

**J. E. ROMBALDI**

**Algorithmique numérique**

**Avec Open Ada sur P. C.**

# TABLE DES MATIERES

<b>Contents</b> .....	xv
<b>Avant-propos</b> .....	xvii
<b>Chapitre 1 — Le langage Ada avec OpenAda</b>	
1. <i>Petit historique de Ada</i> .....	1
2. <i>Une présentation rapide du langage Ada</i> .....	3
2.1 Le manuel de référence .....	3
2.2 Un survol du langage .....	3
2.3 Exemples d'utilisations de Ada .....	6
3. <i>Utilisation du compilateur OpenAda sur micro-ordinateur</i> .....	7
<b>Chapitre 2 — Bibliothèque mathématique</b>	
1. <i>Introduction</i> .....	11
2. <i>Le paquetage CRT</i> .....	12
2.1 Spécifications des paquetages COMMON_CRT et CRT .....	12
2.2 Documentation du paquetage CRT .....	13
2.2.1 La procédure CLREOL .....	13
2.2.2 La procédure CLRSCR .....	13
2.2.3 La procédure GOTOXY .....	13
2.2.4 La fonction READKEY .....	13
2.2.5 Les fonctions WHERE_X et WHERE_Y .....	13
2.2.6 La procédure Cursor .....	13
2.2.7 La procédure PUT_CHAR .....	13
2.2.8 La procédure PUT_STRING .....	13
2.2.9 La procédure PUT_LINE_STRING .....	13
2.2.10 Partie initialisation de CRT .....	13
3. <i>Le paquetage MATH0</i> .....	14
3.1 Spécifications des paquetages COMMON_MATH0 et MATH0 .....	14
3.2 Documentation du paquetage MATH0 .....	15
3.2.1 La procédure GET_TIME .....	15
3.2.2 La procédure RANDOMIZE .....	15
3.2.3 Les fonctions RANDOM .....	15
3.2.4 La fonction MAJUSCULE .....	15
3.2.5 Les procédures GET_LINE .....	16
3.2.6 La fonction LIRE_REPONSE .....	16
3.2.7 La procédure PAUSE .....	16
3.2.8 La procédure BIP .....	16
3.2.9 La procédure CHRONOMETRE .....	16
3.2.10 Les procédures d'entrée des entiers et des réels .....	16
3.2.11 Les fonctions CHAINE_ENTIERE et CHAINE_REELLE .....	16
3.2.12 La procédure LECTURE_REPERTOIRE .....	16
3.2.13 La procédure LIRE_FICHER .....	17
3.2.14 La procédure FICHER_EXISTE .....	17
3.2.15 La procédure CREER_FICHER .....	17
3.2.16 La procédure MODE_AFFICHAGE .....	17
3.2.17 La fonction TRUNC .....	17
3.2.18 La fonction FRAC .....	17

3.2.19	Les procédures de permutation de deux entiers ou deux réels .....	17
3.2.20	Les procédures VAL .....	17
3.2.21	La fonction LONGUEUR_CHAINE .....	18
3.2.22	Un exemple de programme utilisant Math0 .....	18
4.	<i>Le paquetage MATH1</i> .....	19
4.1	Spécifications des paquetages COMMON_MATH1 et MATH1 .....	20
4.2	Documentation du paquetage MATH1 .....	20
4.2.1	Les fonctions Puissances .....	20
4.2.2	Les fonctions ArcTan, ArcSin et ArcCos .....	20
4.2.3	Les fonctions hyperboliques inverses .....	21
4.2.4	Les fonctions Gamma et LnGamma .....	21
5.	<i>Le paquetage Math2</i> .....	21
5.1	Spécifications du paquetage MATH2 .....	22
5.2	Documentation de MATH2 .....	22
5.2.1	Les fonctions EVALUE .....	22
5.2.2	La procédure DETRUIT_FONCTION .....	23
5.2.3	La procédure GET_LINE .....	23
5.2.4	La fonction TEXTE_DE_FONCTION .....	23
6.	<i>Le paquetage GRAPH</i> .....	23
6.1	Spécifications des paquetages COMMON_GRAPH et GRAPH .....	23
6.2	Documentation de GRAPH .....	24
6.2.1	La procédure CLEAR_SCREEN .....	24
6.2.2	La procédure CLOSE_GRAPH .....	24
6.2.3	La procédure INIT_GRAPH .....	24
6.2.4	La procédure MOVE_TO .....	24
6.2.5	La procédure PUT_PIXEL .....	24
6.2.6	La procédure LINE .....	24
6.2.7	La procédure LINE_TO .....	24
6.2.8	La procédure RECTANGLE .....	24
7.	<i>Le paquetage MATH3</i> .....	24
7.1	Spécifications des paquetages COMMON_MATH3 et MATH3 .....	24
7.2	Documentation de MATH3 .....	25
7.2.1	La procédure INIT_GRAPHIQUE .....	25
7.2.2	La procédure MODE_TEXTE .....	25
7.2.3	La procédure MESSAGE_ERREUR_GRAPHIQUE .....	26
7.2.4	La procédure TEST_MODE_GRAPHIQUE .....	26
7.2.5	La procédure TEXT_MOVE_TO .....	26
7.2.6	La procédure PAUSE_GRAPHIQUE .....	26
7.2.7	La procédure FENETRE .....	26
7.2.8	La procédure CADRE_GRAPHIQUE .....	26
7.2.9	La procédure FENETRE_GRAPHIQUE .....	26
7.2.10	La procédure CONVERSION .....	26
7.2.11	La procédure DEPLACE .....	26
7.2.12	La procédure TRACE .....	27
7.2.13	La procédure CROIX .....	27
7.2.14	La procédure TRACE_SEGMENT .....	27
7.2.15	La procédure POINT .....	27
7.2.16	La procédure CERCLE .....	27
7.2.17	La procédure TITRE .....	27
7.2.18	La procédure X_AXE .....	27
7.2.19	La procédure Y_AXE .....	27
7.2.20	La procédure XY_AXES .....	27
7.2.21	La procédure SORTIE_GRAPHIQUE .....	27
7.2.22	Un exemple d'utilisation de MATH2 et de MATH3 .....	27
8.	<i>Exercices</i> .....	32

**Chapitre 3 — Analyse numérique linéaire**

1. <i>Introduction</i> .....	39
1.1 Position des problèmes .....	39
1.2 Notations .....	39
1.3 Remarques .....	39
1.4 Problèmes numériques liés à la résolution des systèmes linéaires .....	40
1.5 Systèmes dégénérés et numériquement dégénérés .....	40
1.6 Problèmes de stabilité numérique .....	42
1.7 Méthodes de résolution des systèmes linéaires .....	43
1.8 Exemples d'applications .....	43
2. <i>Rappels et compléments sur le calcul matriciel</i> .....	44
2.1 Normes vectorielles et matricielles .....	44
2.1.1 Norme matricielle induite par une norme vectorielle .....	44
2.1.2 Norme utilisée dans ce livre .....	44
2.2 Conditionnement d'une matrice .....	45
2.3 Valeurs propres et rayon spectral d'une matrice .....	46
2.3.1 Définitions .....	46
2.3.2 Propriétés du rayon spectral .....	46
2.4 Quelques propriétés des matrices symétriques réelles .....	47
2.4.1 Transposée d'une matrice .....	47
2.4.2 Spectre des matrices symétriques .....	47
2.5 Matrices à diagonale strictement dominante .....	47
3. <i>Résolution numérique des systèmes linéaires. Inversion des matrices</i> .....	48
3.1 Position des problèmes. Notations .....	48
3.2 Une méthode impraticable : la méthode de Cramer .....	49
3.3 Cas des systèmes triangulaires .....	49
3.3.1 Résolution des systèmes triangulaires .....	49
3.3.2 Inversion des matrices triangulaires .....	50
3.4 Méthode des pivots de Gauss .....	51
3.4.1 Principe de la méthode .....	51
3.4.2 Description de la méthode .....	51
3.4.3 Remarque sur le choix des pivots : méthode de Gauss avec pivots partiels .....	52
3.4.4 Nombre d'opérations élémentaires dans la méthode de Gauss .....	52
3.4.5 Programmation structurée .....	53
3.4.6 Remarques .....	54
3.5 Méthode de Gauss-Jordan .....	55
3.5.1 Méthode d'élimination de Gauss-Jordan .....	55
3.5.2 Inversion et calcul du déterminant d'une matrice .....	57
3.6 Interprétation algébrique de la méthode de Gauss. Décomposition L-R ou méthode de Crout .....	58
3.6.1 Condition nécessaire et suffisante pour qu'il n'y ait pas de permutations dans la méthode de Gauss .....	58
3.6.2 Interprétation algébrique de la méthode Gauss .....	59
3.6.3 Détermination pratique de la décomposition L-R .....	60
3.6.4 Résolution d'un système linéaire par la méthode de Crout .....	61
3.6.5 Calcul de l'inverse d'une matrice par la méthode de Crout .....	61
3.7 Cas des matrices tridiagonales .....	61
3.7.1 Notations et hypothèses .....	61
3.7.2 Décomposition L-R d'une matrice tridiagonale .....	61
3.7.3 Méthode de résolution d'un système tridiagonal .....	62
3.7.4 Déterminant d'une matrice tridiagonale .....	63
3.7.5 Exemples d'applications .....	63
3.8 Résolution des systèmes symétriques .....	63

3.8.1	Décomposition $L \cdot D \cdot L^t$ .....	63
3.8.2	Calculs pratiques .....	64
3.8.3	Programmation structurée .....	64
3.8.4	Cas particulier des matrices symétriques définies positives. Décomposition $B \cdot B^t$ de Cholesky .....	65
3.9	Résolution des systèmes linéaires par des méthodes itératives .....	65
3.9.1	Remarques préliminaires .....	65
3.9.2	Principe des méthodes itératives .....	65
3.9.3	Méthode de Jacobi .....	66
3.9.4	Méthode de Gauss-Seidel .....	67
3.9.5	Comparaison des méthodes de Jacobi et de Gauss-Seidel dans le cas des matrices tridiagonales d'ordre $n \geq 3$ .....	69
3.9.6	Méthode de relaxation .....	70
3.9.7	Exemple d'application .....	73
4.	<i>Calcul des valeurs propres et des vecteurs propres de certaines matrices réelles</i> .....	75
4.1	Introduction .....	75
4.2	Calcul du rayon spectral d'une matrice. Méthode de la puissance itérée .....	75
4.2.1	Hypothèses et notations .....	75
4.2.2	Méthode de la puissance itérée pour calculer le rayon spectral .....	76
4.3	Calcul des autres valeurs propres par la méthode de déflation .....	78
4.3.1	Hypothèses et notations .....	78
4.3.2	Un lemme .....	78
4.3.3	Programmation structurée .....	79
4.4	Méthode de Rutishauser .....	79
4.4.1	Un lemme .....	79
4.4.2	Hypothèses .....	80
4.4.3	Principe de la méthode de Rutishauser .....	80
4.4.4	Une condition suffisante de convergence de la méthode de Rutishauser .....	80
4.5	Méthode de Jacobi pour calculer les valeurs et vecteurs propres d'une matrice symétrique .....	81
4.6	Méthodes de calcul du polynôme caractéristique .....	85
4.6.1	Méthode de Souriau .....	85
4.6.2	Méthode de Krylov .....	87
4.6.3	Méthode de Leverrier .....	87
5.	<i>Exercices</i> .....	88
6.	<i>Programmation Ada</i> .....	93
(6.1)	Spécifications des paquetages COMMON_MATRIX et MATRIX .....	93
(6.2)	Spécification du paquetage Algeblin .....	94
(6.3)	Spécification du paquetage Donnees_ALGEBLIN .....	95
(6.4)	Démonstration de la méthode des pivots de Gauss .....	95
(6.5)	Démonstration de la méthode L-R .....	96
(6.6)	Démonstration de la méthode du double balayage de Cholesky .....	97
(6.7)	Démonstration de la méthode de Gauss-Jordan .....	97
(6.8)	Démonstration de la méthode de Cholesky .....	98
(6.9)	Démonstration de la méthode de Jacobi .....	99
(6.10)	Démonstration de la méthode de relaxation .....	100
(6.11)	Démonstration de la décomposition Q-R .....	100
(6.12)	Démonstration du calcul de l'inverse d'une matrice de Van Der Monde .....	102
(6.13)	Spécification du paquetage SPECTRE .....	102
(6.14)	Démonstration de la méthode de déflation .....	103
(6.15)	Démonstration de la méthode de Rutishauser .....	104
(6.16)	Démonstration de la méthode de Jacobi .....	105
(6.17)	Démonstration de la méthode de Souriau .....	106

## Chapitre 4 — Résolution numérique des systèmes non linéaires

1. Introduction .....	107
1.1 Position des problèmes .....	107
1.2 Remarques .....	107
2. Cas des équations numériques .....	108
2.1 Notations .....	108
2.2 La méthode dichotomie ou de bisection .....	108
2.3 La méthode de Newton-Raphson .....	109
3. Cas des équations algébriques .....	111
3.1 Introduction .....	111
3.2 La méthode Newton_Maehly .....	111
4. Résolution des systèmes non linéaires par la méthode de Newton-Raphson .....	113
4.1 Introduction .....	113
4.2 Algorithme de Newton-Raphson .....	114
4.3 Calcul de la matrice jacobienne .....	114
4.4 Programmation structurée .....	114
5. Racines d'un polynôme. Méthode de Bairstow .....	115
5.1 Principe de la méthode .....	115
5.2 Recherche des coefficients p et q .....	115
5.3 Programmation structurée .....	118
6. Exercices .....	119
7. Programmation Ada .....	120
7.1 Spécification du paquetage EQUATIONS_GENERIQUE .....	120
7.2 Démonstration du paquetage EQUATIONS_GENERIQUE .....	120
7.3 Spécification du paquetage SYSTEME_EQUATIONS_GENERIQUE .....	122
7.4 Démonstration du paquetage SYSTEME_EQUATIONS_GENERIQUE .....	123
7.5 Spécifications des paquetages COMMON_POLY et POLY .....	124
7.6 Démonstration du paquetage POLY .....	125

## Chapitre 5 — Approximation et interpolation

1. Introduction .....	127
2. Problèmes d'approximation. Méthode des moindres carrés .....	127
2.1 Introduction .....	127
2.2 Exemples .....	128
2.2.1 La loi d'Ohm .....	128
2.2.2 Elongation d'un ressort .....	128
2.2.3 Problème du fil chaud .....	128
2.2.4 Corrélation statistique .....	128
2.3 Détermination des paramètres .....	128
2.4 Principes des méthodes des moindres carrés .....	129
2.5 Les modèles linéaires .....	130
2.5.1 Généralités .....	130
2.5.2 La régression affine .....	130
2.5.3 Exemple : le problème du fil chaud .....	133
2.5.4 Régression polynomiale .....	133
2.5.5 Cas général .....	136
2.5.6 Régression trigonométrique .....	136
2.6 Les modèles non linéaires .....	137
2.7 Un exemple d'application : détermination du coefficient de traînée d'une particule sphérique .....	138
2.8 Calcul du cercle des moindres carrés .....	146
2.8.1 Notations et hypothèses .....	146
2.8.2 Calcul du cercle des moindres carrés .....	147
2.8.3 Application à un calcul de rayon de courbure .....	148
3. Approximation uniforme des fonctions continues. Courbes de Bernstein, Bézier et B-Splines .....	149

3.1	Position du problème	149
3.2	Les bases de Bernstein. Polynômes de Bernstein	150
3.2.1	Les bases de Bernstein	150
3.2.2	Polynômes de Bernstein associés à une fonction continue	152
3.3	Les courbes de Bézier	152
3.3.1	Données du problème	152
3.3.2	Les courbes de Bézier	153
3.3.3	Algorithme de De Casteljeau	154
3.4	Les surfaces de Bézier	155
3.4.1	Définition	155
3.4.2	Propriétés	156
3.4.3	Algorithme de De Casteljeau	157
3.5	Les courbes B-Splines	157
3.5.1	Introduction	157
3.5.2	Fonctions de base B-Splines	157
3.5.3	Courbes B-Splines	160
4.	Problèmes d'interpolation	162
4.1	Introduction. Position des problèmes	162
4.2	L'interpolation polynomiale de Lagrange	163
4.2.1	L'interpolation linéaire	163
4.2.2	Le théorème d'interpolation de Lagrange	163
4.2.3	Algorithme de Neville pour calculer le polynôme de Lagrange	164
4.2.4	Forme de Newton du polynôme de Lagrange	165
4.2.5	Sur le choix des abscisses d'interpolation. Les polynômes de Tchébychev	167
4.3	Interpolation spline cubique	168
4.3.1	Introduction	168
4.3.2	Position du problème. Notations	168
4.3.3	Calcul des coefficients $d_i$	169
4.3.4	Calcul des $b_i$ et des $a_i$	169
4.3.5	Calcul des $c_i$	170
4.3.6	Calcul des $s_i = x^i$	170
4.3.7	L'interpolation spline naturelle	170
4.3.8	Programmation structurée pour l'interpolation spline naturelle	171
4.3.9	Application à un calcul d'intégrale	172
4.3.10	Cas des fonctions périodiques	173
4.3.11	Interpolation spline pour les courbes fermées	175
4.3.12	Majoration de l'erreur d'interpolation	176
5.	Exercices	177
6.	Programmation Ada	179
6.1	Spécification du paquetage REGRESSIONS	179
6.2	Le paquetage DONNEES_POINTS	179
6.3	Démonstration de la régression affine	179
6.4	Démonstration de la régression polynomiale	180
6.5	Démonstration de la régression circulaire	181
6.6	Spécification du paquetage BEZIER	181
6.7	Démonstration de l'approximation de Bézier	182
6.8	Spécification du paquetage B_SPLINE	183
6.9	Démonstration de l'approximation B_Spline	183
6.10	Spécification du paquetage LAGRANGE	184
6.11	Démonstration de l'interpolation de Lagrange	184
6.12	Spécification du paquetage INTERPOLATION_SPLINE	185
6.13	Démonstration de l'interpolation spline	186
<b>Chapitre 6 — Calcul numérique des intégrales</b>		
1.	Introduction. Position des problèmes	187

1.1 Remarques préliminaires sur le calcul des primitives .....	187
1.1.1 Calcul direct des intégrales.....	187
1.1.2 A propos des primitives élémentaires.....	187
1.2 Elaboration de méthodes numériques.....	188
2. Méthodes de calcul par interpolation polynomiale. Schémas d'intégration classiques .....	189
2.1 Idée des méthodes par interpolation .....	189
2.2 Résolution du problème (2).....	190
2.2.1 Existence et unicité d'une solution .....	190
2.2.2 Détermination pratique des coefficients $a_i$ .....	190
2.2.3 Majoration de l'erreur .....	191
2.3 Cas particulier où les $x_i$ sont équidistants. Méthodes de Newton et Cotes .....	192
2.3.1 Calcul des coefficients $a_i$ .....	192
2.3.2 Majoration de l'erreur dans les méthodes de Newton_Cotes.....	193
2.4 Méthodes classiques d'intégration .....	193
2.4.1 Cas $n = 1$ . Méthode des trapèzes .....	193
2.4.2 Cas $n = 2$ . Méthode de Simpson.....	196
2.4.3 Méthode de Romberg .....	198
3. Utilisation des polynômes orthogonaux. Quadratures de Gauss .....	200
3.1 Introduction. Idées des méthodes .....	200
3.2 Notations. Position du problème .....	201
3.2.1 Notations .....	201
3.2.2 Position du problème.....	201
3.3 Calcul des abscisses $x_i$ et des coefficients de Gauss $a_i$ .....	202
3.3.1 Remarque.....	202
3.3.2 Condition nécessaire et suffisante sur $P$ .....	202
3.3.3 Détermination explicite de $P_n$ .....	202
3.4 Propriétés des polynômes orthogonaux. Calcul des coefficients de Gauss.....	203
3.4.1 Récurrence vérifiée par les polynômes orthogonaux .....	204
3.4.2 Formule de Darboux-Christoffel .....	204
3.4.3 Calcul des coefficients de Gauss $a_i$ .....	204
3.5 Majoration de l'erreur de quadrature .....	205
3.6 Exemples classiques de polynômes orthogonaux et formules de quadrature correspondantes .....	205
3.6.1 Les polynômes de Legendre.....	205
3.6.2 Polynômes de Tchébychev .....	207
3.6.3 Polynômes de Laguerre .....	208
3.6.4 Polynômes d'Hermite.....	209
4. La méthode probabiliste de Monte-Carlo.....	210
4.1 Nombres pseudo-aléatoires .....	210
4.2 La méthode de Monte-Carlo. Première version : tirage par « noir ou blanc » .....	210
4.2.1 Equiprobabilité sur un segment, ou un pavé de $\mathbb{R}^n$ .....	210
4.2.2 Calcul d'une intégrale simple.....	211
4.2.3 Calcul d'une intégrale multiple .....	211
4.2.4 Remarques .....	212
4.3 Méthode de Monte Carlo, deuxième version.....	212
4.3.1 Remarques préliminaires .....	212
4.3.2 Méthode de Monte Carlo avec échantillonnage simple.....	213
4.3.3 Utilisation de transformations antithétiques .....	214
4.3.4 Calcul d'intégrales multiples.....	215
4.3.5 Exemple : calcul des coordonnées du centre de gravité d'un corps.....	215
5. Transformation de Fourier rapide.....	216
5.1 Position du problème. Notations .....	216
5.2 Approximation des $\hat{f}(x_k)$ .....	217
5.3 La transformation de Fourier discrète .....	218



5.3.1 Définition, propriétés.....	218
5.3.2 Calcul direct de la transformée de Fourier discrète .....	219
5.4 L'algorithme F.F.T. de Cooley et Tukey .....	220
5.4.1 Introduction .....	220
5.4.2 Cas particulier $n = 2$ .....	220
5.4.3 Cas particulier $n = 4$ .....	220
5.4.4 Cas général $n = 2^p$ .....	221
5.4.5 Utilisation de $(p)$ pour calculer la transformation de Fourier discrète.....	221
5.4.6 Nombre d'opérations élémentaires dans l'algorithme de Cooley et Tukey.....	222
5.4.7 Programmation structurée .....	223
5.5 Application au calcul des coefficients de Fourier d'une fonction périodique.....	224
6. Exercices.....	225
7. Programmation Ada .....	229
7.1 Spécification du paquetage INTEGRATION_GENERIQUE .....	229
7.2 Démonstration du paquetage INTEGRATION_GENERIQUE .....	229
7.3 Spécifications des paquetages COMMON_FOURIER et FOURIER_GENERIQUE .....	232
7.4 Démonstration du paquetage FOURIER_GENERIQUE .....	233
<b>Chapitre 7 — Résolution numérique des équations différentielles</b>	
1. Introduction. Origines des problèmes d'équations différentielles .....	237
2. Problème de Cauchy.....	238
2.1 Position du problème.....	238
2.2 Problème de l'existence et l'unicité de solutions.....	240
2.3 Approximation de la solution d'un problème de Cauchy par discrétisation .....	242
3. Généralités sur les méthodes d'intégration à un pas .....	243
3.1 Définitions .....	243
3.2 Lien entre l'erreur de consistance et l'erreur de discrétisation .....	243
3.2.1 Hypothèses .....	243
3.2.2 Première évaluation de l'erreur de consistance.....	243
3.2.3 Hypothèses supplémentaires .....	244
3.2.4 Deuxième évaluation de l'erreur de consistance.....	244
3.2.5 Majoration de l'erreur de discrétisation .....	245
3.3 Les méthodes de Runge-Kutta.....	246
3.3.1 Principe des méthodes de Runge-Kutta.....	246
3.3.2 Exemples classiques .....	246
3.3.3 Remarque.....	248
3.3.4 Les schémas RK2 .....	249
3.4 Programmation structurée de la méthode RK4 pour les équations différentielles d'ordre 1 .....	249
3.5 Programmation structurée de la méthode RK4 pour les systèmes de $p$ équations différentielles d'ordre 1 .....	250
3.6 Programmation structurée de la méthode RK4 pour les équations différentielles d'ordre $q$ .....	251
3.7 Programmation structurée de la méthode RK4 pour les systèmes de $p$ équations différentielles d'ordre $q$ .....	252
4. Contrôle du pas d'intégration .....	253
4.1 Position du problème.....	253
4.2 Choix du pas d'intégration à l'étape $i$ du calcul .....	254
4.3 Programmation structurée de la méthode de Runge-Kutta d'ordre 4 avec contrôle du pas .....	255
5. Problèmes avec conditions aux limites. Méthode du tir .....	256
5.1 Introduction : position du problème .....	256
5.2 La méthode du tir .....	257
5.2.1 Rappels .....	257
5.2.2 Programmation structurée de la méthode du tir.....	257

6. Un exemple d'application. Mouvement de translation d'un corps sphérique pesant dans un fluide au repos .....	259
6.1 Position du problème et notations .....	259
6.2 Equations du mouvement de la sphère .....	260
7. Exercices .....	266
8. Programmation Ada .....	268
8.1 Spécification du paquetage EQUADIFF_11_GENERIQUE .....	268
8.2 Spécification du paquetage EQUADIFF_P1_GENERIQUE .....	269
8.3 Spécification du paquetage EQUADIFF_1Q_GENERIQUE .....	269
8.4 Spécification du paquetage EQUADIFF_PQ_GENERIQUE .....	270
8.5 Démonstration du paquetage EQUADIFF_11_GENERIQUE.....	270
8.6 Démonstration du paquetage EQUADIFF_P1_GENERIQUE .....	271
8.7 Démonstration du paquetage EQUADIFF_1Q_GENERIQUE.....	274
8.8 Démonstration du paquetage EQUADIFF_PQ_GENERIQUE.....	275
<b>Chapitre 8 — Méthode des différences finies</b>	
1. Problème de Dirichlet linéaire en dimension un .....	279
1.1 Introduction .....	279
1.2 Théorème d'existence et d'unicité de solutions du problème de Dirichlet linéaire .....	280
1.3 Forme canonique de l'équation de Dirichlet. Le problème de Poisson .....	280
1.4 Résolution approchée du problème de Dirichlet par discrétisation.....	281
2. Approximations des dérivées d'une fonction par différences finies .....	281
2.1 Approximation de $f'(x_0)$ par différence finie centrée .....	281
2.2 Approximation de $f''(x_0)$ par différence finie centrée .....	282
3. Résolution approchée du problème de Dirichlet par la méthode des différences finies .....	282
3.1 Discrétisation du problème de Dirichlet .....	282
3.2 Discrétisation de l'équation de Poisson .....	284
3.3 Majoration de l'erreur .....	285
3.3.1 Majoration de $h(x) = y(x) - \psi(x)$ .....	285
3.3.2 Majoration de $q(x) = y(x) - j(x)$ .....	286
3.3.3 Majoration de $g(x) = y(x) - \phi(x)$ .....	286
4. Exemples d'application .....	287
4.1 Dissipation de la chaleur dans un disque.....	287
4.2 Fléchissement d'une poutre.....	288
5. Résolution approchée d'équations aux dérivées partielles par la méthode des différences finies.....	288
5.1 Intervention des équations aux dérivées partielles en physique .....	288
5.1.1 L'équation des cordes vibrantes ou équation des ondes en dimension un .....	288
5.1.2 L'équation des ondes en dimension deux .....	291
5.1.3 Equation de la chaleur en dimension un.....	292
5.2 Classification des équations aux dérivées partielles d'ordre 2.....	294
5.2.1 Introduction .....	294
5.2.2 Cas particulier des équations à coefficients constants et sans second membre .....	294
5.2.3 Cas général .....	295
5.2.4 Exemple.....	296
5.2.5 Remarque.....	296
5.3 Principe de la méthode des différences finies .....	296
5.4 Cas d'un problème avec conditions au bord : l'équation de Poisson .....	297
5.5 Equations elliptiques avec conditions aux bords sur un rectangle .....	298
5.5.1 Position du problème et notations .....	298
5.5.2 Discrétisation du problème (1) .....	299
5.5.3 Méthode SOR (Simultaneous Over-Relaxation) pour résoudre un système tridiagonal par blocs .....	302
5.5.4 Programmation structurée de la méthode S.O.R avec accélération de Tchébycheff.....	303

5.6 Cas particulier de l'équation de Poisson .....	304
5.7 Exemple d'application : distribution de la température en régime stationnaire dans une cheminée .....	305
5.8 Premier exemple d'équation parabolique : l'équation de la chaleur à une dimension.....	306
5.8.1 Introduction .....	306
5.8.2 Existence et unicité de solutions de (1) .....	306
5.8.3 Résolution de l'équation de la chaleur à une dimension par des méthodes de différences finies .....	308
5.8.4 Problèmes de stabilité .....	310
5.9 Equations paraboliques linéaires à une variable d'espace avec conditions initiales .....	312
6. <i>Exemple d'application : Phénomène de sustentation par utilisation d'une source de pression. Le patin hydrostatique</i> .....	313
6.1 Présentation du problème .....	313
6.2 Détermination du coefficient $K_s$ .....	314
6.3 Détermination du coefficient $K_q$ .....	315
6.4 Raideur du patin .....	315
6.5 Modélisation du champ de pression .....	316
6.5.1 Introduction .....	316
6.5.2 Les équations de laminage.....	316
6.5.3 Résolution dans le cas où les termes d'accélération sont très petits.....	317
7. <i>Exercices</i> .....	317
8. <i>Programmation Ada</i> .....	319
8.1 Spécification du paquetage DIRICHLET .....	319
8.2 Démonstration du paquetage DIRICHLET .....	319
8.3 Spécification du paquetage NIVEAUX.....	321
8.4 Spécification du paquetage EDP_ELLIPTIQUES .....	321
8.5 Démonstration du paquetage EDP_ELLIPTIQUES.....	322
8.6 Spécification du paquetage EDP_PARABOLIQUES .....	326
8.7 Démonstration du paquetage EDP_PARABOLIQUES .....	326
<b>Bibliographie</b> .....	329
<b>Index</b> .....	333

# Avant-propos

Le but de cet ouvrage est de décrire, sans référence à un quelconque langage de programmation, des méthodes classiques d'analyse numérique, et de mettre à la disposition du programmeur scientifique (mathématicien ou utilisateur des mathématiques) une bibliothèque mathématique écrite en Ada et prête à l'emploi sur compatible PC.

Cet ouvrage pourra être très utile aux auditeurs du cours « programmation scientifique » du cycle B du Cnam ainsi qu'à tout étudiant en licence, maîtrise d'analyse numérique ou école d'ingénieurs. De manière plus générale il intéressera tout utilisateur de l'outil mathématique sur ordinateur.

L'informatique et les mathématiques sont ici des outils. C'est un ouvrage de *programmation mathématique* dans la lignée des « Numerical Recipes ».

Cet ouvrage n'a pas la prétention d'être un ouvrage d'informatique. Son but est de montrer que le langage Ada peut être très efficace pour résoudre des problèmes numériques (créneau non encore exploité). J'espère inciter le programmeur scientifique en Pascal ou autre à évoluer vers Ada. Il est vrai que la lecture sera plus aisée pour le lecteur connaissant Pascal (c'est théoriquement le cas pour tout étudiant en école d'ingénieur).

Dans le premier chapitre on fait un rapide tour d'horizon du langage Ada en mettant l'accent sur les grandes qualités de ce langage (essentiellement l'apprentissage d'une « bonne hygiène de programmation »). Le lecteur devra consulter la bibliographie pour plus de détails. Ce chapitre se termine par une brève description du compilateur Open Ada sur PC, distribué en France par Cerus Informatique.

Dans le deuxième chapitre on décrit une série de paquetages de base permettant de manipuler des fonctions d'une ou plusieurs variables réelles et d'exploiter les spécificités du PC pour travailler en mode graphique. Cette bibliothèque est inspirée du logiciel Modulog utilisé par les élèves de classes préparatoires. Elle est adaptée à l'utilisation du compilateur Open Ada sur PC. En général, les ouvrages sur le langage Ada ne se préoccupent pas d'une implémentation particulière d'un compilateur et pourtant il faudra bien travailler sur un type particulier de matériel.

A ma connaissance une telle bibliothèque livrée avec les sources Ada n'est pas disponible dans le domaine public et pourra être très utile au programmeur scientifique.

Dans l'ouvrage, seules les spécifications sont présentées. C'est tout ce que l'utilisateur à besoin de connaître. Le corps de ces paquetages est disponible sur la disquette livrée avec l'ouvrage.

Dans les chapitres 3 à 8 on décrit des méthodes numériques de base.

Pour chaque chapitre j'ai adopté la même philosophie de travail :

- Définir le problème à résoudre.
- Analyser le problème en donnant des résultats théoriques sur l'existence et l'unicité de solutions.
- Décrire les méthodes numériques classiques d'intérêt pédagogique et présenter des méthodes améliorées.

Par exemple, pour le calcul des intégrales on a les méthodes classiques des trapèzes et de Simpson puis celle plus performante de Romberg.

Ces méthodes sont décrites, du point de vue mathématique de façon rigoureuse, mais je n'ai pas jugé utile de réécrire certaines démonstrations que l'on peut trouver dans la littérature. J'utilise des références du type : voir Stoer et Burlisch p. 234.

- Ces méthodes étant destinées à être programmées, l'analyse précédente se termine par une programmation structurée en Français (en *italique* dans le texte) sans référence à un quelconque langage de programmation.

L'intérêt de cette façon de programmer est de ne pas restreindre le public aux seuls programmeurs Ada.

- Une série d'exercice est proposée pour chaque chapitre.
- Enfin le chapitre se termine par la programmation Ada correspondante sous forme de paquetages et d'exemples d'utilisation (en « *courier* » dans le texte). Seules les spécifications des paquetages et des programmes d'exemples figurent sur le manuscrit, l'intégralité se trouvant sur la disquette.

On donne également des exemples d'applications de ces méthodes à des problèmes issus de la mécanique des fluides. Ces exemples sont dus à Guy Aubry, professeur de mécanique des fluides à l'Ensam.

Les sujets traités sont :

*Chapitre 3* — Résolution de systèmes d'équations linéaires et recherche de valeurs propres.

*Chapitre 4* — Résolution de systèmes d'équations non linéaires.

*Chapitre 5* — Mathématiques pour la CAO (Conception Assistée par Ordinateur) : approximations par des courbes de régression, de Bézier et B-Splines ; interpolations de Lagrange et Splines cubiques.

*Chapitre 6* — Calculs d'intégrales et transformée de Fourier rapide (outils de la théorie du signal).

*Chapitre 7* — Résolution d'équations et de systèmes différentiels.

*Chapitre 8* — Résolution d'équations aux dérivées partielles par la méthode des différences finies.

Avec cet ouvrage je pense convaincre le lecteur que le langage Ada est vraiment bien adapté à la programmation mathématique (sécurité nettement supérieure à ce qu'on peut espérer avec les autres langages de programmation, réutilisabilité des composants logiciels, gestion des tableaux très efficace sans avoir à utiliser de pointeurs, ...).

*Remerciements* — Je tiens à remercier le Professeur André Warusfel qui a bien voulu étudier mon manuscrit. C'est pour moi un grand honneur de le publier dans la collection qu'il coordonne : « Logique Mathématiques et Informatique ».

Je remercie également Marcel Nicolas, Directeur de l'ENSAM de Châlons sur Marne qui a mis à ma disposition le matériel nécessaire à la réalisation de ce livre.

Enfin un grand merci à mon ami et collègue Jean Luc Bauchat pour avoir lu et critiqué une première version.

*Contenu de la disquette fournie avec le livre* — Cette disquette est découpée en répertoires comme décrit ci-dessous.

*A: \DIVERS* — Chapitre 2

PREMIERS.ADA	LIST.ADA	POLYGO.ADA	DEM_CART.ADA
DEM_NIV.ADA	DEM_CURV.ADA		

*A: \BIBLIO* — Chapitres 2 à 8

LIO.ADS	CO_CRT.ADS	CRT.ADS	CRT.ADB
CO_MATH0.ADS	MATH0.ADS	MATH0.ADB	CO_MATH1.ADS
MATH1.ADS	MATH1.ADB	LIRE_FCT.ADA	MATH2.ADS
MATH2.ADB	CO_MATH3.ADS	MATH3.ADS	MATH3.ADB
CO_GRAPH.ADS	GRAPH.ADS	GRAPH.ADB	CURVE.ADS
CURVE.ADB	CARTESIE.ADS	CARTESIE.ADB	CO_MATRI.ADS
MATRIX.ADS	MATRIX.ADB	NIVEAUX.ADS	NIVEAUX.ADB
ALGEBLIN.ADS	ALGEBLIN.ADB	SPECTRE.ADS	SPECTRE.ADB
DON_ALGL.ADS	DON_ALGL.ADB	CO_POLY.ADS	POLY.ADS
POLY.ADB	COMPLEXE.ADS	COMPLEXE.ADB	EQUATION.ADS
EQUATION.ADB	SYST_EQU.ADS	SYST_EQU.ADB	REGRESS.ADS
REGRESS.ADB	BEZIER.ADS	BEZIER.ADB	B_SPLINE.ADS
B_SPLINE.ADB	LAGRANGE.ADS	LAGRANGE.ADB	SPLINE.ADS
SPLINE.ADB	DON_PTS.ADS	DON_PTS.ADB	INTEGRAT.ADS
INTEGRAT.ADB	CO_FOURI.ADS	FOURIER.ADS	FOURIER.ADB
CO_QDIF.ADS	QDIF_11.ADS	QDIF_11.ADB	QDIF_P1.ADS
QDIF_P1.ADB	QDIF_1Q.ADS	QDIF_1Q.ADB	QDIF_PQ.ADS
QDIF_PQ.ADB	DIRICHLE.ADS	DIRICHLE.ADB	ELLIPTIC.ADS
ELLIPTIC.ADB	PARABOLI.ADS	PARABOLI.ADB	

*A: \DONNEES* — Données pour les programmes

REGR\_AFF.DAT : Nuage de points pour la régression affine.  
 SPL\_CART.DAT    SPL\_PARA.DAT : Nuages de points pour l'interpolation spline.  
 BILLE.DAT : Données pour BILLE1.ADA et BILLE.ADA.  
 B\_SPLINE.DAT : Nuage de points pour l'interpolation B\_Spline.

*ALGEBLIN* — Chapitre 3

HILBERT.ADA	DEM_PIV.ADA	DEM_LR.ADA	DEM_GJ.ADA
DEM_TRID.ADA	DEM_CHOL.ADA	DEM_JAC.ADA	DEM_REL.ADA
DEM_DEFL.ADA	DEM_VPJA.ADA	DEM_RUTI.ADA	DEM_SOU.ADA
DEM_VAND.ADA	DEM_QR.ADA	DEM_ALGL.ADA	

*EQUATION* — Chapitre 4

DEM_POL.ADA	DEM_EQU.ADA	DEM_SYST.ADA
-------------	-------------	--------------

*INTERPOL* — Chapitre 5

DEM_AFF.ADA	DEM_POL.ADA	DEM_CIRC.ADA	DEM_APPR.ADA
BILLE1.ADA	DEM_BEZI.ADA	DEM_BSPL.ADA	DEM_LAGR.ADA
DEM_SPL.ADA	DEM_APP1.ADA	DEM_APP2.ADA	

*INTEGRAT* — Chapitre 6

DEM_INT.ADA	DEM_FFT.ADA
-------------	-------------

*EQUADIFF* — Chapitre 7

BILLE.ADA	DEM_RK11.ADA	DEM_RKP1.ADA	DEM_RK1Q.ADA
DEM_RKPQ.ADA	DEM_QDIF.ADA		

*EDP* — Chapitre 8

DEM_DIRI.ADA	DEM_ELLI.ADA	DEM_CHAL.ADA	DEM_EDP.ADA
--------------	--------------	--------------	-------------

*Notes* — Les programmes proposés dans ce livre le sont à titre pédagogique. Je ne garantis en aucune manière qu'ils soient exempts d'erreurs. Le lecteur peut les utiliser à ses propres risques et périls. Ni l'auteur ni les éditions Masson ne peuvent être tenus pour responsables des dommages que l'utilisation de ces programmes pourrait occasionner.

En aucun cas ces programmes ne peuvent être utilisés à des fins commerciales ou industrielles.

Open Ada est une marque déposée de Meridian Software.

First Ada est une marque déposée de Alsys.

## CHAPITRE 1

# Le langage Ada avec OpenAda

### 1. Petit historique de Ada

C'est dans les années 1970 que l'on se rend compte des problèmes posés par le trop grand nombre de langages informatiques existants. Cette diversité de langages entraîne une maintenance de plus en plus difficile. En passant d'une machine à l'autre on ne dispose pas nécessairement des mêmes langages et il est parfois nécessaire de tout reprendre à zéro, la fiabilité n'étant pas très bien assurée.

Le langage Ada est donc né du besoin de définir un standard en matière de langage de programmation.

En 1975, le département de la défense Américain (le « *DOD* » = *Department Of Defense*), qui est considéré comme le plus gros consommateur d'informatique au monde, a lancé un concours international en vue de définir un nouveau langage de programmation prenant en compte toutes les améliorations actuelles. Ce langage doit permettre, entre autres, de traiter des problèmes numériques, de gestion, de calculs en temps réel, ainsi que la gestion d'activités parallèles et la programmation de systèmes embarqués (calculateurs sans systèmes d'exploitation, tel un missile guidé). Il doit également permettre la programmation modulaire et la définition de nouveaux types de variables, tout en étant lisible, facile à maintenir et assurer la portabilité maximale.

Les programmes étant de plus en plus ambitieux et nécessitant beaucoup de participants, le travail d'un membre de l'équipe doit pouvoir être lu sans documentation.

En 1979, c'est le projet du Français *Jean Ichbiach*, de chez *CII Honeywell-Bull*, qui a été retenu. Le langage obtenu a reçu le nom de « *ADA* » en hommage à la comtesse de Lovelace, *Adélaïde Augusta Byron* (1815-1851), fille du poète Anglais Lord Byron, qui était l'assistante de *Babbage* et qui est considérée comme le premier programmeur.

Le standard définitif du langage est apparu en 1983 avec la norme *ANSI MIL-STD-1815A* et les premiers compilateurs sont disponibles à partir de 1985 sur différents types de matériels.



C'est la première fois qu'une norme est définie avant l'existence de tout compilateur. Pour les autres langages c'est l'inverse qui s'est produit d'où la diversité des compilateurs C, Fortran, Pascal ...

La validation des compilateurs est très sévère et elle doit être renouvelée tous les 18 mois. Ces restrictions donnent des compilateurs de très bonne qualité.

La norme est révisée tous les dix ans. Le projet actuel de révision est le projet *ADA 9X*. Avec cette nouvelle norme Ada sera véritablement orienté objet. Si Ada est à l'origine de la conception orientée objet, avec le livre de *Grady Booch* : Ingénierie du logiciel avec Ada, ce n'est pas à l'heure actuel un langage orienté objet au même titre que C++, Objective C ou Eiffel.

Actuellement, le département de la défense Américain utilise ce langage de façon quasi exclusive pour tout traitement informatique.

Le langage Ada est fortement inspiré du Pascal et de l'Algol60, mais on y retrouve certains des points positifs d'autres langages tels que Fortran, C, Cobol, Pl1 ou Algol. Il est considéré comme un langage d'usage général.

On ne peut pas vraiment considérer Ada comme un nouveau langage, mais plutôt comme la synthèse des connaissances actuelles sur les langages de programmation. Néanmoins, avec Ada, une nouvelle philosophie de programmation s'est développée .

Si l'apprentissage à un premier niveau est relativement aisé, spécialement pour le programmeur en Pascal, il est beaucoup plus difficile d'en exploiter toutes les subtilités.

Pour la lecture de cet ouvrage une connaissance du langage à un premier niveau est suffisante. Le livre de B. Leguy indiqué en référence est un très bon ouvrage d'initiation au langage Ada. On pourra également consulter le livre de Barnes qui est considéré comme un manuel de base. Dans la bibliographie on trouvera une liste d'ouvrages permettant une étude plus approfondie.

Comme Pascal, Ada est un langage très « typé », c'est-à-dire que toutes les variables utilisées doivent être déclarées au préalable et le programmeur a la possibilité de définir de nouveaux types de variables.

On dispose en Ada d'un grand éventail de types de variables dont les tableaux, les enregistrements et les pointeurs. Mais Ada est encore plus rigoureux que Pascal dans la manipulation des variables. Si, par exemple, en Pascal on peut combiner dans une expression travaillant sur des réels une variable de type INTEGER et une de type REAL, une telle opération n'est pas possible en Ada sans une conversion préalable de la variable de type INTEGER en type FLOAT (c'est le type réel de Ada).

Comme Pascal et C, Ada permet la récursivité.

Le langage Ada assure une meilleure sécurité que le Pascal tout en étant aussi puissant que le C (on a la possibilité de travailler au niveau de la machine).

La gestion des erreurs, grâce à la notion *d'exception*, est beaucoup plus efficace qu'en Pascal ou C.

L'écriture de gros programmes découpés en unités compilées séparément se fait de manière beaucoup plus sûre.

Toutes ces propriétés facilitent de beaucoup la maintenance des programmes qui constituent la partie « chère » des logiciels.

Parmi les nouveautés, on dispose de la notion de « *tâche* ». Un programme Ada peut consister en plusieurs tâches s'effectuant de manière asynchrone et en parallèle. La synchronisation des tâches est possible grâce à la notion de « *rendez-vous* ». Ce qui permet la programmation en « *temps réel* ».

La portabilité du langage, c'est-à-dire la qualité d'un logiciel à transférer son savoir faire sur un autre type de matériel à un coût très inférieur à sa réécriture, est également un des aspects intéressants du langage. Cette portabilité est assurée par les normes ANSI (en 1983) et ISO (en 1987).

## 2. Une présentation rapide du langage Ada

### 2.1 Le manuel de référence

Le langage est défini par son manuel de référence qui est contrôlé par une commission de validation. Tout compilateur doit être en parfait accord avec ce manuel. Ce qui dépend de l'implémentation de Ada sur un type de machine doit être décrit dans l'annexe F de la documentation fournie avec le compilateur.

Le manuel de référence étant une norme, il est identique pour tout type de compilateur. Les remarques particulières doivent être repérables.

Une mention à ce manuel se fera sous la forme [LRM 4.3.2 (11)], pour chapitre 4, section 3, paragraphe 2 et alinéa 11.

Souvent les messages d'erreur d'un compilateur feront une mention de ce type là.

Il existe une version Française de ce manuel, c'est : Manuel de Référence du langage de programmation Ada ALSYS : 29, Av. Lucien René Duschesne, 78170 La Celle Saint-Cloud.

### 2.2 Un survol du langage

Dans le survol qui suit, nous présentons les éléments du langage Ada qui peuvent être très facilement assimilés par le programmeur en Pascal ou C.

Pour une étude plus détaillée on pourra consulter les ouvrages de la bibliographie.

Un des points forts du langage est le *typage des données* qui assure une fiabilité optimale. Selon le cabinet américain de « consulting » Reifer, portant sur 41 projets et 15 millions de lignes de codes, on détecte en moyenne 4 à 13 erreurs pour 1000 lignes avec des langages différents de Ada. Ce taux est ramené à 3 à 5 erreurs pour Ada (01 Informatique, 10/01/92).

Le langage reprend le typage de Pascal, mais avec beaucoup plus de rigueur.

On dispose des *types énumératifs* identiques à ceux de Pascal, par exemple, on peut écrire :

```
type WEEK_END is (SAMEDI, DIMANCHE) ;
type COULEUR is (ORANGE, VERT, ROUGE) ;
type FRUIT is (ORANGE, POMME) ;
```

et il n'y aura pas de confusion entre une variable de type FRUIT prenant la valeur ORANGE et une variable de type COULEUR prenant la valeur ORANGE. Le compilateur vérifie avec beaucoup de soins la cohérence des affectations.

On dispose également de la possibilité de distinguer de façon logique des données ayant la même représentation. Ainsi, en écrivant :

```
type DOLLAR is range 0..100_000 ;
type FRANC is range 0..100_000 ;
```

on ne risque pas d'ajouter des dollars et des francs, une telle opération sera rejetée à la compilation.

On dispose également de la possibilité de préciser de façon très stricte le domaine de validité de certaines variables réelles, on peut par exemple écrire :

```
type DISTANCE is digits 3 range 0.0..3245.0 ;
type VOLT is delta 0.01 range 0.0..220.0 ;
```

pour des distances calculées avec trois décimales et des volts, avec un pas de 0.01. Tout dépassement des valeurs permises déclenchera une erreur.

On a également la possibilité de définir des tableaux homogènes ou hétérogènes, avec une dimension fixée, quelconque ou précisée à l'exécution.

Par exemple, on peut écrire :

```
type VECTEUR is array(POSITIVE range 1..20) of FLOAT ;
```

pour des vecteurs à 20 composantes, ou

```
type VECTEUR is array(POSITIVE range <>) of FLOAT ;
```

pour des tableaux dont la dimension sera précisée dans le programme appelant avec une déclaration du type :

```
v : VECTEUR(2..2) ;
```

Le type *access* permet de déclarer des variables de type *pointeur*. L'adresse d'un objet est contenue dans une variable et cette adresse mènera forcément vers un objet du bon type. On ne peut pas pointer sur n'importe quoi comme en C.

L'utilisation des types permet de masquer les détails d'implémentation, par exemple le concept abstrait de fichier est suffisant. C'est la notion de *types privés* qui permet de gérer cela au mieux. Cette manipulation abstraite des objets permet d'avoir une structure logique très stricte ce qui facilite le côté maintenance.

Un autre exemple significatif est celui des piles. On peut définir un type privé comme suit :

```
type PILE is private ;
```

et l'utilisateur n'a pas besoin de savoir comment est défini ce genre d'objet, s'il sait que c'est défini à l'aide de tableaux on a le risque de le voir trafiquer l'élément numéro 56, alors qu'on ne manipule que le sommet de la pile. Ce genre de comportement est catastrophique à la maintenance. Sachant qu'il dispose d'un type pile, l'utilisateur a seulement besoin de connaître les opérations disponibles sur ces objets, par exemple il peut avoir à sa disposition les procédures suivantes :

```
procedure EMPILER(TOTAL : in FLOAT ; P : in out PILE) ;
procedure DEPILER(TOTAL : out FLOAT ; P : in out PILE) ;
```

le changement d'implémentation de la pile et des procédures associées ne doit pas changer une ligne du programme de l'utilisateur.

Parmi les choses nouvelles, on dispose de la notion de *tâche* pour les applications en temps réel et les activités parallèles. Il n'est pas nécessaire de connaître cette notion pour lire cet ouvrage.

On dispose aussi de la notion de *type dérivé* qui permet de définir des sous-types d'un type donné. Les type « fils » héritant des propriétés des types « parents ».

Le langage est évidemment structuré. Les instructions structurées ayant une indication de fin plus pratique que les BEGIN...END ; de Pascal. Par exemple, on a :

```
loop          if...      case...
...          ...        ...
end loop;    end if;    end case;
```

Les sous-programmes peuvent être traités séparément ou regroupés dans des *paquetages*. Ils doivent être écrits de façon à être facilement réutilisables.

Pour ce qui est des paramètres formels des procédures, on dispose des modes *in*, *out* et *in out*. En mode *in* la variable doit avoir été initialisée et ne peut être modifiée dans la procédure. En mode *out* la variable n'est pas nécessairement initialisée et aura obtenu une valeur en fin de sous-programme et le mode *in out* permet de modifier une variable initialisée.

Pour les fonctions, seul le mode *in* est autorisé.

On dispose de la notion *d'agrégat* qui permet de placer les paramètres effectifs d'une procédure dans n'importe quel ordre sans risque d'erreur. Par exemple si on a la déclaration :

```
procedure P1(n in : POSITIVE ; x in out FLOAT ; y out FLOAT) ;
```

on pourra l'appeler de la façon suivante :

```
P1(y => 6.8, n => 2, x => 8.9) ;
```

cette façon de procéder présente l'avantage de ne pas être obligé de retourner à la documentation pour connaître l'ordre logique des paramètres.

On a également la possibilité de *surcharger les opérateurs existants*. On peut par exemple définir une opération d'addition notée + qui agit sur des objets de type matrice. On écrira :

```
function "+" (A, B : MATRICE) return MATRICE is
```

le type MATRICE ayant été préalablement défini comme un tableau à deux indices.

Quand on définit un paquetage, c'est-à-dire une bibliothèque de sous-programmes, on peut séparer la *spécification*, qui est la section de déclaration, de son corps principal où les procédures et fonctions sont effectivement traitées. Pour l'utilisateur l'interface seule compte, les détails d'implémentation peuvent ne pas l'intéresser. Cela présente aussi l'avantage de pouvoir travailler en équipe pour de gros projets. Si deux équipes différentes ont donné le même nom à des procédures différentes, on peut les utiliser sans accéder au code source, en fonction de la nature des paramètres le compilateur se débrouille pour savoir de quoi il s'agit.

On dispose aussi de la notion de sous-unité qui permet le développement hiérarchique de gros programmes en reportant à plus tard l'écriture de certaines procédures. Pour ce faire on utilise la clause *separate*. La sous-unité sera compilée à part sans avoir à recompiler l'unité mère.

La notion *d'unité générique* permet de définir des procédures agissant sur des objets dont on ne connaît pas la nature. Par exemple il est inutile d'écrire une procédure d'échange de deux variables pour chaque type de variable, il suffit d'écrire, dans la partie spécification de l'unité :

```
generic
  type VARIABLE is private ;
  procedure ECHANGER(x, y : in out VARIABLE) ;
puis dans le corps du paquetage :
  procedure ECHANGER(x, y : in out VARIABLE) is
  AUX : VARIABLE := x ;
  begin
    x := y ;
    y := AUX ;
  end ;
```

On précise alors sur quel type d'objet va agir ECHANGER en l'instanciant de la manière suivante :

```
procedure ECHANGER_REELS is new ECHANGER(FLOAT) ;
procedure ECHANGER_ENTIERS is new ECHANGER(INTEGER) ;
```

On peut contrôler les erreurs d'exécution avec la notion d'exception. Ce contrôle d'erreur augmente la taille du code exécutable de 10 à 20%. Pour des applications critiques du type embarqué où on compte le moindre octet on peut demander au compilateur de ne plus générer le code relatif aux vérifications. On procède de la sorte pour le produit final.

