which he had not observed : he muttered : 'Il faut que j'y songe encore', and put the paper in his pocket. » A De Morgan *Budget of Paradoxes*.

- Born : 25 Jan 1736 in Turin, Sardinia-Piedmont (now Italy).
- Died : 10 April 1813 in Paris, France.



Karl Gustav Jacob Jacobi

Jacobi published three treatises on determinants in 1841. These were important in that for the first time the definition of the determinant was made in an algorithmic way and the entries in the determinant were not specified so his results applied equally well to cases where the entries were numbers

or to where they were functions. These three papers by Jacobi made the idea of a determinant widely known.

- Born : 10 Dec 1804 in Potsdam, Prussia (now Germany).
- Died : 18 Feb 1851 in Berlin, Germany.

Ludwig Philipp von Seidel

Philipp Seidel entered the University of Berlin in 1840 and studied under Dirichlet and Encke. He moved to Königsberg where he studied under Bessel, Jacobi and Franz Neumann. Seidel obtained his doctorate from Munich in 1846 and he went on to become professor there.

Seidel wrote on dioptics and mathematical analysis. His work on lens identified mathematically five coefficients describing the aberration of a lens, now called 'Seidel sums'. These Seidel sums correspond to spherical aberration, coma, astigmatism, Petzval curvature and distortion. He also introduced the concept of nonuniform convergence and applied probability to astronomy.

- Born : 24 Oct 1821 in Zweibröcken, Germany.
- Died : 13 Aug 1896 in Munich, Germany.



James Hardy Wilkinson

Wilkinson won a Foundation Scholarship to Sir Joseph Williamson's Mathematical School, Rochester at the age of 11. In 1940, Wilkinson began war work which involved mathematical and numerical work on ballistics.

Wilkinson continued work becoming more involved in writing many high quality papers on numerical analysis, particularly numerical linear algebra where he developed backward error analysis methods. He worked on numerical methods for solving systems of linear equations and eigenvalue problems.

As well as the large numbers of papers on his theoretical work, Wilkinson developed computer software, working on the production of libraries of numerical routines. The NAG (Numerical Algorithms Group) began work in 1970 and much of the linear algebra routines were due to Wilkinson.

- Born : 27 Sept 1919 in Strood, Kent, England.
- Died : 5 Oct 1986 in London, England.



John William Strutt Lord Rayleigh

His first paper in 1865 was on Maxwell's electromagnetic theory. He worked on propagation of sound and, while on an excursion to Egypt taken for health reasons, Strutt wrote Treatise on Sound (1870-1). In 1879, he wrote a paper on travelling waves, this theory has now developed into the theory of solitons. His theory of scattering (1871) was the first correct explanation of why the sky is blue.

- Born : 12 Nov 1842 in Langford Grove (near Maldon), Essex, England.
- Died : 30 June 1919 in Terling Place, Witham, Essex, England.



Helmut Wielandt

Wielandt entered the University of Berlin in 1929 and there he studied mathematics, physics and philosophy. There he was greatly influenced by Schmidt and Schur. It was on the topic of permutation groups that Wielandt wrote his doctoral dissertation and he was awarded a doctorate in 1935. At the end of World War II Wielandt was appointed associate professor at the University of Mainz, and in 1951, Ordinary Professor at the University of Tübingen. For 20 years beginning in 1952, Wielandt was managing editor of *Mathematische Zeitschrift*.

- Born : 19 Dec 1910 in Niedereggenen, Loerrach, Germany.
- Died : 1984.



Alston Scott Householder

Householder then taught mathematics in a number of different places. He was awarded a Ph.D. by the University of Chicago in 1947 for a thesis on the calculus of variations. However his interests were moving towards applications of mathematics, particularly applications of mathematics to biology.

In 1946, after the end of the war, Householder joined the Mathematics Division of Oak Ridge National Laboratory. Here he changed topic, moving into numerical analysis which was increasing in importance due to the advances in computers. He started publishing on this new topic with *Some numerical methods for solving systems of linear equations* which appeared in 1950.

In 1964, Householder published one of his most important books : *The theory of matrices in numerical analysis*.

- Born : 5 May 1904 in Rockford, Illinois, USA.
- Died : 4 July 1993 in Malibu, California, USA.



Cornelius Lanczos

In 1938 at Purdue, Lanczos published his first work in numerical analysis. Two years later he published a matrix method of calculating Fourier coefficients which, over 25 years later, was recognised as the Fast Fourier Transform algorithm. In 1946, with Boeing, he worked on applications of mathematics to aircraft design and was able to develop new numerical methods to solve the problems. In 1949 he moved to the Institute for Numerical Analysis of the National Bureau of Standards in Los Angeles. Here he worked on developing digital computers and was able to produce versions of the numerical methods he had developed earlier to program on the digital computers.