On a aussi la possibilité de programmer à un bas niveau en gérant les interruptions.

Enfin on dispose, grâce au « *pragma interface* » de la possibilité d'utiliser des sous-programmes écrits dans un autre langage (Assembleur, C, Fortran, Pascal, ...).

### 2.3 Exemples d'utilisations de Ada

Le langage Ada est utilisé de façon systématique par le département de la défense Américain et diverses armées dans le monde.

On trouve des utilisations dans les domaines cités ci dessous :

- Avionique. Avec les réalisations suivantes :

Aux Etats Unis :

Contrôle du trafic aérien, pour le Federal Aviation Authority ;

Contrôle des cockpits des 767 et 777 pour Boeing ;

Simulateur de vol, pour Ford Aerospace.

Au Canada :

Contrôle du trafic aérien, pour CSC/IBM.

En Europe :

Contrôle du trafic aérien en Europe, pour Eurocontrol ;

Contrôle du trafic aérien, pour Thompson SDC ;

Contrôle d'équipements au sol, pour Plessey Avionics ;

.....

- L'espace, avec les applications suivantes :

Logiciel d'environnement et de gestion des données pour la station spatiale Freedom ;

Contrôle de robots pour la NASA ;

Station de transmission pour les communications spatiales pour Jet Propulsion Lab ;

Logiciels pour Ariane V ;

...

- L'industrie, avec les applications suivantes :  
 Processeurs électroniques pour Lockheed ;  
 Logiciel d'aide à l'extraction de pétrole pour BP Oil, Total et Schlumberger Oil ;  
 Exploration géophysique du pétrole pour Shell Oil ;  
 Système de transmission dans le domaine des chemins de fer pour la CSEE ;  
 Systèmes d'instrumentation, pour Schlumberger ;  
 Automate de simulation pour 3IP/Citroën ;

....

On a également des applications dans les domaines de la banque et des assurances, dans le domaine commercial, médical, des télécommunications, de l'électronique, de l'intelligence artificielle et du nucléaire.

### 3. Utilisation du compilateur OpenAda sur micro-ordinateur

Ce compilateur, qui est validé, permet de travailler sur tout micro-ordinateur du type IBM-PC (de préférence un 386 ou 486) sous environnement MS-DOS. De plus, il est livré avec un environnement de programmation à base de menus déroulants, une bibliothèque mathématique et une bibliothèque graphique.

La dernière version permet de travailler en mode réel, 16 bits ou 32 bits. De plus une version sous Windows est disponible. Ces versions sont proposées avec un « debugger ». Dans ce qui suit on s'intéresse seulement à la version DOS.

Il est distribué par Meridian Software, cette société étant représentée en France par Cerus-Informatique.

On peut l'utiliser directement en ligne de commande, sous DOS, ou depuis l'environnement de programmation ACE (Ada Compiler Environnement). Nous décrivons tout d'abord son utilisation en ligne de commande, pour bien comprendre les étapes de création d'une bibliothèque, de compilation et d'édition de liens.

Une fois l'installation réalisée sur le disque C (ou D, ...), on dispose du répertoire nommé Ada (ou autre), contenant les sous-répertoires *ACE*, *BIN*, *PACLIB*, *INASM*, *MSC*, *ADAUTIL*, *DOSENV*, *MATH*, *AGUL*, *SBOOCH* et *TEST*. Pour la version 386 on a de plus les répertoires *ADAUTIL.32*, *DOSENV.32* et *PACLIB.32*. Ils correspondent à la version 32 bits.

Le répertoire ACE contient les fichiers nécessaires à l'environnement de programmation. BIN contient les programmes exécutables indispensables au compilateur. PACLIB contient les bibliothèques et fichiers de configuration. INASM contient des exemples de programmes interfacés avec le langage assembleur ainsi qu'une documentation sur le sujet. MSC contient des exemples de programmes interfacés avec le langage C de Microsoft. Math contient une

bibliothèque mathématique avec un exemple d'utilisation. Agul (Ada Graphic Utility Library) contient une bibliothèque graphique avec un exemple d'utilisation. Et Test est un répertoire contenant des programmes de démonstration.

Avant d'utiliser un compilateur Ada on doit créer une *bibliothèque* (en Anglais *library*). Avec OpenAda, on utilise le programme *NEWLIB.BAT* qui crée dans le répertoire de travail le fichier *ADA.LIB* et le répertoire *ADA.AUX*. C'est dans *ADA.AUX* que seront stockés les codes objets correspondants aux sous-programmes compilés. Cette création se fait une seule fois dans un répertoire donné.

Pour plus de détails sur la gestion d'une telle bibliothèque, on se reportera au chapitre 8 du « Compiler User's Guide » fourni avec la documentation.

Cette bibliothèque est nécessaire pour produire le code objet correspondant à un code source donné. Le processus de compilation est schématisé ci-dessous :

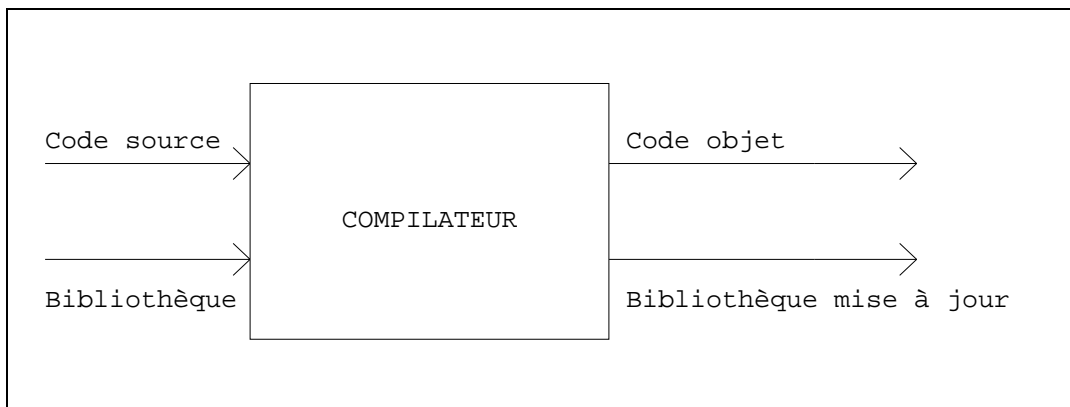


Figure 1.1

Cela permet une compilation séparée et dépendante, c'est-à-dire qu'on garde une trace de toutes les compilations et des liens avec les autres bibliothèques.

Si on compare avec ce qui se passe en Fortran où on a une compilation séparée, mais indépendante, on a une meilleure sécurité. En Fortran on peut écrire dans un programme `CALL Sub(i,j,k)`, alors que la subroutine `Sub`, se trouvant dans un autre fichier, a été déclarée comme agissant seulement sur deux paramètres. Une telle erreur ne sera pas décelée à la compilation. Alors qu'en Ada, si on a compilé la procédure `Sub` agissant sur deux paramètres entiers `i` et `j`, l'appel `Sub(i,j,k)` dans n'importe quelle procédure ou fonction provoquera une erreur à la compilation.

L'avantage de la compilation séparée est de pouvoir traiter et tester séparément les divers morceaux d'un programme.

Pour compiler le programme `Simple.Ada`, il suffit d'écrire : `ADA Simple.Ada`.

Puis l'édition de liens se fait avec : `BAMP Simple`, où `BAMP` signifie *Build Ada Main Program*. Et si tout se passe normalement, le fichier `Simple.Exe` a été créé. `Bamp` doit être suivi du nom du programme principal qui peut être différent de celui du fichier qui le contient.

On dispose avec ADA et BAMP de plusieurs options. On se reportera, pour plus de détails, au fascicule « Compiler User's Guide PC DOS » faisant partie de la documentation.

En fait, il est beaucoup plus agréable d'utiliser l'environnement de programmation ACE, d'utilisation relativement simple.

Pour l'utiliser, il suffit de taper soit ACE, soit ACE suivi du nom du programme à charger.

Cet éditeur peut être configuré à sa convenance. Cette configuration se fait en utilisant le sous-menu *CHANGE BINDINGS* du sous-menu *OPTIONS* du menu principal. On peut y accéder directement avec la combinaison de touches ALT 2.

On a également la possibilité d'utiliser une « souris ». Le bouton droit permet d'accéder au menu principal et ensuite il suffit de savoir lire.

Dans le tableau 1.1 nous décrivons les options par défaut et celles que nous avons choisies.

On peut aussi utiliser le compilateur *FirstAda* de *ALSYS*. Ce compilateur est utilisable sous MS-DOS en mode 32 bits, mais nécessite au minimum un 386 avec 3 Mo. de mémoire vive. Il est livré avec une bibliothèque mathématique mais pas de bibliothèque graphique. Ce compilateur est d'aussi bonne qualité que celui de Meridian Software, mais plus onéreux et de maniement plus complexe.

Les paquetages décrits dans le chapitre suivant et utilisés dans les autres chapitres sont utilisables avec *OpenAda*. Pour les utiliser avec *FirstAda*, il faudra en réécrire certaines parties (tout ce qui est relatif au clavier et à l'écran).



ACTION	TOUCHES PAR DEFAUT	TOUCHES CHOISIES
— Menu Principal	ESC	ESC
— Quitter un menu	ESC	ESC
— Charger un fichier	CTRL N	F3
— Sauvegarder un fichier	ALT W et <ENTER>	F2 et <ENTER>
— Le sauver sous un autre nom	ALT W, Nom et <ENTER>	F2, Nom et <ENTER>
— Créer une bibliothèque		ALT F8
— Compiler le fichier édité		F9
— Faire l'édition de liens		ALT F9
— Lancer l'exécution du programme édité		CTRL F9, <ENTER>
— Exécuter un programme.		CTRL F9, Nom
— Aide sur l'éditeur	ALT H	ALT H
— Retour temporaire au DOS	ALT F2	SHIFT F10
— Chercher une chaîne de caractères et la remplacer par une autre	CTRL R	CTRL R
— Effacer une ligne	CTRL D	CTRL Y
— Quitter sans sauvegarder	ALT Z	ALT X
— Quitter en sauvegardant	CTRL Z	CTRL X
— Marquer le début d'un bloc	F1	F1
— Marquer la fin d'un bloc	Flèches de déplacement	Flèches de déplacement
— Copier le bloc marqué dans un « buffer »	F3	CTRL C
— Insérer à la position du curseur un bloc du « buffer »	F7	F7
— Imprimer le bloc marqué		CTRL P
— Supprimer le bloc marqué	F3	ALT D

Tableau 1.1

## CHAPITRE 2

# Bibliothèque mathématique

### 1. Introduction

La bibliothèque décrite dans ce chapitre est formée des paquetages *CRT*, *MATH0*, *MATH1*, *MATH2*, *GRAPH* et *MATH3*. A chacun de ces paquetages est associé le paquetage *COMMON\_NOM\_DU\_PAQUETAGE* où sont déclarés des constantes, types, variables globales et exceptions.

La conception de cette bibliothèque destinée à faciliter l'écriture de programmes mathématiques est inspirée du logiciel Modulog (écrit en Pascal), qui est produit par l'ALE Sup. sous la responsabilité d'Hervé Lehning et Robert Rolland, et du livre de Ducamp et Reverchon : *Mathématiques en Turbo-Pascal*, Volume 1 : Analyse.

Tout programme utilisant les paquetages *CRT*, *MATH0*, et *COMMON\_MATH1* (par exemple) devra commencer par :

```
with CRT, MATH0, COMMON_MATH1 ;  
use CRT, MATH0, COMMON_MATH1 ;
```

L'utilité de chacun de ces paquetages est décrite ci-dessous.

- Dans le paquetage *CRT*, on trouve des procédures et fonctions destinées à faciliter la gestion du clavier et de l'écran. Il est inspiré de l'unité *CRT* de Turbo-Pascal. Ce paquetage utilise les interruptions DOS.
- Dans le paquetage *MATH0*, on trouve des procédures destinées à faciliter la communication entre un programme et son utilisateur, ainsi que des procédures permettant de gérer les erreurs d'entrées-sorties.
- Dans le paquetage *MATH1*, on définit des fonctions mathématiques qui ne sont pas directement disponibles avec Ada. Ce paquetage utilise le paquetage non générique *MATH\_LIB* fourni avec OpenAda.

On pourrait aussi écrire une version générique qui aurait l'avantage de permettre à l'utilisateur de jouer avec les précisions. Pour la conception d'un tel paquetage on pourra consulter le manuel « Math Library » fourni avec la documentation du compilateur OpenAda. Dans ce manuel on décrit le paquetage

GENERIC\_ELEMENTARY\_FUNCTIONS et des exemples d'utilisation sont donnés dans le répertoire MATH.

Dans MATH1 on travaille directement avec le type FLOAT par souci de simplicité.

- Avec le paquetage *MATH2*, on définit un nouveau type de variable, le type FONCTION, qui permet d'entrer au clavier des fonctions de une à plusieurs variables (ce nombre est défini en constante) et de les évaluer en fonction des paramètres effectifs donnés.
- Dans le paquetage *GRAPH*, on trouve des procédures permettant d'initialiser un mode graphique, de dessiner des points et des segments en mode graphique et de revenir au mode texte. Il est inspiré de l'unité Graph de Turbo-Pascal et utilise les interruptions DOS.
- Avec le paquetage *MATH3*, on dispose de procédures permettant de travailler en mode graphique dans une fenêtre graphique utilisateur en coordonnées réelles. Elle sera très utile pour le tracé de courbes.

Pour utiliser les entiers du type LONG\_INTEGER, on instancie le paquetage INTEGER\_IO comme suit :

```
with TEXT_IO ;
package LIO is new TEXT_IO.INTEGER_IO(LONG_INTEGER) ;
```

Dans les paragraphes qui suivent seules les spécifications des paquetages sont présentées ainsi qu'une documentation. Les codes sources des corps correspondants se trouvant sur la disquette fournie avec l'ouvrage.

## 2. Le paquetage CRT

Ce paquetage permettant la gestion du clavier et de l'écran n'est évidemment pas facilement portable.

La version décrite dans ce paragraphe utilise le paquetage *INTERRUPT* de OpenAda. Ce paquetage permet de manipuler les interruptions BIOS et DOS. L'action de chacune des interruptions utilisée est décrite en commentaire dans le code du corps de CRT.

Pour écrire une version équivalente avec FirstAda, on peut utiliser le paquetage *INTERRUPT\_MANAGER* ou le paquetage *DOS*. Pour plus de détails, on se reportera à la documentation : Application Developer's Guide and Appendix F.

Pour certaines procédures, on peut se passer des interruptions en utilisant les séquences d'échappement ANSI.

### 2.1 Spécifications des paquetages *COMMON\_CRT* et *CRT*

```
package COMMON_CRT is
subtype COLONNE_TEXTE is INTEGER range 0..79 ;
subtype LIGNE_TEXTE is INTEGER range 0..24 ;
end COMMON_CRT ;
```

```
with COMMON_CRT ;
use COMMON_CRT ;
package CRT is
```

```
procedure CURSOR(B : in BOOLEAN) ;
procedure CLREOL ;
procedure CLRSCR ;
```

```
procedure GOTOXY(C : in COLONNE_TEXTE ; L : in LIGNE_TEXTE) ;
function READKEY return CHARACTER ;
procedure PUT_CHAR(C : in CHARACTER) ;
procedure PUT_STRING(Msg : in STRING) ;
procedure PUT_LINE_STRING(Msg : in STRING) ;
function WHERE_X return COLONNE_TEXTE ;
function WHERE_Y return LIGNE_TEXTE ;
end CRT ;
```

## **2.2 Documentation du paquetage CRT**

### **2.2.1 La procédure CLREOL**

Cette procédure permet d'effacer tous les caractères d'une ligne à partir de la position du curseur.

### **2.2.2 La procédure CLRSCR**

Cette procédure permet d'effacer l'écran.

### **2.2.3 La procédure GOTOXY**

Permet de se positionner en colonne C et ligne L.

### **2.2.4 La fonction READKEY**

Permet de lire un caractère au clavier sans écho à l'écran.

### **2.2.5 Les fonctions WHERE\_X et WHERE\_Y**

Ces fonctions renvoient, respectivement, le numéro de colonne et de ligne de la position actuelle du curseur.

### **2.2.6 La procédure CURSOR**

Cette procédure permet d'inhiber (B => FALSE) ou non (B => TRUE) le curseur.

### **2.2.7 La procédure PUT\_CHAR**

Cette procédure permet un affichage plus rapide et plus pratique d'un caractère, que la procédure PUT de TEXT\_IO.

### **2.2.8 La procédure PUT\_STRING**

Permet d'afficher une chaîne de caractères en utilisant la procédure précédente et sans passer à la ligne suivante.

### **2.2.9 La procédure PUT\_LINE\_STRING**

Même effet que la précédente en passant à la ligne suivante après affichage.

### **2.2.10 Partie initialisation de CRT**

```
begin
  CLRSCR ;
  CURSOR(TRUE) ;
end CRT ;
```

## 3. Le paquetage MATH0

### 3.1 Spécifications des paquetages *COMMON\_MATH0* et *MATH0*

La variable IMP, de type fichier texte (i. e. fichier ASCII) est utilisée dans la procédure MODE\_AFFICHAGE.

```

with TEXT_IO ;
package COMMON_MATH0 is
  subtype HEURE is INTEGER range 0..23 ;
  subtype MINUTE is INTEGER range 0..59 ;
  subtype SECONDE is INTEGER range 0..59 ;
  subtype CENTIEME is INTEGER range 0..99 ;
  type TIME is record
    H : HEURE := 0 ;
    M : MINUTE := 0 ;
    S : SECONDE := 0 ;
    C : CENTIEME := 0 ;
  end record ;
  type PERIPHERIQUE_DE_SORTIE is (ECRAN, IMPRIMANTE, FICHER) ;
  CHOIX_AFFICHAGE : PERIPHERIQUE_DE_SORTIE := ECRAN ;
  IMP : TEXT_IO.FILE_TYPE ;
end COMMON_MATH0 ;

with TEXT_IO, COMMON_MATH0 ;
use COMMON_MATH0 ;
package Math0 is

  procedure RANDOMIZE ;
  function RANDOM return FLOAT ;
  function RANDOM(n : in INTEGER) return NATURAL ;
  function MAJUSCULE(C : in CHARACTER) return CHARACTER ;
  function CHAINE_ENTIERE(n : in INTEGER) return STRING ;
  function CHAINE_REELLE(x : in FLOAT) return STRING ;
  procedure GET_LINE(n : out INTEGER) ;
  procedure GET_LINE(x : out FLOAT) ;
  procedure GET_LINE(c : out CHARACTER) ;
  procedure GET_LINE(Msg : out STRING) ;
  procedure PAUSE ;
  procedure BIP ;
  function LIRE_REPONSE(Msg : in STRING) return BOOLEAN ;
  function TRUNC(x : in FLOAT) return INTEGER ;
  function FRAC(x : in FLOAT) return FLOAT ;
  procedure CHRONOMETRE(Top : in INTEGER) ;
  procedure ENTRER_REEL(x : out FLOAT ; Msg : in STRING := " ") ;
  procedure ENTRER_ENTIER(n : out INTEGER ; Msg : in STRING := " ") ;
  procedure ENTRER_ENTIER_LONG(n : out LONG_INTEGER ;
    Msg : in STRING := " ") ;
  procedure ENTRER_REEL_BORNE(rMin, rMax : in FLOAT ; x : out FLOAT ;
    Msg : in STRING := " ") ;
  procedure ENTRER_ENTIER_BORNE(nMin, nMax : in INTEGER ;
    n : out INTEGER ; Msg : in STRING := " ") ;
  procedure ENTRER_ENTIER_LONG_BORNE(nMin, nMax : in LONG_INTEGER ;
    n : out LONG_INTEGER ; Msg : in STRING := " ") ;

```

```

procedure ECHANGER_REELS(x, y : in out FLOAT) ;
procedure ECHANGER_ENTIERS(p, q : in out INTEGER) ;
procedure LECTURE_REPERTOIRE(NOM_REPERTOIRE : in STRING) ;
procedure LIRE_FICHER(FICHER_DONNEES : in out TEXT_IO.FILE_TYPE ;
                     NOM_DE_FICHER : in out STRING ;
                     NOM_DE_REPERTOIRE : in STRING := "A:\DONNEES") ;
function FICHER_EXISTE(NOM_DE_FICHER : in STRING) return BOOLEAN ;
procedure CREER_FICHER(FICHER_DONNEES : in out TEXT_IO.FILE_TYPE ;
                      NOM_DE_FICHER : in out STRING ;
                      NOM_DE_REPERTOIRE : in STRING := "A:\DONNEES") ;
procedure MODE_AFFICHAGE ;
procedure VAL(Msg : in STRING ; DEBUT, FIN : in NATURAL ;
              i : in out NATURAL ; Err : in out NATURAL) ;
procedure VAL(Msg : in STRING ; DEBUT, FIN : in NATURAL ;
              x : in out FLOAT ; Err : in out NATURAL) ;
function LONGUEUR_CHAINE(Msg : in STRING) return NATURAL ;
end MATH0 ;

```

## 3.2 Documentation du paquetage MATH0

Le corps de ce paquetage utilise les paquetages INTERRUPT, ERRORS et FILE\_IO livrés avec OPEN\_ADA.

Dans ERRORS est défini un type énuméré EXTENDED\_ERRORS permettant de gérer des erreurs DOS.

Le paquetage FILE\_IO permet la gestion de fichiers sous DOS.

### 3.2.1 La procédure GET\_TIME

Cette procédure donne l'heure maintenue par le système d'exploitation.

### 3.2.2 La procédure RANDOMIZE

Cette procédure permet d'initialiser un générateur de nombres pseudo-aléatoires. Elle doit être appelée avant les fonctions RANDOM. Elle utilise la procédure GET\_TIME qui permet de connaître l'heure maintenue par le système d'exploitation.

### 3.2.3 Les fonctions RANDOM

RANDOM sans argument retourne un réel aléatoire entre 0.0 et 1.0 et RANDOM(n) retourne un entier aléatoire entre 0 et  $n - 1$ , pour  $n \leq 1$  et 0 pour  $n < 1$ .

Voir le paragraphe (4.1) du chapitre sur le calcul numérique des intégrales pour un algorithme de construction de nombres pseudo-aléatoires.

### 3.2.4 La fonction MAJUSCULE

Permet de convertir un caractère compris entre 'a' et 'z' en majuscule.

### 3.2.5 Les procédures GET\_LINE

Ces procédures permettent l'entrée, au clavier, d'un entier, d'un réel, d'un caractère ou d'une chaîne de caractères. Si Ch est une chaîne de 78 caractères, il suffit d'écrire GET\_LINE(Ch) ; et si à l'exécution on entre abc suivi d'un retour chariot alors Ch est la chaîne formée des caractères a, b et c suivis de 75 blancs. Pour les entiers et les réels on utilisera de préférence les procédures ENTRER\_ENTIER, ... (Cf § 3.2.10), où on gère les erreurs de type et de contrainte.

### 3.2.6 La fonction LIRE\_REPONSE

Cette fonction permet d'associer à la question définie par `Msg` une valeur booléenne qui aura la valeur `TRUE` si la réponse est 'o' ou 'O' et `FALSE` si la réponse est 'n' ou 'N'.

On pourra l'utiliser comme suit :

```
if LIRE_REPONSE("Un autre essai? ") then
    .....
end if ;
```

### 3.2.7 La procédure PAUSE

Cette procédure permet de marquer une pause pendant l'exécution d'un programme, la frappe d'une touche quelconque permet ensuite de continuer.

### 3.2.8 La procédure BIP

Cette procédure pourra être utilisée pour signaler une erreur en émettant un son court.

### 3.2.9 La procédure CHRONOMETRE

Cette procédure permet de calculer le temps nécessaire à la réalisation d'une partie, ou de la totalité, d'un programme. Avant la première instruction relative à cette partie de programme, on écrira `CHRONOMETRE(0)` ; ce qui déclenche le chronomètre, puis après la dernière instruction, on écrira : `CHRONOMETRE(1)` ; ce qui arrête le chronomètre et affiche le temps écoulé en secondes et centièmes de secondes.

### 3.2.10 Les procédures d'entrée des entiers et des réels

Ces procédures remplacent les instructions `GET(n)` et `GET(x)` où `n` est une variable de type entier et `x` de type réel.

L'intérêt est de signaler l'erreur si la variable introduite n'est pas du bon type ou dépasse les bornes autorisées en évitant l'arrêt de l'exécution du programme. Le nom de chacune des ces procédures est assez significatif pour comprendre de quoi il s'agit.

### 3.2.11 Les fonctions CHAINE\_ENTIERE et CHAINE\_REELLE

Ces fonctions permettent de convertir un entier (Resp. un réel) en chaîne de caractères.

### 3.2.12 La procédure LECTURE\_REPERTOIRE

Cette procédure donne une liste de fichiers contenus dans le répertoire dont le nom et le chemin d'accès sont indiqués par la variable `NOM_REPERTOIRE`, de type `STRING`. Elle équivaut à l'utilisation de l'instruction `DIR` du `DOS`. Elle utilise les procédures `FIND_FIRST` et `FIND_NEXT` de `FILE_IO`.

### 3.2.13 La procédure LIRE\_FICHER

Cette procédure permet d'ouvrir en lecture un fichier texte dans lequel sont stockées des données. Ce fichier est recherché dans le sous-répertoire `NOM_DE_REPERTOIRE` (par défaut le répertoire `DONNEES` du disque `A`). A l'exécution on entre `NOM_DE_FICHER` et le fichier cherché aura pour nom, `NOM_DE_REPERTOIRE\NOM_DE_FICHER`. En cas de non existence de ce fichier l'erreur est signalée et on peut recommencer.

Au préalable, la liste de tous les fichiers de `NOM_DE_REPERTOIRE` est affichée, ce qui permet de retrouver plus facilement le fichier cherché.

**3.2.14 La procédure FICHER\_EXISTE**

Cette procédure teste si le fichier que l'on veut créer existe déjà, si c'est le cas on demande confirmation pour en effacer le contenu.

**3.2.15 La procédure CREER\_FICHER**

Cette procédure permet de créer un fichier de type texte.

**3.2.16 La procédure MODE\_AFFICHAGE**

Cette procédure permet d'afficher les résultats d'un programme sur:

- l'écran, pour CHOIX\_AFFICHAGE = ECRAN ;
- l'imprimante, pour CHOIX\_AFFICHAGE = IMPRIMANTE ;
- un fichier, pour CHOIX\_AFFICHAGE = FICHER (à l'exécution, il suffit d'entrer NOM\_DE\_FICHER).

Un programme utilisant cette procédure devra débiter par :

```
MODE_AFFICHAGE ;
```

et se terminer par :

```
CLOSE(IMP) ;
```

Les instructions PUT et PUT\_LINE, seront alors remplacées par : PUT(IMP, et PUT\_LINE(IMP, .

**3.2.17 La fonction TRUNC**

Calcule la partie entière d'un réel.

**3.2.18 La fonction FRAC**

Calcule la partie fractionnaire d'un réel.

**3.2.19 Les procédures de permutation de deux entiers ou deux réels**

Ces procédures permettent de permuter deux entiers ou deux réels.

**3.2.20 Les procédures VAL**

La première procédure permet de rechercher un entier dans une chaîne de caractères Msg. i est l'entier trouvé. Si tout ce passe bien Err vaut 0, sinon il indique la position du premier caractère non entier de la chaîne.

La seconde permet la recherche d'un réel dans une chaîne de caractères. x est le réel trouvé. Si tout ce passe bien Err vaut 0, sinon il indique la position du premier caractère non correct de la chaîne.

**3.2.21 La fonction LONGUEUR\_CHAINE**

Cette procédure permet de calculer la longueur d'une chaîne de caractères, sans compter les derniers caractères blancs.

**3.2.22 Un exemple de programme utilisant MATH0**

Le programme suivant permet de rechercher tous les nombres premiers inférieurs à un nombre donné en utilisant le crible d'Eratosthène.

```
with TEXT_IO, IIO, LIO, CRT, COMMON_MATH0, MATH0 ;
use TEXT_IO, IIO, LIO, CRT, COMMON_MATH0, MATH0 ;
procedure PREMIERS is

type TABLEAU_BOOLEENS is array(LONG_INTEGER range <>) of BOOLEAN ;
nMax : constant LONG_INTEGER := 60_000 ;
n : LONG_INTEGER ;

procedure MENU_PREMIER(n : out LONG_INTEGER) is
```



```
begin
  CLRSCR ;
  PUT_LINE("RECHERCHE DE NOMBRES PREMIERS") ; NEW_LINE ;
  ENTREER_ENTIER_LONG_BORNE(2,nMax,n,
    "Borne superieure de recherche (m > 2) : ") ;
  CLRSCR ;
end MENU_PREMIER ;

procedure CALCUL_PREMIERS(n : in LONG_INTEGER) is
k, m, COMPTEUR : LONG_INTEGER ;
C : TABLEAU_BOOLEENS(1..n) ;
begin
  CHRONOMETRE(0) ;
  k := 2 ; m := k ;
  COMPTEUR := 0 ;
  C := (others => TRUE) ;
  C(1) := FALSE ;
  while (k < n)
  loop
    if c(k) then
      if ( (COMPTEUR mod 9) = 0) then
        PUT_LINE(IMP," ") ;
      end if ;
      COMPTEUR := COMPTEUR + 1 ;
      PUT(IMP,k,8) ;
      m := k ;
      while (m < n - k)
      loop
        m := m + k ;
        c(m) := FALSE ;
      end loop ;
    end if ;
    k := k + 1 ;
  end loop ;
  NEW_LINE(IMP,2) ;
```

```

CHRONOMETRE(1) ;
PUT(IMP,"Nombre d'entiers premiers : ") ;
PUT(IMP,COMPTEUR) ; PUT_LINE(IMP," ") ;
end CALCUL_PREMIERS ;
begin
MODE_AFFICHAGE ;
MENU_PREMIER(n) ;
CALCUL_PREMIERS(n) ;
CLOSE(IMP) ;
end PREMIERS ;

```

## 4. Le paquetage *MATH1*

Dans ce paquetage, on définit les fonctions usuellement utilisées en mathématiques. Ce sont les fonctions :

- SQRT* : pour la racine carrée d'un réel positif ;
- LN* : pour le logarithme népérien d'un réel positif ;
- LOG* : pour le logarithme dans une base quelconque ;
- EXP* : pour l'exponentielle ;
- \*\** : pour la fonction puissance ;
- SIN* : pour le sinus ;
- COS* : pour le cosinus ;
- TAN* : pour la tangente ;
- COT* : pour la cotangente ;
- ARCSIN* : pour l'arcsinus d'un réel entre  $-1.0$  et  $1.0$  ;
- ARCCOS* : pour l'arccosinus d'un réel entre  $-1.0$  et  $1.0$  ;
- ARCTAN* : pour l'arctangente d'un réel ;
- SH* : pour le sinus hyperbolique ;
- CH* : pour le cosinus hyperbolique ;
- TH* : pour la tangente hyperbolique ;
- ARGSH* : pour l'argument sinus hyperbolique ;
- ARGCH* : pour l'argument cosinus hyperbolique d'un réel  $\geq 1.0$  ;
- ARGTH* : pour l'argument tangente hyperbolique d'un réel entre  $-1.0$  et  $1.0$  ;
- FACT* : pour la factorielle d'un entier ;
- GAMMA* : pour la fonction Gamma d'Euler ;
- LNGAMMA* : pour le logarithme népérien de Gamma ;
- MIN* : pour le minimum de deux réels ;
- MAX* : pour le maximum de deux réels ;

L'exception *ERREUR\_ARGUMENT* est levée quand l'argument fourni à une fonction ne vérifie pas les bonnes conditions (Par exemple, *LN(x)*, avec  $x$  négatif ou nul).

Pour ce paquetage, on utilise le paquetage *MATH\_LIB* fourni avec *OPEN\_ADA*. Dans *MATH\_LIB* sont définies les 6 fonctions de base : *SIN*, *COS*, *EXP*, *SQRT*, *LN* et *ATAN* agissant sur des paramètres de type *FLOAT*. Pour une version générique de ce paquetage on consultera le paquetage *GENERIC\_ELEMENTARY\_FUNCTION* fourni avec *OpenAda*.

Le compilateur *FirstAda* de *Alslys* fournit un paquetage générique, non documenté, *MATH\_LIB* qui permet d'écrire une version générique de *MATH1* pour ce compilateur.

### 4.1 *Spécifications*

#### *des paquetages COMMON\_MATH1 et MATH1*

```

package COMMON_MATH1 is
PI: constant FLOAT := 3.141_592_653_589_793 ;
PI_SUR_DEUX: constant FLOAT := 0.5*PI ;
PI_SUR_QUATRE : constant FLOAT := 0.25*PI ;
DEUX_PI : constant FLOAT := 2.0*PI ;
E : constant FLOAT := 2.718_281_828_459_045 ;
ERREUR_ARGUMENT : exception ;
end COMMON_MATH1 ;

package MATH1 is
function Sqrt(x : FLOAT) return FLOAT ;
function LN(x : FLOAT) return FLOAT ;
function LOG(x : FLOAT ; BASE : in FLOAT := 10.0) return FLOAT ;
function EXP(x : FLOAT) return FLOAT ;
function "**"(x : FLOAT ; n : INTEGER) return FLOAT ;
function "**"(x,y : FLOAT) return FLOAT ;
function SIN(x : FLOAT) return FLOAT ;
function COS(x : FLOAT) return FLOAT ;
function TAN(x : FLOAT) return FLOAT ;
function COT(x : FLOAT) return FLOAT ;
function ARCSIN(x : FLOAT) return FLOAT ;
function ARCCOS(x : FLOAT) return FLOAT ;
function ARCTAN(x : FLOAT) return FLOAT ;
function SH(x : FLOAT) return FLOAT ;
function CH(x : FLOAT) return FLOAT ;
function TH(x : FLOAT) return FLOAT ;
function ARGSH(x : FLOAT) return FLOAT ;
function ARGCH(x : FLOAT) return FLOAT ;
function ARGTH(x : FLOAT) return FLOAT ;
function FACT(n : INTEGER) return FLOAT ;
function LNGAMMA(x : FLOAT) return FLOAT ;
function GAMMA(x : FLOAT) return FLOAT ;
function MIN(x, y : FLOAT) return FLOAT ;
function MAX(x, y : FLOAT) return FLOAT ;
end MATH1 ;

```

## 4.2 Documentation du paquetage MATH1

### 4.2.1 Les fonctions Puissances

La première associe à  $x \in \mathbb{R}$  et  $n \in \mathbb{Z}$  la puissance  $n^{\text{ème}}$  de  $x$ ,  $x^n$ , avec  $x \neq 0$  pour  $n < 0$  et la seconde associe à  $x > 0$  et  $y \in \mathbb{R}$  la puissance :  $x^y = e^{y \cdot \text{Ln}(x)}$

### 4.2.2 Les fonctions ArcTan, ArcSin et ArcCos

Si on dispose seulement de la fonction Arctangente, les fonctions ArcSin et ArcCos peuvent être définies par :

$$\text{ArcSin}(x) = \text{ArcTan}\left(\frac{x}{\sqrt{1-x^2}}\right), \text{ pour } |x| < 1.0$$

et

$$\text{ArcSin}(1) = -\text{ArcSin}(-1) = \frac{\pi}{2}$$

puis,

$$\text{ArcCos}(x) = \frac{\pi}{2} - \text{ArcSin}(x)$$

#### 4.2.3 Les fonctions hyperboliques inverses

Elles sont définies par :

$$\text{ArgCh}(x) = \text{Ln}\left(x + \sqrt{x^2 - 1}\right), \text{ pour } x \geq 1$$

$$\text{ArgSh}(x) = \text{Ln}\left(x + \sqrt{x^2 + 1}\right), \text{ pour } x \in \mathbb{R}$$

$$\text{ArgTh}(x) = \frac{1}{2} \cdot \text{Ln}\left(\frac{1+x}{1-x}\right), \text{ pour } |x| < 1$$

#### 4.2.4 Les fonctions Gamma et LnGamma

La fonction  $\Gamma$  est définie pour  $x > 0$  par :

$$\forall x > 0, \Gamma(x) = \int_0^{+\infty} t^{x-1} \cdot e^{-t} dt$$

On a alors la relation fonctionnelle :

$$\Gamma(x+1) = x \cdot \Gamma(x)$$

de laquelle on déduit que:

$$\forall n \in \mathbb{N}, \Gamma(n+1) = n!$$

Il y a plusieurs méthodes de calcul de  $\Gamma$ , l'une des plus rapide étant celle donnée par une formule du type Stirling (Cf. Numerical Recipes p. 157) :

$$\frac{\Gamma(x)}{\sqrt{2 \cdot \pi}} \approx (x+4.5)^{x-0.5} \cdot e^{-(x+4.5)} \cdot \left(1 + \frac{C_1}{x} + \frac{C_2}{x+1} + \dots + \frac{C_6}{x+5}\right)$$

L'erreur commise étant de l'ordre de  $10^{-10}$  pour  $x > 1$ .

Pour éviter les débordements arithmétiques, il est préférable de calculer  $\text{Ln}(\Gamma(x))$ , pour  $x > 1$ , avec la formule ci-dessus. Le calcul pour  $0 < x < 1$  se fait avec la relation fonctionnelle, puis celui de  $\Gamma(x)$  avec  $\text{EXP}(\text{Ln}(\Gamma(x)))$ .

## 5. Le paquetage Math2

Dans ce paquetage, on définit un type FONCTION ainsi que des procédures permettant de manipuler des variables de ce type. Ce paquetage est inspiré du livre de Ducamp et Reverchon : Mathématiques en Turbo-Pascal, Vol. 1.

## 5.1 Spécifications du paquetage *MATH2*

```

package MATH2 is
MAX_VARIABLES : constant INTEGER := 99 ;
subtype INDICE_VARIABLES is INTEGER range 0..MAX_VARIABLES ;
type VARIABLE_2 is array(MAX_VARIABLES range <>) of FLOAT ;
type FONCTION(MAX_CHAINE : NATURAL := 80) is limited private ;
function EVALUE(f : FONCTION ; x : FLOAT ; y : VARIABLE_2)
    return FLOAT ;
function EVALUE(f : FONCTION ; x, y, z, t : FLOAT := 0.0)
    return FLOAT ;
procedure DETRUIT_FONCTION(f : in out FONCTION) ;
procedure GET_LINE(f : in out FONCTION ; Msg : in STRING := " ") ;
function TEXTE_DE_FONCTION(f : in FONCTION) return STRING ;
private
type TYPE_NOEUD is
    (CONSTANTE, VARIABLE, OPERATION_UNAIRE, OPERATION_BINAIRE) ;
type UNE_OPERATION_UNAIRE is
    (FMOINS, FPLUS, FSIN, FCOS, FTG, FARCSIN, FARCCOS,
     FARCTAN, FSH, FCH, FTH, FARGSH, FARGCH, FARGTH,
     FABS, FFRAC, FEXP, FLN, FGAMMA, FLNGAMMA, FSQRT) ;
type UNE_OPERATION_BINAIRE is (FSUM, FSUB, FMUL, FDIV, FPUI) ;
type NOEUD ;
type POINTEUR_NOEUD is access NOEUD ;
type NOEUD is
    record
        GAUCHE, DROITE : POINTEUR_NOEUD ;
        TYPE_DU_NOEUD : TYPE_NOEUD ;
        VALEUR : FLOAT ;
        RANG : INTEGER ;
        OPERATION_1 : UNE_OPERATION_UNAIRE ;
        OPERATION_2 : UNE_OPERATION_BINAIRE ;
    end record ;
type FONCTION(MAX_CHAINE : NATURAL := 80) is
    record
        TEXTE_FONCTION : STRING(1..MAX_CHAINE) := (others => ' ') ;
        ARBRE : NOEUD ;
    end record ;
end MATH2 ;

```

## 5.2 Documentation de *Math2*

### 5.2.1 Les fonctions *EVALUE*

Les fonctions d'évaluation :

```

EVALUE(f : Fonction ; x : FLOAT ; y : VARIABLE_2) return FLOAT ;
EVALUE(f : Fonction ; x, y, z, t : FLOAT := 0.0) return FLOAT ;

```

interprètent la fonction *f* pour retourner respectivement les valeurs réelles  $f(x,y)$  et  $f(x)$  ou  $f(x,y)$  ou  $f(x,y,z)$  ou  $f(x,y,z,t)$ , où  $x$  est réelle et  $y$  est un vecteur ayant au plus  $\text{MAX\_VARIABLES} + 1$  composantes dans le premier cas et  $x, y, z, t$  sont réelles dans le second cas. Pour une fonction d'une variable, on écrira  $r := \text{Evalue}(f,x)$  ; pour une fonction de deux variables,  $r := \text{Evalue}(f,x,y)$  ; ...

### 5.2.2 La procédure DETRUIT\_FONCTION

Cette procédure permet d'éliminer une fonction f en vue de libérer l'espace mémoire occupé.

### 5.2.3 La procédure GET\_LINE

Cette procédure permet de lire au clavier le texte du corps d'une fonction f et de créer l'objet f de type FONCTION.

### 5.2.4 La fonction TEXTE\_DE\_FONCTION

Cette fonction retourne le texte de la fonction f sans les derniers caractères blancs.

## 6. Le paquetage GRAPH

Dans ce paquetage on définit des procédures graphiques de base. Ce paquetage étant lié au matériel n'est pas facilement portable. On utilise les interruptions DOS.

### 6.1 Spécifications des paquetages COMMON\_GRAPH et GRAPH

```

package COMMON_GRAPH is
type CARTE_GRAPHIQUE is (CGA_C4_320x200, CGA_GREY4_320x200,
    CGA_C2_640x200, EGA_C16_320x200, EGA_C16_640x200, EGA_BW_640x350,
    EGA_VGA_C16_320x200, EGA_VGA_C16_640x200, EGA_VGA_BW_640x350,
    EGA_VGA_C16_640x350, MCGA_BW_640x480, MCGA_C256_320x200,
    VGA_C16_640x480, SVGA_C16_800x600, SVGA_C16_1024x768) ;
type COULEUR is (NOIR, BLEUE, VERT, CYAN, ROUGE, MAGENTA, BRUN,
    GRIS_CLAIR, GRIS_FONCE, BLEUE_CLAIR, VERT_CLAIR,
    CYAN_CLAIR, ROUGE_CLAIR, MAGENTA_CLAIR, JAUNE, BLANC) ;
subtype TAILLE_CARACTERE is integer range 8 .. 16 ;
subtype TEXTE_HORIZONTAL is INTEGER range 0..127 ;
subtype TEXTE_VERTICAL is INTEGER range 0..47 ;
CARTE_GRAPHIQUE_UTILISATEUR : CARTE_GRAPHIQUE ;
MAX_LIGNE_TEXTE : TEXTE_VERTICAL ;
MAX_COLONNE_TEXTE : TEXTE_HORIZONTAL ;
CLIP_ON : BOOLEAN ;
GR_X_MIN, GR_X_MAX, GR_Y_MIN, GR_Y_MAX, GET_MAX_X, GET_MAX_Y : NATURAL;
LARGEUR_PIXEL, HAUTEUR_PIXEL : TAILLE_CARACTERE ;
X_POS, Y_POS : INTEGER ;
COULEUR_PIXEL : COULEUR := GRIS_CLAIR ;
end COMMON_GRAPH ;

with COMMON_GRAPH ;
use COMMON_GRAPH ;
package GRAPH is
procedure CLEAR_SCREEN ;
procedure CLOSE_GRAPH ;
procedure INIT_GRAPH(CHOIX : CARTE_GRAPHIQUE) ;
procedure MOVE_TO(x : in INTEGER ; y : in INTEGER) ;
procedure PUT_PIXEL(x : in INTEGER ; y : in INTEGER) ;
procedure LINE(x, x1 : in INTEGER ; y, y1 : in INTEGER) ;

```

```

procedure LINE_TO(x : in INTEGER ; y : in INTEGER) ;
procedure RECTANGLE(x, x1 : in INTEGER ; y, y1 : in INTEGER ;
                   CLIP : in BOOLEAN := FALSE) ;
end GRAPH ;

```

## 6.2 Documentation de GRAPH

### 6.2.1 La procédure CLEAR\_SCREEN

Cette procédure permet d'effacer l'écran en mode graphique.

### 6.2.2 La procédure CLOSE\_GRAPH

Cette procédure permet de revenir en mode texte : 25 lignes et 80 colonnes.

### 6.2.3 La procédure INIT\_GRAPH

Cette procédure permet d'initialiser un mode graphique.

### 6.2.4 La procédure MOVE\_TO

Cette procédure permet de mémoriser la position d'un point.

### 6.2.5 La procédure PUT\_PIXEL

Cette procédure permet de se déplacer en (x,y) en allumant un pixel avec la couleur spécifiée par COULEUR\_PIXEL (par défaut c'est GRIS\_CLAIR).

### 6.2.6 La procédure LINE

Cette procédure permet de tracer le segment d'extrémités (x,y) et (x1,y1). le curseur étant ensuite positionné en (x1,y1).

### 6.2.7 La procédure LINE\_TO

Cette procédure permet de tracer un segment depuis la position actuelle du curseur jusqu'au point de coordonnées (x,y).

### 6.2.8 La procédure RECTANGLE

Cette procédure permet de tracer un rectangle de coin supérieur gauche (x,y) et de coin inférieur droit (x1, y1).

## 7. Le paquetage MATH3

### 7.1 Spécifications

#### *des paquetages COMMON\_MATH3 et MATH3*

La constante CHOIX\_CARTE\_GRAPHIQUE devra être modifiée en fonction du type d'écran sur lequel on travaille (CGA, EGA, ...). On se référera à la spécification de COMMON\_GRAPH pour connaître la liste des choix possibles.

L'exception ERREUR\_CONTRAINTES\_GRAPHIQUES sera levée pour toute erreur de contrainte dans la définition d'une fenêtre ou d'un cadre graphique. L'exception ERREUR\_MODE\_ECRAN sera levée quand on tente d'utiliser une procédure graphique alors qu'on est en mode texte.

On pourra ajouter des procédures permettant d'utiliser une table traçante.

```

with COMMON_GRAPH ; use COMMON_GRAPH ;
package COMMON_MATH3 is
type MODE_ECRAN is (TEXTE, GRAPHIQUE) ;
CHOIX_CARTE_GRAPHIQUE :

```

```

constant GRAPH.CARTE_GRAPHIQUE := VGA_C16_640x480 ;
x_MINIMUM, x_MAXIMUM, y_MINIMUM, y_MAXIMUM : NATURAL ;
HAUTEUR_CAR, LARGEUR_CAR, MARGE : NATURAL ;
GR_LARGEUR_ECRAN, GR_HAUTEUR_ECRAN : NATURAL ;
X_MIN, X_MAX, Y_MIN, Y_MAX, X_GRD, Y_GRD : FLOAT ;
LARGEUR_ECRAN, HAUTEUR_ECRAN : FLOAT ;
X_RAPPORT, Y_RAPPORT : FLOAT ;
INDIC_MODE : MODE_ECRAN := TEXTE ;
ERREUR_CONTRAINTES_GRAPHIQUES : exception ;
ERREUR_MODE_ECRAN : exception ;
end COMMON_MATH3 ;

package MATH3 is
procedure INIT_GRAPHIQUE ;
procedure MODE_TEXTE ;
procedure MESSAGE_ERREUR_GRAPHIQUE(Msg : in STRING) ;
procedure TEST_MODE_GRAPHIQUE(Msg : in STRING) ;
procedure FENETRE(X_FEN_MIN,X_FEN_MAX,
                  Y_FEN_MIN,Y_FEN_MAX : in FLOAT) ;
procedure CADRE_GRAPHIQUE(X_CADRE_MIN, X_CADRE_MAX,
                           Y_CADRE_MIN, Y_CADRE_MAX : in INTEGER) ;
procedure FENETRE_GRAPHIQUE(X_FEN_MIN,X_FEN_MAX,
                             Y_FEN_MIN,Y_FEN_MAX : in FLOAT) ;
procedure CONVERSION(x, y : in FLOAT ;
                     X_INFO : out INTEGER ; Y_INFO : out INTEGER) ;
procedure TEXT_MOVE_TO(x : in INTEGER ; y : in INTEGER) ;
procedure PAUSE_GRAPHIQUE ;
procedure DEPLACE(x, y : in FLOAT) ;
procedure POINT(x, y : in FLOAT) ;
procedure TRACE_SEGMENT(x1, y1, x2, y2 : in FLOAT) ;
procedure TRACE(x, y : in FLOAT) ;
procedure CROIX(x, y : in FLOAT) ;
procedure SORTIE_GRAPHIQUE ;
procedure CERCLE(x, y, r : in FLOAT) ;
procedure X_AXE(X_ORIGINE, Y_ORIGINE, X_UNITE : in FLOAT) ;
procedure Y_AXE(X_ORIGINE, Y_ORIGINE, Y_UNITE : in FLOAT) ;
procedure XY_AXES(X_ORIGINE, Y_ORIGINE, X_UNITE, Y_UNITE : in FLOAT ;
                  XY_GRILLE : in BOOLEAN := TRUE) ;
procedure TITRE(Msg : in STRING) ;
end MATH3 ;

```

## 7.2 Documentation de MATH3

### 7.2.1 La procédure INIT\_GRAPHIQUE

Permet de passer en mode graphique et d'initialiser les variables du système.

### 7.2.2 La procédure MODE\_TEXTE

Quitte le mode graphique et repasse en mode texte.

### 7.2.3 La procédure MESSAGE\_ERREUR\_GRAPHIQUE

En cas d'erreur dans une procédure graphique, cette procédure permet d'arrêter l'exécution du programme en retournant en mode texte et en donnant un message sur l'erreur ayant provoqué cet arrêt.



### 7.2.4 La procédure TEST\_MODE\_GRAPHIQUE

Cette procédure permet de vérifier que l'on est dans le bon mode pour utiliser une procédure graphique. Pour cela, on utilise la variable INDIC\_MODE définie par :

- INDIC\_MODE = GRAPHIQUE, si on est en mode graphique ;
- INDIC\_MODE = TEXTE, si on est en mode texte.

### 7.2.5 La procédure TEXT\_MOVE\_TO

Cette procédure permet de se déplacer en un point de l'écran, en mode graphique, pour y écrire du texte.

### 7.2.6 La procédure PAUSE\_GRAPHIQUE

Cette procédure interrompt momentanément le déroulement d'un programme jusqu'à ce qu'une touche soit pressée.

### 7.2.7 La procédure FENETRE

Initialisation des variables, pour définir une fenêtre graphique, permettant de travailler dans les unités de l'utilisateur. Ainsi, si on a donné l'instruction

```
FENETRE (-15.5, 2.3, 1.2, 1.8) ;
```

le point situé en bas et à gauche a pour coordonnées  $X = -15.5$  et  $Y = 1.2$  et le point situé en haut et à droite a pour coordonnées  $X = 2.3$  et  $Y = 1.8$ .

### 7.2.8 La procédure CADRE\_GRAPHIQUE

Cette procédure qui doit suivre l'appel à FENETRE permet de préciser sur quel portion de l'écran graphique on travaille. Elle permet de travailler sur plusieurs fenêtres graphiques dans un même écran (pour une même fenêtre, on peut utiliser plusieurs cadres).

### 7.2.9 La procédure FENETRE\_GRAPHIQUE

Cette procédure permet de définir une fenêtre graphique de travail avec un cadre par défaut qui est l'écran tout entier privé d'une bande en haut de l'écran et d'une bande en bas.

### 7.2.10 La procédure CONVERSION

Cette procédure permet d'associer aux coordonnées de l'utilisateur  $(x,y)$ , dans un repère orienté de façon classique en mathématiques, les coordonnées écran  $(xInfo,yInfo)$  dans le repère « informatique », c'est-à-dire avec des coordonnées entières et les abscisses allant de gauche à droite et les ordonnées de haut en bas.

Cette procédure est nécessaire pour les procédures DEPLACE, TRACE, CROIX et POINT.

### 7.2.11 La procédure DEPLACE

Mémorise la position du curseur au point de coordonnées  $(X,Y)$  sans tracer.

### 7.2.12 La procédure TRACE

Trace un segment de droite joignant la position actuelle du curseur au point de coordonnées  $(X,Y)$ . Le curseur est déplacé en ce dernier point. Le tracé se fait dans la couleur définie par COULEUR\_PIXEL.

### 7.2.13 La procédure CROIX

Trace une croix au point de coordonnées  $(x,y)$ .

**7.2.14 La procédure TRACE\_SEGMENT**

Permet de tracer le segment d'extrémités  $(x_1, y_1)$  et  $(x_2, y_2)$ . La position du curseur est mémorisée en  $(x_2, y_2)$ . Le tracé se fait dans la couleur définie par COULEUR\_PIXEL.

**7.2.15 La procédure POINT**

Positionne le curseur au point de coordonnées  $X, Y$  et dessine ce point dans la couleur précisée par COULEUR\_PIXEL.

**7.2.16 La procédure CERCLE**

Trace un cercle de centre  $(x, y)$  et de rayon  $r$ . Le tracé se fait dans la couleur définie par COULEUR\_PIXEL.

Le tracé se fait par récurrence de la façon suivante : soit un cercle de centre  $O$  et de rayon  $R$  et  $(x, y)$  un point de ce cercle ; le point déduit par une rotation d'angle  $d\theta$  est donné par:

$$\begin{cases} x_{n+1} = \text{Cos}(d\theta) \cdot x_n - \text{Sin}(d\theta) \cdot y_n \\ y_{n+1} = \text{Sin}(d\theta) \cdot x_n + \text{Cos}(d\theta) \cdot y_n \end{cases}$$

En prenant  $d\theta$  assez petit on peut ainsi tracer rapidement un cercle.

**7.2.17 La procédure TITRE**

Efface le dernier titre en date puis dessine un petit rectangle, en haut de l'écran, dans lequel sera écrit le titre Msg.

**7.2.18 La procédure X\_AXE**

Cette procédure permet de tracer un axe gradué dans la direction  $Ox$  passant par le point de coordonnées  $X\_ORIGINE$  et  $Y\_ORIGINE$  et de graduation  $X\_UNITE$ .

**7.2.19 La procédure Y\_AXE**

Procédure analogue à  $X\_AXE$  concernant la direction  $Oy$ .

**7.2.20 La procédure XY\_AXES**

Cette procédure permet de tracer les axes des  $x$  et  $y$  avec la possibilité de tracer une grille pour mieux visualiser les graduations.

**7.2.21 La procédure SORTIE\_GRAPHIQUE**

Permet de sortir du mode graphique.

**7.2.22 Un exemple d'utilisation de MATH2 et de MATH3**

Dans ce paragraphe on décrit tout d'abord le paquetage CURVE, où on définit un type COURBE avec des procédures permettant le tracé de courbes dans le plan, puis le paquetage CARTESIENNE est utilisé pour le tracé de courbes d'équation cartésienne  $y = f(x)$  sur un intervalle  $[a, b]$ .

— Spécification de CURVE

```
package CURVE is
  NB_POINTS : constant NATURAL := 800 ;
  type COORDONNEE is array(NATURAL range <>) of FLOAT ;
  type COURBE(n : NATURAL) is
    record
      uMin, uMax, vMin, vMax : FLOAT ;
```

```

    u, v : COORDONNEE(0..n) ;
    Rep_Axes, Rep_Grille : BOOLEAN ;
end record ;

```

Le rectangle de travail est  $[u_{\min}, u_{\max}] \times [v_{\min}, v_{\max}]$ , le nombre de points de la courbe est  $u_{\text{LENGTH}} = v_{\text{LENGTH}} = n + 1$ , les abscisses sont consignées dans le vecteur  $u$  et les ordonnées dans  $v$ . Si  $\text{Rep\_Axes} = \text{TRUE}$ , on trace les axes et si  $\text{Rep\_Grille} = \text{TRUE}$ , on trace une grille pour mieux visualiser les graduations.

```

procedure CHOIX_AXES(xOrigine , yOrigine , xUnite , yUnite ,
                    xFenMin , xFenMax , yFenMin , yFenMax : out FLOAT) ;
procedure AXES_PAR_DEFAUT(xMin , yMin , xMax , yMax : in FLOAT ;
                        xOrigine , yOrigine , xUnite , yUnite ,
                        xFenMin , xFenMax , yFenMin , yFenMax : in out FLOAT) ;
procedure SYSTEME_AXES(xMin , yMin , xMax , yMax : in FLOAT ;
                      xOrigine , yOrigine , xUnite , yUnite ,
                      xFenMin , xFenMax , yFenMin , yFenMax : in out FLOAT) ;
procedure INIT_COURBE(C : in out COURBE ; xFenMin, xFenMax, yFenMin,
                    yFenMax, xOrigine, yOrigine, xUnite, yUnite : in out FLOAT) ;
procedure TRACE_COURBE(C : in COURBE) ;
procedure TRACE_UNE_COURBE(C : in out COURBE) ;
end CURVE ;

```

### — Corps de CURVE

```

with TEXT_IO, CRT, MATH0, COMMON_MATH3, MATH3 ;
use TEXT_IO, CRT, MATH0, COMMON_MATH3, MATH3 ;
package body CURVE is
procedure CHOIX_AXES(xOrigine , yOrigine , xUnite , yUnite ,
                    xFenMin , xFenMax , yFenMin , yFenMax : out FLOAT) is
begin
    CLRSCR ;
    PUT_LINE("Coordonnees de l'origine:") ;
    ENTREER_REEL(xOrigine," Abscisse de l'origine : ") ;
    ENTREER_REEL(yOrigine," Ordonnee de l'origine : ") ;
    NEW_LINE ;
    PUT_LINE("Choix des unites sur les axes:") ;
    ENTREER_REEL(X_GRD," Unite de l'echelle en x : ") ;
    xUnite := X_GRD ;
    ENTREER_REEL(Y_GRD," Unite de l'echelle en y : ") ;
    yUnite := Y_GRD ; NEW_LINE ;
    PUT_LINE("Coordonnees de la fenetre graphique:") ;
    ENTREER_REEL(xFenMin," Abscisse coin inferieur gauche : ") ;
    ENTREER_REEL(yFenMin," Ordonnee coin inferieur gauche : ") ;
    ENTREER_REEL(xFenMax," Abscisse coin superieur droit : ") ;
    ENTREER_REEL(yFenMax," Ordonnee coin superieur droit : ") ;
end CHOIX_AXES ;

procedure AXES_PAR_DEFAUT(xMin , yMin , xMax , yMax : in FLOAT ;
                        xOrigine , yOrigine , xUnite , yUnite ,
                        xFenMin , xFenMax , yFenMin , yFenMax : in out FLOAT) is
begin
    xOrigine := (xMax + xMin)/2.0 ; yOrigine := (yMin + yMax)/2.0 ;
    xUnite := (xMax - xMin)/20.0 ; X_GRD := xUnite ;
    yUnite := (yMax - yMin)/20.0 ; Y_GRD := yUnite ;
    xFenMin := xMin - xUnite ; xFenMax := xMax + xUnite ;

```

```

    yFenMin := yMin - yUnite ;
    yFenMax := yMax + yUnite ;
end AXES_PAR_DEFAUT ;

procedure SYSTEME_AXES(xMin , yMin , xMax , yMax : in FLOAT ;
    xOrigine , yOrigine , xUnite , yUnite ,
    xFenMin , xFenMax , yFenMin , yFenMax : in out FLOAT) is
CHOIX : INTEGER ;
begin
    CLRSCR ;
    PUT_LINE(" 1 : Systeme d'axe par default." ) ;
    PUT_LINE(" Si les abscisses varient dans l'intervalle [xMin,xMax]");
    PUT_LINE(" et les ordonnees dans [yMin,yMax], on pose alors:" ) ;
    PUT_LINE(" Abscisse de l'origine = (xMax + xMin)/2.0 ;" ) ;
    PUT_LINE(" Ordonnee de l'origine = (yMax + yMin)/2.0 ;" ) ;
    PUT_LINE(" Unite sur l'axe des x = (xMax - xMin)/20.0 ;" ) ;
    PUT_LINE(" Unite sur l'axe des y = (yMax - yMin)/20.0 ;" ) ;
    PUT("Abscisse du coin inferieur gauche de la fenetre graphique = " );
    PUT_LINE("xMin - xUnite ;" ) ;
    PUT("Ordonnee du coin inferieur gauche de la fenetre graphique = " );
    PUT_LINE("yMin - yUnite ;" ) ;
    PUT("Abscisse du coin superieur droit de la fenetre graphique = " ) ;
    PUT_LINE("xMax + xUnite ;" ) ;
    PUT("Ordonnee du coin superieur droit de la fenetre graphique = " ) ;
    PUT_LINE("yMax + yUnite ." ) ;
    NEW_LINE ;
    PUT_LINE(" 2 : Votre systeme d'axes." ) ;
    NEW_LINE ;
    ENTREER_ENTIER_BORNE(1,2,CHOIX,"Votre choix : " ) ;
    CLRSCR ;
    if CHOIX = 1 then
        AXES_PAR_DEFAUT(xMin,yMin,xMax,yMax,xOrigine,yOrigine,
            xUnite,yUnite, xFenMin,xFenMax,yFenMin,yFenMax) ;
    else
        CHOIX_AXES(xOrigine,yOrigine,xUnite, yUnite,
            xFenMin,xFenMax,yFenMin,yFenMax) ;
    end if ;
end SYSTEME_AXES ;

procedure INIT_COURBE(C : in out COURBE ; xFenMin, xFenMax,
    yFenMin, yFenMax, xOrigine, yOrigine,
    xUnite, yUnite : in out FLOAT) is
begin
    C.Rep_Axes := LIRE_REPONSE("Voulez vous le trace des axes? " ) ;
    if C.Rep_Axes then
        NEW_LINE ;
    C.Rep_Grille := LIRE_REPONSE("Voulez vous le trace d'une grille? " ) ;
        NEW_LINE ;
        SYSTEME_AXES(C.uMin,C.vMin,C.uMax,C.vMax,xOrigine,yOrigine,
            xUnite,yUnite, xFenMin,xFenMax,yFenMin,yFenMax) ;
    else
        xUnite := (C.uMax - C.uMin)/20.0 ;
        yUnite := (C.vMax - C.vMin)/20.0 ;
        xFenMin := C.uMin - xUnite ;
        xFenMax := C.uMax + xUnite ;
        yFenMin := C.vMin - yUnite ;
        yFenMax := C.vMax + yUnite ;
    end if ;
end INIT_COURBE ;

```

```
end INIT_COURBE ;
```

```
procedure TRACE_COURBE(C : in COURBE) is
begin
  TEST_MODE_GRAPHIQUE("TRACE_COURBE") ;
  DEPLACE(C.u(0),C.v(0)) ;
  for i in C.u'range
  loop
    TRACE(C.u(i),C.v(i)) ;
  end loop ;
end TRACE_COURBE ;
```

```
procedure TRACE_UNE_COURBE(C : in out COURBE) is
xOrigine, yOrigine, xUnite, yUnite,
xFenMin, xFenMax, yFenMin, yFenMax : FLOAT ;
begin
  INIT_COURBE(C,xFenMin,xFenMax,yFenMin,yFenMax,xOrigine,yOrigine,
              xUnite,yUnite) ;
  INIT_GRAPHIQUE ;
  FENETRE_GRAPHIQUE(xFenMin,xFenMax,yFenMin,yFenMax) ;
  TRACE_COURBE(C) ;
  if C.Rep_Axes then
    XY_AXES(xOrigine,yOrigine,xUnite,yUnite,C.Rep_Grille) ;
  end if ;
end TRACE_UNE_COURBE ;
end CURVE ;
```

## — Spécification de CARTESIENNE

```
with MATH2, CURVE ;
use MATH2, CURVE ;
package CARTESIENNE is
procedure DONNEES_CARTESIENNE(f : in out FONCTION ; C : in out COURBE) ;
procedure TRACE_CARTESIENNE(f : in out FONCTION ; C : in out COURBE) ;
end CARTESIENNE ;
```

## — Corps de CARTESIENNE

```
with TEXT_IO, FIO, MATH0, MATH1, MATH2, MATH3, CURVE ;
use TEXT_IO, FIO, MATH0, MATH1, MATH2, MATH3, CURVE ;

package body CARTESIENNE is

procedure DONNEES_CARTESIENNE(f : in out FONCTION ;
                              C : in out COURBE) is
Pas : FLOAT ;
begin
  GET_LINE(f,"Entrer la fonction f(x) : ") ;
  ENTRER_REEL(C.uMin,"Valeur minimale de l'abscisse, a = ") ;
  ENTRER_REEL_BORNE(C.uMin,1.0E+6,C.uMax,
                   "Valeur maximale de l'abscisse, b = ") ;
  NEW_LINE ;
  PAS := (C.uMax - C.uMin)/FLOAT(C.n) ;
  C.vMIN := 1.0E+37 ; C.vMAX := -C.vMIN ;
  C.u(0) := C.uMin ; C.v(0) := Evaluate(f,C.u(0)) ;
  for i in 1..C.n
  loop
```

```

    C.u(i) := C.u(i-1) + PAS ;
    C.v(i) := Evaluate(f,C.u(i)) ;
    C.vMIN := MIN(C.vMIN,C.v(i-1)) ;
    C.vMAX := MAX(C.vMAX,C.v(i-1)) ;
  end loop ;
end DONNEES_CARTESIENNE ;

procedure TRACE_CARTESIENNE(f : in out FONCTION ; C : in out COURBE) is
begin
  DONNEES_CARTESIENNE(f,C) ;
  TRACE_UNE_COURBE(C) ;
  TITRE(" Graphe de f(x) = " & TEXTE_DE_FONCTION(f)) ;
  SORTIE_GRAPHIQUE ;
end TRACE_CARTESIENNE ;

end CARTESIENNE ;

```

### — Un exemple d'utilisation

```

with TEXT_IO, MATH0, COMMON_MATH2, CURVE, CARTESIENNE ;
use TEXT_IO, MATH0, COMMON_MATH2, CURVE, CARTESIENNE ;

procedure ESS_CART is
NOMBRE_DE_POINTS : NATURAL ;
begin
  PUT_LINE(" TRACE D'UNE COURBE D'EQUATION Y = F(X) SUR [A,B].") ;
  NEW_LINE(3) ;
  ENTRER_ENTIER_BORNE(2,NB_POINTS,NOMBRE_DE_POINTS,
    "Nombre de points de la courbe : ") ;

```

```

declare
  C : COURBE(NOMBRE_DE_POINTS) ;
  f : FONCTION ;
begin
  TRACE_CARTESIENNE(f,C) ;
end ;
end ESS_CART ;

```

## 8. Exercices

### 8.1 *Calculs sur les nombres rationnels*

Le but est d'écrire des procédures permettant d'effectuer des calculs de façon exacte sur les nombres rationnels. Ces procédures seront stockées dans un paquetage.

On devra donc :

- définir un type rationnel ;
- écrire une procédure de simplification d'un rationnel (i. e. l'écrire sous forme irréductible) ;
- définir les 4 opérations élémentaires sur les rationnels.

A titre d'application on pourra écrire des procédures de manipulation des matrices à coefficients dans  $\mathbb{Q}$ , ainsi qu'une procédure de résolution d'un système linéaire à coefficients dans  $\mathbb{Q}$ .

On utilisera cette procédure pour calculer les coefficients qui interviennent dans les méthodes de Newton-Cotes décrites dans le chapitre sur le calcul numérique des intégrales.

On pourra aussi imaginer d'autres applications.

### 8.2 *Approximation d'un réel par des fractions continues*

Le problème est, étant donné un réel  $x$ , de construire une suite de nombres rationnels qui converge vers  $x$ .

Les fractions continues permettent de construire une telle suite de rationnels  $(r_n) = (p_n/q_n)$ , écrits sous forme irréductible.

Algorithme de calcul :

On définit d'abord, la fonction :

$$f : \mathbb{R} - \mathbb{Z} \rightarrow \mathbb{R}$$

$$x \mapsto f(x) = \frac{1}{x - E[x]}$$

où  $E[x]$  est la partie entière de  $x$ , c'est à dire l'entier qui vérifie  $E[x] \leq x < E[x] + 1$ . Puis on définit les suites  $(x_n)$  et  $(a_n)$  par :

$$\begin{cases} x_0 = x; a_0 = E[x_0]; \\ \text{si } x_n \in \mathbb{Z}, \text{ alors on arrête;} \\ \text{sinon, } x_n = f(x_{n-1}) \text{ et } a_n = E[x_n] \end{cases}$$

La suite d'entiers  $(a_n)$  est formée d'une infinité de termes si, et seulement si  $x$  n'est pas rationnel.

On note  $N_{\max}$ , le dernier indice pour lequel  $a_n$  est défini, soit :

$$N_{\max} = \begin{cases} +\infty & \text{si } x \in \mathbb{R} - \mathbb{Q} \\ < +\infty & \text{si } x \in \mathbb{Q} \end{cases}$$

Pour  $0 \leq n \leq N_{\max}$ , on pose alors :

$$r_n = a_0 + \frac{1}{a_1 + \frac{1}{a_{n-1} + \frac{1}{a_n}}}$$

et on a :  $r_n = \frac{p_n}{q_n}$ , où les entiers  $p$  et  $q$  sont premiers entre eux et vérifient la même récurrence :

$$(R) \quad u_n = a_n \cdot u_{n-1} + u_{n-2}, \text{ pour } n > 1$$

avec les valeurs initiales :

$$\begin{cases} p_0 = a_0; & p_1 = a_0 \cdot a_1 + 1 \\ q_0 = 1; & q_1 = a_1 \end{cases}$$

On a alors :

$$\left| x - \frac{p_n}{q_n} \right| \leq \frac{1}{q_n^2}$$

avec, pour  $N_{\max} = +\infty$ ,  $\lim_{n \rightarrow +\infty} (q_n) = +\infty$ , et  $x = r_{N_{\max}}$  si  $x \in \mathbb{Q}$ .

Ecrire une procédure correspondante à l'algorithme décrit ci-dessus.

### 8.3 Calculs en multiprécision sur les entiers

Le but de l'exercice est de calculer, i.e. faire des additions, soustractions, multiplications et divisions avec des entiers ayant au plus  $N_{\max}$  chiffres où  $N_{\max}$  est aussi grand que possible.

Les procédures et fonctions écrites seront rangées dans un paquetage.

Pour ce faire, une idée est d'écrire un entier positif non nul en base 10, soit :

$$a = \sum_{k=0}^{n(a)} a_k \cdot 10^k, \text{ avec } 0 \leq a_k < 10 \text{ et } a_{n(a)} \neq 0;$$

puis de ranger les  $a$  dans un vecteur :

$$A = (a_0, a_1, \dots, a_{n(a)}, 0, \dots, 0) \in \mathbb{N}^{T_{\max}}$$

où  $T_{\max}$  est le nombre maximum chiffres que l'on a choisi.

Si maintenant, on veut travailler sur  $\mathbb{Z}$ , on aura intérêt à stocker les  $a$  dans un vecteur à  $T_{\max} + 1$  composantes  $A = (A_i)_{0 \leq i \leq T_{\max} + 1}$ , avec :

$$\begin{cases} A_k = a_{k-1} \text{ pour } k \geq 1; \\ |A_0| = \text{Nombre de chiffres de } a, \text{ si } a \neq 0; \\ A_0 \text{ du signe de } a, \text{ si } a \neq 0 \text{ et } A_0 = 0, \text{ si } a = 0 \end{cases}$$

1°) Définir un type Entier correspondant à la description ci-dessus.

2°) Ecrire une procédure de lecture d'un entier.

3°) Ecrire une procédure d'affichage d'un entier.

4°) Définir une fonction Longueur qui donne le nombre de chiffres d'un entier.

5°) Décrire des procédures d'addition, de soustraction et de multiplication qui opèrent sur des variables de type Entier.



6°) Décrire une procédure de division euclidienne qui donne le reste et le quotient dans la division euclidienne de deux entiers.

### 8.4 Opérations sur les polynômes

Définir un type polynôme réel, puis des procédures permettant de faire les calculs usuels : addition, soustraction, multiplication, division euclidienne, division suivant les puissances croissantes, calcul de pgcd et ppcm, évaluation d'un polynôme, dérivation, intégration, composition de deux polynômes, opérations sur les fractions rationnelles, ...

*Applications* : Ecrire des procédures de calcul des polynômes orthogonaux classiques (Cf. le chapitre sur le calcul numérique des intégrales).

### 8.5 Ensembles de Mandelbrot et de Julia

1°) On appelle ensemble de Mandelbrot, une partie de  $\mathbb{C}$  définie par :

$$M = \left\{ z \in \mathbb{C}; (z_n)_{n \in \mathbb{N}} \text{ converge} \right\}$$

où  $(z_n)$  est la suite définie par :

$$\begin{cases} z_0 = z \\ z_{n+1} = z_n^2 + z \quad (n \geq 0) \end{cases}$$

Ecrire une procédure permettant de dessiner cet ensemble.

2°) Si dans l'ensemble ci-dessus, on se fixe  $z$  et on fait varier la valeur initiale  $z_0$ , alors l'ensemble obtenu est appelé ensemble de Julia. Soit :

$$G = \left\{ z_0 \in \mathbb{C}; (z_n)_{n \in \mathbb{N}} \text{ converge} \right\}$$

où  $(z_n)$  est la suite définie par :

$$\begin{cases} z_0 \text{ donné dans } \mathbb{C} \\ z_{n+1} = z_n^2 + z \quad (n \geq 0) \end{cases}$$

$z$  étant donné.

Ecrire une procédure permettant de dessiner cet ensemble.

### 8.6 Courbe de Van Koch

La courbe de Van Koch est une courbe continue dérivable nulle part (fractale), de longueur infinie et bornée.

Le principe de la construction est le suivant :

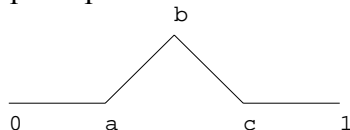


Figure 2.1

Partant de  $[0, 1]$  on remplace ce segment par 4 autres disposés comme sur le schéma ci-contre, puis on recommence indéfiniment avec chacun des nouveaux segments.

1°) La courbe obtenue à l'étape  $n$  est alors polygonale. Calculer le nombre de sommets de cette courbe.

2°) Ecrire une procédure de tracé d'une telle courbe.

3°) Généralisation : au lieu d'un triangle équilatéral, on peut choisir un motif quelconque défini par une courbe polygonale joignant 0 à 1 et on reproduit ce motif sur chacun des segments du motif précédent.

Ecrire une procédure de définition du motif (c'est le générateur de la courbe), puis une procédure de tracé de la fractale obtenue.

### 8.7 Transformations géométriques du plan

Ecrire des procédures de transformations géométriques dans le plan qui permettent, pour une figure donnée  $D$ , d'en déduire les transformées par : une rotation, une homothétie, une symétrie, ..., et toute transformation obtenue par composition.

Une transformation affine du plan peut être décrite par une matrice  $3 \times 3$  :

$$M = \begin{pmatrix} a & b & m \\ c & d & n \\ 0 & 0 & 1 \end{pmatrix}$$

où  $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$  est la matrice de la transformation vectorielle associée et  $\begin{pmatrix} m \\ n \end{pmatrix}$  est le vecteur de translation.

Les coordonnées  $\begin{pmatrix} x' \\ y' \end{pmatrix}$  de l'image de  $\begin{pmatrix} x \\ y \end{pmatrix}$  sont alors données par :

$$\begin{pmatrix} x' \\ y' \\ 1 \end{pmatrix} = M \cdot \begin{pmatrix} x \\ y \\ 1 \end{pmatrix}$$

Par exemple la rotation d'angle  $\theta$  autour du point  $(\alpha, \beta)$  peut être décomposée en une translation de vecteur  $(-\alpha, -\beta)$  (qui nous ramène à l'origine), une rotation de centre  $(0,0)$  et d'angle  $\theta$  et une translation de vecteur  $(\alpha, \beta)$ . Ce qui donne :

$$\begin{pmatrix} x' \\ y' \\ 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & \alpha \\ 0 & 1 & \beta \\ 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} \cos(\theta) & -\sin(\theta) & 0 \\ \sin(\theta) & \cos(\theta) & 0 \\ 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 0 & -\alpha \\ 0 & 1 & -\beta \\ 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} x \\ y \\ 1 \end{pmatrix}$$

La figure que l'on veut déformer sera définie par un fichier de points.

### 8.8 Géométrie dans l'espace

Pour une représentation plane de l'espace, on peut utiliser une perspective cavalière décrite par les formules :

$$\begin{cases} X = X_0 + x \cdot \cos(\alpha) + y \cdot \cos(\beta) \\ Y = Y_0 + x \cdot \sin(\alpha) + y \cdot \sin(\beta) - z \end{cases}$$

où  $(x, y, z)$  sont les coordonnées d'un point de l'espace, avec l'axe  $oz$  vertical,  $\alpha$  est l'angle entre  $OX$  et  $ox$ ,  $\beta$  l'angle entre  $OY$  et  $oy$  et  $(X_0, Y_0)$  sont les coordonnées de  $o$  dans le plan :

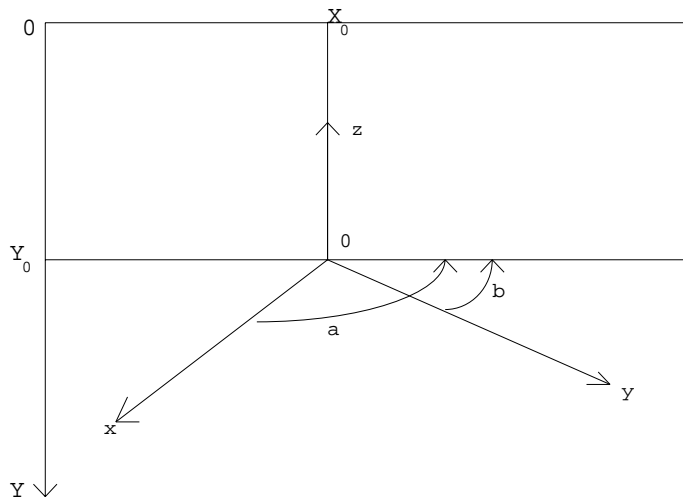


Figure 2.2

On écrira :

- des procédures permettant le tracé de courbes gauches définies par:

$$\begin{cases} x = f(t) \\ y = g(t) \\ z = h(t) \end{cases}$$

- des procédures permettant de tracer les projections de la courbe sur chacun des plans xoy, xoz et yoz ;
- des procédures permettant le tracé « en fil de fer » d'une surface d'équation:

$$z = f(x,y)$$

Si  $x$  varie de  $x_{\text{Min}}$  à  $x_{\text{Max}}$  et  $y$  de  $y_{\text{Min}}$  à  $y_{\text{Max}}$ , on pourra procéder comme suit :

- on coupe la surface par un plan d'équation  $x = h$  et on trace la courbe intersection en faisant varier  $y$  de  $y_{\text{Min}}$  à  $y_{\text{Max}}$  ;
- on coupe la surface par le plan  $y = k$  et on trace la courbe d'intersection en faisant varier  $x$  de  $x_{\text{Min}}$  à  $x_{\text{Max}}$  ;
- des procédures permettant d'éliminer les parties cachées dans le tracé précédent;
- des procédures de tracé de surfaces définies par:

$$\begin{cases} x = f(u, v) \\ y = g(u, v) \\ z = h(u, v) \end{cases}$$

avec et sans parties cachées.

## 8.9 Tracé de courbes paramétriques et polaires

Ecrire des paquetages permettant le tracé de courbes d'équations paramétriques  $x = f(t)$ ,  $y = g(t)$  et d'équation polaire  $r = f(\theta)$ .

## 8.10 Recherche des lignes de niveaux d'une fonction

Soit  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , une fonction continue, il s'agit de tracer les courbes d'équation:

$$f(x,y) = \text{Cste}$$

On définit d'abord un rectangle de travail  $[a,b] \times [c,d]$  que l'on maille en  $m$  lignes et  $n$  colonnes, en posant :

$$hx = \frac{b - a}{n} ; hy = \frac{d - c}{m}$$

$$x_i = a + i \cdot hx ; y_j = c + j \cdot hy \quad (i = 1, \dots, n; j = 1, \dots, m)$$

On pose alors :

$$F_{i,j} = f(x_i, y_j) \quad (i = 1, \dots, n; j = 1, \dots, m)$$

A partir du tableau F, on calcule :

$$f\text{Min} = \text{Min} \{ F_{i,j}, i = 1, \dots, n, j = 1, \dots, m \}$$

$$f\text{Max} = \text{Max} \{ F_{i,j}, i = 1, \dots, n, j = 1, \dots, m \}$$

Chaque petit rectangle est divisé en deux triangles et on trace les lignes de niveaux par triangle.

Soit un de ces triangles :

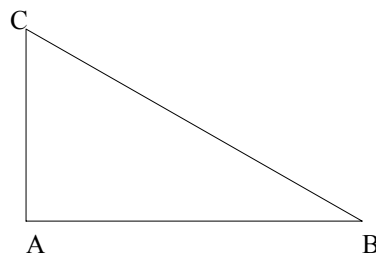


Figure 2.3

La fonction  $f$  sera alors remplacée sur chacun de ces triangles par une fonction affine :

$$f(x,y) \cong \alpha + \beta \cdot x + \gamma \cdot y$$

La fonction  $f$  étant connue en trois points A, B, C, les coefficients  $\alpha$ ,  $\beta$  et  $\gamma$ , s'en déduisent facilement.

Le problème est alors de savoir comment déterminer les morceaux de lignes de niveau qui coupent le triangle.

En programmation linéaire, on montre qu'une fonction affine sur un triangle (de manière plus générale sur un polygone) atteint ses extrema en un sommet.

Connaissant  $f$  en A, B, C, on calcule les extrema et on regarde si la valeur constante cherchée ( $f = \text{Cste}$ ) est entre ses bornes.

En général la courbe de niveau coupe le triangle en deux cotés. On aura des ennuis quand elle passera par l'un des sommets (à détailler).

*Exemples de courbes :*

Ellipse:  $a \cdot x^2 + b \cdot y^2 + \alpha \cdot x + \beta \cdot y + \gamma$  ( $a > 0, b > 0, \gamma - \frac{1}{4} \cdot \left( \frac{\alpha^2}{a} + \frac{\beta^2}{b} \right) < 0$ )

Hyperbole:  $a \cdot x^2 - b \cdot y^2 + \alpha \cdot x + \beta \cdot y + \gamma$  ( $a > 0, b > 0, \gamma - \frac{1}{4} \cdot \left( \frac{\alpha^2}{a} + \frac{\beta^2}{b} \right) < 0$ )

Parabole:  $a \cdot x^2 + \alpha \cdot x + \beta \cdot y + \gamma$  ( $a > 0$ )

? :  $x + y - \text{Ln}(x) - \text{Ln}(y) + C$

## CHAPITRE 3

# Analyse numérique linéaire

### 1. Introduction

#### 1.1 Position des problèmes

Dans ce chapitre, on étudie les trois problèmes suivants :

- la résolution d'un système linéaire  $A \cdot x = b$  ;
- l'inversion d'une matrice  $A$  ;
- la recherche des valeurs propres et des vecteurs propres d'une matrice  $A$ .

#### 1.2 Notations

Une matrice sera notée  $A = ((a_{ij}))_{1 \leq i, j \leq n}$ , où :

$i$  est le numéro de la ligne ;

$j$  est le numéro de la colonne.

Un vecteur sera noté  $v = (v_i)_{1 \leq i \leq n}$ .

#### 1.3 Remarques

(i) Si  $A$  et  $b$  sont à coefficients complexes, en écrivant :  $A = C + i \cdot D$ ,  $b = \alpha + i \cdot \beta$ ,  $x = \xi + i \cdot \eta$ , avec  $C, D, \alpha, \beta, \xi, \eta$  à coefficients réels, la résolution du système de  $n$  équations à  $n$  inconnues  $A \cdot x = b$  dans  $\mathbb{C}^n$  se ramène à un système de  $2 \cdot n$  équations à  $2 \cdot n$  inconnues, à coefficients réels :

$$\begin{pmatrix} C & -D \\ D & C \end{pmatrix} \cdot \begin{pmatrix} \xi \\ \eta \end{pmatrix} = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

Il nous suffit donc de considérer des systèmes à coefficients réels.

(ii) Si  $A$  est inversible, les coefficients de  $A^{-1}$  sont les  $n$  solutions des  $n$  systèmes linéaires :

$$A \cdot x = e_j \quad (j = 1, 2, \dots, n) ;$$

où  $\{e_1, \dots, e_n\}$  désigne la base canonique de  $\mathbb{R}^n$ .

C'est-à-dire que le calcul de l'inverse d'une matrice se ramène à la résolution simultanée de  $n$  systèmes linéaires de même matrice.

#### 1.4 Problèmes numériques liés à la résolution des systèmes linéaires

La précision numérique des ordinateurs n'étant pas infinie, un algorithme de résolution numérique d'un système linéaire ne donnera pas, en général, une solution exacte mais seulement une approximation de cette solution.

D'autre part il sera important de savoir évaluer le temps de calcul nécessaire à un tel algorithme .

Les erreurs numériques sont de trois types :

- les erreurs sur les données (pour des données expérimentales) ;
- les erreurs d'arrondis (calculs en virgule flottante) ;
- les erreurs de troncature.

Les erreurs d'arrondis vont s'accumuler au cours des calculs, il est donc important de connaître le nombre d'opérations élémentaires que nécessite un tel algorithme.

#### 1.5 Systèmes dégénérés et numériquement dégénérés

Un système  $A \cdot x = b$  est dit (théoriquement) dégénéré si  $\text{Dét}(A) = 0$ . Mais vérifier cette condition de manière exacte nécessite une précision infinie, ce qui n'est pas le cas, il faut donc définir un concept de dégénérescence numérique d'un système linéaire.

*Exemple 1* — Considérons le système :

$$\begin{cases} x + y = 2 \\ (1 + a) \cdot x + y = 2 + a \end{cases}$$

où  $a$  est « très petit ».

On a  $\text{Dét}(A) = -a$ , qui est supposé non nul. Mais si  $a = 10^{-12}$ , pour un ordinateur qui travaille avec 11 chiffres significatifs la représentation machine de  $A$  et  $b$  sera :

$$A_{\text{mach}} = \begin{pmatrix} 1.000000000000 & 1.000000000000 \\ 1.000000000000 & 1.000000000000 \end{pmatrix} \text{ et } b_{\text{mach}} = \begin{pmatrix} 2.000000000000 \\ 2.000000000000 \end{pmatrix}$$

et le système apparaîtra comme dégénéré.

*Exemple 2* — Si on considère le système dégénéré :

$$\begin{cases} 3 \cdot x + 5 \cdot y = 4 \\ 3 \cdot x + 5 \cdot y = 4 \end{cases}$$

l'utilisation de la méthode de Gauss (voir paragraphe 3.4) donne  $x = 0.500000000000$  et  $y = 0.500000000000$  comme solutions. Ce qui n'est pas faux, mais ce n'est pas la seule solution, (il y en a une infinité) ce que ne sait pas détecter l'ordinateur.

Un système sera dit numériquement dégénéré, si la valeur de son déterminant calculé en machine est non significative (i.e. trop petite ou trop grande, c'est-à-dire en dehors des limites de la machine).

*Exemple 3* — La matrice de Hilbert d'ordre  $n$  est définie par  $H = ((h_{ij}))_{1 \leq i, j \leq n}$  où  $h_{ij} = \frac{1}{i+j-1}$  pour  $1 \leq i, j \leq n$ . C'est une matrice symétrique définie positive, donc non dégénérée. En fait, on a :

$$\text{Dét}(H) = \frac{\Phi^4(n-1)}{\Phi(2 \cdot n - 1)}, \text{ où } \Phi(n) = \prod_{k=1}^n k!$$

Pour  $n$  grand, on aura  $\text{Dét}(A)_{\text{mach}} = 0$ . Un tel système est donc numériquement dégénéré, mais non dégénéré.

Pour tester la performance d'un algorithme de résolution d'un système linéaire, on pourra l'essayer avec une matrice de Hilbert d'ordre élevé.

Des exemples de calculs de déterminants de matrices de Hilbert, en utilisant, d'une part la décomposition L-R (Cf. paragraphe (3.6)) et d'autre part la formule de Hilbert sont donnés dans le tableau 3.1.

Les résultats sont obtenus en utilisant les réels de base, de type FLOAT, de OpenAda.

Pour la formule de Hilbert, une méthode de calcul est basée sur les récurrences :

$$\Phi(n) = n! \cdot \Phi(n-1)$$

$$\Delta_{n+1} = \frac{\Delta_n}{(2 \cdot n + 1) \cdot (C_{2n}^n)^2}$$

où on a noté  $\Delta_n$  le déterminant de la matrice de Hilbert d'ordre  $n$ .

Si  $u_n = C_{2n}^n$ , on a  $u_{n+1} = \frac{2 \cdot (2 \cdot n + 1)}{n + 1} \cdot u_n$  ce qui donne en définitive l'algorithme de calcul suivant :

D = 1 ;

u = 2 ;

Pour n allant de 1 à nMax faire

Début

$$\Delta = \frac{\Delta}{(2 \cdot n + 1) \cdot u^2} ;$$

$$u = \frac{2 \cdot (2 \cdot n + 1)}{n + 1} \cdot u ;$$

Fin ;

Pour de grandes valeurs de  $n$  la matrice n'apparaît pas numériquement dégénérée mais les résultats ne sont pas fiables du fait des accumulations d'erreurs



d'arrondis, aussi bien pour le déterminant calculé en utilisant la décomposition L-R que pour le calcul utilisant la formule de Hilbert.

n	Déterminant Formule	Déterminant calculé avec L-R
1	1.00000000000000E+00	1.00000000000000E+00
2	8.33333333333333E-02	8.33333333333333E-02
3	4.62962962962963E-04	4.62962962962961E-04
4	1.65343915343915E-07	1.65343915343930E-07
5	3.74929513251509E-12	3.74929513251290E-12
6	5.36729988735869E-18	5.36729988682088E-18
7	4.83580262392612E-25	4.83580261566227E-25
8	2.73705011379151E-33	2.73705023258196E-33
9	9.72023431192500E-43	9.72026335654786E-43
10	2.16417922643149E-53	2.16436997064480E-53
11	3.01909533444935E-65	3.02449714771783E-65
12	2.63778065125355E-78	2.70874882840569E-78
13	1.44289651879114E-92	Pivot trop petit dans LR
14	4.94031491459083E-108	Pivot trop petit dans LR
15	1.05854274306972E-124	Pivot trop petit dans LR
16	1.41913921144317E-142	Pivot trop petit dans LR
17	1.19027124270300E-161	Pivot trop petit dans LR
18	6.24486061412311E-182	Pivot trop petit dans LR
19	2.04934373331812E-203	Pivot trop petit dans LR
20	4.20617895662473E-226	Pivot trop petit dans LR
21	5.39898624190004E-250	Pivot trop petit dans LR

Tableau 3.1

### 1.6 Problèmes de stabilité numérique

Dans la pratique, les coefficients de la matrice et du second membre ne sont connus que de façon approximative. Et un problème important est donc de savoir si de petites variations sur ces coefficients peuvent entraîner de grosses variations sur la solution. La réponse est hélas oui comme le montre l'exemple du système  $A \cdot x = b$ , avec :

$$A = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \text{ et } b = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

dont la solution est  $x = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$ . Si on modifie le second membre en  $b = \begin{pmatrix} 32.1 \\ 22.9 \\ 33.1 \\ 30.9 \end{pmatrix}$ ,

c'est-à-dire que  $b$  est donné avec une erreur relative de 0.3%, alors la solution devient  $x = \begin{pmatrix} 9.2 \\ -12.6 \\ 4.5 \\ -1.1 \end{pmatrix}$ . On a donc une erreur relative sur  $x$  de l'ordre de 1000%.

De même si on perturbe la matrice en prenant  $A = \begin{pmatrix} 10 & 7 & 8.1 & 7.2 \\ 7.08 & 5.04 & 6 & 5 \\ 8 & 5.98 & 9.81 & 9 \\ 6.99 & 4.99 & 9 & 9.98 \end{pmatrix}$

et en gardant le second membre initial, la solution devient  $x = \begin{pmatrix} -8.1 \\ 137 \\ -34 \\ 22 \end{pmatrix}$  Un tel

système est dit mal conditionné.

Le conditionnement d'une matrice n'est pas lié à un déterminant voisin de 0. Pour l'exemple ci-dessus, la matrice est symétrique définie positive de déterminant égal à 1.

Dans le paragraphe 2.2 on définira une quantité permettant de mesurer le bon ou mauvais conditionnement d'un système.

### 1.7 Méthodes de résolution des systèmes linéaires

On a deux catégories de méthodes de résolution d'un système linéaire :

- les méthodes directes ;
- les méthodes itératives.

Si on suppose une précision infinie, les méthodes directes conduisent à la solution exacte du système en un nombre fini d'étapes.

Alors que les méthodes itératives donnent toujours une approximation de la solution, la solution exacte étant obtenue en un nombre infini d'étapes.

Par exemple, la méthode de Cramer est une méthode directe de résolution des systèmes linéaires (mais impraticable).

Comme il n'existe pas d'algorithme direct de résolution d'une équation polynomiale, les méthodes de recherche de valeurs propres seront nécessairement des méthodes itératives.

Les méthodes directes sont plutôt adaptées au « petits systèmes » ( $n < 50$ ) et les méthodes itératives aux « grands systèmes ».

Un aspect intéressant de la matrice  $A$  est le nombre de termes  $a_{ij}$  qui sont nuls. Une matrice avec beaucoup de zéros sera dite creuse. Les méthodes itératives

seront bien adaptées aux matrices creuses et en particulier aux matrices bandes (i.e. telles que  $a_{ij} = 0$  pour  $|i - j| > p$ ) car cette structure bande sera conservée.

Enfin, les problèmes avec des matrices symétriques seront plus faciles à traiter, en particulier les problèmes de valeurs propres.

### ***1.8 Exemples d'applications***

Dans ce cours, après avoir exposé quelques méthodes d'analyse numérique linéaire, nous les utiliseront pour la résolution numérique de certaines équations différentielles (méthode des différences finies et méthode des éléments finis pour des problèmes aux limites en dimension un) ou aux dérivées partielles (méthode

des différences finies et méthode des éléments finis pour des problèmes aux limites en dimension 2) et pour résoudre des problèmes d'interpolation et d'approximation.

En utilisant la méthode de Newton-Raphson, la résolution des systèmes non linéaires se ramène à des résolutions de systèmes linéaires. La méthode de Newton-Raphson sera utilisée avec la méthode du tir qui permet de résoudre des systèmes d'équations différentielles avec conditions aux limites.

## 2. Rappels et compléments sur le calcul matriciel

Dans tout ce paragraphe, on désigne par  $A = ((a_{ij}))_{1 \leq i, j \leq n}$  une matrice dans  $M_n(\mathbb{R}) = \{ \text{Matrices à } n \text{ lignes et } n \text{ colonnes à coefficients réels} \}$ .

### 2.1 Normes vectorielles et matricielles

#### 2.1.1 Norme matricielle induite par une norme vectorielle

*Définition* : Si  $x \mapsto \|x\|$  est une norme sur  $\mathbb{R}^n$ , en posant, pour toute matrice  $A \in M_n(\mathbb{R})$  :

$$\|A\| = \text{Sup} \left\{ \frac{\|A \cdot x\|}{\|x\|} ; x \in \mathbb{R}^n - \{0\} \right\}$$

on définit une norme sur  $M_n(\mathbb{R})$  appelée norme matricielle induite par  $\|\cdot\|$ .

*Remarque 1* — Avec la linéarité de  $A$ , on déduit qu'on a aussi :

$$\|A\| = \text{Sup} \{ \|A \cdot x\| ; x \in \mathbb{R}^n \text{ et } \|x\| = 1 \}$$

*Remarque 2* — De la définition, il résulte que :

$$\forall x \in \mathbb{R}^n, \forall A \in M_n(\mathbb{R}), \|A \cdot x\| \leq \|A\| \cdot \|x\|$$

$$\forall A, B \in M_n(\mathbb{R}), \|A \cdot B\| \leq \|A\| \cdot \|B\|$$

$$\text{Si } I_d \text{ est la matrice identité, alors } \|I_d\| = 1$$

*Remarque 3* — Il existe des normes matricielles qui ne sont pas induites par une norme vectorielle. Un exemple classique est donné par la « norme de Schur » définie par :

$$N_s(A) = \left( \text{Trace}({}^t A \cdot A) \right)^{\frac{1}{2}} = \left( \sum_{i,j=1}^n a_{ij}^2 \right)^{\frac{1}{2}}$$

C'est simplement la norme euclidienne de  $A$  considérée comme vecteur de  $\mathbb{R}^{n^2}$ . Comme  $N_s(I_d) = \sqrt{n}$  cette norme n'est pas induite par une norme vectorielle.

#### 2.1.2 Norme utilisée dans ce livre

Nous prendrons, sauf indication contraire, la norme vectorielle définie par :

$$\|x\| = \text{Max} \{ |x_i| ; 1 \leq i \leq n \}$$

On a alors le résultat suivant qui nous permet de calculer effectivement la norme matricielle induite :

*Théorème* : Pour toute matrice  $A$  dans  $M(\mathbb{R})$ , on a :

$$\|A\| = \text{Max} \left\{ \sum_{j=1}^n |a_{ij}| ; 1 \leq i \leq n \right\}$$

pour la norme matricielle induite par  $\|\cdot\|$  définie ci-dessus.

*Démonstration* — Posons  $\alpha = \text{Max} \left\{ \sum_{j=1}^n |a_{ij}| ; 1 \leq i \leq n \right\}$ .

(a) Pour  $x$  dans  $\mathbb{R}^n$  tel que  $\|x\| = 1$ , on a :

$$\|A \cdot x\| = \text{Max} \left\{ \sum_{j=1}^n |a_{ij} \cdot x_j| ; 1 \leq i \leq n \right\} \leq \alpha$$

(b) Pour montrer que  $\alpha \leq \|A\|$ , il suffit de trouver un vecteur  $x$  tel que  $\|A \cdot x\| = \alpha$ .

Soit  $k$  tel que  $\alpha = \sum_{j=1}^n |a_{kj}|$ . On pose, pour  $j = 1, \dots, n$

$$x_j = \begin{cases} \text{Signe}(a_{kj}) & \text{si } a_{kj} \neq 0 \\ 0 & \text{si } a_{kj} = 0 \end{cases}$$

Et on a alors  $\|A \cdot x\| = \alpha$ .

## 2.2 Conditionnement d'une matrice

Soit  $A$  une matrice inversible et  $x$  solution du système  $A \cdot x = b$ . En notant  $x + \delta x$  la solution du système perturbé  $A \cdot y = b + \delta b$ , on a :

$$\frac{\|\delta x\|}{\|x\|} \leq \|A\| \cdot \|A^{-1}\| \cdot \frac{\|\delta b\|}{\|b\|}$$

de même si  $x + \delta x$  est solution de  $(A + \delta A) \cdot y = b$ , on a :

$$\frac{\|\delta x\|}{\|x + \delta x\|} \leq \|A\| \cdot \|A^{-1}\| \cdot \frac{\|\delta A\|}{\|A\|}$$

ce qui nous amène à poser la définition suivante :

*Définition* : Si  $A$  est une matrice inversible et  $\|\cdot\|$  une norme matricielle induite par une norme vectorielle, alors le conditionnement de  $A$  relativement à cette norme est la quantité :

$$\text{Cond}(A) = \frac{1}{\|A\| \cdot \|A^{-1}\|}$$

Le conditionnement est donc un réel compris entre 0 et 1.

Un système sera bien conditionné si  $\text{Cond}(A)$  est voisin de 1. On a  $\text{Cond}(A) = 1$  pour les matrices unitaires.

Pour l'exemple de matrice mal conditionnée du paragraphe 1.6, on a  $\text{Cond}(A) \cong 1/3000$ .

Pour plus de détails sur ces problèmes de conditionnement, on pourra consulter Théodor et Lascaux, Chap. 3 (p. 167 à 206).

## 2.3 Valeurs propres et rayon spectral d'une matrice

### 2.3.1 Définitions

On dit que  $\mu$  dans  $C$  est *valeur propre* de  $A$  s'il existe un vecteur non nul  $x$  dans  $C^n$  tel que  $A \cdot x = \mu \cdot x$ . On dit alors que  $x$  est un *vecteur propre* associé à  $A$ . Ce qui équivaut à dire que  $A - \mu \cdot I_d$  est non inversible, soit que le *polynôme caractéristique* de  $A$ ,  $P(t) = \text{Dét}(A - t \cdot I_d)$ , est nul pour  $t = \mu$ .

On dit que  $A$  est *diagonalisable*, s'il existe une base de  $R^n$  formée de vecteurs propres de  $A$ .

L'ensemble des valeurs propres de  $A$  est appelé le *spectre de  $A$*  et noté  $\text{Sp}(A)$ . C'est une partie finie de  $C$  ayant au plus  $n$  éléments.

Le *rayon spectral* de  $A$  est alors défini par :

$$\rho(A) = \text{Max} \{ |\mu| ; \mu \in \text{Sp}(A) \}$$

### 2.3.2 Propriétés du rayon spectral

*Théorème* : Si  $A$  est dans  $M_n(R)$ , on a :

$$\rho(A) = \text{Inf} \{ \|A\| ; \|\cdot\| \text{ est une norme matricielle induite par une norme vectorielle} \}$$

*Démonstration* — (a) Pour toute valeur propre  $\mu$  de vecteur propre associé  $x$  et pour toute norme matricielle induite par une norme vectorielle, on a :  $\|A \cdot x\| \leq \|A\| \cdot \|x\|$ , avec  $\|A \cdot x\| = \mu \cdot \|x\|$ , on en déduit donc que  $|\mu| \leq \|A\|$ , pour tout  $\mu$  dans  $\text{Sp}(A)$ , donc  $\rho(A) \leq \|A\|$ .

(b) On peut d'autre part montrer que pour tout  $\tau > 0$ , il existe une norme matricielle induite, telle que  $\|A\| \leq \rho(A) + \tau$  (Voir Ciarlet p.19). Ce qui donne le résultat.

*Corollaire* — Soit  $A$  dans  $M_n(R)$ , les trois conditions suivantes sont équivalentes :

(i)  $\lim_{k \rightarrow +\infty} (A^k) = 0$  ;

(ii) pour toute valeur initiale  $x_0$ , la suite définie par  $x_{k+1} = A \cdot x_k$ , pour  $k \geq 0$ , converge vers le vecteur nul ;

(iii)  $\rho(A) < 1$ .

*Démonstration* — (i)  $\Rightarrow$  (ii). Résulte de :  $\|x_k\| = \|A^k \cdot x_0\| \leq \|A^k\| \cdot \|x_0\|$ .

(ii)  $\Rightarrow$  (iii). Supposons qu'il existe une valeur propre  $\mu$  telle que  $1/\mu \geq 1$ . Alors, si  $x_0$  est un vecteur propre associé à  $\mu$ , en écrivant que  $x_k = A^k \cdot x_0 = \mu^k \cdot x_0$ , on déduit que la suite  $(x_k)$  ne converge pas vers 0.

(iii)  $\Rightarrow$  (i). Soit  $\tau > 0$  tel que  $\rho(A) + \tau < 1$  et  $\|\cdot\|$  une norme matricielle induite telle que  $\|A\| < \rho(A) + \tau$ . Avec  $\|A^k\| \leq \|A\|^k$ , on déduit alors que  $\lim_{k \rightarrow +\infty} (A^k) = 0$ .

*Remarque* — On peut aussi montrer que, pour toute norme matricielle :

$$\rho(A) = \lim_{k \rightarrow +\infty} \left( \|A^k\|^{1/k} \right).$$

(Voir Ciarlet p. 22).

## 2.4 Quelques propriétés des matrices symétriques réelles

Le produit scalaire euclidien de  $\mathbb{R}^n$  est noté :

$$(x, y) \mapsto \langle x | y \rangle = \sum_{i=1}^n x_i \cdot y_i$$

La norme associée est la *norme euclidienne* et notée  $x \mapsto \|x\|_2$ .

### 2.4.1 Transposée d'une matrice

La *transposée* de la matrice  $A = ((a_{ij}))_{1 \leq i, j \leq n}$ , est la matrice :  ${}^tA = ((a_{ji}))_{1 \leq i, j \leq n}$

Elle est aussi définie par :

$$\forall x, y \in \mathbb{R}^n : \langle x | {}^tA \cdot y \rangle = \langle A \cdot x | y \rangle$$

La matrice  $A$  est dite *symétrique*, si  ${}^tA = A$ . On dit qu'elle est *symétrique définie positive* si on a, de plus :  $\forall x \in \mathbb{R}^n - \{0\} : \langle A \cdot x | x \rangle > 0$ .

### 2.4.2 Spectre des matrices symétriques

Pour ce qui est des valeurs propres d'une matrice symétrique réelle, la situation est plus simple que dans le cas général, ce que nous résumons par le :

*Théorème* : (1) Les valeurs propres d'une matrice symétrique réelle sont toutes réelles.  
 (2) Si, de plus, la matrice est définie positive, alors toutes les valeurs propres sont strictement positives.  
 (3) Une matrice symétrique se diagonalise dans une base orthonormale. Ce qui se traduit en disant qu'il existe une matrice orthogonale  $O$  (i.e. telle que  $O^{-1} = {}^tO$ ) telle que  $A = {}^tO \cdot D \cdot O$ , où  $D = \text{Diag}(\mu_1, \dots, \mu_n)$ , les  $\mu_i$  étant les valeurs propres de  $A$ .

## 2.5 Matrices à diagonale strictement dominante

*Définition* — Une matrice  $A$  est à diagonale strictement dominante, si :

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}| \quad (i = 1, \dots, n)$$

*Remarque* — Les matrices à diagonale strictement dominante se rencontrent dans de nombreux problèmes, par exemple dans le problème de l'interpolation par des fonctions splines cubiques ou dans les problèmes de résolutions d'équations aux dérivées partielles par des méthodes de discrétisation par différences finies.

*Théorème* : Une matrice à diagonale strictement dominante est inversible.

*Démonstration* — Montrons que toutes les valeurs propres de  $A$  sont non nulles.

Soit  $\mu$  une valeur propre de  $A$  et  $x$  un vecteur propre associé on a alors :

$$\forall i \in \{1, \dots, n\} \quad (\mu - a_{ii}) \cdot x_i = \sum_{j \neq i} a_{ij} \cdot x_j$$

Et en prenant un indice  $i$  tel que  $|x_i| = \|x\|$ , on déduit que :

$$|\mu - a_{ii}| \leq \sum_{j \neq i} |a_{ij}|.$$

On ne peut donc avoir  $\mu = 0$  si  $A$  est à diagonale strictement dominante. Donc  $A$  est inversible.

*Théorème* : Soit  $A$  une matrice symétrique réelle à diagonale strictement dominante. Alors,  $A$  est définie positive si, et seulement si  $a_{ii} > 0$ , pour tout  $i = 1, \dots, n$ .