- Born : 2 Feb 1893 in Székesfehérvár, Hungary.
- Died : 25 June 1974 in Budapest, Hungary.

Joseph Raphson

Joseph Raphson's life can only be deduced from a number of pointers. It is through the University of Cambridge records that we know that Raphson attended Jesus College Cambridge and graduated with an M.A. in 1692. Rather remarkably Raphson was made a member of the Royal Society in 1691, the year before he graduated. His election to that Society was on the strength of his book *Analysis aequationum universalis* which was published in 1690 contained the Newton method for approximating the roots of an equation.

Raphson's ideas of space and philosophy were based on Cabalist ideas. The Cabala was a Jewish mysticism which was influential from the 12th century on and for which several basic doctrines were strong influences on Raphson's philosophical thinking. The doctrines included the withdrawal of the divine light, thereby creating primordial space, the sinking of luminous particles into matter and a *cosmic restoration*.

- Born : 1648 in Middlesex, England.
- Died : 1715.



Carl David Tolmé Runge

As a German mathematician, physicist, and spectroscopist, he is known to be the co-developer and co-eponym of the Runge-Kutta method in the field of numerical analysis. After taking his secondary school teachers examinations he returned to Berlin where he was influenced by Kronecker. Runge then worked on a procedure for the numerical solution of algebraic equations in which the roots were expressed as infinite series of rational functions of the coefficients.

- Born : 30 Aug 1856 in Bremen, Germany.
- Died : 3 Jan 1927 in Göttingen, Germany.



Sir Isaac Newton

Isaac Newton is perhaps the best known renaissance scientist today, living between 1642 and 1727. We think of gravity, celestial mechanics, and calculus when we think of him. He certainly did develop the calculus by building upon the ideas of Fermat and Barrow (the person whose chair he took when he went to Cambridge). But he was not alone in developing calculus. And his focus was really one of mechanics - how do bodies move. His focus was always on motion and is reflected in the terminology he chose for calculus - what we call *derivatives*, he called *fluxions*. And his integrals were simply *inverse fluxions*. You can see how *flux* was his focus. He also developed the *dot* notation for calculus - one dot meant a first derivative, two dots a second, and so forth. Today this is used, but only with respect to time derivatives. A general derivative still needs to specify the variable of differentiation and more often than not, time derivatives are done in this more general manner.

As indicated earlier, Newton and his followers argued vehemently with Huygens and his followers over the nature of light. Newton subscribed to a *corpuscular* theory, where he en-

visioned light as small compact bodies of energy. Huygens focussed on the wave like nature and developed that theory. The diffraction properties of light were so obvious, that Huygens school eventually won out, and the wave theory of light ruled science for the next three centuries.

- Born : 4 Jan 1643 in Woolsthorpe, Lincolnshire, England.
- Died : 31 March 1727 in London, England.



Leonhard Euler

The publication of many articles and his book *Mechanica* (1736-37), which extensively presented Newtonian dynamics in the form of mathematical analysis for the first time, started Euler on the way to major mathematical work. He integrated Leibniz's differential calculus and Newton's method of fluxions into mathematical analysis. In number theory he stated the prime number theorem and the law of biquadratic reciprocity. Euler made large bounds in modern analytic geometry and trigonometry. He was the most prolific writer of mathematics of all time. His complete works contains 886 books and papers.

- Born : 15 April 1707 in Basel, Switzerland.
- Died : 18 Sept 1783 in St Petersburg, Russia.



Martin Wilhelm Kutta

Kutta was professor at the RWTH Aachen from 1910 to 1911. He then became professor at Stuttgart and remained there until he retired in 1935. In 1901, he had co-developed the Runge-Kutta method, used to solve ordinary differential equations numerically. He is also remembered for the Zhukovsky-Kutta aerofoil, the Kutta-Joukowski theorem and the Kutta condition in aerodynamics.

- Born : 3 Nov 1867 in Pitschen, Upper Silesia (now Byczyna), Poland.
- Died : 25 Dec 1944 in Frstenfeldbruck, Germany.