*Démonstration* — Supposons que  $a_{ii} > 0$  pour tout  $i$ . En reprenant la démonstration du théorème précédent, on a :

$$|\mu - a_{ii}| \leq \sum_{j \neq i} |a_{ij}| < a_{ii}, \text{ pour tout indice } i$$

On en déduit donc que  $\mu > 0$ . La matrice  $A$  a donc toutes ses valeurs propres strictement positives, ce qui équivaut à dire qu'elle est définie positive.

La réciproque provient de l'égalité  $a_{ii} = \langle A \cdot e_i | e_i \rangle$ , où  $e_1, \dots, e_n$  est la base canonique de  $\mathbb{R}^n$

### 3. Résolution numérique des systèmes linéaires.

#### Inversion des matrices

##### 3.1 Position des problèmes. Notations

On se donne une matrice  $A$  dans  $M_n(\mathbb{R})$ , supposée inversible, un vecteur  $b$  dans  $\mathbb{R}^n$  et on veut résoudre le système linéaire :

$$A \cdot x = b \quad (1)$$

L'inverse de  $A$  sera obtenue en prenant pour valeurs successives de  $b$  les éléments de la base canonique de  $\mathbb{R}^n$ .

*Remarque* — Une matrice sera notée  $A = ((a_{ij}))_{1 \leq i, j \leq n}$ , où l'indice  $i$  est le numéro de la ligne et  $j$  celui de la colonne. Mais il faut être conscient du fait que la



machine va stocker un tel objet sur un seul tableau, soit par colonnes (c'est le cas en Fortran) :

$$a_{11}, a_{21}, \dots, a_{n1}, a_{12}, \dots, a_{n2}, \dots, a_{1n}, \dots, a_{nn}$$

soit par lignes (c'est le cas en Ada, C et Pascal) :

$$a_{11}, a_{12}, \dots, a_{1n}, a_{21}, \dots, a_{2n}, \dots, a_{n1}, \dots, a_{nn}$$

Le stockage en colonnes peut être exploité en programmation alors que cela est moins pratique pour le stockage en lignes.

### 3.2 Une méthode impraticable : la méthode de Cramer

Si on note  $A_j$ , la matrice déduite de  $A$  en remplaçant la colonne  $N^o j$  par  $b$ , pour  $j = 1, 2, \dots, n$ , la solution du système (1) est donnée par les formules de Cramer :

$$x_j = \frac{\text{Dét}(A_j)}{\text{Dét}(A)} \quad (j = 1, 2, \dots, n)$$

Si on calcule un déterminant par la formule :

$$\text{Dét}(A) = \sum \text{Sgn}(\sigma) \cdot \prod a_{j,\sigma(j)}$$

où  $\sigma$  décrit l'ensemble des permutations de  $\{1, 2, \dots, n\}$  et  $\text{Sgn}(\sigma)$  est la signature de  $\sigma$ , cela nécessite :

$n! \cdot (n - 1)$  multiplications et  $(n! - 1)$  additions,

donc, environ  $n \cdot n!$  opérations élémentaires.

Comme il y a  $n$  déterminants à calculer puis  $n$  divisions à faire, on aura un total, environ, de  $n^2 \cdot n!$  opérations élémentaires.

En considérant qu'une opération élémentaire demande à un micro-ordinateur une microseconde, pour un système de 20 équations à 20 inconnues, la méthode de Cramer demandera  $10^{15}$  secondes soit plus de  $10^7$  années, c'est-à-dire plus de 10 millions d'années !!!

Cette méthode n'a donc qu'un intérêt théorique, ce qui est loin d'être négligeable.

### 3.3 Cas des systèmes triangulaires

On suppose, ici, que  $A$  est une matrice triangulaire supérieure.

#### 3.3.1 Résolution des systèmes triangulaires

Le système (1) s'écrit alors :

$$a_{ii} \cdot x_i + a_{i,i+1} \cdot x_{i+1} + \dots + a_{i,n} \cdot x_n = b_i \quad (1 \leq i \leq n).$$

La résolution se fait donc « en remontant » :

$$\begin{cases} x_n = \frac{b_n}{a_{nn}} \\ x_i = \frac{\left( b_i - \sum_{j=i+1}^n a_{ij} \cdot x_j \right)}{a_{ii}} \quad (i = n-1, \dots, 1) \end{cases}$$

Ce qui donne la programmation structurée :

*PROCEDURE SystTriSup(Entrée n : Entier ; A : Matrice ; b : Vecteur ; Sortie x : Vecteur) ;*

*Début*

*Pour i allant de n à 1 faire*

*Début*

*S = 0 ;*

*Pour j allant de i + 1 à n faire S = S + a<sub>ij</sub>x<sub>j</sub> ;*

*Si (a<sub>ii</sub> = 0) Alors Stop('Dét(A) = 0) ;*

*x<sub>i</sub> = (b<sub>i</sub> - S)/a<sub>ii</sub> ;*

*Fin ;*

*Fin ;*

Pour un système triangulaire inférieur, on procède de même « en descendant ».

### 3.3.2 Inversion des matrices triangulaires

Les colonnes de  $A^{-1}$  sont les solutions de  $A \cdot x = e_i$ , où  $\{e_1, \dots, e_n\}$  est la base canonique de  $\mathbb{R}^n$ . Ce qui donne, si A est triangulaire inférieure, pour les coefficients de  $U = A^{-1}$  :

$$\begin{cases} u_{ii} = \frac{1}{a_{ii}} \\ u_{ij} = 0 \quad (1 \leq i < j \leq n) \\ u_{ki} = -\frac{\sum_{j=i}^{k-1} a_{kj} \cdot u_{ji}}{a_{kk}} \quad (1 \leq i < k \leq n) \end{cases}$$

Si A est triangulaire supérieure,  ${}^tA$  est triangulaire inférieure et il suffit d'écrire :

$$A^{-1} = {}^t(A^{-1}).$$

Le déterminant d'une telle matrice se calculant facilement par :

$$\text{Dét}(A) = \prod_{i=1}^n a_{ii}$$

On a donc la programmation structurée :

*PROCEDURE InvTriInf(Entrée n : Entier ; A : Matrice ; Sortie U : Matrice) ;*

*Début*

*Pour i allant de 1 à n faire*

*Début*

*Si (a<sub>ii</sub> = 0) Alors Stop('Matrice non inversible') ;*

*u<sub>ii</sub> = 1/a<sub>ii</sub> ;*

*Pour j allant de i + 1 à n faire u<sub>ij</sub> = 0 ;*

*Fin ;*

*Pour i allant de 1 à n - 1 faire*

*Début*

*Pour k allant de i + 1 à n faire*

*Début*

*S = 0 ;*

*Pour j allant de i à k - 1 faire S = S + a<sub>kj</sub>u<sub>ji</sub> ;*

*u<sub>ki</sub> = -S/a<sub>kk</sub> ;*

*Fin ;*

Fin ;  
Fin ;

### 3.4 Méthode des pivots de Gauss

#### 3.4.1 Principe de la méthode

En faisant des opérations élémentaires sur le système (1) (i.e. permuter deux lignes ou ajouter à une ligne un multiple d'une autre, c'est-à-dire des opérations qui ne vont pas changer l'ensemble des solutions de (1)), on transforme ce système en un système triangulaire supérieur :  $R \cdot x = c$ . De plus comme une permutation de ligne change le déterminant de signe et que le fait d'ajouter un multiple d'une ligne ne change pas ce dernier, on aura

$$\text{Dét}(A) = \text{Dét}(R) = \pm \prod_{i=1}^n r_{ii}$$

L'intérêt de cette méthode est surtout pédagogique, dans la pratique on lui préférera la méthode L-R (Cf. paragraphe (3.6)) quand cette dernière peut s'appliquer, ou la méthode de Gaus-Jordan (Cf. paragraphe (3.5)).

#### 3.4.2 Description de la méthode

On note  $L_i$  la ligne N°  $i$  d'un système linéaire.

*Etape 0* — On se ramène à un système tel que  $a_{11}$  non nul.

Si pour tout  $i = 1, 2, \dots, n$ , on a  $a_{i1} = 0$ , alors  $\text{Dét}(A) = 0$  et c'est fini. Sinon, il existe  $i > 1$  tel que  $a_{i1}$  non nul, et en permutant les lignes 1 et  $i$  (si  $i = 1$ , on ne fait rien), on se ramène à un système  $A^{(1)} \cdot x = b^{(1)}$ , avec  $a_{11}^{(1)}$  non nul.

Le coefficient  $a_{11}^{(1)}$  est le premier pivot.

On a alors  $\text{Dét}(A) = \pm \text{Dét}(A^{(1)})$ , avec le signe moins si, et seulement si il y a une permutation des lignes 1 et  $i$  avec  $i > 1$ .

*Etape 1* — Elimination de  $x_1$  dans les équations 2, ..., n.

On fait pour cela :

$$L_i^{(1)} - \frac{a_{i1}^{(1)}}{a_{11}^{(1)}} \cdot L_1^{(1)} \quad (i = 2, 3, \dots, n)$$

et, après une éventuelle permutation des lignes 2 et  $i > 3$ , on a le système :

$$A^{(2)} \cdot x = b^{(2)}, \text{ avec } a_{22}^{(2)} \text{ non nul.}$$

où

$$A^{(2)} = \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \cdot & \cdot & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & \cdot & \cdot & a_{2n}^{(2)} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & a_{n2}^{(2)} & \cdot & \cdot & a_{nn}^{(2)} \end{pmatrix}$$

Le coefficient  $a_{22}^{(2)}$  est le deuxième pivot.

*Etape k* — Elimination de  $x_k$  dans les équations  $k + 1, \dots, n$ .

A la fin de l'étape  $k - 1$ , on a obtenu le système :

$$A^{(k)} \cdot x = b^{(k)}, \text{ avec } a_{kk}^{(k)} \text{ non nul.}$$

et

$$A^{(k)} = \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \cdot & \cdot & \cdot & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & \cdot & \cdot & \cdot & a_{2n}^{(2)} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & 0 & a_{kk}^{(k)} & \cdot & a_{kn}^{(k)} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & 0 & a_{nk}^{(k)} & \cdot & a_{nn}^{(k)} \end{pmatrix}$$

Le coefficient  $a_{kk}^{(k)}$  est le pivot N° k.

On fait alors  $L_i^{(k)} - \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} \cdot L_k^{(k)}$  ( $i = k + 1, \dots, n$ ), puis une éventuelle permutation

des lignes  $k + 1$  et  $j > k + 1$  pour se ramener à  $a_{k+1,k+1}^{(k+1)}$  non nul.

Au bout de  $n - 1$  étapes, on est donc ramené à un système triangulaire supérieur :

$$A^{(n)} \cdot x = b^{(n)}.$$

De plus, on a  $\text{Dét}(A) = (-1)^p \cdot \text{Dét}(A^{(n)})$ , où  $p$  est le nombre de permutations qui ont été nécessaires pour avoir des pivots non nuls et  $\text{Dét}(A^{(n)})$  est le produit des pivots.

### 3.4.3 Remarque sur le choix des pivots : méthode de Gauss avec pivots partiels

Pour éviter de faire une division par un nombre trop petit, dans le choix du pivot, on aura intérêt, à l'étape  $k - 1$ , à permuter la ligne  $k$  avec la ligne  $j > k$  telle que  $|a_{jk}| = \text{Max}\{|a_{ik}|; i = k, \dots, n\}$ , de manière à avoir le pivot le plus grand possible en valeur absolue.

Si ce maximum est trop petit, alors le système est numériquement dégénéré et la méthode ne sera pas fiable.

### 3.4.4 Nombre d'opérations élémentaires dans la méthode de Gauss

Les opérations  $L_i^{(k)} - \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} \cdot L_k^{(k)}$  ( $i = k + 1, \dots, n$ ), à l'étape N° k, demandent :

- $n - k$  divisions ;
- $(n - k) \cdot (n - k)$  multiplications ;
- $(n - k) \cdot (n - k)$  additions .

Soit un total de :

$$\sum_{k=1}^{n-1} n - k = \frac{n \cdot (n - 1)}{2} \text{ divisions}$$

et

$$2 \cdot \sum_{k=1}^{n-1} (n - k)^2 = \frac{2 \cdot n \cdot (n - 1) \cdot (2 \cdot n - 1)}{6} \text{ additions et multiplications}$$

Il nous faut donc  $\frac{n \cdot (n - 1) \cdot (4n + 1)}{6} \approx \frac{2 \cdot n^3}{3}$  opérations élémentaires pour aboutir à un système triangulaire supérieur. Puis la résolution de ce dernier système nécessite  $n$  divisions et  $n \cdot (n - 1)$  additions et multiplications, soit  $n^2$  opérations élémentaires.

En résumé, la résolution d'un système de  $n$  équations à  $n$  inconnues par la méthode de Gauss va nécessiter un nombre d'opérations qui est un  $O(n^3)$ .

Par exemple, un système de 100 équations demandera moins d'une seconde.

### 3.4.5 Programmation structurée

La variable  $L$  utilisée dans la procédure PivotMax est le numéro de la ligne où se trouve le pivot le plus grand en valeur absolue, pour tout  $k = 1, 2, \dots, n$ .

La procédure PermuterLigne, facile à écrire permet de permuter deux lignes du système.

Après exécution de la procédure la matrice  $A$  est devenue triangulaire et le vecteur  $b$  est transformé, de sorte que les valeurs initiales de  $A$  et  $b$  sont perdues. C'est-à-dire que les variables  $A$  et  $b$  sont passées par adresse (c'est le mode Entrée\_Sortie).

*PROCEDURE PivotMax(Entrée  $k, n$  : Entier ; Entrée\_Sortie  $A$  : Matrice ;  $b$  : Vecteur) ;*

*Début*

*$L = k$  ;*

*Pour  $i$  Allant de  $k + 1$  à  $n$  Faire*

*Si  $|a_{ik}| > |a_{Lk}|$  Alors  $L = i$  ;*

*Si  $|a_{Lk}| < \varepsilon$  Alors Stop('Pivot trop petit') ;*

*Si  $L \neq k$  Alors PermuterLignes( $n, k, L, A, b$ ) ;*

*Fin ;*

*PROCEDURE Eliminer(Entrée  $k, n$  : Entier ; Entrée\_Sortie  $A$  : Matrice ;  $b$  : Vecteur) ;*

*Début*

*Pour  $i$  Allant de  $k + 1$  à  $n$  Faire*

*Début*

*$m = \frac{a_{ik}}{a_{kk}}$  ;*

*Pour  $j$  Allant de  $k + 1$  à  $n$  Faire  $a_{ij} = a_{ij} - m \cdot a_{kj}$  ;*

*$a_{ik} = 0$  ;*

*$b_i = b_i - m \cdot b_k$  ;*

*Fin ;*

*Fin ;*

*PROCEDURE Pivotage(Entrée  $n$  : Entier ; Entrée\_Sortie  $A$  : Matrice ;  $b$  : Vecteur) ;*

*Début*

*Pour  $k$  Allant de 1 à  $n - 1$  Faire*

*Début*

*PivotMax( $k, n, A, b$ ) ;*

*Eliminer( $k, n, A, b$ ) ;*

*Fin ;*

*Fin ;*

Ensuite la résolution de  $A \cdot x = b$  se fait en appelant la procédure de résolution d'un système triangulaire supérieur.

### 3.4.6 Remarques

(i) — Dans la programmation structurée ci-dessus, on peut aussi travailler avec une matrice à  $n$  lignes et  $n + 1$  colonnes en prenant  $b$  pour colonne  $n + 1$ .

De manière plus générale, pour résoudre en parallèle  $p$  systèmes linéaires de même matrice  $A$  et de seconds membres respectifs  $b_1, \dots, b_p$ , on travaillera avec une matrice à  $n$  lignes et  $n + p$  colonnes en prenant, pour  $j = 1, \dots, p$ ,  $b_j$  comme colonne  $n + j$ .

En prenant pour valeurs particulières des seconds membres les vecteurs de la base canonique de  $\mathbb{R}^n$ , on a ainsi un procédé de calcul de l'inverse d'une matrice. Mais, en fait, pour calculer l'inverse d'une matrice on préfère utiliser la méthode Gauss-Jordan qui est une variante de la méthode de Gauss (Cf. paragraphe (3.5)).

(ii) *Cas particulier des systèmes linéaires à coefficients entiers* — Si  $A$  et  $b$  sont à coefficients entiers relatifs, en modifiant la méthode de Gauss, on peut calculer la solution  $x$  de façon exacte comme vecteur à coefficients dans  $\mathbb{Q}$ , en utilisant, au besoin, un procédé de calcul en multiprécision sur les entiers.

On procède comme suit :

*Etape 0* — Après une éventuelle permutation de lignes, on a le système, à coefficients dans  $\mathbb{Z}$ ,  $A^{(1)} \cdot x = b^{(1)}$ , avec un premier pivot non nul.

*Etape 1* — Elimination de  $x_1$  dans les équations 2, ...,  $n$ .

Si on fait la division de la ligne 1 par le premier pivot, on perdra le caractère entier du système.

On fait plutôt :

garder la ligne 1 ;  
 $a_{i1}^{(1)} \cdot L_i^{(1)} - a_{i1}^{(1)} \cdot L_1^{(1)}$ , pour  $i = 2, \dots, n$ .

Ce qui donne le système,  $A^{(2)} \cdot x = b^{(2)}$ , avec un deuxième pivot non nul (au prix d'une éventuelle permutation).

*Etape 2* — Elimination de  $x_2$  dans les équations 3, ...,  $n$ .

Si on procède de la même façon pour cette étape, on constate que les coefficients de la matrice  $A^{(3)}$  et du vecteur  $b^{(3)}$  sont divisibles par le premier pivot. On peut donc s'autoriser une division par ce premier pivot pour cette étape 2, ce qui aura l'avantage de diminuer les coefficients obtenus.

*Etape  $k$*  — Elimination de  $x_k$  dans les équations  $k + 1, \dots, n$ .

Cette étape est décrite par les formules :

$$a_{ij}^{(k+1)} = \frac{\left( a_{kk}^{(k)} \cdot a_{ij}^{(k)} - a_{ik}^{(k)} \cdot a_{kj}^{(k)} \right)}{a_{k-1,k-1}^{(k-1)}}$$

pour  $i = k + 1, \dots, n$  et  $j = k + 1, \dots, n + 1$  (en prenant  $b$  comme colonne  $N^\circ n + 1$ ).

Ce qui donne au bout de  $n - 1$  étapes un système triangulaire supérieur à coefficients entiers qui peut se résoudre de façon exacte dans  $\mathbb{Q}^n$ .

*Calcul du déterminant* — Au signe près, on a :

$$\text{Dét}(A^{(2)}) = (a_{11}^{(1)})^{n-1} \cdot \text{Dét}(A)$$

$$\text{Dét}(A^{(3)}) = \left(\frac{a_{22}^{(2)}}{a_{11}^{(1)}}\right)^{n-2} \cdot \text{Dét}(A^{(2)}) = (a_{22}^{(2)})^{n-2} \cdot a_{11}^{(1)} \cdot \text{Dét}(A^{(1)})$$

De proche en proche, on a alors,  $\text{Det}(A^{(n)}) = (-1)^p \cdot \left(\prod_{k=1}^n a_{kk}^{(k)}\right) \cdot \text{Dét}(A)$ , ce qui donne, en tenant compte de l'expression de  $\text{Det}(A^{(n)})$  comme produit des pivots,  $\text{Dét}(A) = (-1)^p \cdot a_{nn}^{(n)}$ .

### 3.5 Méthode de Gauss-Jordan

#### 3.5.1 Méthode d'élimination de Gauss-Jordan

La méthode de Gauss-Jordan est une variante de la méthode de Gauss, qui consiste à diagonaliser la matrice  $A$ , à l'aide de transformations élémentaires (c.a.d. des permutations et des combinaisons linéaires de lignes).

A l'étape  $k$ , on ne se contente pas seulement d'éliminer le coefficient de  $x_k$  dans les lignes  $k + 1$  à  $n$ , mais on l'élimine aussi dans les lignes du dessus soit  $1$  à  $k - 1$ . Ce qui donnera au bout de  $n$  étapes un système diagonal.

Description de la  $k^{\text{ième}}$  étape du calcul ( $1 < k < n$ ) — A la  $k^{\text{ième}}$  étape, on est au départ dans la situation suivante :

$$\left\{ \begin{array}{l} x_1 \qquad \qquad \qquad + a_{1k}^{(k)} \cdot x_k + \dots + a_{1n}^{(k)} \cdot x_n = b_1^{(k)} \\ \dots \\ x_{k-1} + a_{k-1,k}^{(k)} \cdot x_k + \dots + a_{k-1,n}^{(k)} \cdot x_n = b_{k-1}^{(k)} \\ \qquad a_{kk}^{(k)} \cdot x_k + \dots + a_{kn}^{(k)} \cdot x_n = b_k^{(k)} \\ \qquad a_{k+1,k}^{(k)} \cdot x_k + \dots + a_{k+1,n}^{(k)} \cdot x_n = b_{k+1}^{(k)} \\ \dots \\ \qquad a_{nk}^{(k)} \cdot x_k + \dots + a_{nn}^{(k)} \cdot x_n = b_n^{(k)} \end{array} \right.$$

On peut supposer que  $a_{kk}^{(k)}$  est le pivot maximum, sinon, il suffit de permuter avec une des lignes suivantes.

On commence par tirer  $x_k$  de la  $k^{\text{ième}}$  équation :

$$x_k = \frac{(b_k^{(k)} - a_{k,k+1}^{(k)} \cdot x_{k+1} - \dots - a_{kn}^{(k)} \cdot x_n)}{a_{k,k}^{(k)}}$$

que l'on reporte ensuite dans les autres équations, ce qui donne :

$$\left\{ \begin{array}{l} x_1 + a_{1,k+1}^{(k+1)} \cdot x_{k+1} + \dots + a_{1n}^{(k+1)} \cdot x_n = b_1^{(k+1)} \\ \dots \\ x_{k-1} + a_{k-1,k+1}^{(k+1)} \cdot x_{k+1} + \dots + a_{k-1,n}^{(k+1)} \cdot x_n = b_{k-1}^{(k+1)} \\ x_k + a_{k,k+1}^{(k+1)} \cdot x_{k+1} + \dots + a_{kn}^{(k+1)} \cdot x_n = b_k^{(k+1)} \\ a_{k+1,k+1}^{(k+1)} \cdot x_{k+1} + \dots + a_{k+1,n}^{(k+1)} \cdot x_n = b_{k+1}^{(k+1)} \\ \dots \\ a_{n,k+1}^{(k+1)} \cdot x_{k+1} + \dots + a_{nn}^{(k+1)} \cdot x_n = b_n^{(k+1)} \end{array} \right.$$

avec, pour  $j = k + 1, \dots, n$  :

$$a_{kj}^{(k+1)} = \frac{a_{kj}^{(k)}}{a_{kk}^{(k)}} \quad \text{et} \quad b_k^{(k+1)} = \frac{b_k^{(k)}}{a_{kk}^{(k)}}$$

et, pour  $i$  dans  $\{1, \dots, n\} - \{k\}$ ,  $j$  dans  $\{k+1, \dots, n\}$  :

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - a_{ik}^{(k)} \cdot a_{kj}^{(k+1)} \quad \text{et} \quad b_i^{(k+1)} = b_i^{(k)} - a_{ik}^{(k)} \cdot b_k^{(k+1)}.$$

On obtient alors directement la solution, après la  $n^{\text{ème}}$  étape :

$$\left\{ \begin{array}{l} x_1 = b_1^{(n)} \\ \cdot \\ \cdot \\ x_n = b_n^{(n)} \end{array} \right.$$

Pour la recherche du pivot maximum, on utilisera la procédure PivotMax du paragraphe (3.4.5).

Ce qui donne la programmation structurée suivante, le vecteur  $b$  contenant en fin d'opération la solution du système.

*PROCEDURE GaussJordan(Entrée  $n$  : Entier ; Entrée\_Sortie  $A$  : Matrice ;  $b$  : Vecteur) ;*

*PROCEDURE EliminerGaussJordan(Entrée  $k, n$  : Entier ;*

*Entrée\_Sortie  $A$  : Matrice ;  $b$  : Vecteur) ;*

*Début*

*Pour  $j$  Allant de  $k + 1$  à  $n$  Faire  $a_{kj} = \frac{a_{kj}}{a_{kk}}$  ;*

*$b_k = \frac{b_k}{a_{kk}}$  ;  $a_{kk} = 1$  ;*

*Pour  $i$  Allant de  $1$  à  $k - 1$  Faire*

*Début*

*Pour  $j$  Allant de  $k + 1$  à  $n$  Faire  $a_{ij} = a_{ij} - a_{ik} * a_{kj}$  ;*

*$b_i = b_i - a_{ik} * b_k$  ;*



```

     $a_{ik} = 0$  ;
  Fin ;
  Pour  $i$  Allant de  $k + 1$  à  $n$  Faire
  Début
    Pour  $j$  Allant de  $k + 1$  à  $n$  Faire  $a_{ij} = a_{ij} - a_{ik} * a_{kj}$  ;
     $b_i = b_i - a_{ik} * b_k$  ;
     $a_{ik} = 0$  ;
  Fin ;
Fin ;

```

```

Début
  Pour  $k$  Allant de  $1$  à  $n$  Faire
  Début
    PivotMax( $k, n, A, b$ ) ;
    EliminerGaussJordan( $k, n, A, b$ ) ;
  Fin ;
Fin ;

```

### 3.5.2 Inversion et calcul du déterminant d'une matrice

On adaptera ce qui précède en remplaçant  $b$  par la matrice identité, ce qui revient à résoudre en parallèle les  $n$  systèmes linéaires de même matrice  $A$  et de seconds membres respectifs  $e_1, \dots, e_n$  formant la base canonique de  $\mathbb{R}^n$ .

On obtient alors la programmation structurée suivante :

*PROCEDURE InversionGaussJordan(Entrée  $n$  : Entier ; Entrée\_Sortie  $A, Inv$  : Matrice) ;*

*PROCEDURE InitInverse(Entrée :  $n$  : Entier ; Sortie  $Inv$  : Matrice) ;*

```

Début
  Pour  $i$  Allant de  $1$  à  $n$  Faire
  Début
    Pour  $j$  Allant de  $1$  à  $n$  Faire
    Début
      Si  $j \neq i$ 
      Alors  $Inv_{ij} = 0$ 
      Sinon  $Inv_{ij} = 1$ 
    Fin ;
  Fin ;
Fin ;

```

*PROCEDURE PivotMaxInv(Entrées  $k, n$  : Entiers ; Entrée\_Sortie  $A, Inv$  : Matrice) ;*

```

Début
   $L = k$  ;
  Pour  $i$  Allant de  $k + 1$  à  $n$  Faire
  Si  $|a_{ik}| > |a_{Lk}|$ 
  Alors  $L = i$  ;
  Si  $|a_{Lk}| < \epsilon$ 
  Alors Arrêter('Pivot trop petit') ;
  Si  $L \neq k$ 
  Alors Début
    Pour  $j$  Allant de  $k$  à  $n$  Echanger( $a_{kj}, a_{Lj}$ ) ;
    Pour  $j$  Allant de  $1$  à  $n$  Echanger( $Inv_{kj}, Inv_{Lj}$ ) ;
  Fin ;
Fin ;

```

*PROCEDURE EliminerInvGaussJordan(Entrées  $k, n$  : Entiers ; Entrée\_Sortie  $A, Inv$  : Matrice) ;*

*Début*

*Pour  $j$  Allant de  $k + 1$  à  $n$  Faire  $a_{kj} = a_{kj}/a_{kk}$  ;*

*Pour  $j$  Allant de  $1$  à  $n$  Faire  $Inv_{kj} = Inv_{kj}/a_{kk}$  ;*

*$a_{kk} = 1$  ;*

*Pour  $i$  Allant de  $1$  à  $k - 1$  Faire*

*Début*

*Pour  $j$  Allant de  $k + 1$  à  $n$  Faire  $a_{ij} = a_{ij} - a_{ik} * a_{kj}$  ;*

*Pour  $j$  Allant de  $1$  à  $n$  Faire  $Inv_{ij} = Inv_{ij} - a_{ik} * Inv_{kj}$  ;*

*$a_{ik} = 0$  ;*

*Fin*

*Pour  $i$  Allant de  $k + 1$  à  $n$  Faire*

*Début*

*Pour  $j$  Allant de  $k + 1$  à  $n$  Faire  $a_{ij} = a_{ij} - a_{ik} * a_{kj}$  ;*

*Pour  $j$  Allant de  $1$  à  $n$  Faire  $Inv_{ij} = Inv_{ij} - a_{ik} * Inv_{kj}$  ;*

*$a_{ik} = 0$  ;*

*Fin*

*Fin ;*

*Début*

*InitInverse( $n, Inv$ ) ;*

*Pour  $k$  Allant de  $1$  à  $n$  Faire*

*Début*

*PivotMaxInv( $k, n, A, Inv$ ) ;*

*EliminerInvGaussJordan( $k, n, A, Inv$ ) ;*

*Fin ;*

*Fin ;*

La matrice  $A$  étant passée par adresse, on a perdu sa valeur initiale en fin d'opération.

### 3.6 Interprétation algébrique de la méthode de Gauss.

#### Décomposition $L \cdot R$ ou méthode de Crout

##### 3.6.1 Condition nécessaire et suffisante pour qu'il n'y ait pas de permutations dans la méthode de Gauss

*Définition* : Les sous-matrices principales de la matrice  $A$  sont les matrices :

$$A_k = \left( (a_{ij})_{1 \leq i, j \leq k} \right), \quad (k = 1, \dots, n).$$

Et les déterminants principaux sont les  $\Delta_k = \text{Dét}(A_k)$ .

Si tous les déterminants principaux de  $A$  sont non nuls alors, dans la méthode de Gauss, tous les pivots seront non nuls et il ne sera pas nécessaire de faire des permutations (on ne s'occupe pas de la taille des pivots). Réciproquement, si tous les pivots sont non nuls, alors tous les déterminants principaux sont non nuls.

*Exemple 1* — Si  $A$  est à diagonale strictement dominante alors toutes les sous-matrices principales de  $A$  sont aussi à diagonale strictement dominante et donc

inversibles. On en déduit donc que pour  $A$  à diagonale strictement dominante, la méthode de Gauss ne nécessite pas de permutations.

*Exemple 2* — Si  $A$  est symétrique définie positive, il en est de même de toutes les sous-matrices principales de  $A$ . Donc, pour  $A$  symétrique définie positive, la méthode de Gauss ne nécessite pas de permutations.

### 3.6.2 Interprétation algébrique de la méthode Gauss

On suppose que  $A$  est une matrice ayant tous ses déterminants principaux non nuls.

Pour  $k = 1, \dots, n$ , on note  $E_k$  la matrice déduite de la matrice identité en remplaçant la colonne  $N^\circ k$  par :

$$C_k = \begin{pmatrix} 0 \\ \cdot \\ 0 \\ 1 \\ c_{k+1,k} \\ \cdot \\ C_{n,k} \end{pmatrix}, \text{ avec } c_{ik} = -\frac{a_{ik}^{(k)}}{a_{kk}^{(k)}}, \text{ pour } k < i < n.$$

en gardant les notations de la méthode de Gauss.

On vérifie facilement que :

$$A^{(k)} = E_{k-1} \cdot A^{(k-1)} = E_{k-1} \cdot E_{k-2} \cdot \dots \cdot E_1 \cdot A, \text{ pour } k = 2, \dots, n-1,$$

c'est-à-dire que l'opération d'élimination de  $x_{k-1}$  dans les équations  $k, \dots, n$  revient à multiplier à gauche  $A^{(k-1)}$  par  $E_{k-1}$ .

Au bout des  $n - 1$  étapes, on a donc  $A^{(n)} = (E_{n-1} \cdot \dots \cdot E_1) \cdot A$  qui est triangulaire supérieure.

Ce qui donne la décomposition de  $A : A = L \cdot R$ , où  $L$  est triangulaire inférieure de diagonale unité et  $R$  est triangulaire supérieure ( $R = A^{(n)}$ ) et  $L$  est l'inverse du produit des  $E_k$ .

De plus une telle décomposition est unique : en effet si  $A = L \cdot R$  et  $A = L' \cdot R'$  avec  $L$  et  $L'$  triangulaires inférieures de diagonale unité et  $R$  et  $R'$  triangulaires supérieures, alors  $L^{-1} \cdot L = R' \cdot R^{-1}$  est à la fois triangulaire inférieure et triangulaire supérieure de diagonale unité, c'est nécessairement l'identité.

On a donc le :

*Théorème* : La matrice  $A$  possède une décomposition  $A = L \cdot R$ , avec  $L$  triangulaire inférieure à diagonale unité et  $R$  triangulaire supérieure si, et seulement si tous les déterminants principaux de  $A$  sont non nuls. Dans ce cas une telle décomposition est unique et les coefficients diagonaux de  $R$  sont donnés par :

$$\begin{cases} r_{11} = a_{11} \\ r_{kk} = \frac{\text{Dét}(A^{(k)})}{\text{Dét}(A^{(k-1)})}, \text{ pour } k = 2, \dots, n \end{cases}$$

*Démonstration* — Voir Lascaux et Théodor, Vol. 1 p. 228.

### 3.6.3 Détermination pratique de la décomposition L·R

On procède par coefficients indéterminés, c'est-à-dire qu'on écrit  $A = L \cdot R$ , on effectue le produit et on identifie, ce qui donne, en effectuant les calculs dans l'ordre indiqué :

$$\begin{cases} r_{1j} = a_{1j} & (j = 1, \dots, n) \\ L_{i1} = \frac{a_{i1}}{a_{11}} & (i = 2, \dots, n) \end{cases}$$

puis, pour  $k = 2, \dots, n$

$$r_{kj} = a_{kj} - \sum_{i=1}^{k-1} L_{ki} \cdot r_{ij} \quad (j = k, \dots, n)$$

et

$$L_{ik} = \frac{\left( a_{ik} - \sum_{j=1}^{k-1} L_{ij} \cdot r_{jk} \right)}{r_{kk}} \quad (i = k+1, \dots, n)$$

avec :

$$L_{ii} = 1$$

Ce qui donne la programmation structurée :

*PROCEDURE L·R(Entrée n : Entier ; A : Matrice ; Sortie : L, R : Matrice) ;*

*Début*

*Pour i allant de 1 à n faire*

*Début*

*$L_{ii} = 1$  ;*

*Pour j allant de i + 1 à n faire  $L_{ij} = 0$  ;*

*Pour j allant de 1 à i - 1 faire  $r_{ij} = 0$  ;*

*Fin ;*

*Pour j allant de 1 à n faire  $r_{1j} = a_{1j}$  ;*

*Pour i allant de 2 à n faire  $L_{i1} = \frac{a_{i1}}{a_{11}}$  ;*

*Pour k allant de 2 à n faire*

*Début*

*Pour j allant de k à n faire*

*Début*

*$S = 0$  ;*

*Pour i allant de 1 à k - 1 faire  $S = S + r_{ij} \cdot L_{ki}$  ;*

*$r_{kj} = a_{kj} - S$  ;*

*Fin ;*

*Si  $k \neq n$*

*Alors Début*

*Pour i allant de k + 1 à n faire*

*Début*

*$S = 0$  ;*

*Pour j allant de 1 à k - 1 faire  $S = S + L_{ij} \cdot r_{jk}$  ;*

*$L_{ik} = (a_{ik} - S) / r_{kk}$  ;*

*Fin ;*

Fin ;

Fin ;

Fin ;

### 3.6.4 Résolution d'un système linéaire par la méthode de Crout

Comme  $A = L \cdot R$ , l'équation  $A \cdot x = b$  équivaut aux deux systèmes triangulaires :

$$\begin{cases} L \cdot y = b & \text{(triangulaire inférieur)} \\ R \cdot x = y & \text{(triangulaire supérieur)} \end{cases}$$

Cette méthode est intéressante pour résoudre en parallèle plusieurs systèmes linéaires de même matrice  $A$ , car la décomposition  $L \cdot R$  ne fait pas intervenir le second membre.

### 3.6.5 Calcul de l'inverse d'une matrice par la méthode de Crout

Avec  $A = L \cdot R$ , on a  $A^{-1} = R^{-1} \cdot L^{-1}$ , et il suffit alors d'utiliser les procédures d'inversion des matrices triangulaires.

## 3.7 Cas des matrices tridiagonales

### 3.7.1 Notations et hypothèses

On suppose que  $A$  est tridiagonale, inversible et telle qu'il n'y ait pas de permutations dans la méthode de Gauss (par exemple si  $A$  est à diagonale strictement dominante ou symétrique définie positive).

On note :

$$A = \begin{pmatrix} a_1 & c_1 & 0 & . & . & 0 \\ b_2 & a_2 & c_2 & 0 & . & 0 \\ . & . & . & . & . & . \\ 0 & . & 0 & b_{n-1} & a_{n-1} & c_{n-1} \\ 0 & . & . & 0 & b_n & a_n \end{pmatrix}$$

### 3.7.2 Décomposition L·R d'une matrice tridiagonale

On cherche les matrices  $L$  et  $R$  sous la forme bidiagonale, soit :

$$L = \begin{pmatrix} 1 & 0 & 0 & . & . & 0 \\ L_2 & 1 & 0 & 0 & . & 0 \\ . & . & . & . & . & . \\ 0 & . & 0 & L_{n-1} & 1 & 0 \\ 0 & . & . & 0 & L_n & 1 \end{pmatrix} \text{ et } R = \begin{pmatrix} d_1 & r_1 & 0 & . & . & 0 \\ 0 & d_2 & r_2 & 0 & . & 0 \\ . & . & . & . & . & . \\ 0 & . & 0 & 0 & d_{n-1} & r_{n-1} \\ 0 & . & . & 0 & . & d_n \end{pmatrix}$$

On effectuant le produit  $L \cdot R$  et en identifiant avec les coefficients de  $A$ , on obtient alors :

$$\begin{cases} d_1 = a_1 \\ r_j = c_j \quad (j = 1, \dots, n - 1) \end{cases}$$

puis :

$$\begin{cases} L_{j+1} = \frac{b_{j+1}}{d_j} \\ d_{j+1} = a_{j+1} - r_j \cdot L_{j+1} \end{cases} \quad (j = 1, \dots, n - 1)$$

La programmation structurée est alors évidente.

### 3.7.3 Méthode de résolution d'un système tridiagonal

Soit à résoudre  $A \cdot x = e$ , pour  $A$  tridiagonale, vérifiant les hypothèses précédentes.

*Première méthode : utilisation de la décomposition L·R* — Tout d'abord, on résout  $L \cdot y = e$ , ce qui donne :

$$\begin{cases} y_1 = e_1 \\ y_j = e_j - L_j \cdot y_{j-1} \quad (j = 2, \dots, n) \end{cases}$$

Puis on résout  $R \cdot x = y$ , ce qui donne :

$$\begin{cases} x_n = \frac{y_n}{d_n} \\ x_j = \frac{(y_j - r_j \cdot x_{j+1})}{d_j} \quad (j = n - 1, \dots, 1) \end{cases}$$

Là encore la préparation structurée est évidente.

*Deuxième méthode : utilisation de la méthode de Gauss* — On obtient alors une méthode aussi appelée méthode du double balayage de Cholesky.

En supposant qu'il n'y a pas de permutations dans la méthode de Gauss (une telle permutation ferait perdre le caractère tridiagonal), on n'aura à chaque étape qu'une ligne à traiter et pour chaque ligne seulement deux opérations.

À l'étape  $k$  de la méthode, la ligne  $L_k$  du système devient :

$$L_k - m \cdot L_{k-1} \quad (k = 2, \dots, n)$$

avec :

$$m = \frac{b_k}{a_{k-1}}$$

Ce qui donne les formules de transformations :

$$\begin{aligned} b_k &= 0 ; \\ a_k &= a_k - m \cdot c_{k-1} \\ c_k &\text{ inchangé} \\ e_k &= e_k - m \cdot e_{k-1} \end{aligned}$$

Le système triangulaire supérieur obtenu sera bidiagonal et aura pour solution :

$$\begin{cases} x_n = \frac{e_n}{a_n} \\ x_i = \frac{(e_i - c_i \cdot x_{i+1})}{a_i} \quad (i = n - 1, \dots, 1) \end{cases}$$

*Remarque* — Le produit des  $a_i$  donnera, en fin d'opération, le déterminant de A.

Pour la programmation structurée, en vue d'économiser la place mémoire, on aura intérêt à stocker une matrice tridiagonale dans trois vecteurs. Une telle matrice sera donc notée :

$$T = \left( T_{ij} \right)_{\substack{1 \leq i \leq n \\ 1 \leq j \leq 3}}$$

où  $i$  est le numéro de la ligne et  $j$  celui de la diagonale, c'est-à-dire que  $(T_{i1})_{1 \leq i \leq n}$  représente la diagonale inférieure,  $(T_{i2})_{1 \leq i \leq n}$  la diagonale et  $(T_{i3})_{1 \leq i \leq n}$  la diagonale supérieure.

Ce qui donne la procédure :

*PROCEDURE SystTridiag*(Entrée  $n$  : Entier ;  
Entrée\_Sortie  $T$  : MatriceTridiagonale ;  $b$  : Vecteur) ;

*Début*

*Pour i* Allant de 2 à  $n$  *Faire*

*Début*

*Si*  $(T_{i-1,2} = 0)$

*Alors Arrêter*('Déterminant nul') ;

$m = T_{i,1}/T_{i-1,2}$  ;

$T_{i,2} = T_{i,2} - m * T_{i-1,3}$  ;

$b_i = b_i - m * b_{i-1}$  ;

*Fin* ;

*Si*  $(T_{n,2} = 0)$

*Alors Arrêter*('Déterminant nul') ;

$b_n = b_n / T_{n,2}$  ;

*Pour i* Allant de  $n-1$  à 1 *Faire*  $b_i = (b_i - T_{i,3} * b_{i+1}) / T_{i,2}$  ;

*Fin* ;

*Remarque* — En fin d'opération les valeurs initiales de la matrice  $T$  et du second membre  $b$  sont perdues, le vecteur  $b$  contenant la solution du système.

### 3.7.4 Déterminant d'une matrice tridiagonale

On reprend les notations du paragraphe précédent, en notant  $D_n$  le déterminant d'une matrice tridiagonale d'ordre  $n$ .

En développant suivant la dernière ligne, on a :

$$D_n = a_n \cdot D_{n-1} - b_n \cdot c_{n-1} \cdot D_{n-2}$$

Ce qui donne, avec les valeurs initiales  $D_0 = 1$  et  $D_1 = a_1$ , un algorithme très simple de calcul.

### 3.7.5 Exemples d'applications

- Interpolation spline cubique.
- Résolution de certaines équations différentielles linéaires par la méthode des différences finies (problème de Dirichlet linéaire en dimension 1).
- Résolution de l'équation de la chaleur à une dimension par la méthode des différences finies.

### 3.8 Résolution des systèmes symétriques

On suppose que  $A$  est symétrique et qu'il n'y a pas de permutations dans la méthode de Gauss.

#### 3.8.1 Décomposition $L \cdot D \cdot {}^tL$

Dans la décomposition  $L \cdot R$  de  $A$ , on a  $\text{Det}(R) = \text{Det}(A)$  qui est non nul, donc tous les termes diagonaux de  $R$  sont non nuls et on peut écrire  $R$  sous la forme :  $R = D \cdot R'$ , où  $D$  est diagonale et  $R'$  est triangulaire supérieure à diagonale unité (il suffit de diviser chaque ligne de  $R$  par son terme diagonal). On a donc  $A = L \cdot D \cdot R'$ , puis en écrivant que  ${}^tA = A$  et en utilisant le fait que la décomposition  $L \cdot R$  est unique, on déduit que  $R' = {}^tL$ .

On a donc pour toute matrice symétrique  $A$ , dont tous les déterminants principaux sont non nuls, la décomposition unique  $A = L \cdot D \cdot {}^tL$ , la matrice  $L$  étant triangulaire inférieure et à diagonale unité et la matrice  $D$  étant diagonale.

*Remarque : Calcul de la signature d'une matrice symétrique* — On sait qu'une matrice symétrique a toutes ses valeurs propres réelles et que, pour  $A$  inversible, le couple d'entiers  $(p, q)$  formé du nombre  $p$  de valeurs propres strictement positives et du nombre  $q = n - p$  de valeurs propres strictement négatives est uniquement déterminé par  $A$ .

Ce couple d'entiers est la signature de  $A$  et la décomposition ci-dessus nous permet de la calculer.

On dit que deux matrices  $A$  et  $B$  sont congruentes, s'il existe une matrice inversible  $L$  telle que  $A = L \cdot B \cdot {}^tL$ . Et le théorème de Sylvester nous dit que deux matrices congruentes ont même signature.

Avec les notations ci-dessus, on a donc :

$$\text{Signature}(A) = \text{Signature}(D) = (p, q)$$

où  $p$  est le nombre de termes strictement positifs de la diagonale de  $D$ .

De plus  $A$  est définie positive, si et seulement si, tous les coefficients de  $D$  sont strictement positifs.

#### 3.8.2 Calculs pratiques

Comme pour la décomposition  $L \cdot R$ , on trouve les coefficients de  $L$  et  $D$  par identification, ce qui donne,  $d_1 = a_{11}$  et pour  $i = 2, \dots, n$  :

$$\left\{ \begin{array}{l} L_{ij} = \frac{\left( a_{ij} - \sum_{k=1}^{j-1} L_{ik} \cdot d_k \cdot L_{jk} \right)}{d_j} \quad (j = 1, \dots, i - 1) \\ d_i = \left( a_{ii} - \sum_{k=1}^{i-1} L_{ik}^2 \cdot d_k \right) \end{array} \right.$$

#### 3.8.3 Programmation structurée

*PROCEDURE L-R\_Sym(Entrée  $n$  : Entier ;  $A$  : Matrice ; Sortie :  $L$  : Matrice ;  $D$  : Vecteur) ;*  
Début



```

Pour i allant de 1 à n faire
Début
  Pour j allant de 1 à i - 1 faire
  Début
    S = 0 ;
    Pour k allant de 1 à j - 1 faire S = S + Lik·dk·Ljk ;
    Lij = (aij - S)/dj ;
  Fin ;
  S = 0 ;
  Pour k allant de 1 à i - 1 faire S = S + Lik2 · dk ;
  di = aii - S ;
Fin ;

```

### 3.8.4 Cas particulier des matrices symétriques définies positives. Décomposition B<sup>t</sup>B de Cholesky

Si A est symétrique définie positive, dans la décomposition précédente, on a alors d<sub>i</sub> > 0 pour tout i = 1, ..., n. On peut donc écrire D = D<sup>2</sup>, et en posant : B = L·D', on a A = B<sup>t</sup>B avec B triangulaire inférieure, la diagonale de B étant formée des ±√d<sub>i</sub>. Si on impose la positivité des coefficients diagonaux de B une telle décomposition est unique. on a donc le :

*Théorème* : Une matrice A est symétrique définie positive si, et seulement si il existe une matrice B triangulaire inférieure inversible telle que A = B<sup>t</sup>B. De plus une telle décomposition est unique si on impose la positivité des coefficients diagonaux de B.

Le calcul effectif des coefficients de B se fait par identification ce qui donne :

$$\left\{ \begin{array}{l} b_{ij} = \frac{\left( a_{ij} - \sum_{k=1}^{j-1} b_{ik} \cdot b_{jk} \right)}{b_{jj}} \quad (j = 1, \dots, i - 1) \\ b_{ii}^2 = a_{ii} - \sum_{k=1}^{i-1} b_{ik}^2 \end{array} \right. \quad (i = 1, \dots, n)$$

La programmation structurée étant évidente.

Le nombre d'opérations, dans la décomposition de Cholesky est alors un O(n<sup>3</sup>).

Pour calculer l'inverse de A il suffit d'inverser la matrice triangulaire B.

## 3.9 Résolution des systèmes linéaires par des méthodes itératives

### 3.9.1 Remarques préliminaires

(i) Les méthodes directes pour résoudre les systèmes linéaires ont un caractère exact, mais pour les « grands systèmes », (i.e. pour n > 100) la propagation des erreurs d'arrondis en diminue l'efficacité.

(ii) Les méthodes itératives sont bien adaptées au cas des matrices creuses (i.e. avec beaucoup de zéros), car ces méthodes ne transforment pas la matrice de

départ, contrairement à la méthode de Gauss qui peut faire apparaître des termes non nuls à la place de zéros.

### 3.9.2 Principe des méthodes itératives

L'idée des méthodes itératives est la suivante : construire une suite  $(x^{(k)})$  d'éléments de  $\mathbb{R}^n$  qui va converger vers  $x$  solution de  $A \cdot x = b$ , où  $\text{Dét}(A)$  est non nul ; le calcul de chaque  $x^{(k)}$  devant être plus facile que la résolution directe de  $A \cdot x = b$ .

On écrira la matrice  $A$  sous la forme  $A = M - N$ , où  $M$  est « facilement inversible », et la résolution de  $A \cdot x = b$  est ramenée au « problème de point fixe » :

$$\text{trouver } x \text{ dans } \mathbb{R}^n \text{ solution de : } x = M^{-1} \cdot N \cdot x + M^{-1} \cdot b.$$

Pour résoudre ce problème, on utilise la « méthode des approximations successives », c'est-à-dire qu'on considère la suite  $(x^{(k)})$  définie par :

$$(1) \begin{cases} x^{(0)} \text{ donné dans } \mathbb{R}^n \\ x^{(k+1)} = M^{-1} \cdot N \cdot x^{(k)} + M^{-1} \cdot b \end{cases}$$

Si cette suite converge, c'est nécessairement vers la solution de  $A \cdot x = b$ .

*Problèmes posés par les méthodes itératives*

- A quelles conditions sur la matrice  $M$ , la suite définie par (1) converge ?
- Comment choisir  $M$  facilement inversible ?
- Quelle est la rapidité de la convergence ?

*Convergence de la suite définie par (1)*

On suppose que  $A$  est inversible et  $x$  désigne la solution de  $A \cdot x = b$ .

Si  $A = M - N$ , avec  $M$  inversible, on pose  $B = M^{-1} \cdot N$ . On a alors le :

*Théorème* : La suite définie par (1) est convergente, quelle que soit la valeur initiale  $x^{(0)}$ , si et seulement si,  $\rho(B) < 1$ , où  $\rho(B)$  désigne le rayon spectral de  $B$ .

*Démonstration* — En écrivant que  $A \cdot x = b$  équivaut à  $x = B \cdot x + M^{-1} \cdot b$ , on déduit que :

$$x^{(k+1)} - x = B \cdot (x^{(k)} - x) = B^{k+1} \cdot (x^{(0)} - x), \text{ pour tout } k \in \mathbb{N},$$

et donc  $(x^{(k)})$  va converger vers  $x$  pour toute valeur initiale si, et seulement si  $(B^{(k)})$  tend vers 0, ce qui équivaut à  $\rho(B) < 1$ .

### 3.9.3 Méthode de Jacobi

On prend  $M = D$ , matrice diagonale définie par  $d_{ii} = a_{ii}$  pour tout  $i$ . Si on veut que  $M$  soit inversible il faut donc supposer que tous les coefficients diagonaux de  $A$  sont non nuls.

L'algorithme de construction des  $x^{(k)}$ , pour un tel choix de  $M$  est alors le suivant :

$$a_{ii} \cdot x_i^{(k+1)} = -\left(a_{i1} \cdot x_1^{(k)} + \dots + a_{i,i-1} \cdot x_{i-1}^{(k)}\right) - \left(a_{i,i+1} \cdot x_{i+1}^{(k)} + \dots + a_{in} \cdot x_n^{(k)}\right) + b_i \quad (i = 1, \dots, n)$$

*Remarques* — (i) Le calcul des composantes de  $x^{(k+1)}$  nécessite de garder en mémoire le vecteur  $x^{(k)}$ . Une itération va donc immobiliser  $2 \cdot n$  cases mémoires.

(ii) On peut décider d'arrêter les itérations à un rang *MaxIter* donné et lorsque  $\|x^{(k+1)} - x^{(k)}\| < \varepsilon \cdot \|x^{(k+1)}\|$  ( $\varepsilon$  une précision Donnée).

(iii) Comme valeur initiale, on peut prendre  $x^{(0)} = \left(\left(\frac{b_i}{a_{ii}}\right)\right)$ .

Ce qui donne la programmation structurée :

*PROCEDURE Jacobi*(Entrée  $n$  : Entier ;  $A$  : Matrice ;  $b$  : Vecteur ; Sortie  $x$  : Vecteur) ;

*Début*

*Pour*  $i$  allant de 1 à  $n$  faire

*Début*

*Si*  $a_{ii} = 0$  *Alors* Stop *Sinon*  $x_i = b_i/a_{ii}$  ;

*Fin* ;

$p = 0$  ; *Continuer* = *Vrai* ;

*Tant que* ( $p < \text{MaxIter}$ ) *et* (*Continuer*) faire

*Début*

$p = p + 1$  ;

*Pour*  $i$  allant de 1 à  $n$  faire

*Début*

$S = b_i$  ;

*Pour*  $j$  allant de 1 à  $n$  faire

*Début*

*Si*  $j \neq i$  *Alors*  $S = S - a_{ij} \cdot x_j$  ;

*Fin* ;

$y_i = S/a_{ii}$  ;

*Fin* ;

*Continuer* =  $\|y - x\| > \varepsilon \cdot \|y\|$  ;

$x = y$  ;

*Fin* ;

*Si Non* *Continuer*

*Alors afficher* ('Solution approchée = ',  $x$ )

*Sinon arrêter* (convergence trop lente ou divergence) ;

*Fin* ;

### Cas des matrices à diagonale strictement dominante

*Théorème* : Si  $A$  est à diagonale strictement dominante, alors la méthode de Jacobi converge.

*Démonstration* — Si  $J = D^{-1} \cdot N = ((b_{ij}))$ , avec  $A = D - N$ , on a :

$$\|J\| = \text{Max} \left\{ \sum_{j=1}^n |b_{ij}|, i = 1, \dots, n \right\} = \text{Max} \left\{ \sum_{j \neq i} \left| \frac{a_{ij}}{a_{ii}} \right|, i = 1, \dots, n \right\} < 1$$

c'est-à-dire que  $\rho(J) < 1$  et la méthode de Jacobi est convergente.

### 3.9.4 Méthode de Gauss-Seidel

On prend  $M$  définie par :

$$\begin{cases} m_{ij} = 0, & \text{pour } 1 \leq i < j \leq n; \\ m_{ij} = a_{ij}, & \text{pour } 1 \leq j \leq i \leq n; \end{cases}$$

En notant  $A = M - N$ ,  $x^{(k+1)}$  est solution du système triangulaire inférieur  $M \cdot x^{(k+1)} = N \cdot x^{(k)} + b$  qui se résout en « descente », d'où l'algorithme de Gauss-Seidel :

$$a_{ii} \cdot x_i^{(k+1)} = -\left(a_{i1} \cdot x_1^{(k+1)} + \dots + a_{i,i-1} \cdot x_{i-1}^{(k+1)}\right) - \left(a_{i,i+1} \cdot x_{i+1}^{(k)} + \dots + a_{i,n} \cdot x_n^{(k)}\right) + b_i \quad (i = 1, \dots, n)$$

*Remarque 1* — Cet algorithme est en fait une amélioration de l'algorithme de Jacobi : dans le calcul de  $x_i^{(k+1)}$ , on utilise les composantes  $1, \dots, i-1$  de  $x^{(k+1)}$  (alors que dans la méthode de Jacobi ce sont celles de  $x^{(k)}$  qui sont utilisées) et les composantes  $i+1, \dots, n$  de  $x^{(k)}$  (comme dans la méthode de Jacobi). Cet algorithme sera donc en général plus performant que celui de Jacobi.

*Remarque 2* — Pour une itération on ne garde donc que  $n$  termes en mémoire.

Ce qui donne la programmation structurée :

*PROCEDURE Gauss-Seidel(Entrée  $n$  : Entier ;  $A$  : Matrice ;  $b$  : Vecteur ; Sortie  $x$  : Vecteur) ;*

*Début*

*Continuer = Vrai ;  $p = 0$  ;*

*Pour  $i$  allant de 1 à  $n$  faire*

*Début*

*Si  $a_{ii} = 0$  Alors Stop Sinon  $x_i = b_i/a_{ii}$  ;*

*Fin ;*

*Tant que ( $p < \text{MaxIter}$ ) et (Continuer) faire*

*Début*

*$p = p + 1$  ;*

*DeltaMax = 0 ;*

*Pour  $i$  allant de 1 à  $n$  faire*

*Début*

*$S = b_i$  ;*

*Pour  $j$  allant de 1 à  $n$  faire*

*Début*

*Si  $j \neq i$  Alors  $S = S - a_{ij}x_j$  ;*

*Fin ;*

*$S = S/a_{ii}$  ;*

*Si  $|x_i - S| > \text{DeltaMax}$*

*Alors DeltaMax =  $|x_i - S|$  ;*

*$x_i = S$  ;*

*Fin ;*

*Continuer = DeltaMax > Epsilon ;*

*Fin ;*

*Si Non Continuer*

*Alors Afficher('Solution approchée', $x$ )*

*Sinon Arrêter('convergence trop lente ou divergence') ;*

*Fin ;*

*Théorème* : La méthode de Gauss-Seidel converge si A vérifie l'une des conditions suivantes :

- (i) A à diagonale strictement dominante ;
- (ii) A symétrique définie positive.

*Démonstration* — (i) Soit A à diagonale strictement dominante,  $G = M^{-1} \cdot N$ ,  $\mu$  une valeur propre de G et  $x$ , non nul, un vecteur propre associé. On a alors :  $N \cdot x = \mu \cdot M \cdot x$ , c'est-à-dire :

$$\begin{cases} -\sum_{j=i+1}^n a_{ij} \cdot x_j = \mu \cdot \sum_{j=1}^i a_{ij} \cdot x_j & (i = 1, \dots, n - 1) \\ 0 = \mu \cdot \sum_{j=1}^n a_{nj} \cdot x_j \end{cases}$$

Si on suppose que  $|\mu| \geq 1$ , en prenant  $i$  tel que  $\|x\| = |x_i|$ , on déduit que :  $|\mu| \cdot |a_{ii}| \leq \sum_{j \neq i} |a_{ij}| \cdot |\mu|$ , ce qui contredit A à diagonale strictement dominante.

(ii) Soit A symétrique définie positive,  $\mu$  un vecteur propre de G et  $x$ , non nul, un vecteur propre associé. En écrivant  $A = D + E + F$ , où D est la diagonale de A, E le triangle inférieur strict et F le triangle supérieur strict, on en déduit alors que :

$$\begin{aligned} \langle x | A \cdot x \rangle &= (1 - \mu) \cdot \langle x | D \cdot x \rangle + (1 - \mu) \cdot \langle x | E \cdot x \rangle \\ \langle x | A \cdot x \rangle &= (1 - \mu) \cdot \langle x | D \cdot x \rangle + (1 - \mu) \cdot \langle x | F \cdot x \rangle \\ \langle x | A \cdot x \rangle &= \langle x | D \cdot x \rangle + \langle x | E \cdot x \rangle + \langle x | F \cdot x \rangle \end{aligned}$$

D'où :

$$(1 - |\mu|^2) \cdot \langle x | A \cdot x \rangle = |1 - \mu|^2 \cdot \langle x | D \cdot x \rangle$$

avec  $\langle x | A \cdot x \rangle > 0$  et  $\langle x | D \cdot x \rangle > 0$  car A et D sont définies positives, donc  $|\mu| < 1$ . C'est-à-dire que  $\rho(G) < 1$  et la méthode est convergente.

*Remarque* — Quand les méthodes de Jacobi et de Gauss-Seidel sont convergentes il vaut mieux choisir celle de Gauss-Seidel. Mais il se peut que la méthode de Gauss-Seidel diverge alors que celle de Jacobi converge comme le montre l'exemple suivant :

$$A = \begin{pmatrix} 1 & 2 & -2 \\ 1 & 1 & 1 \\ 2 & 2 & 1 \end{pmatrix}$$

### 3.9.5 Comparaison des méthodes de Jacobi et de Gauss-Seidel dans le cas des matrices tridiagonales d'ordre $n \geq 3$

On se donne :

$$A = \begin{pmatrix} a_1 & c_1 & 0 & \cdot & \cdot & 0 \\ b_2 & a_2 & c_2 & 0 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & 0 & b_{n-1} & a_{n-1} & c_{n-1} \\ 0 & \cdot & \cdot & 0 & b_n & a_n \end{pmatrix}$$

avec  $n \geq 3$  et tous les  $a_i$  non nuls.

On a alors le

*Lemme* — 1°) Si  $\mu$  est valeur propre de  $J$ , alors  $-\mu$  est aussi valeur propre de  $J$ .  
 2°) Le polynôme caractéristique de  $J$  s'écrit :  $P_J(x) = P(x^2) \cdot x^q$  où  $P$  est un polynôme de degré  $p$  et  $2 \cdot p + q = n$ .  
 Si  $P_G$  est le polynôme caractéristique de  $G$ , alors :  $P_G(x^2) = x \cdot P_J(x)$ .

On en déduit donc le

*Théorème* : Si  $A$  est tridiagonale d'ordre  $n > 3$  alors :  $\rho(G) = (\rho(J))^2$ .  
 Donc les deux méthodes convergent ou divergent simultanément, et dans le cas de la convergence, c'est la méthode de Gauss-Seidel qui est la plus rapide.

### 3.9.6 Méthode de relaxation

En vue d'accélérer la convergence de la méthode de Gauss-Seidel, on introduit un paramètre dans la matrice  $G = M^{-1} \cdot N$ . C'est-à-dire qu'on pose  $M_\omega = ((m_{ij}))$ , avec :

$$\begin{cases} m_{ij} = 0, & \text{pour } 1 \leq i < j \leq n \\ m_{ii} = \frac{a_{ii}}{\omega}, & \text{pour } i = 1, \dots, n, \text{ avec :} \\ m_{ij} = a_{ij}, & \text{pour } 1 \leq j < i \leq n \end{cases}$$

où  $\omega$  est non nul à préciser.

On pose  $L_\omega = M^{-1} \cdot N_\omega$ , si  $A = M_\omega - N_\omega$ .

Si  $D$  est la diagonale de  $A$ ,  $E$  le triangle inférieur strict et  $F$  le triangle supérieur strict, on a :

$$M_\omega = \omega^{-1} \cdot D + E \text{ et } A = D + E + F$$

On a alors :

$$L_\omega = (D + \omega \cdot E)^{-1} \cdot ((1 - \omega) \cdot D - \omega \cdot F)$$

Le vecteur  $x^{(k+1)}$  est alors solution du système triangulaire :

$$(D + \omega \cdot E) \cdot x^{(k+1)} = ((1 - \omega) \cdot D - \omega \cdot F) \cdot x^{(k)} + \omega \cdot b$$

et l'algorithme de calcul des  $x^{(k)}$  est :

$$x_i^{(k+1)} = (1 - \omega) \cdot x_i^{(k)} + \omega \cdot \hat{x}_i^{(k+1)} \quad (i = 1, \dots, n)$$

où  $\hat{x}_i^{(k+1)}$  est donné par les formules de Gauss-seidel en fonction des composantes  $1, \dots, i-1$  de  $x^{(k+1)}$  et des composantes  $i+1, \dots, n$  de  $x^{(k)}$ , soit :

$$a_{ii} \cdot \hat{x}_i^{(k+1)} = -\left(a_{i1} \cdot x_1^{(k+1)} + \dots + a_{i,i-1} \cdot x_{i-1}^{(k+1)}\right) - \left(a_{i,i+1} \cdot x_{i+1}^{(k)} + \dots + a_{in} \cdot x_n^{(k)}\right) + b_i \quad (i = 1, \dots, n)$$

Ce qui peut encore s'écrire :

$$x^{(k+1)} = x^{(k)} - \omega \cdot \xi^{(k+1)} \quad (k \geq 0)$$

où  $\xi^{(k+1)}$  est le « vecteur résidu » défini par :

$$\xi_i^{(k+1)} = \frac{\left(\sum_{j=1}^{i-1} a_{ij} \cdot x_j^{(k+1)} + \sum_{j=i}^n a_{ij} \cdot x_j^{(k)} - b_i\right)}{a_{ii}} \quad (i = 1, \dots, n)$$

Un test de convergence sera alors  $\|\xi\| < \varepsilon$ , où  $\varepsilon$  est une précision donnée.

La programmation structurée est alors la suivante :

*PROCEDURE Relaxation(Entrée n : Entier ; A : Matrice ; b : Vecteur ;  $\omega$  : Réel ;  
Sortie : x : Vecteur) ;*

*Début*

*p = 0 ;*

*Continuer = Vrai ;*

*Pour i Allant de 1 à n Faire*

*Début*

*Si  $a_{ii} = 0$  Alors Stop Sinon  $x_i = b_i/a_{ii}$  ;*

*Fin ;*

*Tant que (p < MaxIter) et (Continuer) Faire*

*Début*

*p = p + 1 ;*

*NormeResidu = 0 ;*

*Pour i Allant de 1 à n Faire*

*Début*

*Residu = -b<sub>i</sub> ;*

*Pour j Allant de 1 à n Faire*

*Residu = Residu + a<sub>ij</sub>·x<sub>j</sub> ;*

*Residu = Residu/a<sub>ii</sub> ;*

*x<sub>i</sub> = x<sub>i</sub> -  $\omega$ ·Residu ;*

*Si |Residu| > NormeResidu*

*Alors NormeResidu = |Residu| ;*

*Fin ;*

*Continuer = (NormeResidu >  $\varepsilon$ ) ;*

*Fin ;*

*Si (p = MaxIter)*

*Alors Arrêter('convergence trop lente ou divergence') ;*

*Fin ;*

Une condition nécessaire sur  $\omega$  pour que la méthode converge est donnée par le

*Théorème* : Si,  $\mu_1, \dots, \mu_n$  sont les valeurs propres de  $L_\omega$ , on a alors  $\prod_{i=1}^n \mu_i = (1 - \omega)^n$ . Donc  $\rho(L_\omega) \geq |\omega - 1|$ . On en déduit donc que la méthode de relaxation ne peut converger que si  $\omega \in ]0, 2[$ .

*Démonstration* — Voir Stoer et Burlisch p. 547.

*Cas des matrices symétriques définies positives* — Dans ce cas la condition précédente est aussi suffisante :

*Théorème (Ostrowski, Reich)* :

(a) Pour toute valeur propre  $\mu$  de  $L$ , on a :  $\mu \neq 1$  et :

$$\frac{(1 - |\mu|^2)}{|1 - \mu|^2} \cdot \langle x | A \cdot x \rangle = (2 - \omega) \cdot \frac{\langle x | A \cdot x \rangle}{\omega}$$

où  $x$  est un vecteur propre associé à  $\mu$ .

(b) Donc une condition nécessaire et suffisante de convergence de la méthode de relaxation, pour une matrice symétrique définie positive, est  $\omega \in ]0, 2[$ .

*Démonstration* — Voir Stoer et Burlisch p. 547.

*Remarque* — Si de plus  $A$  est tridiagonale, on a le

*Théorème (Young, Varga)* : Si  $A$  est symétrique définie positive et tridiagonale, alors :

- (i) les méthodes de Jacobi et de relaxation pour  $\omega \in ]0, 2[$  sont convergentes ;  
 (ii) la fonction  $\omega \mapsto \rho(L_\omega)$  a l'allure indiquée ci-dessous :

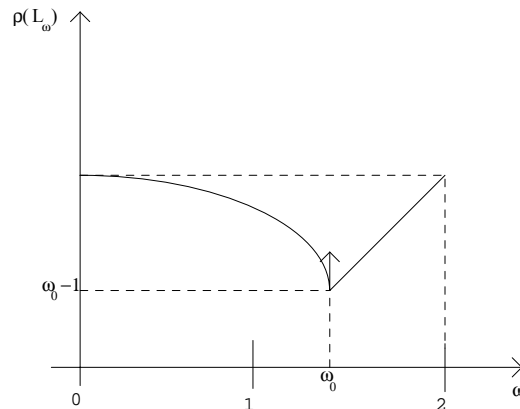


Figure 3.1

La valeur du paramètre de relaxation optimal étant :

$$\omega_0 = \frac{2}{1 + \sqrt{1 - \rho(J)^2}}$$

avec



$$\rho(L_{\omega_0}) = \omega_0 - 1.$$

*Démonstration* — Voir Stoer et Burlisch p.553.

*Remarque 1 : Procédé d'accélération de la convergence de Tchébycheff* — Si on sait calculer la valeur optimale de  $\omega$ , ce choix n'est pas judicieux dès la première itération. L'idée de Tchébycheff est de changer de paramètre  $\omega$  à chaque itération de la manière suivante : au départ, on prend  $\omega_1 = 1$ , puis à l'étape suivante on prend

$$\omega_2 = \frac{1}{1 - \frac{\rho(J)^2}{2}}$$

et aux étapes suivantes on prendra :

$$\omega_{k+1} = \frac{1}{1 - \frac{\rho(J)^2 \cdot \omega_k}{4}}$$

La suite  $(\omega_k)$  converge en fait vers  $\omega_0$  (la valeur optimale) et le procédé de Tchébycheff va diminuer, en général, le nombre d'itérations.

*Remarque 2 : Méthode de relaxation par blocs* — Lorsque la dimension du système est très grande, on aura intérêt à partitionner la matrice A en blocs, ce découpage étant adapté à la forme particulière de A (Pour la résolution de certaines équations aux dérivées partielles, on doit résoudre des systèmes tridiagonaux par blocs.).

On a donc :

$$A = \begin{pmatrix} A_{11} & A_{12} & \cdot & \cdot & \cdot & A_{1m} \\ A_{21} & A_{22} & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ A_{m1} & A_{m2} & \cdot & \cdot & \cdot & A_{mm} \end{pmatrix}$$

où  $A_{ij}$  est une matrice à  $n_i$  lignes et  $n_j$  colonnes, avec  $\sum_{i=1}^m n_i = n$ .

De plus, les matrices  $A_{ii}$  sont supposées inversibles.

Le second membre b et l'inconnue x sont alors partitionnés de manière analogue, soit :

$$b = \begin{pmatrix} B_1 \\ B_2 \\ \cdot \\ \cdot \\ B_m \end{pmatrix} \text{ et } x = \begin{pmatrix} X_1 \\ X_2 \\ \cdot \\ \cdot \\ X_m \end{pmatrix}$$

Le système à résoudre peut alors s'écrire par blocs :

$$\sum_{j=1}^m A_{ij} \cdot X_j = B_i \quad (i = 1, \dots, m)$$

Et les formules pour la méthode de relaxation s'écrivent :

$$X^{(k+1)} = X^{(k)} - \omega \cdot \Xi^{(k+1)} \quad (k \geq 0)$$

où  $\Xi^{(k+1)}$  est le « vecteur résidu » défini par :

$$\Xi^{(k+1)} = A_{ii}^{-1} \cdot \left( \sum_{j=1}^{i-1} A_{ij} \cdot X_j^{(k+1)} + \sum_{j=i+1}^n A_{ij} \cdot X_j^{(k)} - B_i \right) \quad (i = 1, \dots, m)$$

Cette méthode sera intéressante si les  $A_{ii}$  sont facilement inversibles.

### 3.9.7 Exemple d'application

La méthode de relaxation sera utilisée pour la résolution d'une équation aux dérivées partielles de type elliptique par la méthode des différences finies : Cf. le chapitre sur la résolution des équations aux dérivées partielles paragraphe (5.5). On aura à résoudre des systèmes tridiagonaux par blocs, les matrices blocs de la diagonale étant tridiagonales et les autres diagonales.

Dans le cas particulier de l'équation de Poisson,  $\Delta \cdot u = f$ , on est ramené à résoudre un système tridiagonal par blocs :

$$T \cdot u = v \quad (1)$$

où la matrice  $T$ , tridiagonale par blocs est définie par :

$$T = \begin{pmatrix} A & I_d & 0 & \cdot & \cdot & 0 \\ I_d & A & I_d & 0 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot & 0 & I_d & A \end{pmatrix}$$

où  $I_d$  désigne la matrice identité de  $\mathbb{R}^n$  et  $A$ , d'ordre  $n$  est définie par :

$$A = \begin{pmatrix} -4 & 1 & 0 & \cdot & \cdot & 0 \\ 1 & -4 & 1 & 0 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot & 0 & 1 & -4 \end{pmatrix}$$

les vecteurs étant notés :

$$u = \begin{pmatrix} u_1 \\ u_2 \\ \cdot \\ \cdot \\ u_m \end{pmatrix}, \text{ avec } u_j = \begin{pmatrix} u_{1j} \\ u_{2j} \\ \cdot \\ \cdot \\ u_{nj} \end{pmatrix}, \text{ pour } j = 1, \dots, m$$

La composante vectorielle  $N^{\circ j}$  de  $M \cdot u$  est alors :

$$u_{j-1} + A \cdot u_j + u_{j+1} \quad (j = 1, \dots, m)$$

avec la convention  $u_0 = u_{m+1} = 0$ .

C'est-à-dire que la composante  $N^{\circ (i,j)}$  de  $M \cdot u$  est :

$$u_{i,j-1} + u_{i-1,j} - 4 \cdot u_{i,j} + u_{i+1,j} + u_{i,j+1} \quad (i = 1, \dots, n)$$

avec les conventions  $u_{0,j} = u_{n+1,j} = 0$ , pour tout  $j = 1, \dots, m$ .

La suite  $(u^{(k)})_{k \in \mathbb{N}}$  construite par la méthode de relaxation, avec paramètre  $\omega$  dans  $]0,2[$  est alors définie par :

$$\begin{cases} u^{(0)} \text{ donné dans } \mathbb{R}^{n \cdot m} \\ u^{(k+1)} = u^{(k)} - \omega \cdot \xi^{(k+1)} \quad (k \geq 0) \end{cases}$$

où le vecteur résidu  $\xi^{(k+1)}$  est défini par :

$$\xi_{ij}^{(k+1)} = \frac{1}{-4} \cdot (u_{i,j-1}^{(k+1)} + u_{i-1,j}^{(k+1)} - 4 \cdot u_{i,j}^{(k)} + u_{i+1,j}^{(k)} + u_{i,j+1}^{(k)} - v_{i,j}) \quad (1 \leq i \leq n, 1 \leq j \leq m)$$

pour  $1 < i < n$  et  $1 < j < m$ .

Dans ce cas on connaît explicitement le rayon spectral de la matrice de Jacobi :

$$\rho(J) = \frac{\cos\left(\frac{\pi}{n+1}\right) + \cos\left(\frac{\pi}{m+1}\right)}{2}$$

ce qui permet d'appliquer le procédé d'accélération de Tchébycheff.

Pour la programmation structurée, on se reportera au paragraphe (5.5.4) du chapitre sur la résolution des équations aux dérivées partielles par la méthode des différences finies.

## 4. Calcul des valeurs propres et des vecteurs propres de certaines matrices réelles

### 4.1 Introduction

Le problème de la détermination des valeurs propres d'une matrice est, en général, beaucoup plus difficile que celui de la résolution des systèmes linéaires.

Tout d'abord nous allons voir la méthode de la puissance itérée qui permet de calculer la valeur propre de plus grand module sous certaines hypothèses. Puis en itérant ce procédé on peut en déduire toutes les autres quand elles sont toutes distinctes en modules. Cette méthode, peu performante, est à utiliser pour les matrices ayant toutes les valeurs propres distinctes et si on en cherche seulement quelques-unes.

D'autres méthodes sont les méthodes de Rutishauser, Givens et Householder et de Jacobi qui reposent sur le principe suivant : on construit une suite  $(A_k)$  de matrices semblables à  $A$ , donc ayant les mêmes valeurs propres, et qui « converge » vers une matrice plus simple (diagonale, pour Jacobi, ou tridiagonale, pour Givens et Householder, ou triangulaire ; pour Rutishauser) encore semblable à  $A$ . Les valeurs propres cherchées sont alors les termes diagonaux de la matrice limite, dans le cas des méthodes de Jacobi et de Rutishauser.

Les problèmes posés par ces méthodes sont :

- (i) savoir décrire le passage de  $A_k$  à  $A_{k+1}$  aussi simplement que possible c'est-à-dire avec un coût de calcul peu élevé ;
- (ii) les valeurs propres doivent être bien conservées (stabilité numérique de la méthode).

Enfin une autre idée est d'essayer de déterminer le polynôme caractéristique de  $A$  puis d'utiliser un procédé numérique de recherche des racines d'une équation polynomiale. On verra, par exemple, les méthodes de Leverrier, de Krylov et de Souriau. Mais de telles méthodes sont déconseillées pour des grandes matrices à cause des problèmes d'instabilité numérique (évaluation du polynôme, extraction des racines d'un polynôme mal déterminé qui peut être de degré élevé ...).

### 4.2 Calcul du rayon spectral d'une matrice. Méthode de la puissance itérée

#### 4.2.1 Hypothèses et notations

On se donne une matrice  $A$  à coefficients réels telle que la valeur propre dominante soit unique, c'est-à-dire que si  $\mu_1, \dots, \mu_n$  sont les valeurs propres de  $A$  dans  $\mathbb{C}$ , alors on suppose que :

$$|\mu_1| > |\mu_2| \geq |\mu_3| \geq \dots \geq |\mu_n|.$$

Dans ce cas la valeur propre  $\mu_1$  est nécessairement réelle et simple. De plus, le rayon spectral de  $A$  est donné par :  $\rho(A) = |\mu_1|$ .

*Remarque* — Il existe un théorème qui donne une condition suffisante sur les coefficients de la matrice  $A$ , pour que la valeur propre dominante soit unique :

*Théorème* : Si  $a_{ij} > 0$  pour tous  $i, j$  alors  $\rho(A) = \mu$ , avec  $\mu > 0$  unique valeur propre dominante de la matrice  $A$ .

#### 4.2.2 Méthode de la puissance itérée pour calculer le rayon spectral

On note  $e_1$  un vecteur propre, de norme 1, associé à  $\mu_1$  et on définit la suite de vecteurs  $(v_p)$  de la façon suivante :

$$\begin{cases} v_0 \text{ donné tel que } \|v_0\| = 1 \\ v_p = \frac{1}{\|A \cdot v_{p-1}\|} \cdot A \cdot v_{p-1}, \text{ pour } p > 0 \end{cases}$$

*Lemme* : Pour tout  $p \geq 0$ , on a :

$$v_p = \frac{1}{\|A^p \cdot v_0\|} \cdot A^p \cdot v_0$$

*Démonstration* — Pour  $p = 1$  le résultat est vrai et en le supposant vrai pour  $p$ , on a :

$$v_{p+1} = \frac{1}{\|A \cdot v_p\|} \cdot A \cdot v_p,$$

avec :

$$A \cdot v_p = A \left( \frac{1}{\|A^p \cdot v_0\|} \cdot A^p \cdot v_0 \right) = \frac{1}{\|A^p \cdot v_0\|} \cdot A^{p+1} \cdot v_0$$

d'où le résultat.

*Théorème* : Si  $v_0$  n'est pas orthogonal à  $e_1$  (i.e. si la composante de  $v_0$  suivant  $e_1$  dans une base de Jordan associée à  $A$  est non nulle) alors la suite  $(v_p)$  converge vers un vecteur propre associé à la valeur propre dominante  $\mu_1$  et la suite  $(\|A \cdot v_p\|)$  converge vers le rayon spectral  $\rho(A)$ .

*Démonstration* — Avec le théorème de Jordan, on peut écrire que  $A = P^{-1} \cdot A' \cdot P$ , où :

$$A' = \begin{pmatrix} \mu_1 & 0 & \dots & 0 \\ 0 & \left( \begin{array}{c} \\ \\ B \\ \end{array} \right) \\ \vdots & \\ 0 & \end{pmatrix}$$

la matrice  $B$  étant triangulaire supérieure de valeurs propres  $\mu_2, \dots, \mu_n$ . En écrivant  $B = \mu_1 \cdot B'$ , où  $B'$  a pour valeurs propres les  $\mu_j/\mu_1$  pour  $j = 2, \dots, n$ , on a donc :

$$A^p \cdot v_0 = \mu_1^p \cdot \begin{pmatrix} 1 & & 0 \dots 0 \\ 0 & & \\ \vdots & & \\ 0 & & (B')^p \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ v'_0 \end{pmatrix} = \mu_1^p \cdot \begin{pmatrix} x_1 \\ u_p \end{pmatrix}$$

si  $v_0 = \begin{pmatrix} x_1 \\ v_0 \end{pmatrix}$ , avec  $u_p \rightarrow 0$ . On a alors :

$$v_p = \text{Signe}(x_1) \cdot (\text{Signe}(\mu_1))^p \cdot \begin{pmatrix} 1 \\ u_p \end{pmatrix}.$$

Donc :

$$v_p \xrightarrow{p \rightarrow +\infty} (-\text{Signe}(x_1))^p \cdot e_1$$

et

$$\|A \cdot v_p\| \xrightarrow{p \rightarrow +\infty} |\mu_1|.$$

*Remarque 1* — En général on n'a pas d'informations sur le sous-espace propre associé à  $\mu_1$  et en prenant  $v_0$  quelconque dans  $\mathbb{R}^n$  on pourrait avoir  $x_1 = 0$  et, a priori la méthode décrite ci-dessus ne donnera pas le résultat espéré. En fait  $V = \{ v \in \mathbb{R}^n ; x_1 = 0 \}$  est un hyperplan de  $\mathbb{R}^n$ , donc de mesure nulle, et en prenant  $v_0$  quelconque dans  $\mathbb{R}^n$  on a très peu de chances de tomber sur  $V$ . Et même si  $v_0$  est dans  $V$ , à cause des erreurs d'arrondis, au bout de quelques itérations on sortira de l'hyperplan et on aura les bonnes hypothèses mais la convergence peut être lente. Une autre façon d'éviter les problèmes est de prendre deux valeurs différentes de  $v_0$  et de garder celle qui donne la convergence la plus rapide.

*Remarque 2* — La méthode de la puissance itérée est une adaptation de la méthode de Bernoulli pour calculer la racine de plus grand module d'un polynôme, quand cette dernière est unique.

*Remarque 3* — En appliquant la méthode de la puissance itérée à la matrice  $A^{-1}$ , on a un moyen de calculer la valeur propre de plus petit module de  $A$  (quand cette dernière est unique) car  $\rho(A^{-1}) = \frac{1}{|\mu_n|}$  où  $\mu_n$  est la valeur propre de plus petit

module de  $A$ . La méthode obtenue est aussi appelée méthode de la puissance inverse.

*Remarque 4* — Quand la méthode converge, pour  $p$  grand,  $v_p$  peut être considéré comme un vecteur propre de  $A$  associé à  $\mu_1$ , de sorte que  $A \cdot v_p \approx \mu_1 \cdot v_p$ , ce qui permet de calculer  $\mu_1$  avec :

$$\mu_1 \approx \frac{(A \cdot v_p)(i)}{v_p(i)}, \text{ où } i \text{ est tel que } |v_p(i)| = \|v_p\|.$$

### 4.3 Calcul des autres valeurs propres par la méthode de déflation

#### 4.3.1 Hypothèses et notations

On suppose maintenant que les valeurs propres de  $A$  vérifient :

$$|\mu_1| > |\mu_2| > |\mu_3| > \dots > |\mu_n|.$$

Soit  $X$  un vecteur propre associé à  $\mu_1$  de norme euclidienne égale à 1.

On va alors vérifier que les valeurs propres de la matrice

$$B = A - \mu_1 \cdot X \cdot {}^t X = ((a_{ij} - \mu_1 \cdot x_i \cdot x_j))$$

sont  $0, \mu_2, \dots, \mu_n$ . On appliquera alors la méthode de la puissance itérée à cette nouvelle matrice pour avoir  $\mu_2$ , puis on continue ainsi de suite.

#### 4.3.2 Un lemme

*Lemme* : Soit  $\alpha$  une valeur propre de  $A$  différente de  $\mu_1$ , alors

- (i)  $\alpha$  est aussi valeur propre de  ${}^t A$  ;
- (ii) si  $Y$  est un vecteur propre de  ${}^t A$  associé à  $\alpha$  alors  $Y$  est orthogonal à  $X$ , c'est-à-dire  $\langle X | Y \rangle = 0$  ;
- (iii)  $Y$  est aussi valeur propre de  $B = A - \mu_1 \cdot X \cdot {}^t X$  ;
- (iv)  $0$  est valeur propre de  $B$ .

*Démonstration* — Il suffit de vérifier.

*Remarque* — Pour calculer la deuxième valeur propre de  $B$  on peut se passer de calculer la matrice  $B$ .

Le procédé de calcul est le suivant :

$$u_p = B \cdot v_{p-1} ;$$

puis

$$v_p = \frac{1}{\|u_p\|} \cdot u_p.$$

Mais on a :

$$u_p = A \cdot v_{p-1} - \mu_1 \cdot X \cdot {}^t X \cdot v_{p-1} = A \cdot (v_{p-1} - \langle X | v_{p-1} \rangle \cdot X) ;$$

car  $\mu_1 \cdot X = A \cdot X$ .

L'algorithme de calcul est donc :

$$\text{calcul de } \beta_p = \langle X | v_{p-1} \rangle ;$$

$$\text{calcul de } u_p = A \cdot (v_{p-1} - \beta_p \cdot X) ;$$

$$\text{calcul de } v_p = \frac{1}{\|u_p\|} \cdot u_p.$$

Mais pour calculer toutes les valeurs propres il est préférable d'utiliser une procédure de calcul de la matrice  $B$ .

### 4.3.3 Programmation structurée

Les valeurs propres peuvent être stockées dans un vecteur Valeurs\_Propres et les vecteurs propres dans une matrice Vecteurs\_Propres.

La procédure ci-dessous calcule les p premières valeurs propres et les vecteurs propres associés pour  $p = 1, \dots, n$ . Les valeurs propres sont calculées dans l'ordre des valeurs absolues décroissantes, ce qui suppose qu'elles sont toutes distinctes.

*PROCEDURE Déflation(Entrée  $n, p$  : Entier ;  $A$  : Matrice ;  
Sortie Valeurs\_Propres : Vecteur ; Vecteurs\_Propres : Matrice) ;*

*Début*

*Pour k allant de 1 à p faire*

*Début*

*Se donner  $V_k$  aléatoire ;*

*Iter = 0 ;*

*Répéter*

*Iter = Iter + 1 ;*

*$v = A \cdot V_k$  ;*

*$i_0 = 1$  ;*

*Pour i allant de 2 à n faire*

*Début*

*Si  $|V_{k_i}| > |V_{k_{i_0}}|$*

*Alors  $i_0 = i$  ;*

*Fin ;*

*Valeurs\_Propres(k) =  $v_{i_0} / V_{k_{i_0}}$  ;*

*Aux = NormeEuclidienne(v) ;*

*$v = (1/Aux) \cdot v$  ;*

*$u_{Ancien} = V_k$  ;*

*$V_k = v$  ;*

*Jusqu'à ( $\|V_k - u_{Ancien}\| < \epsilon$ ) ou (Iter = MaxIter) ;*

*Pour i allant de 1 à n faire Vecteurs\_Propres<sub>ik</sub> =  $v_i$  ;*

*Si Iter = MaxIter*

*Alors Arrêter('Convergence trop lente ou divergence') ;*

*Pour i allant de 1 à n faire*

*Début*

*Pour j allant de 1 à n faire*

*Début*

*$a_{ij} = a_{ij} - Valeurs\_Propres(k) \cdot v_i \cdot v_j$  ;*

*Fin ;*

*Fin ;*

*Fin ;*

*Fin ;*

## 4.4 Méthode de Rutishauser

### 4.4.1 Un lemme

*Lemme* : Si  $A$  admet une décomposition L·R, alors la matrice  $A_1 = R \cdot L$  est semblable à  $A$ , donc ces deux matrices ont les mêmes valeurs propres.



*Démonstration* — On a :  $R \cdot L = R \cdot (L \cdot R) \cdot R^{-1}$ .

#### 4.4.2 Hypothèses

On suppose que la matrice  $A$  admet une décomposition  $L \cdot R$  et que toutes ses valeurs propres sont réelles. Un exemple de telle matrice est donné par  $A$  symétrique définie positive.

#### 4.4.3 Principe de la méthode de Rutishauser

*Étape 1* — A partir de la décomposition  $A = L_1 \cdot R_1$ , on pose  $A_1 = R_1 \cdot L_1$  ;

*Étape 2* — on écrit la décomposition  $L \cdot R$  de  $A_1$ , soit  $A_1 = L_2 \cdot R_2$ , puis on pose  $A_2 = R_2 \cdot L_2$  ;

*Étape p* — ayant obtenu  $A_{p-1} = R_{p-1} \cdot L_{p-1}$ , on écrit sa décomposition  $L \cdot R$ , soit  $A_{p-1} = L_p \cdot R_p$ , puis on pose  $A_p = R_p \cdot L_p$ .

Si à chaque étape la décomposition  $L \cdot R$  est possible, alors  $A_p$  est semblable à  $A$ .

Sous certaines conditions, la suite de matrices ainsi obtenue va converger vers une matrice  $A'$  qui est triangulaire supérieure et semblable à  $A$ . Donc les valeurs propres de  $A$  seront les termes diagonaux de  $A'$ .

#### 4.4.4 Une condition suffisante de convergence de la méthode de Rutishauser

*Théorème* : Si la matrice  $A$  est diagonalisable avec ses valeurs propres réelles et si les matrices  $A_p$  admettent toutes des décompositions  $L \cdot R$ , alors la suite  $(A_p)$  va converger vers une matrice triangulaire supérieure  $A'$  dont les termes diagonaux sont les valeurs propres de  $A$ .

*Démonstration* — Voir Stoer et Burlisch Chap. 6.

*Remarque 1* — Si on ajoute  $\alpha \cdot I_d$  à la matrice  $A$ , du point de vue spectral, cela revient à ajouter  $\alpha$  aux valeurs propres de  $A$ .

Donc, dans le cas où des pivots trop petits apparaissent à l'étape  $p$ , avant de factoriser la matrice  $A_p$  on lui ajoute  $\alpha_p \cdot I_d$  avec  $\alpha_p$  choisi tel que  $A_p + \alpha \cdot I_d$  soit à diagonale strictement dominante. Ensuite on effectue la décomposition  $L \cdot R$ . A la fin il ne faudra pas oublier de retrancher des valeurs propres obtenues, la somme de tous les  $\alpha_p$  qu'il a fallu ajouter.

*Remarque 2* — Comme test d'arrêt on pourra choisir de s'arrêter pour  $p < \text{MaxIter}$  et quand :

$$\sum_{i=1}^n |a_{ii}^{(p+1)} - a_{ii}^{(p)}| < \varepsilon \text{ et } \sum_{1 \leq j < i \leq n} |a_{ij}^{(p+1)}| < \varepsilon.$$

La programmation structurée s'écrit alors :

*PROCEDURE* Rutishauser(Entrée  $n$  : Entier ;  $A$  : Matrice ; Sortie Valeurs\_Propres : Vecteur) ;

*Début*

Iter = 0 ;

Répéter

$L \cdot R(n, A, L, R)$  ;

Pour  $i$  allant de 1 à  $n$  Faire  $Valeurs\_Propres(i) = a_{ii}$  ;  
 $A = R \cdot L$  ;  $c = 0$  ;  $d = 0$  ;  
 Pour  $i$  allant de 1 à  $n$  Faire  
 Début  
 $c = c + |a_{ii} - Valeurs\_Propres(i)|$  ;  
 Pour  $j$  allant de 1 à  $i - 1$  Faire  $d = d + |a_{ij}|$  ;  
 Fin ;  
 $Iter = Iter + 1$  ;  
 Jusqu'à  $(c < \varepsilon)$  et  $(d < \varepsilon)$  ou  $(Iter = MaxIter)$  ;  
 Fin ;

#### 4.5 Méthode de Jacobi pour calculer les valeurs et vecteurs propres d'une matrice symétrique

On se donne, dans ce paragraphe, une matrice  $A$  symétrique réelle. Ses valeurs propres sont donc toutes réelles.

Le principe de la méthode de Jacobi consiste à construire une suite de matrices de rotations planes ( $R_k$ ) telle que pour  $k$  tendant vers l'infini, la suite de matrices ( $A_k$ ) définie par  $A_0 = A$  et  $A_k = {}^tR_k \cdot A_{k-1} \cdot R_k$  tend vers une matrice diagonale. Comme chacune des matrices  $A_k$  est semblable à  $A$  (puisque  ${}^tR_k = R_k^{-1}$ ), on en déduit que les valeurs propres de  $A$  sont les termes diagonaux de la matrice limite.

Cette méthode est bien adaptée aux matrices symétriques de taille au plus 10. Pour les matrices de taille supérieure, on lui préférera la méthode Q-R (Cf. Exercice 9).

Pour tous  $\theta$  dans  $[0, 2\pi[$  et  $1 \leq p < q \leq n$ , on note  $R_{p,q}(\theta)$  la matrice de rotation d'angle  $\theta$  dans le plan défini par les vecteurs  $e_p$  et  $e_q$  de la base canonique de  $\mathbb{R}^n$ . Elle est donc définie par :

$$R_{p,q}(\theta) = \begin{pmatrix} 1 & 0 & \cdot & \cdot & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & c & \cdot & -s & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & s & \cdot & c & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot & \cdot & \cdot & 0 & 1 \end{pmatrix} \begin{array}{l} \leftarrow p \\ \leftarrow q \end{array}$$

$$\begin{array}{cc} \uparrow & \uparrow \\ p & q \end{array}$$

où  $c = \cos(\theta)$  et  $s = \sin(\theta)$ .

Pour toute matrice symétrique  $A$ , le produit  $A' = R_{p,q}^{-1}(\theta) \cdot A \cdot R_{p,q}(\theta)$  est alors facile à calculer. La multiplication à gauche par  $R^{-1} = {}^tR$  modifie seulement les

lignes  $p$  et  $q$  de  $A$  et la multiplication à droite par  $R$  change seulement les colonnes  $p$  et  $q$  de  $A$ . Ce qui donne pour les coefficients de  $A'$  :

$$\begin{aligned} a'_{ij} &= a_{ij} \quad (i \neq p, i \neq q, j \neq p, j \neq q) \\ a'_{ip} &= c \cdot a_{ip} + s \cdot a_{iq} \quad (i \neq p, i \neq q) \\ a'_{pp} &= c^2 \cdot a_{pp} + s^2 \cdot a_{qq} + 2 \cdot s \cdot c \cdot a_{pq} \\ a'_{iq} &= c \cdot a_{iq} - s \cdot a_{ip} \quad (i \neq p, i \neq q) \\ a'_{qq} &= s^2 \cdot a_{pp} + c^2 \cdot a_{qq} - 2 \cdot s \cdot c \cdot a_{pq} \\ a'_{pq} &= (c^2 - s^2) \cdot a_{pq} - s \cdot c \cdot (a_{pp} - a_{qq}) \end{aligned}$$

la matrice  $A'$  étant symétrique.

L'idée de Jacobi est alors de déterminer  $\theta$  de manière à annuler les coefficients  $a'_{pq}$  et  $a'_{qp}$  de  $A'$ , ce qui donne pour l'angle  $\theta$  :

$$2 \cdot \cos(2 \cdot \theta) \cdot a_{pq} - \sin(2 \cdot \theta) \cdot (a_{pp} - a_{qq}) = 0$$

pour  $a_{pq} \neq 0$ , on a nécessairement  $\sin(2 \cdot \theta) \neq 0$  et :

$$\text{Cotg}(2 \cdot \theta) = \frac{a_{pp} - a_{qq}}{2 \cdot a_{pq}} = b$$

(si  $a_{pq} = 0$ , le problème ne se pose pas).

Le calcul de  $c = \text{Cos}(\theta)$  et  $s = \text{Sin}(\theta)$  peut se faire de façon purement algébrique, ce qui est plus rapide que de faire appel aux fonctions trigonométriques.

En posant  $t = \text{tg}(\theta)$ , on a :

$$b = \frac{1}{2} \cdot \left( \frac{1}{\text{tg}(\theta)} - \text{tg}(\theta) \right)$$

c'est-à-dire que  $t$  est solution de l'équation du second degré :

$$t^2 + 2 \cdot b \cdot t - 1 = 0$$

Cette équation admet deux racines réelles :

$$t = -b + \sqrt{1 + b^2} = \frac{1}{b + \sqrt{1 + b^2}} \quad \text{et} \quad t' = -\frac{1}{t} = \frac{-1}{\sqrt{1 + b^2} - b}$$

Comme  $|t \cdot t'| = 1$ , l'une de ces racines est de valeur absolue inférieure ou égale à 1. Pour minimiser les problèmes d'erreurs d'arrondis, on choisit la racine de valeur absolue inférieure à 1, ce qui donne un angle  $\theta$  entre  $-\frac{\pi}{4}$  et  $\frac{\pi}{4}$ . Cette racine est donnée par :

$$t = \frac{\text{signe}(b)}{\sqrt{b^2 + 1} + |b|}$$

Avec  $\cos(\theta) = \frac{1}{\sqrt{1 + \text{tg}^2(\theta)}}$  et  $\sin(\theta) = \cos(\theta) \cdot \text{tg}(\theta)$ , on déduit que :

$$c = \frac{1}{\sqrt{1 + t^2}} \quad \text{et} \quad s = t \cdot c$$

Les formules donnant les coefficients de  $A'$  se simplifient alors pour donner finalement :

$$a'_{pq} = 0$$

$$\begin{aligned} a'_{pp} &= a_{pp} + t \cdot a_{pq} \\ a'_{qq} &= a_{qq} - t \cdot a_{pq} \\ a'_{ip} &= a_{ip} + s \cdot (a_{iq} - \tau \cdot a_{ip}) \quad (i \neq p, i \neq q) \\ a'_{iq} &= a_{iq} - s \cdot (a_{ip} + \tau \cdot a_{iq}) \quad (i \neq p, i \neq q) \end{aligned}$$

où  $\tau$  est donné par :

$$\tau = \operatorname{tg}\left(\frac{\theta}{2}\right) = \frac{s}{1+c}$$

La méthode de Jacobi classique consiste à partir de la matrice  $A$  à construire la matrice  $A^{(1)}$  par le procédé décrit ci-dessus où les indices  $p$  et  $q$  sont définis par :

$$|a_{pq}| = \operatorname{Max}\{|a_{ij}|; 1 \leq i < j \leq n\}$$

puis on recommence avec  $A^{(1)}$  pour construire  $A^{(2)}$  et en continuant ainsi de suite on construit une suite de matrices  $A^{(k)}$  toutes semblables à  $A$ .

On peut alors montrer le :

*Théorème* : Pour toute matrice symétrique réelle, la suite  $(A^{(k)})$  converge vers une matrice  $D$  diagonale et semblable à  $A$ .  
 Les valeurs propres de  $A$  sont donc les éléments diagonaux de  $D$ .  
 De plus, si on écrit  $D = {}^tV \cdot A \cdot V$ , alors pour tout  $j = 1, \dots, n$ , la colonne  $j$  de  $V$  est un vecteur propre associé à la valeur propre  $D_{jj}$ .

*Démonstration* — Voir Lascaux et Théodor p. 697.

*Remarque 1* — D'un point de vue numérique, le calcul du maximum des termes non diagonaux de  $A^{(k)}$  à chaque étape n'est pas intéressant car il augmente le temps de calcul. On préfère procéder de la manière suivante : à l'étape  $k$  du calcul, la matrice  $A^{(k-1)}$  étant construite on construit la matrice  $A^{(k)}$  en effectuant  $\frac{n \cdot (n-1)}{2}$  transformations de Jacobi en prenant pour valeurs successives de  $(p,q)$ , les valeurs  $(1,2), (1,3), \dots, (1,n); (2,3), \dots, (2,n); \dots; (n-1,n)$ . Un tel calcul est appelé un « balayage ».

En notant :

$$S_k = \sum_{1 \leq i < j \leq n} |a_{ij}|$$

on arrête les itérations quand  $S < \varepsilon$ , où  $\varepsilon$  est une précision donnée.

On peut alors montrer que la convergence devient très rapidement quadratique.

*Remarque 2* — En pratique on n'effectuera pas de transformation de Jacobi si  $|a_{pq}| < \sigma$ , où  $\sigma$  est un seuil donné. Usuellement, on prend pour valeur de  $\sigma$  :

$\sigma = \frac{S_k}{5 \cdot n^2}$  seulement pour les trois premiers balayages, puis à partir du quatrième,

on décide  $a_{pq} = 0$  si  $|a_{pq}| < \varepsilon' \cdot |a_{pp}|$  et  $|a_{pq}| < \varepsilon' \cdot |a_{qq}|$  avec  $\varepsilon'$  assez petit.

*Remarque 3* — La matrice  $V$ , donnant les vecteurs propres de  $A$ , est donnée par :

$$V = R_1 \cdot R_2 \cdot \dots$$

où les  $R_i$  sont les rotations successives. Cette matrice est donc définie par la récurrence :

$$V' = V \cdot R_i$$

si  $A' = {}^tR_i \cdot A \cdot R_i$  est une étape de calcul. Ce qui donne les relations :

$$\begin{aligned} v'_{ij} &= v_{ij} \quad (j \neq p, j \neq q) \\ v'_{ip} &= v_{ip} + s \cdot (v_{iq} - \tau \cdot v_{ip}) \quad (i = 1, \dots, n) \\ v'_{iq} &= v_{iq} - s \cdot (v_{ip} + \tau \cdot v_{iq}) \quad (i = 1, \dots, n) \end{aligned}$$

Ce qui nous donne en définitive la programmation structurée suivante, où  $D$  est un vecteur qui contiendra les valeurs propres et  $V$  une matrice dont les colonnes seront des vecteurs propres associés.

*PROCEDURE Jacobi*(Entrée  $n$  : Entier ;  $A$  : Matrice ; Sorties  $D$  : Vecteur ;  $V$  : Matrice) ;

*Début*

$V = Id$  ;  $Correct = Faux$  ;  $Iteration = 0$  ;  $D = 0$  ;

Pour  $i$  allant de 1 à  $n$  Faire  $d_i = a_{ii}$  ;

*Répéter*

$Iteration = Iteration + 1$  ;

$Sk = 0$  ;

Pour  $i$  allant de 1 à  $n - 1$  Faire

*Début*

Pour  $j$  allant de  $i + 1$  à  $n$  Faire

*Début*

$$Sk = Sk + |a_{ij}| ;$$

*Fin* ;

*Fin* ;

$Correct = (Sk < \epsilon)$  Ou ( $Iteration = MaxIteration$ ) ;

Si  $Iteration < 4$

Alors  $Seuil = Sk / (5 \cdot n^2)$

Sinon  $Seuil = \epsilon^2$  ;

Pour  $p$  allant de 1 à  $n - 1$  Faire

*Début*

Pour  $q$  allant de  $p + 1$  à  $n$  Faire

*Début*

$$\text{Si } (Iteration \geq 4) \text{ Et } (10^8 \cdot |a_{pq}| < |d_p|) \text{ Et } (10^8 \cdot |a_{pq}| < |d_q|)$$

Alors  $a_{pq} = 0$

Sinon Si  $|a_{pq}| > Seuil$

Alors *Début*

$$b = 0.5 \cdot (d_p - d_q) / a_{pq} ;$$

$$t = \frac{1}{\sqrt{b^2 + 1 + |b|}} ;$$

Si  $b < 0$  Alors  $t = -t$  ;

$$c = \frac{1}{\sqrt{1 + t^2}} ;$$

$$s = t \cdot c ;$$

```

 $\tau = s/(1 + c)$  ;
 $Aux = t \cdot a_{pq}$  ;
 $d_p = d_p + Aux$  ;
 $d_q = d_q - Aux$  ;
 $a_{pq} = 0$  ;
Pour j Allant de 1 à p - 1 Faire
Début
   $Aux = a_{jp}$  ;
   $a_{jp} = a_{jp} + s \cdot (a_{jq} - \tau a_{jp})$  ;
   $a_{jq} = a_{jq} - s \cdot (Aux + \tau a_{jq})$  ;
Fin ;
Pour j Allant de p + 1 à q - 1 Faire
Début
   $Aux = a_{pj}$  ;
   $a_{pj} = a_{pj} + s \cdot (a_{jq} - \tau a_{pj})$  ;
   $a_{jq} = a_{jq} - s \cdot (Aux + \tau a_{jq})$  ;
Fin ;
Pour j Allant de q + 1 à n Faire
Début
   $Aux = a_{pj}$  ;
   $a_{pj} = a_{pj} + s \cdot (a_{qj} - \tau a_{pj})$  ;
   $a_{qj} = a_{qj} - s \cdot (Aux + \tau a_{qj})$  ;
Fin ;
Pour j Allant de 1 à n Faire
Début
   $Aux = v_{jp}$  ;
   $v_{jp} = v_{jp} + s \cdot (v_{jq} - \tau v_{jp})$  ;
   $v_{jq} = v_{jq} - s \cdot (Aux + \tau v_{jq})$  ;
Fin ;
Fin ;
Fin ;
Jusqu'à Correct ;
Si Iteration = MaxIteration
Alors Arrêter('Nombre d'itérations trop grand') ;
Fin ;

```

## 4.6 Méthodes de calcul du polynôme caractéristique

Ces méthodes peuvent être utilisées pour des matrices de petite taille en relation avec des méthodes de résolution d'équations algébriques. On pourra, par exemple, utiliser la méthode de Bairstow pour chercher les racines du polynôme caractéristique.

### 4.6.1 Méthode de Souriau

Le polynôme caractéristique de A est noté ici :

$$P(x) = \text{Dét}(x \cdot I_d - A) = x^n + p_1 \cdot x^{n-1} + \dots + p_{n-1} \cdot x + p_n$$

avec les coefficients  $p_k$  dans  $\mathbb{R}$  inconnus a priori.

On note  $\text{Tr}(A)$  la trace de la matrice A, c'est-à-dire la somme de tous les termes diagonaux.

On définit les suites  $(A_p)$ ,  $(t_p)$  et  $(B_p)$  de la façon suivante :

$$\begin{cases} A_1 = A & t_1 = -\text{Tr}(A_1) & B_1 = A_1 + t_1 \cdot I_d \\ A_2 = B_1 \cdot A & t_2 = -\frac{\text{Tr}(A_2)}{2} & B_2 = A_2 + t_2 \cdot I_d \\ \cdot & \cdot & \cdot \\ A_n = B_{n-1} \cdot A & t_n = -\frac{\text{Tr}(A_n)}{n} & B_n = A_n + t_n \cdot I_d \end{cases}$$

Si  $\lambda_1, \lambda_2, \dots, \lambda_n$  sont les valeurs propres, dans  $C$ , de la matrice  $A$ , en notant :

$$S_p = \sum_{i=1}^n \lambda_i^p \quad (p \geq 1)$$

on a alors les *formules de Newton* :

$$(1) S_k + p_1 \cdot S_{k-1} + \dots + p_{k-1} \cdot S_1 + p_k \cdot k = 0 \quad (1 \leq k \leq n)$$

D'où le résultat suivant qui permet de calculer très simplement le polynôme caractéristique de  $A$  :

*Théorème* : Les coefficients  $t_k$ ,  $k = 1, 2, \dots, n$ , définis ci-dessus sont les coefficients du polynôme caractéristique de  $A$ .

*Démonstration* — On procède par récurrence sur  $k$ .