INDEX

- A**
accélération
 d'Aitken 20, 77
 méthode de l'accélération linéaire 56, 60
 méthode de l'accélération moyenne 56
 schéma de Newmark 55
Aitken 79
 technique d'accélération 20, 77
- B**
Broyden
 Méthode BFGS 41
- C**
Cauchy
 problème de 43, 52
Cholesky 79
 factorisation de 12, 13, 23, 25, 73
coût
 CPU 9
 d'un algorithme . 8–10, 12, 15, 33, 52, 54
 d'une méthode 17
 fonction 40
 fonction (optimisation) 38
complexité 9, 12, 13, 32
conditionnement 6, 7, 24, 26
convergence 17–20
 asymptotique 18
 cubique 30
 gradient conjugué préconditionné 23
 linéaire 33
 méthode de Jacobi 28
 méthode de Lanczos 35
 méthode des itérations inverses .. 30
 méthode des puissances 28
 méthode du gradient 22
 quadratique 64
 schéma d'Euler explicite 44
 taux de 18, 19, 77
 valeurs propres multiples 29
Cramer
 méthode de 8
Crout
 factorisation 13
 factorisation de 13–15
Cuthill-MacKee 9
- E**
éléments finis 1, 5, 7, 13, 37, 75, 76
 code 31
 conditionnement 7
 conditions aux limites 14
 convergence 19
 discrétisation 8, 25
 problème de mémoire 41
 super-élément 13
 système différentiel du premier ordre 53
- équation
 caractéristique 26, 56
 d'équilibre 66
 d'état 63
 d'admissibilité cinématique 37
 d'admissibilité statique 37
 d'Euler 14–16
 de comportement 66
 de la dynamique 67
 différentielle 67
 différentielle non-linéaire 62
 du mouvement 67
Euler 84
 équation d' 14–16
 schéma d'
 explicite, 44, 46, 48
 implicite, 47, 50, 62
- F**
factorisation
 de Cholesky 12, 13, 23, 25, 73
 de Crout 13–15
 de Gauss 13, 15
 LDL^T 13, 31
 LDM^T 14
 LU 11–13
 QR 31, 32
- G**
Gauss 79
 élimination de 10
 algorithme de Gauss-Seidel 19
 factorisation de 13, 15
 méthode de Gauss-Seidel 18, 19
- H**
Hessenberg
 forme de 33
Householder 82
 algorithme de Householder-Businger-Golub 32
 factorisation QR 32
 méthode de 26
 transformation de 31, 34
- J**
Jacobi 80, 81

- algorithme de 27
 méthode de 18–20, 26, 27, 29
 convergence, 28
- K**
- Kahan
 théorème de 20
- Kutta 84
 Kutta-Joukowski 84
 méthode de Runge-Kutta . 51–53, 55,
 62
 Runge-Kutta 83
 Zhukovsky-Kutta 84
- L**
- Lagrange 80
 multiplicateurs de 15, 20
- Lanczos 82
 itérations de 36
 méthode de 26, 35, 36
- Leibniz 84
- linéaire(s)
 équation(s) 46, 54
 combinaison(s) 10, 17, 29
 convergence(s) 33
 méthode des accélérations 56, 60
 récurrence(s) 18
 relation(s) 16
 système(s) 5, 18, 26, 32, 39, 54, 55
- M**
- matrice 5
 bande 15
 booléenne 14
 condensée 11
 d'amortissement 25, 55
 d'assemblage 75
 d'itération 19, 56
 décalée 30
 de capacité 53
 de conduction 53
 de masse 54, 55
 de raideur 16
 de rigidité 13, 14, 38, 55
 tangente, 64
 de rotation 27
 jacobienne 38–40, 63
 pleine 12
 semi-définie 17
 sous-matrice 11, 12
 symétrique 13
 symétrique définie positive 15
 triangulaire supérieure 11
 valeurs propres 20, 25
- N**
- Newmark
 méthode de 55, 57, 61
- schéma de 55, 57, 58
 stabilité 56
- Newton 80, 83
 fluxions 84
 loi de 66
 méthode de 37, 38, 63
 modèle de 68
 Newton-Raphson 39, 64
 quasi-newton 40
- norme 6–8, 18, 22, 24, 29, 30, 36
 infinie 29
 matricielle 6
- normer 28, 35
- O**
- orthogonalisation
 d'un sous-espace 12
 de Gram-Schmidt 73
 ré-orthogonalisation de grand problème
 35
- P**
- pivot 17
 Crout 15
 facrotisation LU 11
 partiel 11
 stratégie 11
 total 11
- préconditionnement 23
 gradient conjugué préconditionné 23,
 41
 préconditionneur 23
- R**
- Raphson 83
 Newton-Raphson 39, 64
- Rayleigh 81
 amotissement de 57
 itérations inverses de 30
 quotient de 26, 75
- Runge 83
 méthode de Runge-Kutta . 51–53, 55,
 62, 84
 Runge-Kutta 83
- S**
- Schur 80, 82
 complément de 11
 condensation de 13
- Seidel 81
 algorithme de Gauss-Seidel 19
 méthode de Gauss-Seidel 18, 19
- Southwell 19
- système
 aux valeurs propres 25, 26, 75, 76
 d'exploitation 10
 différentiel
 du premier ordre, 67

non régulier, 67	
dynamique	67
linéaire	5, 8, 26, 32, 39
creux, 8	
triangulaire, 10, 12	
non-linéaire	18, 37
T	
trapèze(s)	
méthode des	47, 48, 61, 67
schéma des	48
W	
Wielandt	82
méthode de	30
Wilkinson	81
test de	36