Pour  $k = 1$  :

$$t_1 = -\text{Tr}(A_1) = -\sum_{i=1}^n \lambda_i = p_1$$

Supposons le résultat vrai pour  $j = 1, \dots, k-1$ , avec  $k > 1$  : par définition, on a :

$$A_k = B_{k-1} \cdot A = A_{k-1} \cdot A + t_{k-1} \cdot A$$

et par récurrence, il vient :

$$A_k = A^k + t_1 \cdot A^{k-1} + \dots + t_{k-1} \cdot A$$

donc :

$$(1) A_k = A^k + p_1 \cdot A^{k-1} + \dots + p_{k-1} \cdot A \quad (k = 2, \dots, n)$$

On déduit alors que :

$$\text{Tr}(A_k) = \text{Tr}(A^k) + p_1 \cdot \text{Tr}(A^{k-1}) + \dots + t_{k-1} \cdot \text{Tr}(A)$$

C'est-à-dire :

$$t_k = -\frac{1}{k} \cdot (S_k + p_1 \cdot S_{k-1} + \dots + p_{k-1} \cdot S_1)$$

Ce qui s'écrit, en tenant compte des formules de Newton :

$$t_k = p_k$$

*Remarque* : *Calcul de l'inverse de  $A$*  — En faisant  $k = n$  dans l'identité (1) ci-dessus, on déduit que :

$$B_n = A_n + t_n \cdot I_d = P(A) = 0$$

d'après le théorème de Cayley-Hamilton.

Ce qui permet d'en déduire l'inverse de  $A$  :

$$A^{-1} = -\frac{1}{t_n} \cdot B_{n-1}$$

Ce qui donne, pour la programmation structurée, en définissant un type polynôme comme un tableau de réels numérotés de 0 à DimMax, où DimMax est la dimension maximale des matrices :

*PROCEDURE Souriau(Entrée n : Entier ; A : Matrice ;  
Sortie t : Polynome ; Inverse : Matrice) ;*

*Début*

*Pour i Allant de 1 à n Faire*

*Début*

*Pour j Allant de 1 à n Faire  $B_{ij} = 0$  ;*

*$B_{ii} = 1$  ;*

*Fin ;*

*Pour i Allant de 1 à n Faire*

*Début*

*$B = B \cdot A$  ;*

*$t_i = -\text{Trace}(B)/i$  ;*

*Pour j Allant de 1 à n Faire  $B_{jj} := B_{jj} + t_i$  ;*

*Si  $i = n - 1$  Alors  $\text{Inverse} = B$  ;*

*Fin ;*

*Pour i Allant de 1 à n Faire*

*Pour j Allant de 1 à n Faire  $\text{Inverse}_{ij} := -\text{Inverse}_{ij}/t_n$  ;*

*Fin ;*

#### 4.6.2 Méthode de Krylov

Avec les notations de (4.6.1), en notant, pour tout vecteur  $u$  de  $\mathbb{R}^n$  :  $u_k = A^k \cdot u$  ( $k \in \mathbb{N}$ ) on déduit du théorème de Cayley-Hamilton que les coefficients  $p_i$  sont solutions du système linéaire :

$$p_1 \cdot u_{n-1} + p_2 \cdot u_{n-2} + \dots + p_n \cdot u = -u_n$$

En prenant  $u$  au hasard, on a toutes les chances que le système ci-dessus soit non dégénéré de sorte que les coefficients du polynôme caractéristique s'obtiennent comme solution d'un système linéaire que l'on peut résoudre par la méthode de Gauss-Jordan (par exemple).

La programmation structurée étant alors évidente.

On notera « Krylov(Entrée n : Entier ; A : Matrice ; Sortie P : Vecteur) ; » une telle procédure de calcul.

De plus, toujours avec le théorème de Cayley-Hamilton, on déduit que l'inverse de  $A$  est donné par :

$$A^{-1} = -\frac{1}{p_n} \cdot (A^{n-1} + p_1 \cdot A^{n-2} + \dots + p_{n-1} \cdot I_d)$$

Ce qui donne, en s'inspirant de l'algorithme de Horner pour l'évaluation d'un polynôme, la procédure suivante de calcul de l'inverse :

*PROCEDURE InversionKrylov(Entrée n : Entier ; A : Matrice ; Sortie Inverse : Matrice) ;*

*Début*

*Krylov(n,A,P) ;*

*Inv = A + p<sub>1</sub> · I<sub>d</sub> ;*

*Pour i Allant de 2 à n - 1 Faire Inverse = A · Inverse + p<sub>i</sub> · I<sub>d</sub> ;*

*Inverse =  $-\frac{1}{p_n} \cdot \text{Inverse}$  ;*



*Fin ;*

#### **4.6.3 Méthode de Leverrier**

On conserve toujours les notations de (4.6.3).

En utilisant les formules de Newton, les coefficients  $p_k$  apparaissent comme solutions d'un système triangulaire inférieur.

La matrice de ce système se calcule facilement avec les  $S_k = \text{Tr}(A^k)$  (Cf Exercice 5).

## 5. Exercices

### 5.1 Etude d'un système tridiagonal

On se donne un entier  $k \geq 1$  et un réel  $a$  tel que  $|a| > 2$ .

Si  $n = 2^k - 1$ , on considère le système linéaire de  $n$  équations à  $n$  inconnues :

$$A \cdot x = b \quad (1)$$

avec :

$$\begin{cases} a_{ii} = a \quad (i = 1, \dots, n) \\ a_{i,i-1} = a_{i,i+1} = 1 \quad (i = 2, \dots, n-1) \\ a_{ij} = 0 \quad \text{dans les autres cas} \end{cases}$$

Un tel système s'écrit donc :

$$x_{i-1} + a \cdot x_i + x_{i+1} = b_i \quad (i = 1, \dots, n)$$

avec la convention :  $x_0 = x_{n+1} = 0$ .

On va étudier une méthode de résolution qui va consister à diviser  $n$  par 2 à chaque étape.

1°) Cas particulier  $k = 2$ .

(a) Ecrire le système obtenu.

(b) Comment calculer  $x_1$  et  $x_3$  connaissant  $x_2$ ?

(c) Montrer que  $x_2$  est solution d'un système  $1 \times 1$  :  $a^{(1)} \cdot x_2 = b_2^{(1)}$ .

2°) Cas particulier  $k = 4$ ,  $b_i = 1$  ( $i = 1, \dots, 15$ ). On note  $A^{(0)} = A$ ,  $b^{(0)} = b$ .

(a) Comment déterminer les  $x_{2^{p+1}}$  connaissant les  $x_{2^j}$ , pour  $p = 1, \dots, 7$ .

(b) On pose  $y_p = x_{2^p}$ , pour  $p = 1, \dots, 7$ . Et  $y$  est le vecteur de  $\mathbb{R}^7$  de composantes  $y_p$ .

Montrer que  $y$  est solution d'un système tridiagonal :  $A^{(1)} \cdot y = b^{(1)}$ .

(c) On pose  $z_q = y_{2^q}$ , pour  $q = 1, 2, 3$ . Et  $z$  est le vecteur de  $\mathbb{R}^3$  de composantes  $z_q$ .

Montrer que  $z$  est solution d'un système tridiagonal :  $A^{(2)} \cdot y = b^{(2)}$ .

On suppose de plus que  $a = 4$ .

(d) Calculer  $w_1, z_2, y_4, x_8$ .

(e) En déduire les valeurs de  $z_1$  et  $z_3$ .

(f) En déduire celles de  $y_1, y_3, y_5$  et  $y_7$ .

(g) Donner la solution du système (1).

3°) En s'inspirant du 2°) écrire les procédures permettant de résoudre (1) dans le cas général.

### 5.2 Valeurs propres d'une matrice tridiagonale

Soient  $a, b$  dans  $\mathbb{R}$ . Trouver les valeurs propres et les vecteurs propres de la matrice :

$$A = \begin{pmatrix} a & b & 0 & \cdot & \cdot & 0 \\ b & a & b & 0 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & 0 & b & a & b \\ 0 & \cdot & \cdot & 0 & b & a \end{pmatrix}$$

Cas particulier :  $a = 2$  et  $b = -1$ .

### 5.3 Matrices de Vandermonde

Si  $t_1, t_2, \dots, t_n$  sont des réels, on note :

$$V(t_1, t_2, \dots, t_n) = \left( (a_{ij}) \right)_{1 \leq i, j \leq n}, \text{ avec } a_{ij} = t_i^{j-1}$$

une matrice de Vandermonde.

On veut résoudre le système  $A \cdot x = b$ , où  $A$  est de Vandermonde. Ce problème étant lié à l'interpolation polynomiale, on introduit les notations suivantes :

$$\begin{cases} P_i(t) \text{ est le polynôme de degré } n-1 \text{ dans } \mathbb{R}[t] \text{ défini par } P_i(t_j) = \delta_{ij} \\ P_i(t) = \sum_{j=1}^n c_{ij} \cdot t^{j-1} \end{cases}$$

1°) Calculer le déterminant d'une matrice de Vandermonde.

2°) Montrer que, si les  $t_i$  sont deux à deux distincts, alors :

$$V(t_1, \dots, t_n)^{-1} = {}^t(c_{ij})$$

3°) En déduire une procédure de résolution d'un système de Vandermonde.

### 5.4 Polynôme caractéristique d'une matrice tridiagonale

$$\text{Soit } A = \begin{pmatrix} d_1 & e_1 & 0 & \cdot & \cdot & 0 \\ c_2 & d_2 & e_2 & 0 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & 0 & c_{n-1} & d_{n-1} & e_{n-1} \\ 0 & \cdot & \cdot & 0 & c_n & d_n \end{pmatrix} \text{ une matrice tridiagonale de}$$

polynôme caractéristique :

$$P(t) = \alpha_1 \cdot t^n + \alpha_2 \cdot t^{n-1} + \dots + \alpha_n \cdot t + \alpha_{n+1}$$

Pour tout  $k = 1, 2, \dots, n-1$ , on pose  $b_k = c_{k+1} \cdot e_k$ , on note  $A^{(k)}$  la  $k^{\text{ième}}$  sous-matrice principale de  $A$  et son polynôme caractéristique est noté :

$$P_k(t) = \alpha_1^{(k)} \cdot t^k + \dots + \alpha_{k+1}^{(k)}$$

1°) Trouver une relation de récurrence sur les polynômes  $P_k$ .

2°) En déduire un algorithme qui permet de calculer les coefficients de  $P$ .

### 5.5 Méthode de Levérier

#### et calcul du rayon spectral d'une matrice

Soit  $A$  une matrice à  $n$  lignes et  $n$  colonnes à coefficients réels de valeurs propres  $\lambda_1, \lambda_2, \dots, \lambda_n$  dans  $\mathbb{C}$ . On note :

$$S_p = \sum_{i=1}^n \lambda_i^p \quad (p \geq 1)$$

et :

$$P(x) = \text{Dét}(A - x \cdot I_d) = p_0 \cdot x^n + p_1 \cdot x^{n-1} + \dots + p_{n-1} \cdot x + p_n$$

le polynôme caractéristique de A, avec  $p_0 = (-1)^n$ .

On a alors les *formules de Newton* :

$$(1) \quad p_0 \cdot S_k + p_1 \cdot S_{k-1} + \dots + p_{k-1} \cdot S_1 + p_k \cdot k = 0 \quad (1 \leq k \leq n)$$

1°) Dédire avec (1), une procédure de calcul des coefficients  $p_k$ .

2°) On suppose que A admet des valeurs propres toutes distinctes dans R, avec :

$$0 < |\lambda_1| < |\lambda_2| < \dots < |\lambda_n|.$$

(a) Calculer  $\lim_{k \rightarrow +\infty} \left( \frac{S_k}{S_{k-1}} \right)$ .

(b) Calculer  $\lim_{k \rightarrow +\infty} \left( \sqrt[k]{|S_k|} \right)$ .

3°) Montrer que pour  $k \geq n$  :

$$p_0 \cdot S_k + p_1 \cdot S_{k-1} + \dots + p_n \cdot S_{k-n} = 0$$

4°) Donner une procédure de calcul de  $(S_k)$  pour  $k$  dans N et en déduire une procédure de calcul de  $\lambda_n$ .

### 5.6 Tridiagonalisation d'une matrice symétrique.

On se donne une matrice symétrique A.

On pose  $A_1 = A$  et on suppose qu'on a construit la matrice  $A_k$ , pour  $k = 1, \dots, n - 1$ , de la forme :

$$A_k = \begin{pmatrix} * & * & & & & & & \\ * & * & * & & & & & \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ & & * & * & \cdot & \cdot & * & \\ & & & * & \cdot & \cdot & \cdot & \\ & & & & \cdot & \cdot & \cdot & \cdot \\ & & & & * & \cdot & \cdot & * \end{pmatrix} \leftarrow \text{Ligne } k$$

↑  
colonne k

où les éléments non nuls (a priori) sont repérés par \*.

Montrer qu'il existe un réel r, de même signe que  $a_{k,k+1}^{(k)}$  tel que si  $v_k$  est le vecteur de coordonnées  $(0, \dots, 0, r + a_{k,k+1}^{(k)}, a_{k,k+2}^{(k)}, \dots, a_{k,n}^{(k)})$  et  $H_k$  la matrice de Householder définie par :

$$H_k = I_d - \frac{2}{\|v_k\|^2} \cdot v_k \cdot {}^t v_k$$

alors la matrice  $A_{k+1} = {}^t H_k \cdot A_k \cdot H_k$  vérifie bien la condition voulue.

En déduire que  $A$  est semblable à une matrice tridiagonale.

Ecrire une procédure de tridiagonalisation d'une matrice symétrique, puis une procédure de calcul du polynôme caractéristique.

Proposer une méthode de recherche des valeurs propres d'une matrice symétrique.

### 5.7 Méthode de Givens et Householder pour les matrices tridiagonales

$$\text{Soit } A = \begin{pmatrix} a_1 & b_1 & 0 & \cdot & \cdot & 0 \\ b_1 & a_2 & b_2 & 0 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & 0 & b_{n-2} & a_{n-1} & b_{n-1} \\ 0 & \cdot & \cdot & 0 & b_{n-1} & a_n \end{pmatrix} \text{ une matrice tridiagonale et}$$

symétrique. Les valeurs propres de  $A$  sont alors réelles et notées :

$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

On suppose, en outre, que les  $b_i$  sont tous non nuls (sinon on découpe  $A$  en deux blocs et on travaille séparément sur chacun d'eux).

On construit la suite de polynômes  $(P_k)_{0 \leq k \leq n}$  de la manière suivante :

$$\begin{cases} P_0 = 1; P_1 = a_1 - x; \\ P_k = (a_k - x) \cdot P_{k-1} - (b_{k-1})^2 \cdot P_{k-2} \quad (k = 2, \dots, n) \end{cases}$$

1°) Montrer que la suite  $(P_k)$  vérifie les propriétés suivantes :

(a)  $P_k$  est de degré égal à  $k$  ;

(b) Pour  $k = 1, \dots, n$  :

$$\lim_{k \rightarrow \infty} P_k(x) = +\infty$$

(c)  $P_k$  possède  $k$  racines réelles distinctes séparées par les  $k - 1$  racines de  $P_{k-1}$  ;

(d)  $P_k$  est le polynôme caractéristique de la matrice  $A_k$  déduite de  $A$  en supprimant les  $n - k$  dernières lignes et colonnes.

2°) On note, pour tout  $x$  réel,  $s_k(x) = \pm 1 =$  le signe de  $P_k(x)$  (avec la convention  $s_k(x) = s_{k-1}(x)$ , si  $P_k(x) = 0$ ).  $N(x)$  désigne le nombre de changements de signes entre deux termes consécutifs de la suite  $(s_0(x), s_1(x), \dots, s_n(x))$ .

(a) Montrer que  $N(x)$  est égal au nombre de racines de  $P_n$  qui sont inférieures à  $x$ .

(b) Ecrire une procédure de calcul de  $N(x)$ .

3°) Toutes les valeurs propres de  $A$  sont dans l'intervalle  $[-M, M]$ , où

$$M = \text{Max} \left\{ \sum_{i=1}^n |a_{ij}| ; 1 \leq j \leq n \right\}$$

(a) Ecrire une procédure de calcul de  $M$ .

(b) Ecrire une procédure de calcul de la  $k^{\text{ème}}$  valeur propre de  $A$ , pour tout  $k = 1, 2, \dots, n$ .

### 5.8 Factorisation $Q \cdot R$

Si  $v$  est un vecteur non nul, on définit la *matrice de Householder* par :

$$H(v) = I_d - \frac{2}{\|v\|} \cdot v \cdot {}^t v$$

1°) Calculer le déterminant et l'inverse de  $H(v)$ .

2°) Soit  $a$  un vecteur de  $\mathbb{R}^n$  non colinéaire à  $e_1$  (premier vecteur de base). Montrer qu'il existe un unique vecteur  $v$  tel que :

- (i)  $v = a + \lambda \cdot e_1$  ;
- (ii)  $H(v) \cdot a = \mu \cdot e_1$  ;
- (iii)  $\|v\| \geq \sqrt{2} \cdot \|a\|$

Ce vecteur sera noté  $v = J_n(a)$ . Si  $a$  est proportionnel à  $e_1$ , on pose  $J_n(a) = 0$ .

3°) On construit la suite de matrices  $(A_k)$ , de la manière suivante :

$$A_1 = A ;$$

Pour  $k = 1, \dots, n$ , on note  $A_k = \left( (a_{ij}^{(k)}) \right)_{1 \leq i, j \leq n}$  et  $a'_k$  est le vecteur de  $\mathbb{R}^{n-k+1}$  de

composantes  $a_{ik}^{(k)}$  avec  $i = k, \dots, n$ .

On pose alors  $v'_k = J_{n-k+1}(a'_k)$  et  $v_k$  est le vecteur de  $\mathbb{R}^n$  dont les  $k - 1$  premières composantes sont nulles, les autres étant égales à celles de  $v'_k$ .

Enfin  $H_k = H(v_k)$ .

On pose alors  $H_{k+1} = H_k \cdot A_k$ , pour  $k = 1, \dots, n - 1$ .

(a) Montrer que  $A_n$  est triangulaire supérieure.

(b) Montrer qu'il existe une matrice diagonale  $D$  telle que si on pose :

$$Q = H_1 \cdot \dots \cdot H_{n-1} \cdot D \text{ et } R = D \cdot A_n$$

alors  $Q$  est une matrice orthogonale,  $R$  est triangulaire supérieure à coefficients diagonaux positifs et :

$$A = Q \cdot R$$

4°) Ecrire une procédure permettant de réaliser la factorisation  $Q \cdot R$  d'une matrice réelle inversible. En déduire un moyen de calculer l'inverse.

5°) En déduire une procédure de résolution du système  $A \cdot x = b$ .

### 5.9 Matrices complexes

Ecrire un paquetage permettant de manipuler des variables de type complexe. Puis un paquetage permettant de manipuler des matrices à coefficients complexes (voir le paragraphe 6.1).

### 5.10 Paquetage générique de manipulation des matrices

Ecrire un paquetage générique de manipulation des matrices (voir le paragraphe 6.1). On l'instanciera ensuite sur les réels, les rationnels et les complexes.

## 6. Programmation Ada

### 6.1 Spécifications des paquetages *COMMON\_MATRIX* et *MATRIX*

Ce paquetage est formé d'utilitaires pour travailler avec des matrices carrées et des vecteurs de dimension au plus DIM\_MAX déclarée en constante.

Le nom de chaque procédure est assez explicite pour comprendre de quoi il s'agit.

```

package COMMON_MATRIX is
DIM_MAX : constant NATURAL := 80 ;
type MATRICE is array(NATURAL range <>, NATURAL range <>) of FLOAT ;
type VECTEUR is array(NATURAL range <>) of FLOAT ;
type MATRICE_TRIDIAGONALE is array(NATURAL range <>) of VECTEUR(1..3) ;
ERREUR_DIMENSION : exception ;
ERREUR_DETERMINANT_NUL : exception ;
end COMMON_MATRIX ;

with TEXT_IO, COMMON_MATRIX ;
use COMMON_MATRIX ;
package MATRIX is
procedure GET(A : out MATRICE) ;
procedure MATRICE_ALEATOIRE(A : out MATRICE) ;
procedure MATRICE_DU_FICHER(A : out MATRICE ;
    NOM_DE_REPERTOIRE : in STRING := "A :\DONNEES") ;
procedure MATRICE_HILBERT(A : out MATRICE) ;
procedure MATRICE_TEST(A : in out MATRICE) ;
procedure MATRICE_VAN_DER_MONDE(v : in VECTEUR ; A : in out MATRICE) ;
procedure GET(A : out MATRICE_TRIDIAGONALE) ;
procedure MATRICE_TRIDIAG_ALEATOIRE(A : out MATRICE_TRIDIAGONALE) ;
procedure MATRICE_TRIDIAG_DU_FICHER(A : out MATRICE_TRIDIAGONALE ;
    NOM_DE_REPERTOIRE : in STRING := "A :\DONNEES") ;
procedure MATRICE_SYMETRIQUE(A : in out MATRICE) ;
procedure MATRICE_SYM_ALEATOIRE(A : in out MATRICE) ;
procedure MATRICE_SYM_DU_FICHER(A : in out MATRICE ;
    NOM_DE_REPERTOIRE : in STRING := "A :\DONNEES") ;
procedure PUT_LINE(A : in MATRICE ; FORE : in TEXT_IO.FIELD := 6 ;
    AFT : in TEXT_IO.FIELD := 4 ; EXP : in TEXT_IO.FIELD := 0) ;
procedure PUT_LINE(A : in MATRICE_TRIDIAGONALE ;
    FORE : in TEXT_IO.FIELD := 6 ; AFT : in TEXT_IO.FIELD := 4 ;
    EXP : in TEXT_IO.FIELD := 0) ;
procedure GET(v : out VECTEUR) ;
procedure VECTEUR_ALEATOIRE(v : out VECTEUR) ;
procedure VECTEUR_DU_FICHER(v : out VECTEUR ;
    NOM_DE_REPERTOIRE : in STRING := "A :\DONNEES") ;
procedure PUT_LINE(v : in VECTEUR ; FORE : in TEXT_IO.FIELD := 6 ;
    AFT : in TEXT_IO.FIELD := 4 ; EXP : in TEXT_IO.FIELD := 0) ;
function TRACE(A : in MATRICE) return FLOAT ;
function TRANSPOSE(A : in MATRICE) return MATRICE ;
function "+"(A, B : in MATRICE) return MATRICE ;

```

```

function "-"(A, B : in MATRICE) return MATRICE ;
function "*" (x : in FLOAT ; A : in MATRICE) return MATRICE ;
function "*" (A : in MATRICE ; v : in VECTEUR) return VECTEUR ;
function "*" (A : in MATRICE_TRIDIAGONALE ; v : in VECTEUR)
    return VECTEUR ;
function "+"(u, v : in VECTEUR) return VECTEUR ;
function "-"(u, v : in VECTEUR) return VECTEUR ;
function "*" (x : in FLOAT ; u : in VECTEUR) return VECTEUR ;
function NORME_DU_SUP(v : in VECTEUR) return FLOAT ;
function NORME_EUCLIDIENNE(v : in VECTEUR) return FLOAT ;
function PRODUIT_SCALAIRE(u, v : in VECTEUR) return FLOAT ;
end MATRIX ;

```

## 6.2 Spécification du paquetage ALGEBLIN

Dans ce paquetage, on trouvera des procédures permettant de résoudre des systèmes linéaires par des méthodes directes et itératives.

En commentaire, on indique à quel paragraphe de ce chapitre on devra se reporter.

```

with COMMON_MATRIX, MATRIX ;
use COMMON_MATRIX, MATRIX ;
package ALGEBLIN is
procedure SYSTEME_TRIANGULAIRE_SUPERIEUR(A : in MATRICE ;
    b : in VECTEUR ; x : out VECTEUR) ; -- § 3.3
procedure PIVOTAGE(A : in out MATRICE ; b : in out VECTEUR) ; -- § 3.4
function DET_MATRICE_TRIANGULAIRE(A : in MATRICE)
    return FLOAT ; -- § 3.3
function INV_MATRICE_TRIANG_INF(A : in MATRICE)
    return MATRICE ; -- § 3.3
function INV_MATRICE_TRIANG_SUP(A : in MATRICE)
    return MATRICE ; -- § 3.3
procedure DECOMPOSITION_LR(A : in MATRICE ; L, R : in out MATRICE) ;
    -- § 3.6
function DETERMINANT_LR(A : in MATRICE) return FLOAT ; -- § 3.6
function INVERSE_LR(A : in MATRICE) return MATRICE ; -- § 3.6
procedure GAUSS_JORDAN(A : in out MATRICE ; b : in out VECTEUR) ;
    -- § 3.5
function INVERSE_GAUSS_JORDAN(A : in MATRICE) return MATRICE ; -- § 3.5
procedure SYSTEME_TRIDIAGONAL(A : in out MATRICE_TRIDIAGONALE ;
    b : in out VECTEUR) ; -- § 3.7
procedure CHOLESKY(A : in MATRICE ; C : in out MATRICE) ; -- § 3.8
function INVERSE_VAN_DER_MONDE(x : in VECTEUR) return MATRICE ;
    -- Ex. 3
procedure DECOMPOSITION_QR(A : in MATRICE ; Q, R : in out MATRICE) ;
    -- Ex. 8
procedure JACOBI(A : in MATRICE ; b : in VECTEUR ; x : in out VECTEUR) ;
    -- § 3.9.3
procedure RELAXATION(A : in MATRICE ; b : in VECTEUR ;
    OMEGA : in FLOAT := 1.0 ; x : in out VECTEUR) ; -- § 3.9.6
ERREUR_DIAGONALE_NON_DOMINANTE : exception
ERREUR_NON_CONVERGENCE : exception ;

```



end ALGEBLIN ;

### 6.3 Spécification du paquetage *DONNEES\_ALGEBLIN*

Dans ce paquetage on trouvera des procédures facilitant l'entrée de matrices et vecteurs de divers types. Il est utilisé dans les procédures de démonstrations du paragraphe suivant.

```
with COMMON_MATRIX ;
use COMMON_MATRIX ;
package DONNEES_ALGEBLIN is
procedure DONNEE_MATRICE(A : in out MATRICE) ;
procedure DONNEE_VECTEUR(b : in out VECTEUR) ;
procedure DONNEE_MATRICE_SYMETRIQUE(A : in out MATRICE) ;
procedure DONNEE_MATRICE_TRIDIAGONALE(A : in out MATRICE_TRIDIAGONALE);
end DONNEES_ALGEBLIN ;
```

### 6.4 Démonstration de la méthode des pivots de Gauss

```
with TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
ALGEBLIN, DONNEES_ALGEBLIN ;
use TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
ALGEBLIN, DONNEES_ALGEBLIN ;
procedure DEM_PIV is
n : NATURAL ;
begin
MODE_AFFICHAGE ;
PUT_LINE(IMP, " RESOLUTION D'UN SYSTEME LINEAIRE" ) ;
PUT_LINE(IMP, " METHODE DES PIVOTS DE GAUSS" ) ;
NEW_LINE(IMP) ;
ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Dimension du systeme : " ) ;
declare
A : MATRICE(1..n,1..n) ;
b : VECTEUR(1..n) ;
begin
DONNEE_MATRICE(A) ; DONNEE_VECTEUR(b) ;
PUT_LINE(IMP, " SYSTEME AVANT PIVOTAGE" ) ;
PUT_LINE(IMP, "MATRICE A : " ) ; PUT_LINE(A) ;
PAUSE ;
PUT_LINE(IMP, "VECTEUR b" ) ; PUT_LINE(b) ;
PAUSE ;
PIVOTAGE(A,b) ;
CLRSCR ;
PUT_LINE(IMP, " SYSTEME APRES PIVOTAGE" ) ;
PUT_LINE(IMP, "MATRICE A : " ) ; PUT_LINE(A) ;
PAUSE ;
PUT_LINE(IMP, "VECTEUR b" ) ; PUT_LINE(b) ;
PAUSE ;
CLRSCR ;
PUT_LINE(IMP, " SOLUTION DU SYSTEME" ) ;
SYSTEME_TRIANGULAIRE_SUPERIEUR(A,b,b) ;
PUT_LINE(b) ;
CLOSE(IMP) ;
end ;
```

```
end DEM_PIV ;
```

## 6.5 Démonstration de la méthode L·R

```
with TEXT_IO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
     MATRIX, ALGEBLIN, DONNEES_ALGEBLIN ;
use TEXT_IO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
     MATRIX, ALGEBLIN, DONNEES_ALGEBLIN ;
procedure DEM_LR is
n : NATURAL ;
begin
  MODE_AFFICHAGE ;
  PUT_LINE(IMP, "    DECOMPOSITION A = L.R DE CROUT") ;
  PUT_LINE(IMP, "    CALCUL DE L'INVERSE DE A") ;
  NEW_LINE(IMP) ;
  PUT_LINE("Attention : La decomposition L.R est possible s.s.i ") ;
  PUT_LINE("il n'y a pas de permutations dans la methode de GAUSS") ;
  NEW_LINE ;
  ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Dimension de la matrice : ") ;
  declare
    A, L, R, Inv : MATRICE(1..n,1..n) ;
    DET : FLOAT ;
  begin
    DONNEE_MATRICE(A) ;
    DECOMPOSITION_LR(A,L,R) ;
    PUT_LINE(IMP, "    DECOMPOSITION L.R") ;
    PUT_LINE(IMP, "MATRICE L :") ; PUT_LINE(L) ;
    PAUSE ;
    CLRSCR ;
    PUT_LINE(IMP, "MATRICE R") ; PUT_LINE(R) ;
    PAUSE ;
    CLRSCR ;
    PUT_LINE(IMP, "    VERIFICATION : L*R = ") ;
    A := L*R ;
    PUT_LINE(A) ;
    PAUSE ;
    CLRSCR ;
    PUT_LINE(IMP, "    CALCUL DU DETERMINANT") ;
    PUT("Det(A) = ") ;
    DET := DET_MATRICE_TRIANGULAIRE(R) ;
    FIO.PUT(DET,6,4,0) ; NEW_LINE ;
    PAUSE ;
    CLRSCR ;
    PUT_LINE(IMP, "    CALCUL DE L'INVERSE DE A") ;
    INV := INV_MATRICE_TRIANG_SUP(R)*INV_MATRICE_TRIANG_INF(L) ;
    PUT_LINE(INV) ;
    PAUSE ; CLRSCR ;
    PUT_LINE(IMP, "    VERIFICATION : A*INV(A) =") ;
    A := A*INV ;
    PUT_LINE(A) ;
    PAUSE ;
    CLRSCR ;
    CLOSE(IMP) ;
  end ;
end ;
```

```
end DEM_LR ;
```

## 6.6 Démonstration de la méthode du double balayage de Cholesky

```
with TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
     ALGEBLIN, DONNEES_ALGEBLIN ;
use TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
     ALGEBLIN, DONNEES_ALGEBLIN ;

procedure DEM_TRID is
n : NATURAL ;
begin
  MODE_AFFICHAGE ;
  PUT_LINE(IMP, "      RESOLUTION D'UN SYSTEME LINEAIRE TRIDIAGONAL") ;
  PUT_LINE(IMP, "      METHODE DES PIVOTS DE GAUSS") ;
  NEW_LINE(IMP) ;
  PUT_LINE("Attention : la methode est valable si il n'y a pas de") ;
  PUT_LINE("permutations dans la methode de Gauss.") ;
  NEW_LINE ;
  ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Dimension du systeme : ") ;
  declare
    A : MATRICE_TRIDIAGONALE(1..n) ;
    b : VECTEUR(1..n) ;
  begin
    DONNEE_MATRICE_TRIDIAGONALE(A) ;
    DONNEE_VECTEUR(b) ;
    PUT_LINE(IMP, "      SYSTEME AVANT PIVOTAGE") ;
    PUT_LINE("MATRICE A :") ; PUT_LINE(A) ;
    PAUSE ;
    PUT_LINE(IMP, "VECTEUR b") ; PUT_LINE(b) ;
    PAUSE ;
    CLRSCR ;
    SYSTEME_TRIDIAGONAL(A,b) ;
    PUT_LINE(IMP, "      SOLUTION DU SYSTEME") ;
    PUT_LINE(b) ;
    CLOSE(IMP) ;
  end ;
end DEM_TRID ;
```

## 6.7 Démonstration de la méthode de Gauss-Jordan

```
with TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
     ALGEBLIN, DONNEES_ALGEBLIN ;
use TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
     ALGEBLIN, DONNEES_ALGEBLIN ;
procedure DEM_GJ is
n : NATURAL ;
begin
  MODE_AFFICHAGE ;
  PUT_LINE(IMP, "      Resolution d'un systeme lineaire.") ;
  PUT_LINE(IMP, "      Calcul de l'inverse d'une matrice.") ;
  PUT_LINE(IMP, "      Methode de Gauss-Jordan.") ;
  NEW_LINE(IMP) ;
  ENTRER_ENTIER_BORNE(1,DIM_MAX,n,"Dimension du systeme : ") ;
```

```

declare
  A, INV : MATRICE(1..n,1..n) ;
  b : VECTEUR(1..n) ;
begin
  DONNEE_MATRICE(A) ;
  DONNEE_VECTEUR(b) ;
  PUT_LINE(IMP,"RESOLUTION DE A.x = b, AVEC :") ;
  PUT_LINE(IMP," A = " ) ; PUT_LINE(A) ;
  PAUSE ; CLRSCR ;
  PUT_LINE(IMP," et b = " ) ; PUT_LINE(b) ;
  PAUSE ; CLRSCR ;
  GAUSS_JORDAN(A,b) ;
  PUT_LINE(IMP,"          Solution, de A.x = b") ;
  NEW_LINE(IMP) ;
  PUT_LINE(b) ;
  PAUSE ; CLRSCR ;
  PUT_LINE(IMP,"Inversion d'une matrice") ;
  NEW_LINE(IMP) ;
  DONNEE_MATRICE(A) ;
  PUT_LINE(IMP,"INVERSION DE LA MATRICE A = " ) ;
  NEW_LINE(IMP) ;
  PUT_LINE(A) ;
  PAUSE ; CLRSCR ;
  PUT_LINE(IMP,"Inv(A) = " ) ;
  NEW_LINE(IMP) ;
  PUT_LINE(INVERSE_GAUSS_JORDAN(A)) ;
  PAUSE ; CLRSCR ;
  CLOSE(IMP) ;
end ;
end DEM_GJ ;

```

## 6.8 Démonstration de la méthode de Cholesky

```

with TEXT_IO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
  MATRIX, ALGEBLIN, DONNEES_ALGEBLIN ;
use TEXT_IO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
  MATRIX, ALGEBLIN, DONNEES_ALGEBLIN ;

procedure DEM_CHOL is
  n : NATURAL ;
begin
  MODE_AFFICHAGE ;
  PUT_LINE(IMP,"          DECOMPOSITION DE CHOLESKY A = B.Transpose(B)" ) ;
  NEW_LINE(IMP) ;
  PUT_LINE("Attention : La decomposition est valable pour" ) ;
  PUT_LINE("A symetrique definie positive." ) ;
  NEW_LINE ;
  ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Dimension de la matrice : " ) ;
  declare
    A, B, INV : MATRICE(1..n,1..n) ;
    DET : FLOAT ;
  begin
    DONNEE_MATRICE(A) ;

```

```

    CHOLESKY(A,B) ;
    PUT_LINE(IMP,"    MATRICE B") ;
    PUT_LINE(B) ;
    PAUSE ; CLRSCR ;
    PUT_LINE(IMP,"    VERIFICATION : B*Transpose(B) = ") ;
    A := B*TRANSPPOSE(B) ;
    PUT_LINE(A) ;
    PAUSE ; CLRSCR ;
    PUT_LINE(IMP,"    CALCUL DU DETERMINANT") ;
    PUT(IMP,"Det(A) = ") ;
    DET := DET_MATRICE_TRIANGULAIRE(B) ;
    FIO.PUT(IMP,DET,6,4,0) ;
    NEW_LINE(IMP) ;
    PAUSE ; CLRSCR ;
    PUT_LINE(IMP,"    CALCUL DE L'INVERSE DE A") ;
    INV :=
INV_MATRICE_TRIANG_SUP(TRANSPPOSE(B))*INV_MATRICE_TRIANG_INF(B) ;
    PUT_LINE(INV) ;
    PAUSE ; CLRSCR ;
    PUT_LINE(IMP,"    VERIFICATION : A*INV(A) =") ;
    A := A*INV ;
    PUT_LINE(A) ;
    PAUSE ; CLRSCR ;
    CLOSE(IMP) ;
end ;
end DEM_CHOL ;

```

## 6.9 Démonstration de la méthode de Jacobi

```

with TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
    ALGEBLIN, DONNEES_ALGEBLIN ;
use TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
    ALGEBLIN, DONNEES_ALGEBLIN ;

procedure DEM_JAC is
n : NATURAL ;
begin
    MODE_AFFICHAGE ;
    TEXT_IO.PUT_LINE(IMP,"Resolution d'un systeme lineaire.") ;
    TEXT_IO.PUT_LINE(IMP,"Methode iterative de Jacobi.") ;
    TEXT_IO.NEW_LINE(IMP) ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Dimension du systeme : ") ;
    declare
        A : MATRICE(1..n,1..n) ;
        b, x : VECTEUR(1..n) ;
    begin
        DONNEE_MATRICE(A) ;
        DONNEE_VECTEUR(b) ;
        TEXT_IO.PUT_LINE(IMP,"RESOLUTION DE A.x = b, AVEC :") ;
        TEXT_IO.PUT_LINE(IMP," A = ") ; PUT_LINE(A) ;
        PAUSE ; CLRSCR ;
        TEXT_IO.PUT_LINE(IMP," et b = ") ; PUT_LINE(b) ;
        PAUSE ; CLRSCR ;
    end ;
end ;

```

```

    CHRONOMETRE(0) ;
    JACOBI(A,b,x) ;
    CHRONOMETRE(1) ;
    TEXT_IO.PUT_LINE(IMP,"La solution est :") ; PUT_LINE(x) ;
    PAUSE ; CLRSCR ;
    CLOSE(IMP) ;
end ;
end DEM_JAC ;

```

## 6.10 Démonstration de la méthode de relaxation

```

with TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
    ALGEBLIN, DONNEES_ALGEBLIN ;
use TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
    ALGEBLIN, DONNEES_ALGEBLIN ;

procedure DEM_REL is
n : NATURAL ;
begin
    MODE_AFFICHAGE ;
    TEXT_IO.PUT_LINE(IMP,"Resolution d'un systeme lineaire.") ;
    TEXT_IO.PUT_LINE(IMP,"Methode iterative de Relaxation.") ;
    TEXT_IO.NEW_LINE(IMP) ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Dimension du systeme : ") ;
    declare
        A : MATRICE(1..n,1..n) ;
        b, x : VECTEUR(1..n) ;
        Omega : FLOAT ;
    begin
        DONNEE_MATRICE(A) ;
        DONNEE_VECTEUR(b) ;
        TEXT_IO.PUT_LINE(IMP,"RESOLUTION DE A.x = b, avec :") ;
        TEXT_IO.PUT_LINE(IMP," A = ") ; PUT_LINE(A) ;
        PAUSE ; CLRSCR ;
        TEXT_IO.PUT_LINE(IMP," et b = ") ; PUT_LINE(b) ;
        PAUSE ; CLRSCR ;
        loop
            ENTRER_REEL_BORNE(0.0001,1.9999,OMEGA,
                "Parametre de relaxation (entre 0 et 2) : ") ;
            CHRONOMETRE(0) ;
            RELAXATION(A,b,Omega,x) ;
            CHRONOMETRE(1) ;
            TEXT_IO.PUT_LINE(IMP,"La solution est :") ;
            PUT_LINE(x) ;
            PAUSE ; CLRSCR ;
            exit when (not LIRE_REPONSE("Un autre essai : ")) ;
        end loop ;
        CLOSE(IMP) ;
    end ;
end DEM_REL ;

```

## 6.11 Démonstration de la décomposition $Q \cdot R$

```

with TEXT_IO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
     MATRIX, ALGEBLIN, DONNEES_ALGEBLIN ;
use TEXT_IO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
     ALGEBLIN, DONNEES_ALGEBLIN ;

procedure DEM_QR is
n : NATURAL ;
begin
  MODE_AFFICHAGE ;
  TEXT_IO.PUT_LINE(IMP,"Decomposition A = Q*R, avec :
                    Q orthogonale R triangulaire superieure") ;
  NEW_LINE(IMP) ;
  ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Entrer la dimension de la matrice") ;
  declare
    A, Q, R : MATRICE(1..n,1..n) ;
    b : VECTEUR(1..n) ;
  begin
    DONNEE_MATRICE(A) ;
    TEXT_IO.PUT_LINE(IMP,"Matrice A :") ; PUT_LINE(A) ;
    PAUSE ; CLRSCR ;
    TEXT_IO.PUT_LINE(IMP,"Calcul de la decomposition Q*R ...") ;
    DECOMPOSITION_QR(A,Q,R) ;
    TEXT_IO.PUT_LINE(IMP,"Matrice Q :") ; PUT_LINE(Q) ;
    PAUSE ; CLRSCR ;
    TEXT_IO.PUT_LINE(IMP,"Matrice R :") ; PUT_LINE(R) ;
    PAUSE ; CLRSCR ;
    TEXT_IO.PUT_LINE(IMP,"Determinant de A.") ;
    PUT(IMP,"Det(A) = Det(R) = ") ;
    FIO.PUT(IMP,DET_MATRICE_TRIANGULAIRE(R),6,4,0) ;
    PAUSE ; CLRSCR ;
    TEXT_IO.PUT_LINE(IMP,"Resolution de A*x = b.") ;
    DONNEE_VECTEUR(b) ;
    b := TRANSPOSE(Q)*b ;
    SYSTEME_TRIANGULAIRE_SUPERIEUR(R,b,b) ;
    TEXT_IO.PUT_LINE(IMP,"Solution x :") ;
    PUT_LINE(b) ;
    PAUSE ; CLRSCR ;
    TEXT_IO.PUT_LINE(IMP,"Verification A*x = ") ;
    b := A*b ;
    PUT_LINE(b) ;
    PAUSE ; CLRSCR ;
    TEXT_IO.PUT_LINE(IMP,"Inverse de A") ;
    R := INV_MATRICE_TRIANG_SUP(R) ;
    Q := R*TRANSPOSE(Q) ;
    PUT_LINE(Q) ;
    PAUSE ; CLRSCR ;
    TEXT_IO.PUT_LINE(IMP,"Verification A*Inv(A) = ") ;
    A := A*Q ;
    PUT_LINE(A) ;
    PAUSE ; CLRSCR ;
    CLOSE(IMP) ;
  end ;
end DEM_QR ;

```

```

end ;
end DEM_QR ;

```

## 6.12 Démonstration du calcul de l'inverse d'une matrice de Van Der Monde

```

with TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
     ALGEBLIN, DONNEES_ALGEBLIN ;
use TEXT_IO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
     ALGEBLIN, DONNEES_ALGEBLIN ;

procedure DEM_VAND is
n : NATURAL ;
begin
  MODE_AFFICHAGE ;
  PUT_LINE(IMP, "      INVERSE D'UNE MATRICE DE VAN DER MONDE") ;
  NEW_LINE(IMP, 5) ;
  ENTRER_ENTIER_BORNE(2, 10,
    "Dimension du vecteur definissant la matrice (< 10) : ") ;
  declare
    A, INV : MATRICE(1..n, 1..n) ;
    x : VECTEUR(1..n) ;
  begin
    DONNEE_VECTEUR(x) ;
    MATRICE_VAN_DER_MONDE(x, A) ;
    PUT_LINE(IMP, "      MATRICE DE VAN DER MONDE ASSOCIEE") ;
    PUT_LINE(A) ;
    PAUSE ; CLRSCR ;
    PUT_LINE(IMP, "      INVERSE DE LA MATRICE") ;
    INV := INVERSE_VAN_DER_MONDE(x) ;
    PUT_LINE(INV) ;
    PAUSE ; CLRSCR ;
    PUT_LINE(IMP, "VERIFICATION : A*INV = ") ;
    PUT_LINE(A*INV) ;
    PAUSE ; CLRSCR ;
  end ;
  CLOSE(IMP) ;
end DEM_VAND ;

```

## 6.13 Spécification du paquetage SPECTRE

Dans ce paquetage, on trouvera des procédures permettant de calculer les éléments propres d'une matrice.

En commentaire, on indique à quel paragraphe de ce chapitre on devra se reporter.

```

with COMMON_MATRIX ;
use COMMON_MATRIX ;
package SPECTRE is
procedure DEFLATION(A : in MATRICE ; VALEURS_PROPRES : in out VECTEUR ;
  VECTEURS_PROPRES : in out MATRICE) ; -- § 4.3
procedure RUTISHAUSER(A : in MATRICE ;
  VALEURS_PROPRES : in out VECTEUR) ; -- § 4.4

```



```

procedure JACOBI(A : in MATRICE ; D : in out VECTEUR ;
                V : in out MATRICE) ; -- § 4.5
procedure SOURIAU(A : in MATRICE ; T : in out VECTEUR ;
                 INVERSE : in out MATRICE) ; -- 4.6.1
end SPECTRE ;

```

## 6.14 Démonstration de la méthode de déflation

```

with TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
     MATRIX, DONNEES_ALGEBLIN, SPECTRE ;
use TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
     MATRIX, DONNEES_ALGEBLIN, SPECTRE ;

procedure DEM_DEFL is
n, p : NATURAL ;
begin
    MODE_AFFICHAGE ;
    PUT_LINE(IMP,"Valeurs propres d'une matrice : Methode de deflation")
;
    NEW_LINE(IMP) ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Dimension de la matrice : ") ;
    ENTRER_ENTIER_BORNE(1,n,p,"Nombre de valeurs propres a calculer : ") ;
    declare
        A : MATRICE(1..n,1..n) ;
        VALEURS_PROPRES : VECTEUR(1..p) ;
        VECTEURS_PROPRES : MATRICE(1..n,1..p) ;
        S : FLOAT ;
    begin
        DONNEE_MATRICE(A) ;
        S := TRACE(A) ;
        PUT_LINE(IMP,"Valeurs propres de la matrice A = ") ;
        PUT_LINE(A) ; NEW_LINE(IMP) ;
        PAUSE ; CLRSCR ;
        CHRONOMETRE(0) ;
        DEFLATION(A,VALEURS_PROPRES,VECTEURS_PROPRES) ;
        CHRONOMETRE(1) ; NEW_LINE(IMP) ;
        for j in 1..p
        loop
            PUT(IMP,"Valeur propre Numero ") ;
            PUT(IMP,j) ; PUT(IMP," = ") ;
            PUT(IMP,VALEURS_PROPRES(j),5,5,0) ; NEW_LINE(IMP) ;
            PUT_LINE(IMP,"Vecteur propre associe :") ;
            PUT(IMP,"(") ;
            for i in 1..n - 1
            loop
                PUT(IMP,VECTEURS_PROPRES(i,j),5,5,0) ;
                PUT(IMP,",") ;
            end loop ;
            PUT(IMP,VECTEURS_PROPRES(n,j),5,5,0) ; PUT_LINE(IMP,")") ;
            if (j mod 5 = 0) then
                PAUSE ;
            end if ;
        end loop ;
    end loop ;
end DEM_DEFL ;

```

```

if (p = n) then
  PAUSE ; CLRSCR ;
  PUT_LINE(IMP,"Verification") ;
  NEW_LINE(IMP) ;
  PUT(IMP,"Trace de A = ") ;
  PUT(IMP,S,5,5,0) ; NEW_LINE(IMP) ;
  PUT(IMP,"Somme des valeurs propres = ") ;
  S := 0.0 ;
  for i in 1..n
  loop
    S := S + VALEURS_PROPRES(i) ;
  end loop ;
  PUT(IMP,S,5,5,0) ; NEW_LINE(IMP) ;
end if ;
end ;
CLOSE(IMP) ;
end DEM_DEFL ;

```

### ***6.15 Démonstration de la méthode de Rutishauser***

```

with TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
  MATRIX, DONNEES_ALGEBLIN, SPECTRE ;
use TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
  MATRIX, DONNEES_ALGEBLIN, SPECTRE ;

```

```

procedure DEM_RUTI is
n : NATURAL ;
begin
  MODE_AFFICHAGE ;
  PUT_LINE(IMP,
    "Valeurs propres d'une matrice : Methode de Rutishauser") ;
  NEW_LINE(IMP) ;
  ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Dimension de la matrice : ") ;
  declare
    A : MATRICE(1..n,1..n) ;
    V : VECTEUR(1..n) ;
    S : FLOAT ;
  begin
    DONNEE_MATRICE(A) ;
    S := TRACE(A) ;
    PUT_LINE(IMP,"Valeurs propres de la matrice A = ") ;
    PUT_LINE(A) ; NEW_LINE(IMP) ;
    PAUSE ; CLRSCR ;
    CHRONOMETRE(0) ;
    RUTISHAUSER(A,V) ;
    CHRONOMETRE(1) ; NEW_LINE(IMP) ;
    PUT_LINE(IMP,"Les valeurs propres sont : ") ;
    PUT_LINE(V) ;
    PAUSE ; CLRSCR ;
    PUT_LINE(IMP,"Verification") ;
    NEW_LINE(IMP) ;
    PUT(IMP,"Trace de A = ") ;
    PUT(IMP,S,5,5,0) ; NEW_LINE(IMP) ;
  end ;
end ;

```

```

    PUT(IMP,"Somme des valeurs propres = ") ;
    S := 0.0 ;
    for i in V'range
    loop
        S := S + V(i) ;
    end loop ;
    PUT(IMP,S,5,5,0) ;
    NEW_LINE(IMP) ;
end ;
CLOSE(IMP) ;
end DEM_RUTI ;

```

## 6.16 Démonstration de la méthode de Jacobi

```

with TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
    MATRIX, DONNEES_ALGEBLIN, SPECTRE ;
use TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
    MATRIX, DONNEES_ALGEBLIN, SPECTRE ;

procedure DEM_VPJA is
n : NATURAL ;
begin
    MODE_AFFICHAGE ;
    PUT_LINE(IMP,
        "Valeurs propres et vecteurs propres : Methode de Jacobi") ;
    NEW_LINE(IMP) ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Dimension de la matrice : ") ;
    declare
        A, V : MATRICE(1..n,1..n) ;
        D : VECTEUR(1..n) ;
    begin
        DONNEE_MATRICE(A) ;
        PUT_LINE(IMP,"Spectre de la matrice A = ") ;
        PUT_LINE(A) ;
        NEW_LINE(IMP) ;
        PAUSE ; CLRSCR ;
        CHRONOMETRE(0) ;
        JACOBI(A,D,V) ;
        CHRONOMETRE(1) ;
        NEW_LINE(IMP) ;
        for j in D'range loop
            PUT(IMP,"Valeur propre Numero ") ;
            PUT(IMP,j) ;
            PUT(IMP," = ") ;
            PUT(IMP,D(j),5,6,0) ;
            NEW_LINE(IMP) ;
            PUT_LINE(IMP,"Vecteur propre associe") ;
            PUT(IMP,"(") ;
            for i in V'range(1)
            loop
                PUT(IMP,V(i,j),5,6,0) ;
                if (i < D'LAST) then
                    PUT(IMP,",") ;

```

```

        end if ;
    end loop ;
    PUT_LINE(IMP,"") ;
    if (j mod 5 = 0) then
        PAUSE ;
    end if ;
    NEW_LINE(IMP) ;
end loop ;
end ;
CLOSE(IMP) ;
end DEM_VPJA ;

```

### 6.17 Démonstration de la méthode de Souriau

```

with TEXT_IO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX,
    MATRIX, DONNEES_ALGEBLIN, SPECTRE ;
use TEXT_IO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_MATRIX, MATRIX,
    DONNEES_ALGEBLIN, SPECTRE ;

procedure DEM_SOU is
n : NATURAL ;
begin
    MODE_AFFICHAGE ;
    PUT_LINE(IMP,"Calcul du polynome caracteristique : Methode de
Souriau") ;
    NEW_LINE(IMP) ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Dimension de la matrice : ") ;
    declare
        A, INVERSE : MATRICE(1..n,1..n) ;
        t : VECTEUR(0..n) ;
    begin
        DONNEE_MATRICE(A) ;
        PUT(IMP,"Polynome caracteristique (puissances decroissantes)") ;
        PUT_LINE(IMP," de la matrice A = ") ;
        PUT_LINE(A) ;
        NEW_LINE(IMP) ;
        PAUSE ; CLRSCR ;
        CHRONOMETRE(0) ;
        SOURIAU(A,t,INVERSE) ;
        CHRONOMETRE(1) ;
        PUT(IMP,"P(x) = ") ;
        PUT_LINE(T) ;
        NEW_LINE(IMP) ;
        PAUSE ;
        PUT_LINE(IMP,"Inverse de A : ") ;
        PUT_LINE(INVERSE) ;
    end ;
    CLOSE(IMP) ;
end DEM_SOU ;

```

## CHAPITRE 4

# Résolution numérique des systèmes non linéaires

## 1. Introduction

### 1.1 Position des problèmes

Dans ce chapitre, on s'intéresse au problème de la résolution d'un système d'équations non linéaires :

$$\begin{cases} f_1(x_1, \dots, x_n) = 0 \\ \dots\dots\dots \\ f_n(x_1, \dots, x_n) = 0 \end{cases}$$

où les  $f_j$  sont des fonctions à valeurs réelles définies sur une partie de  $\mathbb{R}^n$ .

Le système précédent sera noté de façon plus compacte :

$$(1) \quad f(x) = 0 \quad (x \in \Omega)$$

où  $\Omega$  est une partie de  $\mathbb{R}^n$  et  $f : \Omega \rightarrow \mathbb{R}^n$ .

On étudiera tout d'abord le cas particulier des équations numériques, c'est-à-dire le cas  $n = 1$ , puis le cas des systèmes d'équations ( $n > 1$ ).

### 1.2 Remarques

(i) Dans le cas des systèmes linéaires ( $f(x) = A \cdot x - b$ ), on a distingué deux types de méthodes : les méthodes directes et les méthodes itératives.

Dans le cas des systèmes non linéaires, il y a peu d'espoir d'aboutir avec des méthodes directes. Les méthodes de résolution seront donc directes, c'est-à-dire qu'on cherchera à construire une suite  $(x^{(k)})_{k \in \mathbb{N}}$  qui converge vers la solution quand cette dernière est unique.

(ii) Deux problèmes importants vont alors se poser :

- (i) le problème de l'existence et de l'unicité d'une solution ;  
 (ii) le problème de la localisation de la racine dans un domaine  $\Omega$  donné. Cette localisation est nécessaire pour démarrer une méthode itérative.

Dans le cas des fonctions continues d'une variable réelle, on dispose du théorème des valeurs intermédiaires qui permet d'apporter une réponse à cette question, mais dans le cas des fonctions de plusieurs variables la situation est plus délicate.

Considérons, par exemple, un système de deux équations à deux inconnues :

$$\begin{cases} f(x, y) = 0 \\ g(x, y) = 0 \end{cases} \quad (x, y) \in \Omega \subset \mathbb{R}^2$$

les solutions sont les points d'intersection des courbes  $f(x, y) = 0$  et  $g(x, y) = 0$ . Mais la connaissance des signes de  $f$  et de  $g$  n'est pas d'un grand secours pour localiser les racines :

Pour  $n \geq 2$  la situation est encore plus compliquée, il s'agit de chercher les points d'intersection de  $n$  hypersurfaces.

## 2. Cas des équations numériques

### 2.1 Notations

On se donne un intervalle réel  $I$ , une fonction  $f : I \rightarrow \mathbb{R}$  et on cherche des solutions de  $f(x) = 0$  sur  $I$ .

Dans le cas général on cherchera une solution réelle dans un intervalle donné. Mais dans le cas particulier des polynômes on pourra chercher toutes les solutions réelles ou complexes.

### 2.2 La méthode dichotomie ou de bisection

Si la fonction  $f : [a, b] \rightarrow \mathbb{R}$  est continue et strictement monotone avec  $f(a) \cdot f(b) < 0$ , le théorème des valeurs intermédiaires nous dit que l'équation  $f(x) = 0$  admet une unique solution dans cet intervalle.

L'idée de la méthode de dichotomie est de découper l'intervalle en deux parties égales et de garder celle qui contient la solution. On répète ensuite le processus.

Avec les hypothèses ci-dessus la méthode converge toujours, mais cette convergence peut être lente. De plus elle ne se généralise pas au cas des systèmes.

L'algorithme est donc le suivant :

$$[a_0, b_0] = [a, b]$$

Si  $[a_k, b_k]$  contient la racine de  $f(x) = 0$ , on pose  $c_k = \frac{a_k + b_k}{2}$  et

$$\text{si } f(a_k) \cdot f(c_k) \leq 0$$

$$\text{alors } [a_{k+1}, b_{k+1}] = [a_k, c_k]$$

$$\text{sinon } [a_{k+1}, b_{k+1}] = [c_k, b_k].$$

Les suites  $(a_k)$  et  $(b_k)$  sont donc adjacentes et convergent vers la solution  $x$  de  $f(x) = 0$ .

De plus, pour les erreurs d'approximation, on a :

$$|x - a_k| \leq \frac{b-a}{2^k}, \quad |x - b_k| \leq \frac{b-a}{2^k} \quad (k \geq 0)$$

ce qui permet de connaître le nombre d'itération suffisant pour atteindre une précision  $\varepsilon$  donnée. Ce nombre est  $k \geq \log_2 \left( \frac{b-a}{\varepsilon} \right)$ .

Le test d'arrêt que nous utiliserons est  $|a_k - b_k| < \varepsilon$ .

La programmation structurée est alors la suivante, en supposant  $f(a) \cdot f(b) \leq 0$ .

*PROCEDURE Dichotomie*(Entrée  $f$ : Fonction ;  $a, b, \varepsilon$ : Réel ; Sortie  $x$ : Réel) ;

*Début*

$y = f(a)$  ;

*Répéter*

$$c = \frac{a + b}{2} ;$$

$z = f(c)$  ;

*Si*  $y \cdot z \leq 0$  *Alors*

$b = c$  ;

*Sinon Début*

$a = c$  ;

$y = z$  ;

*Fin* ;

*Jusqu'à*  $(|b - a| < \varepsilon)$  *ou*  $(z = 0)$  ;

*Remarque* — Pour trouver toutes les racines sur un intervalle  $[a, b]$ , on peut le découper en un nombre suffisamment grand de sous-intervalles  $[a_i, b_i]$  et appliquer la procédure précédente sur chacun de ces intervalles.

### 2.3 La méthode de Newton-Raphson

L'idée de cette méthode est de se ramener à un problème de point fixe qui peut se traiter par la méthode des approximations successives.

Soit  $x$  la solution de  $f(x) = 0$  et  $x_k$  une approximation de  $x$ , on cherche  $h > 0$  tel que  $x_k + h$  soit une meilleure approximation. En supposant que  $f$  est continûment dérivable, on peut écrire :

$$f(x_k + h) = f(x_k) + h \cdot f'(x_k) + O(h^2)$$

et pour  $h$  assez petit,  $f(x_k + h) \cong 0$  entraîne :

$$f(x_k) + h \cdot f'(x_k) \cong 0$$

soit :

$$h = -\frac{f(x_k)}{f'(x_k)}$$

Ces considérations heuristiques nous amènent à considérer la suite définie par :

$$(1) \quad \begin{cases} x_0 \in I \\ x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \end{cases} \quad (k \geq 0)$$



ce qui peut s'illustrer comme ci-dessous :

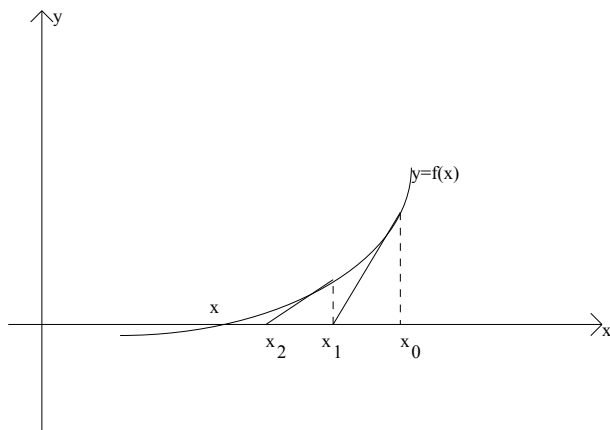


Figure 4.1

Cette méthode s'applique si  $x$  est racine simple de  $f$ , c'est-à-dire si  $f(x) = 0$  et  $f'(x) \neq 0$ . Dans ce cas on aura  $f'(t) \neq 0$  pour  $t$  voisin de la racine  $x$ . De manière plus précise, une condition suffisante de convergence est donnée par le :

*Théorème* : Si  $f : [a,b] \rightarrow \mathbb{R}$  est de classe  $C^2$  avec  $f(a) \cdot f(b) \leq 0$ ,  $f$  strictement monotone (soit  $f'$  de signe constant) et  $f''(x) \neq 0$  sur  $[a,b]$  (concavité dans le même sens), alors pour tout  $x_0$  dans  $[a,b]$  tel que  $f(x_0) \cdot f''(x_0) > 0$ , la suite définie par (1) converge vers l'unique solution de  $f(x) = 0$ .

Une majoration de l'erreur est donnée par :

$$|x_k - x| \leq |x_0 - x|^{2^n} \cdot \left( \frac{M_2}{2 \cdot m_1} \right)^{2^n - 1}$$

où :

$$m_1 = \inf \{ |f'(x)|; x \in [a, b] \}; M_2 = \sup \{ |f''(x)|; x \in [a, b] \}$$

*Démonstration* — Cf. Théodor : P. 26.

*Remarques* — (i) Dans la pratique, la méthode de Newton est très efficace si le point de départ est choisi proche de la solution.

(ii) Cette méthode se généralise au cas des systèmes d'équations.

(iii) L'évaluation de la dérivée peut se faire en utilisant l'approximation par différence finie centrée soit :

$$f'(x) \cong \frac{f(x+h) - f(x-h)}{2 \cdot h}$$

ce qui donne une approximation de l'ordre de  $h^2$ .

(iv) Comme test d'arrêt, on prendra  $|f(x_k)| < \varepsilon$  où  $\varepsilon$  est une précision donnée.

La programmation structurée est alors la suivante :

*PROCEDURE* Newton\_Raphson(Entrée  $f$ : Fonction ;  $x_0$ ,  $\varepsilon$ : Réel ; Sortie  $x$ : Réel) ;

*Début*

$x = x_0$  ;  $k = 0$  ;  $h = 10^{-3}$  ;

*Répéter*  
 $k = k + 1 ;$

$$df = \frac{f(x+h) - f(x-h)}{2 \cdot h};$$

$$\delta = \frac{f(x)}{df};$$

$$x = x - \delta;$$

Jusqu'à ( $k = \text{MaxIter}$ ) ou ( $|\delta| < \epsilon$ );

Fin;

### 3. Cas des équations algébriques

#### 3.1 Introduction

Dans ce paragraphe, on considère le cas où  $f$  est un polynôme de degré  $n$ . On le notera :

$$P(x) = a_0 + a_1 \cdot x + \dots + a_n \cdot x^n$$

où  $a_n \neq 0$ .

Le théorème de D'Alembert-Gauss nous dit que l'équation  $P(x) = 0$  admet  $n$  solutions dans  $\mathbb{C}$ . Mais ce théorème n'est pas constructif.

Pour certaines applications, on sait que toutes les racines de  $P$  sont réelles distinctes et on sait les localiser. C'est le cas pour les polynômes orthogonaux classiques (Cf. § (3.6) du chapitre sur le calcul numérique des intégrales).

Dans ce paragraphe nous décrivons la méthode de Newton-Maehly qui est une adaptation de celle de Newton. Ensuite, après avoir décrit la méthode de Newton pour les systèmes non linéaires, nous l'utiliserons pour décrire la méthode de Bairstow qui permet d'obtenir toutes les racines réelles ou complexes d'un polynôme (Cf. § 5).

#### 3.2 La méthode Newton\_Maehly

On suppose dans ce paragraphe que toutes les racines du polynôme  $P$  sont réelles. C'est le cas pour les polynômes orthogonaux classiques ou pour le polynôme caractéristique d'une matrice symétrique.

La méthode de Newton-Maehly permet d'obtenir toutes les racines de  $P$  avec leur multiplicité.

On note  $\alpha_1, \alpha_2, \dots, \alpha_p$ , les racines de  $P$  de multiplicités respectives  $m_1, m_2, \dots, m_p \geq 1$ .

On a alors le résultat suivant :

*Théorème* : Soit  $(x_k)_{k \in \mathbb{N}}$  la suite définie par :

$$\begin{cases} x_0 > \alpha_p \\ x_{k+1} = x_k - \frac{P(x_k)}{P'(x_k)} \end{cases}$$

est strictement décroissante et  $\lim_{k \rightarrow +\infty} (x_k) = \alpha_p$ . Une majoration de l'erreur étant

$$\text{donnée par :}$$

$$|x_k - \alpha_p| \leq \left(1 - \frac{1}{n}\right)^k \cdot (x_0 - \alpha_p)$$

*Démonstration* — D. Monasse : RMS N° 5, Janvier 1987. Stoer et Burlisch : p. 272.

*Remarque 1* — La convergence est très rapide dans le cas de racines simples.

*Remarque 2* — Pour avoir une méthode convenable, on devra prendre  $x_0 > \alpha_p$ , en étant assez proche de  $\alpha_p$ . Un tel choix est donné par :

$$x_0 = \text{Max} \left\{ \frac{|a_0|}{|a_n|}; 1 + \frac{|a_i|}{|a_n|}; 1 \leq i \leq n-1 \right\}$$

(Cf. Stoer et Burlisch, p. 273).

Le calcul des autres racines peut se faire en appliquant ce qui précède à  $P_1(x) = \frac{P(x)}{(x - \alpha_p)}$  et en remarquant que :

$$\frac{P_1(x)}{P_1'(x)} = \frac{P(x)}{P'(x) - \frac{P(x)}{x - \alpha_p}}$$

La limite de la nouvelle suite  $(x_k)_{k \in \mathbb{N}}$  sera encore  $\alpha_p$  si  $m_p > 1$  et  $\alpha_{p-1}$  si  $m_p = 1$ .

L'algorithme de calcul de toutes les racines est alors le suivant :

*Etape m* — On dispose des racines  $x_1 \geq \dots \geq x_{m-1}$  et la racine  $x_m$  est donnée comme limite de la suite définie par :

$$\begin{cases} x_0^{(m)} = x_0 \\ x_{k+1}^{(m)} = x_k^{(m)} - \frac{P(x_k^{(m)})}{P'(x_k^{(m)}) - P(x_k^{(m)}) \cdot \sum_{i=1}^{m-1} \frac{1}{x_k^{(m)} - x_i}} \end{cases}$$

Comme test d'arrêt, on pourra prendre  $|P(x_k^{(m)})| < \varepsilon$ .

*Remarque 3* — Accélération de la convergence (Cf. Stoer et Burlisch, p. 277).

On considère la suite définie par :

$$\begin{cases} y_0 > \alpha_p \\ y_{k+1} = y_k - 2 \cdot \frac{P(y_k)}{P'(y_k)} \end{cases}$$

et on calcule les  $y_k$  tant que  $y_{k+1} < y_k$  (la suite est théoriquement croissante, mais dès qu'on est très proche de la solution cela peut ne plus être vrai à cause des erreurs d'arrondis), puis dès que  $y_{k0+1} \geq y_{k0}$ , on utilise la suite  $(x_k)$  avec  $x_0 = y_{k0}$ .

On utilise la même idée pour les autres racines.

On a donc en définitive la programmation structurée suivante :

*PROCEDURE* Newton\_Maehly(*Entrée*  $n$  : Entier ;  $P$  : Polynôme ; *Sortie*  $x$  : Vecteur) ;

Début

$$x_0 = \text{Max} \left\{ \frac{|a_0|}{|a_n|} ; 1 + \frac{|a_i|}{|a_n|} ; 1 \leq i \leq n-1 \right\}$$

Pour  $m$  Allant de 1 à  $n$  Faire

Début

$xkP1 = x_0$  ;

Pour  $j$  Allant de 2 à 1 Faire

Début

Répéter

$xk = xkP1$  ;  $U = P(xk)$  ;  $S = 0$  ;

Si  $|U| > \varepsilon$  Alors

Début

$V = P'(xk)$  ;

Pour  $i$  Allant de 1 à  $m-1$  Faire  $S = S + \frac{U}{xk - x_i}$  ;

$V = \frac{U}{V - U \cdot S}$  ;  $xkP1 = xk - j \cdot V$  ;

Fin ;

Jusqu'à ( $xkP1 \geq xk$ ) ou ( $|U| < \varepsilon$ ) ;

Fin ;

$x_m = xk$  ;

Fin ;

Fin ;

## 4. Résolution des systèmes non linéaires par la méthode de Newton-Raphson

### 4.1 Introduction

On suppose ici que  $f : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}$  est de classe  $C^1$  et que l'équation  $f(x) = 0$  admet une solution  $\alpha$  dans  $\mathbb{R}^n$ .

On suppose de plus que, dans un voisinage de  $\alpha$  la *matrice jacobienne* de  $f$  en  $x$  :

$$df(x) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1}(x) & \dots & \frac{\partial f_1}{\partial x_n}(x) \\ \dots & \dots & \dots \\ \frac{\partial f_n}{\partial x_1}(x) & \dots & \frac{\partial f_n}{\partial x_n}(x) \end{pmatrix}$$

est inversible.

## 4.2 Algorithme de Newton-Raphson

La méthode de Newton étudiée dans le cas  $n = 1$  se généralise en considérant la suite  $(x^k)_{k \in \mathbb{N}}$  définie par :

$$(1) \quad \begin{cases} x^{(0)} \in \Omega \\ x^{(k+1)} = x^{(k)} - \delta^{(k)} \quad (k \geq 0) \end{cases}$$

où  $\delta^k$  est solution du système linéaire :  $f(x^{(k)}) \cdot \delta^k = f(x^{(k)})$ .

En général, si  $x^{(0)}$  est choisi assez près de la solution  $\alpha$  alors la méthode converge de façon quadratique. De manière plus précise, on a le

*Théorème* : Si  $f$  est de classe  $C^1$  sur  $\Omega$  convexe avec  $\text{Dét}(df(\alpha)) \neq 0$  et s'il existe une constante  $L > 0$  telle que  $\|df(x) - df(\alpha)\| \leq L \cdot \|x - \alpha\|$  pour tout  $x$  dans  $\Omega$ , alors la suite définie par (1) converge vers la solution  $\alpha$  de  $f(x) = 0$ . De plus la convergence est quadratique, c'est-à-dire que :

$$\|x^{(k+1)} - \alpha\| \leq \text{Cste} \cdot \|x^{(k)} - \alpha\|^2$$

*Démonstration* — Ortega et Rheinboldt : p. 312.

*Remarque* — La convergence est locale, c'est-à-dire que le point de départ de la suite doit être choisi assez proche de la solution.

Par exemple, si on considère le système :

$$\begin{cases} x - \text{tg}(y) = 0 \\ x - y = 0 \end{cases}$$

de solution (4.49, 4.49), en prenant  $(x^{(0)}, y^{(0)}) = (4.1, 4.0)$ , la suite diverge.

## 4.3 Calcul de la matrice jacobienne

Le calcul de la matrice jacobienne en  $x^{(k)}$  se faisant par les *différences finies centrées* ; soit :

$$\frac{\partial f_i}{\partial x_j}(x^{(k)}) \cong \frac{f_i(\beta^{k,j}) - f_i(\gamma^{k,j})}{2 \cdot h}$$

où :  $\beta^{k,j} = (x_1^{(k)}, \dots, x_{j-1}^{(k)}, x_j^{(k)} + h, x_{j+1}^{(k)}, \dots, x_n^{(k)})$ ,  
 $\gamma^{k,j} = (x_1^{(k)}, \dots, x_{j-1}^{(k)}, x_j^{(k)} - h, x_{j+1}^{(k)}, \dots, x_n^{(k)})$  et  $h$  est un pas assez petit.

Comme test d'arrêt dans les itérations, on pourra prendre  $\|f(x^{(k)})\| \leq \varepsilon$ , où  $\varepsilon$  est une précision donnée.

Enfin le système linéaire donnant  $\delta^k$  pourra être résolu par la méthode des pivots de Gauss-Jordan.

## 4.4 Programmation structurée

On écrit tout d'abord une procédure de calcul de la matrice jacobienne.

*PROCEDURE* MatriceJacobienne(Entrée  $n$  : Entier ;  $f$  : Tableau\_Fonctions ;  $x$  : Vecteur ;  
Sortie  $df$  : Matrice) ;

Constante  $h = 10^{-3}$  ;

Début

Pour  $i$  Allant de 1 à  $n$  Faire

Début

Pour  $j$  Allant de 1 à  $n$  Faire

Début

$$x_j = x_j + h ;$$

$$a = f'_1(x) ;$$

$$x_j = x_j - 2 \cdot h ;$$

$$b = f'_1(x) ;$$

$$x_j = x_j + h ;$$

$$df(i, j) = \frac{a - b}{2 \cdot h} ;$$

Fin ;

Fin ;

Fin ;

Puis la procédure principale peut s'écrire :

*PROCEDURE* Newton\_Raphson(Entrée  $n$  : Entier ;  $f$  : Tableau\_Fonctions ;  
Entrée\_Sortie  $x$  : Vecteur) ;

Début

$$k = 0 ;$$

Répéter

$$k = k + 1 ;$$

MatriceJacobienne( $n, f, x, A$ ) ;

$$b = f(x) ;$$

Gauss\_Jordan( $n, A, b$ ) ;

$$x = x - b ;$$

Jusqu'à ( $k = \text{MaxIter}$ ) ou ( $\|f(x)\| < \epsilon$ ) ;

Fin ;

## 5. Racines d'un polynôme. Méthode de Bairstow

### 5.1 Principe de la méthode

L'idée est d'écrire le polynôme  $P(x) = a_0 + a_1 \cdot x + \dots + a_n \cdot x^n$  comme un produit de polynômes de degré au plus 2.

La première étape consiste donc à trouver  $p$  et  $q$  dans  $\mathbb{R}$  tels que :

$$(1) \quad P(x) = Q(x) \cdot (x^2 + p \cdot x + q)$$

On itère ensuite le procédé.

### 5.2 Recherche des coefficients $p$ et $q$

Pour  $p, q$  donnés dans  $\mathbb{R}$ , on peut écrire par division euclidienne :

$$(2) \quad P(x) = B(x) \cdot (x^2 + p \cdot x + q) + R \cdot x + S$$

les coefficients R et S dépendants de p et q. Il s'agit donc de déterminer p et q solutions du système de deux équations à deux inconnues :

$$(3) \quad \begin{cases} R(p, q) = 0 \\ S(p, q) = 0 \end{cases}$$

Le premier problème est donc de calculer R et S en fonction de p, q et des coefficients du polynôme P. Ensuite le système (3) sera résolu par la méthode de Newton-Raphson.

Pour ce faire, on pose :

$$\begin{cases} B(x) = b_0 + b_1 \cdot x + \dots + b_{n-2} \cdot x^{n-2} \\ R(x) = R \cdot x + S = b_{-1} \cdot (x + p) + b_{-2} \end{cases}$$

On va d'abord calculer les  $b_i$  pour  $i = -2, -1, \dots, n-2$ .

Par identification, on a :

$$\begin{cases} a_n = b_{n-2} \\ a_{n-1} = p \cdot b_{n-2} + b_{n-3} \\ a_k = b_{k-2} + p \cdot b_{k-1} + q \cdot b_k \quad (k = n-2, \dots, 0) \end{cases}$$

d'où :

$$(R1) \quad \begin{cases} b_n = b_{n-1} = 0 \\ b_{k-2} = a_k - p \cdot b_{k-1} - q \cdot b_k \quad (k = n, \dots, 0) \end{cases}$$

en particulier on a  $b_{-1}$  et  $b_{-2}$  donc R et S.

Le deuxième problème est de calculer les dérivées partielles de R et S par rapport à p et q.

On a :

$$\begin{cases} R = b_{-1} \\ S = b_{-1} \cdot p + b_{-2} \end{cases}$$

donc :

$$\begin{cases} \frac{\partial R}{\partial p} = \frac{\partial b_{-1}}{\partial p} \\ \frac{\partial R}{\partial q} = \frac{\partial b_{-1}}{\partial q} \\ \frac{\partial S}{\partial p} = b_{-1} + p \cdot \frac{\partial b_{-1}}{\partial p} + \frac{\partial b_{-2}}{\partial p} \\ \frac{\partial S}{\partial q} = p \cdot \frac{\partial b_{-1}}{\partial q} + \frac{\partial b_{-2}}{\partial q} \end{cases}$$



tout revient à calculer les  $\frac{\partial b_k}{\partial p}$  et  $\frac{\partial b_k}{\partial q}$  pour  $k = -2$  et  $k = -1$ . Ce qui se fait en dérivant les relations (R1). On obtient alors la récurrence :

$$\begin{cases} \frac{\partial b_n}{\partial p} = \frac{\partial b_{n-1}}{\partial p} = \frac{\partial b_n}{\partial q} = \frac{\partial b_{n-1}}{\partial q} = 0 \\ \frac{\partial b_{k-2}}{\partial p} = -b_{k-1} - p \cdot \frac{\partial b_{k-1}}{\partial p} - q \cdot \frac{\partial b_k}{\partial p} \quad (k = n, \dots, 0) \\ \frac{\partial b_{k-2}}{\partial q} = -p \cdot \frac{\partial b_{k-1}}{\partial q} - b_k - q \cdot \frac{\partial b_k}{\partial q} \end{cases}$$

En posant :

$$\begin{cases} c_k = -\frac{\partial b_{k-1}}{\partial p} \\ c'_k = -\frac{\partial b_{k-2}}{\partial q} \end{cases}$$

on constate que ces deux suites vérifient la même récurrence :

$$(R2) \begin{cases} u_{n+1} = u_n = 0 \\ u_{k-1} = b_{k-1} + p \cdot u_k + q \cdot u_{k+1} \quad (k = n, \dots, 0) \end{cases}$$

et en particulier  $c_k = c'_k$ .

Le troisième problème est alors de résoudre le système linéaire intervenant dans la méthode de Newton-Raphson.

A l'étape  $k + 1$  du calcul, on a :

$$(p_{k+1}, q_{k+1}) = (p_k, q_k) + (\delta p, \delta q)$$

où  $(\delta p, \delta q)$  est solution du système de deux équations à deux inconnues :

$$d(R, S)(p_k, q_k) \cdot (\delta p, \delta q) = -(R(p_k, q_k), S(p_k, q_k))$$

soit :

$$\begin{cases} \frac{\partial b_{-1}}{\partial p} \cdot \delta p + \frac{\partial b_{-1}}{\partial q} \cdot \delta q = -b_{-1} \\ \left( b_{-1} + p_k \cdot \frac{\partial b_{-1}}{\partial p} + \frac{\partial b_{-2}}{\partial p} \right) \cdot \delta p + \left( p_k \cdot \frac{\partial b_{-1}}{\partial q} + \frac{\partial b_{-2}}{\partial q} \right) \cdot \delta q = -(b_{-1} \cdot p_k + b_{-2}) \end{cases}$$

avec les notations précédentes, ce système va s'écrire :

$$\begin{cases} c_0 \cdot \delta p + c_1 \cdot \delta q = b_{-1} \\ (-b_{-1} + p_k \cdot c_0 + c_{-1}) \cdot \delta p + (p_k \cdot c_1 + c_0) \cdot \delta q = (b_{-1} \cdot p_k + b_{-2}) \end{cases}$$

la solution est :

$$\begin{cases} \delta p = \frac{U}{D} \\ \delta q = \frac{V}{D} \end{cases}$$

où on a noté :

$$\begin{cases} D = c_0^2 - c_1 \cdot c_{-1} + c_1 \cdot b_{-1} \\ U = c_0 \cdot b_{-1} - c_1 \cdot b_{-2} \\ V = c_0 \cdot b_{-2} + b_{-1}^2 - b_{-1} \cdot c_{-1} \end{cases}$$

### 5.3 Programmation structurée

On écrit tout d'abord la procédure de division qui donne les coefficients  $p$  et  $q$  et le polynôme  $B$  tels que  $P(x) = Q(x) \cdot (x^2 + p \cdot x + q)$ .

*PROCEDURE Division(Entrée  $n$  : Entier ;  $P$  : Polynôme ; Sortie  $B$  : Vecteur ;  $p, q$  : Réels) ;*

*Début*

*Se donner  $p_0$  et  $q_0$  ;*

*$k = 0$  ;*

*$p = p_0$  ;  $q = q_0$  ;*

*Répéter*

*Calcule\_Coeff\_bk( $n, P, p, q, B$ ) ; { Récurrence R1 }*

*Calcule\_Coeff\_ck( $n, B, p, q, C$ ) ; { Récurrence R2 }*

*$D = c_0^2 - c_1 \cdot c_{-1} + c_1 \cdot b_{-1}$  ;*

*$U = c_0 \cdot b_{-1} - c_1 \cdot b_{-2}$  ;*

*$V = c_0 \cdot b_{-2} + b_{-1}^2 - b_{-1} \cdot c_{-1}$  ;*

*$\delta p = \frac{U}{D}$  ;  $\delta q = \frac{V}{D}$  ;*

*$p = p + \delta p$  ;  $q = q + \delta q$  ;*

*Jusqu'à ( $|\delta p| < \varepsilon$ ) et ( $|\delta q| < \varepsilon$ ) ;*

*Fin ;*

Les procédures de calcul des coefficients  $b_k$  et  $c_k$  s'écrivent :

*PROCEDURE Calcule\_Coeff\_bk(Entrée  $n$  : Entier ;  $P$  : Polynôme ;  $p, q$  : Réel ;  
Sortie  $B$  : Vecteur) ;*

*Début*

*$b_n = 0$  ;  $b_{n-1} = 0$  ;*

*Pour  $k$  allant de  $n$  à 0 faire  $b_{k-2} = a_k - p \cdot b_{k-1} - q \cdot b_k$  ;*

*Fin ;*

*PROCEDURE Calcule\_Coeff\_ck(Entrée  $n$  : Entier ;  $B$  : Vecteur ;  $p, q$  : Réel ;  
Sortie  $C$  : Vecteur) ;*

*Début*

*$c_{n+1} = 0$  ;  $c_n = 0$  ;*

*Pour  $k$  allant de  $n$  à 0 faire  $c_{k-1} = b_{k-1} + p \cdot c_k + q \cdot c_{k+1}$  ;*

*Fin ;*

La procédure principale peut alors s'écrire :

*PROCEDURE Bairstow(Entrée  $n$  : Entier ;  $P$  : Polynôme ; Sortie  $x$  : VecteurComplexe) ;*

*Début*

*Pour  $i$  allant de 0 à  $n$  faire  $a_i = a_i / a_n$  ;*

*$k = 0$  ;*

*Répéter*

*Division( $n, P, B, p, q$ ) ;*

$n = n - 2 ;$   
*Pour  $i$  allant de 0 à  $n$  faire  $a_i = b_i ;$*   
*RésolEquDegre2( $p, q, x1, x2$ ) ;*  
 $k = k + 1 ; x_k = x1 ;$   
 $k = k + 1 ; x_k = x2 ;$

Jusqu'à ( $n \leq 2$ ) ;

Si  $n = 2$  Alors

Début

$RésolEquDegre2(a_1, a_0, x_1, x_2)$  ;

$k = k + 1 ; x_k = x_1$  ;

$k = k + 1 ; x_k = x_2$  ;

Fin ;

Si  $n = 1$  Alors

Début

$k = k + 1 ; x_k = -a_0/a_1$  ;

Fin ;

Fin ;

## 6. Exercices

### 6.1 Méthode des approximations successives

On veut résoudre l'équation  $g(x) = 0$  sur  $[a, b]$ , où  $g$  est continûment dérivable et admet une unique racine dans  $[a, b]$ . Cette équation équivaut à  $f_\alpha(x) = x - \alpha \cdot g(x) = 0$ , où  $\alpha$  est un réel non nul à choisir. Pour approcher la solution de cette dernière équation, on considère la suite des approximations successives définie par :

$$\begin{cases} x_0 \in [a, b] \\ x_{n+1} = f_\alpha(x_n) \quad (n \geq 0) \end{cases}$$

1°) Montrer que si  $|f'_\alpha(x)| < 1$  sur  $[a, b]$ , alors la suite  $(x_k)_{k \in \mathbb{N}}$  converge vers l'unique solution de  $f(x) = 0$ .

2°) On pose  $m = \inf\{g'(x); x \in [a, b]\}$ ,  $M = \sup\{g'(x); x \in [a, b]\}$  et on suppose que  $m > 0$ .

(a) Déterminer  $\alpha$  tel que  $|f'_\alpha(x)| \leq \frac{M-m}{M+m}$  sur  $[a, b]$ .

(b) En déduire que pour un tel choix de  $\alpha$  la méthode est convergente.

3°) Soit à résoudre  $g(x) = \tan(x) - x = 0$  sur  $\left] \frac{3 \cdot \pi}{2}, \frac{5 \cdot \pi}{2} \right[$ . On pose  $a = 7.65$ ,

$b = 7.75$ .

(a) Montrer que  $g$  admet une unique racine dans  $[a, b]$ .

(b) Déterminer le coefficient  $\alpha$  du 2°) (a).

(c) Décrire un algorithme correspondant à la méthode décrite ci-dessus.

### 6.2 Méthode de Bernoulli

Soit  $P(x) = x^n + a_{n-1} \cdot x^{n-1} + \dots + a_1 \cdot x + a_0$  un polynôme à coefficients réels admettant une unique racine de module maximal  $\lambda_1$ .

On définit la suite  $(x_k)_{k \in \mathbb{N}}$  par :

$$x_k + a_{n-1} \cdot x_{k-1} + \dots + a_0 \cdot x_{k-n} = 0 \quad (k \geq n)$$

les valeurs initiales  $x_0, \dots, x_{n-1}$  étant arbitraires.

1°) Montrer que  $\lim_{k \rightarrow +\infty} \left( \frac{x_{k+1}}{x_k} \right) = \lambda_1$ .

2°) Montrer que  $\lim_{k \rightarrow +\infty} \left( \frac{\frac{x_{k+1} - \lambda_1}{x_k}}{\frac{x_{k1} - \lambda_1}{x_{k-1}}} \right) = \frac{\lambda_2}{\lambda_1}$ .

3°) En déduire l'évaluation de l'erreur :  $\left| \frac{x_{k+1}}{x_k} - \lambda_1 \right| \cong \left| \frac{\lambda_2}{\lambda_1} \right|^n$ .

4°) Comment procéder pour déterminer les autres racines (méthode de déflation).

### 6.3 Méthode de la fausse position

Décrire la méthode obtenue, si dans la méthode de Newton-Raphson on utilise comme approximation de la dérivée en  $x_k$  la quantité  $f'(x_k) \cong \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}$ .

Donner une interprétation géométrique.

### 6.4 Méthode de Muller

Dans la méthode de la fausse position, on fait une interpolation de  $y = f(x)$  par la droite passant par  $(x_{k-1}, f(x_{k-1}))$  et  $(x_k, f(x_k))$ .

En faisant une interpolation par un polynôme de degré 2 passant par  $(x_{k-2}, f(x_{k-2}))$ ,  $(x_{k-1}, f(x_{k-1}))$  et  $(x_k, f(x_k))$  et en prenant comme point d'intersection des deux courbes celui qui donne la valeur la plus proche de  $x_k$ , donner la valeur du point  $x_{k+1}$  obtenu.

## 7. Programmation Ada

### 7.1 Spécification du paquetage EQUATIONS\_GENERIQUE

```
generic
  with function f(x : in FLOAT) return FLOAT ;
package EQUATIONS_GENERIQUE is
  procedure DICHOTOMIE(a, b : in FLOAT ; x : out FLOAT ;
    Eps : in FLOAT := 1.0E-6) ; -- § (2.2)
  procedure NEWTON_RAPHSON(x0 : in FLOAT ; x : out FLOAT ;
    Eps : in FLOAT := 1.0E-6) ; -- § (2.3)
  ERREUR_CONVERGENCE : exception ;
end EQUATIONS_GENERIQUE ;
```

### 7.2 Démonstration du paquetage EQUATIONS\_GENERIQUE

```
with TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH1, MATH2,
  EQUATIONS_GENERIQUE ;
use TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH1, MATH2 ;
```

```

procedure DEM_EQU is

FF : FONCTION ;
Choix : INTEGER range 0..2 ;

function f(x : in FLOAT) return FLOAT is
begin
    return EVALUE(FF,x) ;
end f ;

package EQUATIONS is new EQUATIONS_GENERIQUE(f) ;
use EQUATIONS ;

procedure DEMO_DICHOTOMIE is
a, b, ak, bk, x, y, z, h : FLOAT ;
n, k : INTEGER ;
begin
    MODE_AFFICHAGE ;
    PUT_LINE(
        "RACINES DE f(x) = 0 SUR [a,b] PAR LA METHODE DE DICHOTOMIE." ) ;
    NEW_LINE ;
    PUT_LINE("L'intervalle [a,b] est decoupe en n sous-intervalle") ;
    PUT_LINE(
        "et on cherche une racine dans chaque [ai,bi] si f(ai)*f(bi) <= 0" ) ;
    NEW_LINE ;
    GET_LINE(FF,"Entrez l'equation : f(x) = " ) ; NEW_LINE ;
    ENTRER_REEL(a,"Valeur de a : " ) ;
    ENTRER_REEL_BORNE(a,1.0E+6,b,"Valeur de b : " ) ; NEW_LINE ;
    ENTRER_ENTIER(n,"Nombre de sous-intervalles : " ) ; NEW_LINE ;
    h := (b - a)/FLOAT(n) ;
    ak := a ;
    k := 0 ;
    for i in 0..n - 1
    loop
        y := f(ak) ;
        z := f(ak + h) ;
        if (y = 0.0) then
            k := k + 1 ;
            PUT(IMP,"Solution ") ; PUT(IMP,k,2) ; PUT(IMP," = ") ;
            PUT(IMP,ak,5,6,0) ; NEW_LINE(IMP) ;
        end if ;
        if (y*z < 0.0) then
            DICHOTOMIE(ak,ak + h,x) ;
            k := k + 1 ;
            PUT(IMP,"Solution ") ; PUT(IMP,k,2) ; PUT(IMP," = ") ;
            PUT(IMP,x,5,6,0) ; NEW_LINE(IMP) ;
        end if ;
        ak := ak + h ;
    end loop ;
    CLOSE(IMP) ;
end DEMO_DICHOTOMIE ;

```

```

procedure DEMO_NEWTON_RAPHSON is
x0, x: FLOAT ;
k : INTEGER ;
begin
  MODE_AFFICHAGE ;
  PUT_LINE("RACINE DE f(x) = 0 [a,b] PAR LA METHODE DE NEWTON-RAPHSON.");
  NEW_LINE ;
  GET_LINE(FF,"Entrez l'equation : f(x) = ") ;
  NEW_LINE ;
  ENTRER_REEL(x0,"Valeur initiale x0 : ") ;
  NEW_LINE ;
  NEWTON_RAPHSON(x0,x) ;
  PUT("Solution x = ") ; PUT(x,5,6) ;
  NEW_LINE ;
  CLOSE(IMP) ;
end DEMO_NEWTON_RAPHSON ;

begin
  loop
    CLRSCR ;
    PUT_LINE("-0- FIN") ;
    PUT_LINE("-1- Recherches des racines de f(x) = 0 par dichotomie") ;
    PUT_LINE(
"-2- Recherche d'une racine de f(x) = 0 par la methode de Newton") ;
    NEW_LINE ;
    ENTRER_ENTIER_BORNE(0,2,Choix,"Votre choix : ") ;
    case Choix is
      when 0 => exit ;
      when 1 => DEMO_DICHOTOMIE ; PAUSE ;
      when 2 => DEMO_NEWTON_RAPHSON ; PAUSE ;
    end case ;
  end loop ;
end DEM_EQU ;

```

### 7.3 Spécification du paquetage **SYSTEME\_EQUATIONS\_GENERIQUE**

```

with COMMON_MATRIX, MATH2 ;
use COMMON_MATRIX, MATH2 ;

generic
  with function f(j : in POSITIVE ; x : in VARIABLE_2) return FLOAT ;

package SYSTEME_EQUATIONS_GENERIQUE is

  procedure MATRICE_JACOBIENNE(x : in VARIABLE_2 ; df : out MATRICE) ;
    -- § (4.3)
  procedure NEWTON_RAPHSON(x : in out VARIABLE_2) ; -- § (4.2)

  ERREUR_CONVERGENCE : exception ;

end SYSTEME_EQUATIONS_GENERIQUE ;

```

## 7.4 Démonstration du paquetage *SYSTEME\_EQUATIONS\_GENERIQUE*

```

with TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH2,
     SYSTEME_EQUATIONS_GENERIQUE ;
use TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH2 ;

procedure DEM_SYST is
MAX_FONCTIONS : constant INTEGER := 10 ;
subtype INDICE_FCT is INTEGER range 1..MAX_FONCTIONS ;
type SYSTEME_EQUATIONS is array(INDICE_FCT range <>) of FONCTION ;
n : INTEGER ;
begin
  MODE_AFFICHAGE ;
  PUT_LINE ("          SYSTEME DE N EQUATIONS NON LINEAIRES") ;
  NEW_LINE ;
  PUT_LINE ("          METHODE DE NEWTON_RAPHSON") ;
  NEW_LINE ;
  ENTRER_ENTIER_BORNE(1,MAX_FONCTIONS,n,"Nombre d'equations : ") ;
  NEW_LINE ;
  declare
    FF : SYSTEME_EQUATIONS(1..n) ;
    x : VARIABLE_2(1..n) ;
    function f(j : in POSITIVE ; x : in VARIABLE_2) return FLOAT is
    begin
      return EVALUE(FF(j),0.0,x) ;
    end f ;
    package SYSTEME_FF is new SYSTEME_EQUATIONS_GENERIQUE(F) ;
    use SYSTEME_FF ;
  begin
    for j in x'range
    loop
      PUT("Equation Numero ") ; PUT(j,2) ; NEW_LINE ;
      PUT("f(") ;
      for k in 1..n - 1
      loop
        PUT("y") ; PUT(k,2) ; PUT(",") ;
      end loop ;
      PUT("y") ; PUT(n,2) ; PUT_LINE(") = ") ;
      GET_LINE(ff(j)," ") ;
    end loop ;
    NEW_LINE ;
    PUT_LINE(" Valeur initiale x") ;
    NEW_LINE ;
    for j in x'range
    loop
      PUT("Valeur initiale x") ; PUT(j,2) ; NEW_LINE ;
      ENTRER_REEL(x(j)," x = ") ;
    end loop ;
    NEW_LINE ;
    NEWTON_RAPHSON(x) ;
    PUT_LINE("Solution") ;
    for i in x'range

```



```

    loop
      PUT(IMP,"x") ; PUT(IMP,i,2) ; PUT(IMP," = ") ;
      PUT(IMP,x(i),5,6,0) ;
      NEW_LINE(IMP) ;
    end loop ;
  end ;
  CLOSE(IMP) ;
end DEM_SYST ;

```

### 7.5 Spécifications des paquetages *COMMON\_POLY* et *POLY*

Dans ces paquetages on définit un type POLYNOME avec des procédures d'entrée-sortie, une fonction d'évaluation (algorithme de Horner), une procédure de tracé d'une courbe polynomiale et des procédures de recherche des racines d'un polynôme.

On utilise le paquetage NOMBRES\_COMPLEXES de spécification :

```

with TEXT_IO ;
use TEXT_IO ;
package NOMBRES_COMPLEXES is
type INDICE is (REEL, IMAGINAIRE) ;
type COMPLEXE is array(INDICE) of FLOAT ;

function "+"(A, B : COMPLEXE) return COMPLEXE ;
function "-"(A, B : COMPLEXE) return COMPLEXE ;
function "*" (A, B : COMPLEXE) return COMPLEXE ;
function "*" (A : FLOAT ; B : COMPLEXE) return COMPLEXE ;
function "/"(A, B : COMPLEXE) return COMPLEXE ;
function MODULE(A : COMPLEXE) return FLOAT ;
procedure PUT(z : in COMPLEXE ; FORE : in INTEGER := 3 ;
  AFT : in INTEGER := 3 ; EXP : in INTEGER := 0) ;
procedure PUT(FILE : in FILE_TYPE ; z : in COMPLEXE ;
  FORE : in INTEGER := 3 ; AFT : in INTEGER := 3 ;
  EXP : in INTEGER := 0) ;
end NOMBRES_COMPLEXES ;

package COMMON_POLY is
MAX_DEGRE : constant INTEGER := 30 ;
subtype INDICE_COEFFICIENTS is INTEGER range 0..MAX_DEGRE ;
type POLYNOME is array(INDICE_COEFFICIENTS range <>) of FLOAT ;
end COMMON_POLY ;

with TEXT_IO, COMMON_POLY, COMMON_MATRIX, NOMBRES_COMPLEXES ;
use COMMON_POLY, COMMON_MATRIX, NOMBRES_COMPLEXES ;
package POLY is
type VECTEUR_COMPLEXE is array(INDICE_COEFFICIENTS range <>)
  of COMPLEXE ;

procedure GET(P : out POLYNOME) ;
procedure POLYNOME_ALEATOIRE(P : out POLYNOME) ;
procedure PUT_LINE(P : in POLYNOME ; FORE : in TEXT_IO.FIELD := 6 ;
  AFT : in TEXT_IO.FIELD := 4 ;

```

```

        EXP : in TEXT_IO.FIELD := 0) ;
function EVALUE(P : in POLYNOME ; x : in FLOAT) return FLOAT ;
function EVALUE_DERIVEE(P : in POLYNOME ; x : in FLOAT) return FLOAT ;
procedure TRACE_POLYNOME(P : in POLYNOME ; xMin, xMax : in FLOAT) ;
procedure NEWTON_MAEHLY(P : in POLYNOME ; x : in out VECTEUR) ;
        -- § (3.2)
procedure BAIRSTOW(P : in POLYNOME ; x : out VECTEUR_COMPLEXE) ; -- § 5
ERREUR_CONVERGENCE : exception ;
end POLY ;

```

## 7.6 Démonstration du paquetage POLY

```

with TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_POLY, POLY,
     NOMBRES_COMPLEXES, COMMON_MATRIX ;
use TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, COMMON_POLY, POLY,
     NOMBRES_COMPLEXES, COMMON_MATRIX ;

procedure DEM_POL is
Choix : INTEGER range 0..2 ;
procedure DEMO_NEWTON_MAEHLY is
n : INDICE_COEFFICIENTS ;
begin
    MODE_AFFICHAGE ;
    PUT_LINE("Calcul de toutes les racines reelles d'un polynome") ;
    PUT_LINE("Methode de Newton-Maehly") ;
    PUT_LINE(
        "On suppose que toutes les racines du polynome sont reelles") ;
    NEW_LINE ;
    ENTRER_ENTIER_BORNE(1,MAX_DEGRE,n,"Degre du polynome : ") ;
    declare
        P : POLYNOME(0..n) ;
        x : VECTEUR(1..n) ;
    begin
        PUT_LINE("Entrer le polynome P") ;
        GET(P) ;
        NEW_LINE ;
        NEWTON_MAEHLY(P,x) ;
        for i in x'range
        loop
            PUT(IMP,"Solution ") ; PUT(IMP,i,3) ; PUT(IMP," = ") ;
            PUT(IMP,x(i),5,6,0) ; NEW_LINE(IMP) ;
        end loop ;
    end ;
    CLOSE(IMP) ;
end DEMO_NEWTON_MAEHLY ;

procedure DEMO_BAIRSTOW is
n : INDICE_COEFFICIENTS ;
begin
    MODE_AFFICHAGE ;
    PUT_LINE("Calcul de toutes les racines complexes d'un polynome") ;
    PUT_LINE("Methode de Bairstow") ;
    NEW_LINE ;

```

```

ENTRER_ENTIER_BORNE(1,MAX_DEGRE,n,"Degre du polynome : ") ;
declare
  P : POLYNOME(0..n) ;
  x : VECTEUR_COMPLEXE(1..n) ;
begin
  PUT_LINE("Entrer le polynome P") ;
  GET(P) ;
  NEW_LINE ;
  BAIRSTOW(P,x) ;
  for i in x'range
  loop
    PUT(IMP,"Solution ") ; PUT(IMP,i,3) ; PUT(IMP," = ") ;
    PUT(IMP,x(i),5,6,0) ; NEW_LINE(IMP) ;
  end loop ;
end ;
CLOSE(IMP) ;
end DEMO_BAIRSTOW ;

begin
  loop
    CLRSCR ;
    PUT_LINE("-0- FIN") ;
    PUT_LINE("-1- Recherches des racines reelles d'un polynome.") ;
    PUT_LINE("    Methode de Newton-Maehly") ;
    PUT_LINE("-2- Recherche des racines complexes d'un polynome") ;
    PUT_LINE("    Methode de Bairstow") ;
    NEW_LINE ;
    ENTRER_ENTIER_BORNE(0,2,Choix,"Votre choix : ") ;
    case Choix is
      when 0 => exit ;
      when 1 => DEMO_NEWTON_MAEHLY ; PAUSE ;
      when 2 => DEMO_BAIRSTOW ; PAUSE ;
    end case ;
  end loop ;
end DEM_POL ;

```

## CHAPITRE 5

# Approximation et interpolation

### 1. Introduction

Etant donné un nuage de points,  $S = \{ (x_i, y_i) ; 1 \leq i \leq n \}$ , on cherche une courbe qui passe, « au mieux », par ces points.

On peut demander à cette courbe de passer exactement par ces points, (c'est le cas, par exemple, pour l'interpolation polynomiale de Lagrange ou l'interpolation spline) ou alors on peut demander à cette courbe d'approcher, dans un sens à préciser, le nuage de points, (c'est le cas, par exemple, pour les méthodes des moindres carrés ou pour l'approximation polynomiale uniforme).

Dans un premier temps on étudiera dans ce chapitre les problèmes d'approximation avec les méthodes de moindres carrés, les approximations de Bézier et B-Splines. Puis on s'intéressera aux problèmes d'interpolation avec l'interpolation polynomiale de Lagrange et l'interpolation spline cubique.

Les idées sur l'interpolation polynomiale seront exploitées dans le chapitre sur le calcul numérique des intégrales.

### 2. Problèmes d'approximation.

#### Méthode des moindres carrés

#### 2.1 Introduction

Dans le cas de nuages de points expérimentaux, on n'a aucune raison d'exiger que la courbe approximative passe exactement par ces points. En effet, on peut savoir que la courbe en question sera une droite, un cercle, ... et à cause des erreurs de mesures, aucune courbe de ce type ne pourra passer par tous les points du nuage.

Un modèle de fonction étant choisi, dépendant d'un nombre fini de paramètres, il s'agira alors de déterminer ces derniers. Usuellement cette détermination se ramène à calculer le minimum d'une fonction de plusieurs variables : il s'agira de

minimiser la distance du nuage de points à la courbe définie par le modèle de fonction choisi.

Dans certains cas, comme pour les exemples du paragraphe (2.2), on a une idée du modèle qui va décrire la situation. Mais le problème est beaucoup plus délicat si on n'a aucune idée a priori de la loi décrivant le phénomène étudié.

Cet aspect des choses a été étudié par *Rufener* en 1945 dans son livre « *Mise en équation des résultats d'expériences* ».

L'idée de *Rufener* est, à partir d'une famille de fonctions données a priori, de rechercher, pour un nuage de points donné, le modèle qui décrira au mieux la situation. Cette recherche se faisant à partir d'une série de critères et en utilisant des tests statistiques. Une fois le modèle trouvé il s'agit de déterminer les coefficients qui donneront la courbe de meilleure approximation.

## 2.2 Exemples

Dans les exemples qui suivent des lois de la physique nous disent quel est le modèle à rechercher.

### 2.2.1 La loi d'Ohm

Si  $x$  est une intensité aux bornes d'un circuit et  $y$  la tension correspondante, on aura une relation du type  $y = R \cdot x$ , où la résistance  $R$  est à déterminer. Le nuage de points est alors une série de mesures  $(x_i, y_i)$  des intensités et tensions.

### 2.2.2 Elongation d'un ressort

Si  $x$  représente le temps et  $y$  l'élongation d'un ressort à l'instant  $x$ , on aura une relation du type  $y = A \cdot \cos(\omega \cdot t + \phi)$  où  $A$ ,  $\omega$  et  $\phi$  sont à déterminer.

### 2.2.3 Problème du fil chaud

Pour mesurer la vitesse d'écoulement d'un fluide, on utilise la « méthode du fil chaud » due à *L. V. King* en 1914.

La loi de *King*, qui relie la tension (en Volts) appliquée au fil à la vitesse d'écoulement du fluide (en m/s), est donnée par :

$$T^2 = a \cdot \sqrt{v} + b$$

A partir de mesures expérimentales de  $v$  et  $T$ , on peut déduire les coefficients  $a$  et  $b$  en se ramenant à un modèle affine. En effet, en posant  $y = T^2$  et  $x = \sqrt{v}$ , la loi de *King* s'écrit  $y = a \cdot x + b$ .

### 2.2.4 Corrélation statistique

Il se peut aussi qu'il n'y ait aucune relation déterministe entre  $x$  et  $y$ , mais seulement une corrélation statistique.

Par exemple, en économie, si  $x$  représente une année et  $y$  le nombre de Kw/h consommés cette année-là, il n'y aura aucune relation simple entre  $x$  et  $y$ .

## 2.3 Détermination des paramètres

Pour chaque problème, on commence d'abord par choisir le modèle qui doit rendre compte des observations, ce qui mettra en évidence les paramètres à déterminer.

Un tel modèle de fonction sera noté  $y = f(x, a_1, \dots, a_p)$ , où les  $a_i$  sont les paramètres inconnus.

Le modèle étant choisi, on définit le *résidu* en  $x_i$ , comme l'écart entre l'observation et le modèle, soit :

$$e_i = \left| y_i - f(x_i, a_1, \dots, a_p) \right|, \text{ pour } i = 1, \dots, n.$$

L'idée est alors de trouver les paramètres qui vont minimiser ces résidus. Pour ce faire, on doit travailler de façon globale, c'est-à-dire définir un écart entre le nuage de points S et la fonction f choisie comme modèle. Cet écart doit tenir compte de tous les  $e_i$ . Les choix les plus simples sont :

$$E_\infty = \text{Max}\{|e_i|; 1 \leq i \leq n\}$$

$$E_1 = \sum_{i=1}^n |e_i|$$

$$E_2 = \left( \sum_{i=1}^n e_i^2 \right)^{\frac{1}{2}}$$

qui correspondent aux trois normes classiques de  $\mathbb{R}^n$ .

Le problème étant de minimiser l'écart global E, la première idée est d'utiliser les dérivées partielles, puisqu'une condition nécessaire d'extremum pour une fonction dérivable est la nullité de ces dernières. Cette idée va exclure les quantités  $E_\infty$  et  $E_1$  qui ont peu de chances d'être dérivables, ce qui nous conduit à utiliser  $E_2$ .

Ce choix de  $E_2$  est donc lié à des considérations pratiques et non à des considérations théoriques.

## 2.4 Principes des méthodes des moindres carrés

L'idée des méthodes de moindres carrés est de trouver les paramètres  $a_i$  qui minimisent *l'écart quadratique* :

$$E(a_1, \dots, a_p) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n \left( y_i - f(x_i, a_1, \dots, a_p) \right)^2$$

Si f est dérivable, une condition nécessaire d'extremum est donc donnée par le système de p équations à p inconnues :

$$\frac{\partial E}{\partial a_i}(a_1, \dots, a_p) = 0 \quad (i = 1, \dots, p) \quad (1)$$

Mais il se peut que ce système n'ait pas de solution ou qu'il ait une solution qui ne corresponde pas à un extremum pour E.

Le cas agréable sera celui où on a un système linéaire. Les dérivées partielles donneront des expressions linéaires si E est une expression quadratique des  $a_i$ , ce qui sera réalisé si les  $e_i$  sont des expressions linéaires des  $a_i$ . On parlera alors de modèles linéaires.

De manière plus générale, on peut être amené à introduire des coefficients de pondération liés aux incertitudes sur les mesures des  $x_i$  et des  $y_i$ , ce qui conduit à minimiser :

$$E = \sum_{i=1}^n w_i \cdot e_i^2$$

## 2.5 Les modèles linéaires

### 2.5.1 Généralités

Dans ces cas, on cherche une fonction  $y$  qui est combinaison linéaire de fonctions  $y_1, \dots, y_p$  données, soit :

$$y(x) = \sum_{k=1}^p a_k \cdot y_k(x)$$

Les  $y_k$  pouvant être des monômes ( $y$  est alors un polynôme de degré  $p$ ), des sinus et cosinus ( $y$  est alors une série harmonique), ...

### 2.5.2 La régression affine

On cherche, ici, le modèle sous la forme :

$$y(x) = f(x, a, b) = a \cdot x + b$$

avec le nuage de points  $S = \{ (x_i, y_i) ; 1 \leq i \leq n \}$ .

L'un des tests de Ruffner, pour rechercher un tel modèle, est de dire que nécessairement les quotients  $\frac{y_{i+1} - y_i}{x_{i+1} - x_i}$  devraient être sensiblement constants.

*Remarque* — Si les quotients ci-dessus sont très grands, il sera plus judicieux de chercher une droite presque verticale, c'est-à-dire chercher un modèle du type  $x = \alpha \cdot y + \beta$ .

Il s'agit donc de minimiser :

$$E(a, b) = \sum_{i=1}^n (a \cdot x_i + b - y_i)^2$$

Une condition nécessaire d'extremum est alors :

$$\begin{cases} \frac{\partial E}{\partial a} = 0 \\ \frac{\partial E}{\partial b} = 0 \end{cases} \quad (2)$$

Soit le système :

$$\begin{cases} \sum_{i=1}^n (a \cdot x_i + b - y_i) \cdot x_i = 0 \\ \sum_{i=1}^n (a \cdot x_i + b - y_i) = 0 \end{cases} \quad (3)$$

ou encore :

$$\begin{cases} \left( \sum_{i=1}^n x_i^2 \right) \cdot a + \left( \sum_{i=1}^n x_i \right) \cdot b = \sum_{i=1}^n x_i \cdot y_i \\ \left( \sum_{i=1}^n x_i \right) \cdot a + n \cdot b = \sum_{i=1}^n y_i \end{cases} \quad (4)$$



En notant  $X$  la variable aléatoire prenant les valeurs  $x_i$  de manière équiprobable et  $Y$  celle prenant les valeurs  $y_i$ , le système (4) s'écrit :

$$\begin{cases} E(X^2) \cdot a + E(X) \cdot b = E(X \cdot Y) \\ E(X) \cdot a + b = E(Y) \end{cases} \quad (5)$$

où on a noté :

$$E(X^k) = \frac{1}{n} \cdot \sum_{i=1}^n x_i^k \quad (\text{moment d'ordre } k)$$

$$E(X \cdot Y) = \frac{1}{n} \cdot \sum_{i=1}^n x_i \cdot y_i$$

Le déterminant de ce système est :

$$V(X) = E(X^2) - E(X)^2 = E((X - E(X))^2) \quad (\text{variance de } X)$$

Si  $V(X) = 0$ , tous les  $x_i$  valent  $E(X)$  et les points sont sur une droite verticale.

En général on a donc  $V(X) \neq 0$  et la solution de (5) est :

$$\begin{cases} a = \frac{\sigma_{XY}}{\sigma_X^2} \\ b = E(Y) - a \cdot E(X) \end{cases} \quad (6)$$

où on a noté :

$$V(X) = \sigma_X^2 \quad (\text{écart type de } X)$$

$$\sigma_{XY} = E(X \cdot Y) - E(X) \cdot E(Y) = \frac{1}{n} \cdot \sum_{i=1}^n (x_i - E(X)) \cdot (y_i - E(Y))$$

(Covariance des séries statistiques  $X$  et  $Y$ )

*Remarque* — On a  $E(a', b') = E(a, b) + \sum_{i=1}^n ((a' - a) \cdot x_i + (b' - b))^2$ , pour tous  $a', b'$

dans  $R$  et on a donc bien un minimum.

En considérant le modèle  $x = \alpha \cdot y + \beta$ , on aurait de la même façon :

$$\begin{cases} \alpha = \frac{\sigma_{XY}}{\sigma_Y^2} \\ \beta = E(X) - \alpha \cdot E(Y) \end{cases} \quad (7)$$

Les points seront effectivement alignés si ces deux droites sont identiques c'est-à-dire si, et seulement si, elles ont la même pente, ce qui se traduit par  $\alpha =$

$\frac{1}{a}$ , ou encore :

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \cdot \sigma_Y} = 1$$

$\rho_{XY}$  est le *coefficient de corrélation* de  $X$  et  $Y$ , son importance étant résumée par le :

*Théorème* : (i) Si  $X$  et  $Y$  sont indépendantes, alors  $\rho_{XY} = 0$ .  
 (ii) On a toujours  $|\rho_{XY}| \leq 1$ , avec égalité si, et seulement si, les points sont alignés.

En pratique, si  $|\rho_{XY}| > 0.75$  on recherche un alignement statistique.

*Programmation structurée* — On a :

$$a = \frac{1}{S(X)} \cdot \sum_{i=1}^n (x_i - E(X)) \cdot (y_i - E(Y))$$

où :

$$S(x) = \sum_{i=1}^n (x_i - E(X))^2$$

et en posant  $t_i = x_i - E(X)$ , on a aussi :

$$a = \frac{1}{S(X)} \cdot \sum_{i=1}^n t_i \cdot (y_i - E(Y))$$

et avec  $\sum_{i=1}^n t_i = n \cdot E(X - E(X)) = 0$ , on déduit que :

$$a = \frac{1}{S(X)} \cdot \sum_{i=1}^n t_i \cdot y_i$$

Et de la même façon, on a :

$$\alpha = \frac{1}{S(Y)} \cdot \sum_{i=1}^n t_i' \cdot x_i$$

où  $t_i' = y_i - E(Y)$  et  $S(Y) = \sum_{i=1}^n (y_i - E(Y))^2$ .

Avec les relations  $a = \frac{\sigma_{XY}}{\sigma_X^2}$  et  $\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \cdot \sigma_Y}$ , on déduit que

$$\sigma_{XY} = \frac{\sigma_X}{\sigma_Y} \cdot a.$$

*PROCEDURE Regression\_Affine(Entrée : n : Entier ; x, y : Vecteur ; Sortie : a, b, r : Réel) ;*

*Début*

*a = Ex = Ey = Sx = Sy = 0 ;*

*Pour i allant de 1 à n faire*

*Début*

*Ex = Ex + x<sub>i</sub> ; Ey = Ey + y<sub>i</sub> ;*

*Fin ;*

*Ex = Ex/n ;*

*Ey = Ey/n ;*

*Pour i allant de 1 à n faire*

*Début*

*t = x<sub>i</sub> - Ex ;*

*Sx = Sx + t·x<sub>i</sub> ;*

*a = a + t·y<sub>i</sub> ;*

*t = y<sub>i</sub> - Ey ;*

*Sy = Sy + t·y<sub>i</sub> ;*

*Fin ;*

*a = a/Sx ;*

$$r = a \cdot \sqrt{\frac{Sx}{Sy}} ;$$

$$b = Ey - a \cdot Ex ;$$

$$\text{Si } |r| < 0.75$$

Alors Arrêter ('Modèle affine mal adapté') ;

$$\{ \alpha = r^2/a ; \beta = Ex - \alpha \cdot Ey \}$$

Fin ;

### 2.5.3 Exemple : le problème du fil chaud

Les mesures sont contenues dans  $S = \{ (v_i, T_i) ; 1 \leq i \leq 14 \}$ , avec les vitesses et tensions données par :

$v$ : 5.00; 5.91; 7.07; 8.36; 9.78; 11.54; 13.47; 15.05; 16.98; 19.14; 20.81; 23.45; 26.23; 27.68. $T$ : 4.11; 4.22; 4.34; 4.42; 4.52; 4.63; 4.73; 4.81; 4.89; 4.98; 5.03 ; 5.12 ; 5.20 ; 5.25.
--

Tableau 5.1

La loi cherchée est de la forme :  $T^2 = a \cdot \sqrt{v} + b$ . En posant  $y = T^2$  et  $x = \sqrt{v}$ , on se ramène à un modèle affine et on trouve :

Coefficients de régression :  $a = 3.46$  ;  $b = 9.52$

Coefficient de corrélation :  $\rho = 0.9988$

### 2.5.4 Régression polynomiale

*Position du problème* — De manière plus générale, on peut parfois être amené à chercher un modèle polynomial, soit :

$$f(t, a_0, \dots, a_p) = a_0 + a_1 \cdot t + \dots + a_p \cdot t^p$$

Un critère pour chercher un tel modèle, avec  $p = 2$  et dans le cas de points  $x_i$  équidistants, est donné par Rufener en considérant les *différences finies d'ordre 2*.

En effet, si les points sont près d'une courbe du second degré d'équation  $y = a \cdot x^2 + b \cdot x + c$ , nécessairement on aura  $y_{i-1} - 2 \cdot y_i + y_{i+1} = a \cdot (2 \cdot dx)^2$ , où  $dx = x_{i+1} - x_i$  ne dépend pas de  $i$ , c'est-à-dire que toutes ces différences finies d'ordre 2 devraient être sensiblement constantes.

De manière générale, pour un modèle polynomiale de degré  $p$ , une condition nécessaire pour que le modèle soit acceptable est que les différences finies d'ordre  $p$  soient sensiblement constantes.

On suppose donc une telle condition vérifiée et il s'agit donc de minimiser :

$$E(a_0, \dots, a_p) = \sum_{i=1}^n \left( y_i - (a_0 + a_1 \cdot x_i + \dots + a_p \cdot x_i^p) \right)^2$$

Si on note  $y$ , le vecteur de  $\mathbb{R}^n$  de coordonnées  $y_i$ ,  $v$  un vecteur de coordonnées  $a_i$  dans  $\mathbb{R}^{p+1}$  et  $B$  la matrice à  $n$  lignes et  $p + 1$  colonnes définie par  $b_{ij} = x_i^{j-1}$ , pour  $1 \leq i \leq n$  et  $1 \leq j \leq p + 1$ , alors :

$$E(a_0, \dots, a_p) = \|y - B \cdot v\|^2$$

où  $\|z\|^2 = \sum_{i=1}^n z_i^2$  est la norme euclidienne de  $\mathbb{R}^n$ .

On doit donc minimiser la fonction :

$$E: \mathbb{R}^{p+1} \rightarrow \mathbb{R}^+ \\ v \mapsto E(v) = \|y - B \cdot v\|^2$$

*Recherche du minimum* — Toujours avec la même condition nécessaire d'extremum, on aboutit au système linéaire de  $p + 1$  équations à  $p + 1$  inconnues :

$$\frac{\partial E}{\partial a_i} = 2 \cdot \sum_{k=1}^n (a_j \cdot x_k^j - y_k) \cdot x_k^i = 0 \quad (i = 0, 1, \dots, p) \quad (1)$$

soit, en permutant l'ordre des sommations :

$$\sum_{j=1}^{p+1} S(X^{j+i-2}) \cdot a_{j-1} = S(X^{i-1} \cdot Y) \quad (i = 1, 2, \dots, p + 1) \quad (2)$$

en notant :

$$S(X^{j+i-2}) = \sum_{k=1}^n x_k^{j+i-2} \quad \text{et} \quad S(X^{i-1} \cdot Y) = \sum_{k=1}^n x_k^{i-1} \cdot y_k$$

Le système (2) est appelé *système d'équations normales*.

En notant :

$$\beta_{ij} = S(X^{j+i-2}) = \sum_{k=1}^n x_k^{j+i-2} ; \quad \gamma_i = S(X^{i-1} \cdot Y) = \sum_{k=1}^n x_k^{i-1} \cdot y_k \quad (1 \leq i, j \leq p + 1)$$

le système (2) s'écrit :

$$\beta \cdot v = \gamma$$

On peut en fait montrer que la condition (2) est aussi suffisante. Pour ce faire, on remarque que :

$$(3) \beta = {}^t B \cdot B \quad \text{et} \quad \gamma = {}^t B \cdot y$$

avec :

$$B = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^p \\ 1 & x_2 & x_2^2 & \dots & x_2^p \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & x_n & x_n^2 & \dots & x_n^p \end{pmatrix} ; \quad {}^t B = \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ x_1 & x_2 & x_3 & \dots & x_n \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ x_1^p & x_2^p & x_3^p & \dots & x_n^p \end{pmatrix}$$

puis en écrivant que pour tous  $u, v$  dans  $\mathbb{R}^{p+1}$ , on a :

$$E(u + v) = E(v) + 2 \cdot \langle {}^t B \cdot B \cdot v - {}^t B \cdot y \mid u \rangle + \|B \cdot u\|^2$$

on déduit que si  $v$  est solution de (2), alors :

$$\forall u \in \mathbb{R}^{p+1}: E(u + v) = E(v) + \|B \cdot u\|^2 \geq E(v)$$

l'égalité étant réalisée si, et seulement si,  $u = 0$ . Ce qui prouve que  $E$  admet un minimum en  $v$ .

Donc (2) est bien une condition nécessaire et suffisante d'extremum, ce qui se résume par le :

*Théorème* : Un vecteur  $v$  de  $\mathbb{R}^{p+1}$  est solution du problème de minimisation de  $E$ , si et seulement si, il est solution du système linéaire :  ${}^t B \cdot B \cdot v = {}^t B \cdot y$ .

*Remarque 1* — On peut vérifier que, si  $n > p$  et au moins  $p + 1$   $x_i$  sont distincts, alors la matrice  ${}^tB \cdot B$  est symétrique définie positive (cf. [Ciarlet] p. 69) et donc le problème (1) admet une unique solution.

*Remarque 2* — Souvent on constate que  $\text{Dét}({}^tB \cdot B)$  est proche de 0, ce qui signifie que le système (1) est mal conditionné, c'est-à-dire qu'une petite variation des coefficients de  $B$  peut fortement perturber la solution numérique trouvée. Par exemple, si les  $x_i$  sont régulièrement espacés sur  $[0,1]$  alors la matrice  ${}^tB \cdot B$  est proportionnelle à une matrice de Hilbert d'ordre  $p + 1$ , dont le déterminant décroît très vite vers zéro quand  $n$  augmente.

De manière pratique, le système donnera des résultats médiocres dès que  $p > 7$ . On aura donc intérêt à travailler en double précision.

*Remarque 3* — Si  $p + 1 = n$  et si  $B$  est inversible, alors les équations normales vont s'écrire  $B \cdot u = y$ .

*Résolution du système d'équations normales* — Pour résoudre le système d'équations normales on pourra utiliser la méthode des pivots de Gauss ou la méthode de Cholesky vues dans le chapitre « Analyse Numérique Linéaire ».

Ce qui donne la procédure ci-dessous où  $v$  est le vecteur contenant les coefficients du polynôme de régression. La matrice du système d'équations normales est notée  $\beta$  et le second membre  $\gamma$ .

Tout d'abord, on a une procédure de construction du système linéaire :

*PROCEDURE SystemeRegression(Entrée  $n, p$  : Entiers ;  $x, y$  : Vecteurs ;  
Sortie  $\beta$  : Matrice ;  $\gamma$  : Vecteur) ;*

*Début*

*Pour  $i$  allant de 1 à  $p + 1$  faire*

*Début*

*Pour  $j$  allant de 1 à  $p + 1$  Faire  $\beta_{ij} = 0$  ;*

*$\gamma_i = 0$  ;*

*Fin ;*

*$\beta_{11} = n$  ;*

*$\gamma_1 = y_1$  ;*

*Pour  $k$  allant de 2 à  $n$  Faire  $\gamma_1 = \gamma_1 + y_k$  ;*

*Pour  $i$  allant de 2 à  $p + 1$  Faire*

*Début*

*$\gamma_i = y_i * x_1^{i-1}$  ;*

*Pour  $k$  allant de 2 à  $n$  Faire  $\gamma_i = \gamma_i + y_k * x_k^{i-1}$  ;*

*Pour  $j$  allant de 1 à  $i$  Faire*

*Début*

*$\beta_{ij} = x_1^{i+j-2}$  ;*

*Pour  $k$  allant de 2 à  $n$  Faire  $\beta_{ij} = \beta_{ij} + x_k^{i+j-2}$  ;*

*Fin ;*

*Fin ;*

*Pour  $i$  allant de 1 à  $p$  Faire*

*Début*

*Pour  $j$  allant de  $i + 1$  à  $p + 1$  Faire  $\beta_{ij} = \beta_{ji}$  ;*

*Fin ;*

Fin ;

PROCEDURE RegressionPolynomiale(Entrée :  $p, n$  : Entier ;  $x, y$  : Vecteur ;  
Sortie  $v$  : Vecteur) ;

Début

SystemeRegression( $n, p, x, y, \beta, \gamma$ ) ;

Pivotage( $p + 1, \beta, \gamma$ ) ;

SystTriSup( $p + 1, \beta, \gamma, v$ ) ;

Fin ;

Il suffit ensuite de définir une fonction polynôme à partir de  $u$ , en posant :

$$P(t) = u_1 + u_2 \cdot t + \dots + u_{p+1} \cdot t^p$$

### 2.5.5 Cas général

Les fonctions de base, linéairement indépendantes, étant données, on cherche à minimiser  $E$ . De la même façon que dans le cas polynomiale, on peut vérifier qu'un vecteur  $u$  de composantes  $(a_1, \dots, a_p)$  est solution du problème de minimisation si, et seulement si, il est solution du système d'équations normales :

$${}^t B \cdot u = {}^t B \cdot y \quad \text{où} \\ B = \left( (f_j(x_i))_{\substack{1 \leq i \leq n \\ 1 \leq j \leq p}} \right), \text{ ce qui donne :}$$

$$\begin{cases} \beta = {}^t B \cdot B = ((\beta_{ij}))_{1 \leq i, j \leq p} \\ \gamma = {}^t B \cdot y = ((\gamma_i))_{1 \leq i \leq p} \end{cases} \text{ avec : } \begin{cases} \beta_{ij} = \sum_{k=1}^n f_j(x_k) \cdot f_i(x_k) \\ \gamma_i = \sum_{k=1}^n y_k \cdot f_i(x_k) \end{cases}$$

La matrice  $\beta$  étant symétrique. Si  $B$  est de rang  $p$ , alors  $\beta$  est définie positive et le système d'équations normales a une solution et une seule. Dans tous les cas, le système a au moins une solution.

On a le résultat suivant :

**Théorème :** Si les fonctions de base  $(y_i)_{1 \leq i \leq p}$  sont linéairement indépendantes et s'il y a au moins  $p$  points distincts dans les  $x_i$  alors  $B$  est de rang  $p$  et le système d'équations normales admet une unique solution  $v$ . C'est-à-dire que le problème d'extremum admet une unique solution.

*Remarque* — S'il y a  $p$  points distincts parmi les  $n$ , il peut exister plusieurs points de même abscisse.

Pour la programmation structurée, on procédera alors comme dans le cas polynomiale en remplaçant les  $x^k$ , pour  $k = 0, \dots, p$  par les  $y_k(x)$ , pour  $k = 1, \dots, p$ .

### 2.5.6 Régression trigonométrique

On considère ici le modèle :

$$f(x, a_0, a_1, b_1, \dots, a_p, b_p) = \\ a_0 + a_1 \cdot \cos(\omega \cdot x) + b_1 \cdot \sin(\omega \cdot x) + \dots + a_p \cdot \cos(p \cdot \omega \cdot x) + b_p \cdot \sin(p \cdot \omega \cdot x)$$

où  $\omega$  est connu (i.e. la période  $T = \frac{2 \cdot \pi}{\omega}$  est connue).

La méthode décrite précédemment donne donc un moyen de calculer les coefficients de Fourier d'une fonction périodique échantillonnée.

Dans le cas particulier où les  $x_i$  sont équidistants, on peut déterminer assez facilement les  $a_k$  et  $b_k$ .

Par changement de variable, on peut toujours se ramener à une fonction  $2\cdot\pi$ -périodique, c'est-à-dire à  $\omega = 1$ .

On suppose donc que les  $x_k$  sont donnés par :

$$x_k = \frac{k \cdot \pi}{n} \quad (k = -n + 1, \dots, n)$$

avec  $n \geq p$  (ce sont donc des points équidistants sur  $]-\pi, \pi]$ ).

Il s'agit donc de minimiser :

$$E(a_0, \dots, a_p, b_1, \dots, b_p) = \sum_{i=-n+1}^n \left( y_i - a_0 - \sum_{k=1}^p (a_k \cdot \cos(k \cdot x_i) + b_k \cdot \sin(k \cdot x_i)) \right)^2$$

En annulant les dérivées partielles, on aboutit alors au système linéaire de  $2 \cdot n + 1$  équations à  $2 \cdot n + 1$  inconnues :

$$\sum_{i=-n+1}^n \left( y_i - a_0 - \sum_{k=1}^p (a_k \cdot \cos(k \cdot x_i) + b_k \cdot \sin(k \cdot x_i)) \right) \cdot \cos(j \cdot x_i) = 0$$

$$\sum_{i=-n+1}^n \left( y_i - a_0 - \sum_{k=1}^p (a_k \cdot \cos(k \cdot x_i) + b_k \cdot \sin(k \cdot x_i)) \right) \cdot \sin(j \cdot x_i) = 0$$

avec  $j$  variant de 0 à  $p$  dans les premières équations et de 1 à  $p$  dans les dernières.

En exploitant les propriétés des fonctions trigonométriques, on déduit alors que :

$$a_0 = \frac{1}{2 \cdot n} \cdot \sum_{i=-n+1}^n y_i$$

$$a_k = \frac{1}{n} \cdot \sum_{i=-n+1}^n y_i \cdot \cos(k \cdot x_i) \quad (k = 1, \dots, n)$$

$$b_k = \frac{1}{n} \cdot \sum_{i=-n+1}^n y_i \cdot \sin(k \cdot x_i) \quad (k = 1, \dots, n)$$

(Cf. Pichat, ..., p.118).

Ces formules sont simplement les approximations des coefficients de Fourier en utilisant la méthode des rectangles à droite pour le calcul des intégrales. On pourrait utiliser la méthode de Simpson ou de Romberg (Cf. le chapitre « Calcul Numérique des Intégrales ») pour avoir une meilleure précision, mais cela augmentera le temps de calcul qui est déjà important du fait de la présence des fonctions trigonométriques.

En écrivant les choses sous forme complexe, on aura un algorithme beaucoup plus performant de calcul, c'est l'algorithme FFT (pour Fast Fourier Transform) de Cooley et Tukey (Cf. le chapitre sur le calcul numérique des intégrales).

## 2.6 Les modèles non linéaires

Ce sont les modèles qui vont se ramener au cas linéaire par changement de variable. On a déjà vu un exemple avec le problème du fil chaud.

*Exemple 1* —  $f(x) = a \cdot \cos(\omega \cdot x + \phi) + b$ , où  $\omega$  est connue et  $a$ ,  $\phi$  et  $b$  sont à déterminer.

En écrivant  $f(x) = a \cdot \cos(\phi) \cdot \cos(\omega \cdot x) - a \cdot \sin(\phi) \cdot \sin(\omega \cdot x) + b$ , on se ramène à la régression trigonométrique.

*Exemple 2* —  $f(x) = a \cdot e^{-\alpha \cdot t}$ , en utilisant le logarithme népérien on se ramène au cas linéaire :

$$\text{Ln}(f(x)) = -\alpha \cdot t + \text{Ln}(a)$$

Pour les modèles non linéaires qui résistent à ce genre d'astuce, il existe des méthodes de linéarisation. On écrit le système  $\frac{\partial E}{\partial a_i} = 0$ , que l'on essaye de résoudre par une méthode de type Newton.

## 2.7 Un exemple d'application : détermination du coefficient de traînée d'une particule sphérique

On considère une particule sphérique de masse  $m$ , de diamètre  $d$  et de densité volumique  $\rho$  qu'on laisse tomber, à partir de  $t = 0$  et sans vitesse initiale, dans un fluide au repos de densité volumique  $\rho_f$  et de viscosité cinématique  $\nu$ .

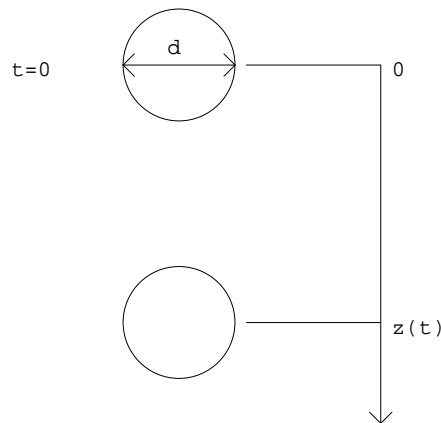


Figure 5.1

La force de réaction, due à la viscosité du fluide, est alors donnée par :

$$F_r = -C_d(v(t)) \cdot \frac{1}{2} \cdot \rho_f \cdot v(t) \cdot |v(t)| \cdot A$$

où,  $v(t)$  est la vitesse de la particule à l'instant  $t$ ,  $A = \pi \cdot \left(\frac{d}{2}\right)^2$  et  $C_d(v(t))$  est le coefficient de traînée de la particule.



Le coefficient de traînée peut s'exprimer en fonction du nombre de Reynolds donné, dans le cas d'une particule sphérique, par :

$$R_e(t) = \frac{v(t) \cdot d}{\nu}$$

L'expression de  $C_d$  en fonction de  $R_e$  est connue seulement de façon empirique, avec des formules valables seulement dans des régions déterminées des  $R_e$ .

Dans le cas d'une sphère, on a les lois classiques suivantes :

(i) loi de Stokes —  $C_d = \frac{24}{R_e}$  ( $0 \leq R_e \leq 0.1$ ) (certains auteurs prennent  $R_e$  dans  $[0,1]$ ) ;

(ii) loi d'Olson —  $C_d = \frac{24}{R_e} \cdot \left(1 + \frac{3}{16} \cdot R_e\right)^{\frac{1}{2}}$  ( $0.1 \leq R_e \leq 100$ )

De manière un peu plus générale, Morsi propose des formules du type :

$$C_d = K_0 + \frac{K_1}{R_e} + \frac{K_2}{R_e^2} \quad (a \leq R_e \leq b)$$

avec des valeurs spécifiques des coefficients  $K_i$  dans les intervalles suivants :

$[0,0.1]$ ,  $[0.1,1]$ ,  $[1,10]$ ,  $[10,100]$ ,  $[100,1000]$ ,  $[1000,5000]$ ,  $[5000,10000]$  et  $[10000,50000]$ .

Une idée pour trouver ces coefficients est de partir d'un tableau de mesures expérimentales du type ci-dessous et de procéder par interpolations polynomiales.

$R_e$ :	0.1	0.2	0.3	0.5	0.7	1
$C_d$ :	240	120	80	49	36.5	26.5
$R_e$ :	1	2	3	5	7	10
$C_d$ :	26.5	14.4	10.4	6.9	5.4	4.1
$R_e$ :	10	20	30	50	70	100
$C_d$ :	4.1	2.55	2	1.5	1.27	1.07
$R_e$ :	100	200	300	500	700	1000
$C_d$ :	1.07	0.77	0.65	0.55	0.5	0.46
$R_e$ :	1000	2000	3000	5000		
$C_d$ :	0.46	0.42	0.4	0.385		
$R_e$ :	5000	7000	10000			
$C_d$ :	0.385	0.39	0.405			
$R_e$ :	10000	20000	30000	50000		
$C_d$ :	0.405	0.45	0.47	0.49		

Tableau 5.2

Dans le programme ci-dessous, on adopte les notations suivantes :

Dans l'intervalle N° j, les mesures de Re et Cd sont notées Re(j,i) et Cd(j,i) où  $1 \leq j \leq \text{Nombre\_Intervalles} = 7$  et  $1 \leq i \leq n(j)$ , avec  $n(j)$  égal au nombre de mesures dans cet intervalle et  $n(j) \leq \text{Nombre\_Mesures} = 6$ .

Dans l'intervalle N° j, les coefficients de Morsi sont notés K(j,i), où  $i = 1, 2, 3$ .

La procédure DONNEES\_DU\_FICHER permet de lire les données du problème depuis un fichier.

On notant, pour j donné,  $u(i) = 1/\text{Re}(j,i)$  et  $v(i) = \text{Cd}(j,i)$ , avec  $1 \leq i \leq n(j) = n$ , on est amené à chercher un polynôme de régression de degré au plus 2 défini par les vecteurs u et v.

La formation du système d'équations normales est réalisée par la procédure Equations\_Normales. Ce système est défini par :

$$M = \begin{pmatrix} n & a & b \\ a & b & c \\ b & c & d \end{pmatrix} \text{ et } \gamma = \begin{pmatrix} \gamma_1 \\ \gamma_2 \\ \gamma_3 \end{pmatrix}$$

avec :

$$a = \sum_{k=1}^n u_k; \quad b = \sum_{k=1}^n u_k^2; \quad c = \sum_{k=1}^n u_k^3; \quad d = \sum_{k=1}^n u_k^4;$$

$$\gamma_1 = \sum_{k=1}^n v_k; \quad \gamma_2 = \sum_{k=1}^n u_k \cdot v_k; \quad \gamma_3 = \sum_{k=1}^n u_k^2 \cdot v_k;$$

La solution de ce système est alors donnée par :

$$w_1 = \frac{\delta_1}{\delta}; \quad w_2 = \frac{\delta_2}{\delta}; \quad w_3 = \frac{\delta_3}{\delta};$$

avec :

$$\delta = n \cdot (b \cdot d - c^2) + 2 \cdot a \cdot b \cdot c - b^3 - a^2 \cdot d$$

$$\delta_1 = \gamma_1 \cdot (b \cdot d - c^2) + \gamma_2 \cdot (b \cdot c - a \cdot d) + \gamma_3 \cdot (a \cdot c - b^2)$$

$$\delta_2 = \gamma_1 \cdot (b \cdot c - a \cdot d) + \gamma_2 \cdot (n \cdot d - b^2) + \gamma_3 \cdot (a \cdot b - n \cdot c)$$

$$\delta_3 = \gamma_1 \cdot (a \cdot c - b^2) + \gamma_2 \cdot (b \cdot a - n \cdot c) + \gamma_3 \cdot (n \cdot b - a^2)$$

Pour ce qui est du tracé graphique, étant donné que  $R_e$  prendra de très grandes valeurs et  $C_d$  des valeurs beaucoup plus petites, il est préférable de tracer le graphe de  $\text{Ln}(C_d)$  en fonction de  $\text{Ln}(R_e)$ .

```
with TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH1, MATH3, CURVE ;
procedure BILLE1 is
use TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH1, MATH3, CURVE ;
RE_MAX : constant FLOAT := 50000.0 ;
NOMBRE_INTERVALLES : constant NATURAL := 7 ;
NOMBRE_MESURES : constant NATURAL := 6 ;
type MATRICE_MESURES is array( POSITIVE range 1..NBR_INTERVALLES,
                               POSITIVE range 1..NBR_MESURES) of FLOAT ;
type VECTEUR_COEFF is array(POSITIVE range 1..NBR_INTERVALLES,
                             NATURAL range 0..2) of FLOAT ;
type VECTEUR_ENTIER is array(POSITIVE range 1..NBR_INTERVALLES)
                             of NATURAL ;
type VECTEUR is array(POSITIVE range 1..NBR_MESURES) of FLOAT ;
```

```

type VECTEUR_3 is array(POSITIVE range 1..3) of FLOAT ;
Re, Cd : MATRICE_MESURES ;
nn : VECTEUR_ENTIER ;
K : VECTEUR_COEFF ;
Rho, Rhof, g, Nu, d, h, t0, tMax, z0, v0 : FLOAT ;

procedure DONNEES_DU_FICHIER(t0,tMax,z0,v0,Rho,Rhof,Nu,g,
                             d : out FLOAT ;   nn : in out VECTEUR_ENTIER ;
                             Re, Cd : out MATRICE_MESURES) is
Nom : STRING(1..50) ;
FichierDonnees : FILE_TYPE ;
begin
  PUT_LINE("  Donnees du fichier") ;
  LIRE_FICHIER(FichierDonnees,Nom) ;
  GET(FichierDonnees,t0) ;
  GET(FichierDonnees,tMax) ;
  GET(FichierDonnees,z0) ;
  GET(FichierDonnees,v0) ;
  GET(FichierDonnees,Rho) ;
  GET(FichierDonnees,Rhof) ;
  GET(FichierDonnees,Nu) ;
  GET(FichierDonnees,g) ;
  GET(FichierDonnees,d) ;
  SKIP_LINE(FichierDonnees) ;
  for j in 1..NOMBRE_INTERVALLES loop
    GET(FichierDonnees,nn(j)) ;
    SKIP_LINE(FichierDonnees) ;
    for i in 1..nn(j) loop
      GET(FichierDonnees,Re(j,i)) ;
    end loop ;
    SKIP_LINE(FichierDonnees) ;
    for i in 1..nn(j) loop
      GET(FichierDonnees,Cd(j,i)) ;
    end loop ;
    SKIP_LINE(FichierDonnees) ;
  end loop ;
  CLOSE(FichierDonnees) ;
end DONNEES_DU_FICHIER ;

procedure AFFICHAGE_DONNEES is
begin
  CLRSCR ;
  PUT_LINE(IMP,"          Donnees du probleme") ; NEW_LINE(IMP) ;
  PUT_LINE(IMP,"Conditions initiales :") ;
  PUT(IMP,"          t0 = ") ;
  PUT(IMP,t0,6,3,0) ; NEW_LINE(IMP) ;
  PUT(IMP,"          tMax = ") ;
  PUT(IMP,tMax,6,3,0) ; NEW_LINE(IMP) ;
  PUT(IMP,"          z0 = ") ;
  PUT(IMP,z0,6,3,0) ; NEW_LINE(IMP) ;
  PUT(IMP,"          v0 = ") ;

```



```

DD, D1, D2, D3 : FLOAT ;
begin
  DD := FLOAT(n)*(b*d - c*c) + 2.0*a*b*c - b*b*b - a*a*d ;
  D1 := Gamma(1)*(b*d - c*c) + Gamma(2)*(b*c - a*d)
      + Gamma(3)*(a*c - b*b) ;
  D2 := Gamma(1)*(b*c - a*d) + Gamma(2)*(FLOAT(n)*d - b*b)
      + Gamma(3)*(a*b - FLOAT(n)*c) ;
  D3 := Gamma(1)*(a*c - b*b) + Gamma(2)*(a*b - FLOAT(n)*c)
      + Gamma(3)*(FLOAT(n)*b - a*a) ;
  Gamma(1) := D1/DD ; Gamma(2) := D2/DD ; Gamma(3) := D3/DD ;
end RESOLUTION_EQU_NORMALES ;

procedure CALCUL_COEFF_MORSI(Re, Cd: in MATRICE_MESURES ;
                             K : out VECTEUR_COEFF) is
  a, b, c, d : FLOAT ;
  Gamma : VECTEUR_3 ;
  u, v : VECTEUR ;
begin
  for j in 1..NOMBRE_INTERVALLES loop
    for i in 1..nn(j) loop
      u(i) := 1.0/Re(j,i) ;
      v(i) := Cd(j,i) ;
    end loop ;
    EQUATIONS_NORMALES(nn(j),u,v,a,b,c,d,Gamma) ;
    RESOLUTION_EQU_NORMALES(nn(j),a,b,c,d,Gamma) ;
    for i in 0..2 loop
      K(j,i) := Gamma(i + 1) ;
    end loop ;
  end loop ;
end CALCUL_COEFF_MORSI ;

procedure AFICCHE_COEFF_MORSI(Re, Cd : in MATRICE_MESURES ;
                              K : in VECTEUR_COEFF) is
begin
  for j in 1..NOMBRE_INTERVALLES
  loop
    PUT_LINE(IMP,"Intervalle Numero " & INTEGER'IMAGE(j)) ;
    for i in 0..2 loop
    PUT(IMP,"K(" & INTEGER'IMAGE(j) & "," & INTEGER'IMAGE(i) & ") = " ) ;
      PUT(IMP,k(j,i),3,3,0) ; NEW_LINE(IMP) ;
    end loop ;
    PUT_LINE(IMP,"Re          Cd (Experimental)   Cd (Calcule)") ;
    for i in 1..nn(j) loop
      PUT(IMP,Re(j,i),6,3,0) ;
      PUT(IMP,"          ") ;
      PUT(IMP,Cd(j,i),6,3,0) ;
      PUT(IMP,"          ") ;
      PUT(IMP,Cd_j(K,j,1.0/Re(j,i)),6,3,0) ; NEW_LINE(IMP) ;
    end loop ;
    Pause ;
  end loop ;
end AFICCHE_COEFF_MORSI ;

```

```

procedure COURBE_CD(Re, Cd : in MATRICE_MESURES ;
                   K : in VECTEUR_COEFF ; CC : in out Courbe) is
-- Definition de la courbe Ln(Cd) en fonction de Ln(Re).
Pas, Pas_r, r : FLOAT ;
i : INTEGER ;
begin
  i := 0 ;
  r := 0.01 ;
  CC.uMin := Ln(r) ;
  CC.uMax := Ln(RE_MAX) ;
  CC.vMin := 1.0E+37 ;
  CC.vMax := -1.0E+37 ;
  CC.u(0) := CC.uMin ;
  CC.v(0) := Ln(Cd_j(K,1,1.0/r)) ;
  Pas := (CC.uMax - CC.uMin)/FLOAT(CC.n) ;
  Pas_r := Exp(Pas) ;
  r := r*Pas_r ;
  for j in 1..NOMBRE_INTERVALLES loop
    while (r <= Re(j,nn(j))) loop
      i := i + 1 ;
      CC.u(i) := CC.u(i-1) + Pas ;
      CC.v(i) := Ln(Cd_j(K,j,1.0/r)) ;
      CC.vMin := Min(CC.vMin,CC.v(i)) ;
      CC.vMax := Max(CC.vMax,CC.v(i)) ;
      r := r*Pas_r ;
    end loop ;
  end loop ;
end COURBE_CD ;

procedure TRACE_CD(Re, Cd : in MATRICE_MESURES) is
CC : Courbe(200) ;
begin
  COURBE_CD(Re,Cd,K,CC) ;
  Trace_Une_Courbe(CC) ;
  for j in 1..NOMBRE_INTERVALLES loop
    for i in 1..nn(j) loop
      CROIX(Ln(Re(j,i)),Ln(Cd(j,i))) ;
    end loop ;
  end loop ;
  SORTIE_GRAPHIQUE ;
end TRACE_CD ;

begin
  MODE_AFFICHAGE ;
  DONNEES_DU_FICHER(t0,tMax,z0,v0,Rho,Rhof,Nu,g,d,nn,Re,Cd) ;
  AFFICHAGE_DONNEES ;
  PUT_LINE("Calcul des coefficients de regression") ;
  CALCUL_COEFF_MORSI(Re,Cd,K) ;
  AFICCHE_COEFF_MORSI(Re,Cd,K) ;
  CLRSCR ;
  PUT_LINE(" Trace de la courbe Ln(Cd) en fonction de Ln(Re)") ;
  PUT_LINE(" Patience je calcul ...") ;
  TRACE_CD(Re,Cd) ;

```

```
CLOSE(IMP) ;
end BILLE1 ;
```

Les résultats obtenus sont alors les suivants :

*Conditions initiales* —

t0	=	0.000
tMax	=	10.000
z0	=	0.000
v0	=	0.000

*Caractéristiques de la bille et du fluide* —

Rho	=	8000.000
Rhof	=	1.220
Nu	=	0.000015
g	=	9.800
d	=	0.010

*Calculs des coefficients de régression*

*Intervalle Numero 1* —  $K(1,0) = 3.486$   $K(1,1) = 22.782$   $K(1,2) = 0.087$

Re	Cd (Expérimental)	Cd (Calculé)
0.100	240.000	240.051
0.200	120.000	119.583
0.300	80.000	80.398
0.500	49.000	49.400
0.700	36.500	36.211
1.000	26.500	26.356

*Intervalle Numero 2* —  $K(2,0) = 1.538$   $K(2,1) = 26.943$   $K(2,2) = -1.995$

Re	Cd (Expérimental)	Cd (Calculé)
1.000	26.500	26.486
2.000	14.400	14.511
3.000	10.400	10.297
5.000	6.900	6.847
7.000	5.400	5.346
10.000	4.100	4.213

*Intervalle Numero 3* —  $K(3,0) = 0.687$   $K(3,1) = 41.253$   $K(3,2) = -71.552$

Re	Cd (Expérimental)	Cd (Calculé)
10.000	4.100	4.097
20.000	2.550	2.571
30.000	2.000	1.983
50.000	1.500	1.484
70.000	1.270	1.262
100.000	1.070	1.093

*Intervalle Numero 4* —  $K(4,0) = 0.377$   $K(4,1) = 88.621$   $K(4,2) = -1929.682$

Re	Cd (Expérimental)	Cd (Calculé)
100.000	1.070	1.070
200.000	0.770	0.771
300.000	0.650	0.650
500.000	0.550	0.546
700.000	0.500	0.499
1000.000	0.460	0.463

*Intervalle Numero 5* —  $K(5,0) = 0.357$   $K(5,1) = 146.799$   
 $K(5,2) = -43665.282$

Re	Cd (Expérimental)	Cd (Calculé)
1000.000	0.460	0.460
2000.000	0.420	0.419
3000.000	0.400	0.401
5000.000	0.385	0.385

*Intervalle Numero 6* —  $K(6,0) = 0.478$   $K(6,1) = -987.500$   
 $K(6,2) = 2625000.000$

Re	Cd (Expérimental)	Cd (Calculé)
5000.000	0.385	0.385
7000.000	0.390	0.390
10000.000	0.405	0.405

*Intervalle Numero 7* —  $K(7,0) = 0.522$   $K(7,1) = -1725.126$   
 $K(7,2) = 5556952.339$

Re	Cd (Expérimental)	Cd (Calculé)
10000.000	0.405	0.405
20000.000	0.450	0.450
30000.000	0.470	0.471
50000.000	0.490	0.490

Enfin le graphe de  $\ln(C)$  en fonction de  $\ln(R)$  est donné sur la figure 5.2 :



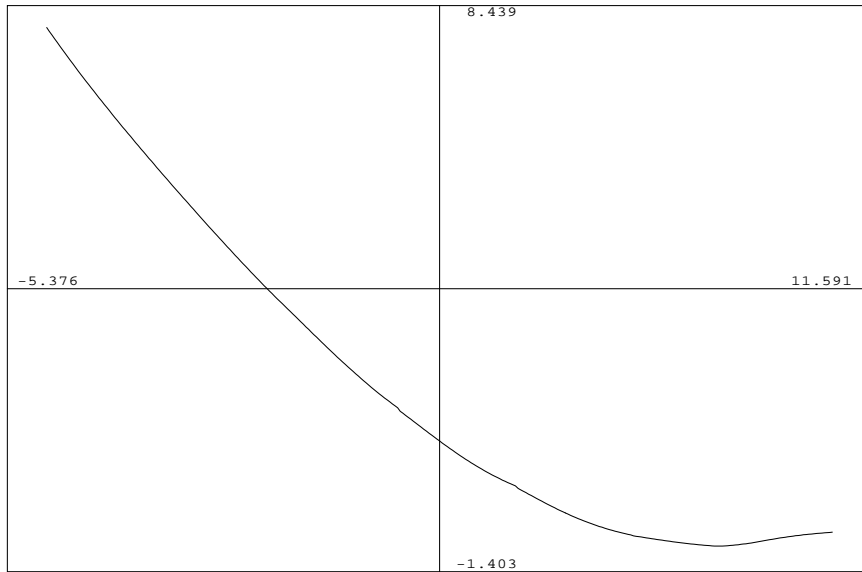


Figure 5.2

## 2.8 Calcul du cercle des moindres carrés

### 2.8.1 Notations et hypothèses

On dispose d'un nuage de points expérimentaux :

$$S = \{ M_i = (x_i, y_i) ; 1 \leq i \leq n \}$$

qui sont supposés être sur un cercle de centre et de rayon inconnus.

Un cercle sera décrit par une équation implicite de la forme :

$$(1) \quad x^2 + y^2 + a \cdot x + b \cdot y + c = 0$$

C'est le cercle de centre  $\Omega$  et de rayon  $R$ , avec :

$$\Omega = \left( \frac{-a}{2}, \frac{-b}{2} \right) \text{ et } R = \sqrt{\frac{a^2 + b^2}{4} - c}$$

On prendra pour écart entre le nuage  $S$  et un cercle d'équation (1), la quantité :

$$E(a, b, c) = \sum_{i=1}^n (x_i^2 + y_i^2 + a \cdot x_i + b \cdot y_i + c)^2$$

### 2.8.2 Calcul du cercle des moindres carrés

En écrivant les conditions nécessaires d'extremum, on est alors amené à résoudre le système linéaire :

$$\begin{cases} E(X^2) \cdot a + E(X \cdot Y) \cdot b + E(X) \cdot c = -E(X^3) - E(X \cdot Y^2) \\ E(X \cdot Y) \cdot a + E(Y^2) \cdot b + E(Y) \cdot c = -E(X^2 \cdot Y) - E(Y^3) \\ E(X) \cdot a + E(Y) \cdot b + c = -E(X^2) - E(Y^2) \end{cases}$$

où on a noté :

$$E(X^p \cdot Y^q) = \frac{1}{n} \cdot \sum_{i=1}^n x_i^p \cdot y_i^q \quad (p, q \geq 0)$$

En retranchant à la ligne 1, la ligne 3 multipliée par  $E(X)$  et à la ligne 2, la ligne 3 multipliée par  $E(Y)$ , on obtient le système :

$$\begin{cases} V(X) \cdot a + \text{Cov}(X, Y) \cdot b = -\text{Cov}(X, X^2 + Y^2) \\ \text{Cov}(X, Y) \cdot a + V(Y) \cdot b = -\text{Cov}(X^2 + Y^2, Y) \end{cases}$$

le déterminant de ce système étant :

$$\Delta = V(X) \cdot V(Y) \cdot (1 - \rho(X, Y)^2)$$

où  $\rho(X, Y)$  est le coefficient de corrélation de  $X$  et  $Y$ .

Ce déterminant est donc nul si, et seulement si, les points sont alignés, ce qui n'est pas le cas en général pour ce type de problème.

Notre système aura donc une unique solution donnée par :

$$\begin{cases} a = \frac{\alpha}{\Delta} \\ b = \frac{\beta}{\Delta} \\ c = -(E(X) \cdot a + E(Y) \cdot b + E(X^2 + Y^2)) \end{cases}$$

où :

$$\begin{cases} \alpha = \text{Cov}(X, Y) \cdot \text{Cov}(X^2 + Y^2, Y) - V(Y) \cdot \text{Cov}(X^2 + Y^2, X) \\ \beta = \text{Cov}(X, Y) \cdot \text{Cov}(X^2 + Y^2, X) - V(X) \cdot \text{Cov}(X^2 + Y^2, Y) \end{cases}$$

*Programmation structurée* — En reprenant les idées utilisées dans le cas de la régression affine, on aboutit à la procédure suivante,  $E$  est la covariance de  $X$  et  $X^2 + Y^2$ ,  $F$  celle de  $Y$  et  $X^2 + Y^2$  et  $G$  est la moyenne de  $X^2 + Y^2$  :

*PROCEDURE Regression\_Cercle(Entrée  $n$  : Entier ;  $x, y$  : Vecteur ; Sortie  $a, b, R$  : Réels) ;*

*Début*

*Ex = Ey = Vx = Vy = Cov\_xy = E = F = G = 0 ;*

*Pour i allant de 1 à n Faire*

*Début*

*Ex = Ex + x<sub>i</sub> ; Ey = Ey + y<sub>i</sub> ;*

*Fin ;*

*Ex = Ex/n ; Ey = Ey/n ;*

*Pour i allant de 1 à n Faire*

*Début*

*t = x<sub>i</sub> - Ex ; s = x<sub>i</sub><sup>2</sup> + y<sub>i</sub><sup>2</sup> ;*

*Vx = Vx + t·x<sub>i</sub> ;*

*Cov\_xy = Cov\_xy + t·y<sub>i</sub> ;*

*E = E + t·s ;*

*t = y<sub>i</sub> - Ey ;*

*Vy = Vy + t·y<sub>i</sub> ;*

*F = F + t·s ;*

*G = G + s ;*

*Fin ;*

*Δ = Vx·Vy - (Cov\_xy)<sup>2</sup> ;*

*G = G/n ;*

*α = Cov\_xy·F - Vy·E ;*

$$\begin{aligned} \beta &= \text{Cov}_{xy} \cdot E - V_x \cdot F ; \\ a &= \alpha/\Delta ; b = \beta/\Delta ; R = G - a \cdot E_x - b \cdot E_y ; \\ a &= -a/2 ; b = -b/2 ; R = \sqrt{a^2 + b^2} - R ; \\ \text{Fin ;} \end{aligned}$$

### 2.8.3 Application à un calcul de rayon de courbure

On se donne, ici, un nuage de points :

$$S = \{ M_i = (x_i, y_i) ; -n \leq i \leq n \}$$

et on voudrait calculer une approximation du rayon de courbure en  $M_0$ .

Une première idée serait d'utiliser la formule :

$$R_0 = \frac{(1 + y_0'^2)^{\frac{3}{2}}}{y_0''}$$

et d'utiliser des approximations des dérivées par différences finies centrées.

Cette idée donnera des résultats médiocres si les points  $M_i$  sont éloignés pour  $i = -1, 0$  et  $1$ .

Une autre idée, plus efficace en général, est de prendre pour approximation de  $R_0$  le rayon du cercle des moindres carrés défini par les  $M_i$ , avec  $-p \leq i \leq p$ , où  $p$  est à choisir, en fonction de la forme de  $S$ , entre  $1$  et  $n$ . Pour favoriser les points proches de  $M_0$ , on cherchera plutôt à minimiser la quantité :

$$E(a,b,c) = \sum_{i=-p}^p \frac{1}{i^2} \cdot (x_i^2 + y_i^2 + a \cdot x_i + b \cdot y_i + c)^2$$

où on pose  $1/0 = 1$ .

### 3. Approximation uniforme des fonctions continues. Courbes de Bernstein, Bézier et B-Splines

#### 3.1 Position du problème

Etant donnée une fonction continue,  $f : [a,b] \rightarrow \mathbb{R}$ , dont on connaît les valeurs en un nombre fini de points  $x_0, x_1, \dots, x_n$ , on voudrait calculer, pour tout  $x$  dans  $[a,b]$ , une valeur approchée de  $f(x)$ . De plus on voudrait que la précision de cette approximation augmente avec  $n$ .

Le *théorème de Weirstrass* (1855), qui nous dit que toute fonction continue sur un intervalle fermé borné peut être approchée uniformément par des polynômes, nous permet de répondre à la question. Une démonstration de ce théorème peut se faire de façon constructive en utilisant les *polynômes de Bernstein* (Cf § 3.2).

Les restrictions de cette construction sont, d'une part que le modèle de courbe est de type cartésien, i. e. de la forme  $y = f(x)$ , ce qui exclut les courbes fermées, et d'autre part qu'il est nécessaire de connaître  $f$  en des points équidistants.

Si on s'intéresse au cas, plus général, des courbes définies par des équations paramétriques dans  $\mathbb{R}^2$  ou  $\mathbb{R}^3$ , on aboutit aux procédés d'approximations de Bézier et aux approximants B-Splines.

L'idée de base de ces méthodes est, étant donné un nuage de points :

$S = \{M_i = (x_i, y_i) ; 0 \leq i \leq n\}$ , de chercher des approximations polynomiales d'une paramétrisation  $(x(t), y(t))$  en prenant  $[0,1]$  comme espace des paramètres, avec :

$$\begin{cases} x\left(\frac{i}{n}\right) = x_i \\ y\left(\frac{i}{n}\right) = y_i \end{cases} \quad (0 \leq i \leq n)$$

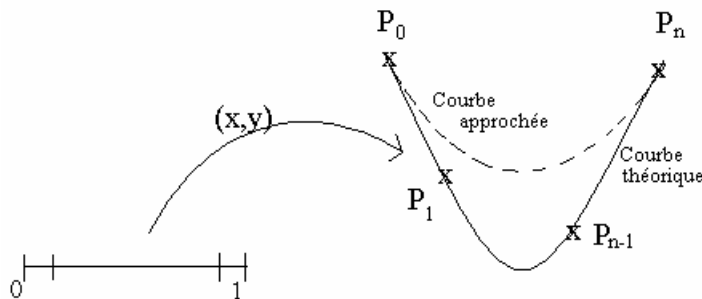


Figure 5.3

Ces procédés vont aussi s'appliquer à la définition de surfaces à partir d'une série de points de  $\mathbb{R}^3$  donnés par l'utilisateur en fonction d'expérimentations ou toute autre considération.

Ces idées sont utilisées en *conception assistée par ordinateur (CAO)* pour définir des carrosseries de voiture ou des fuselages d'avions.

En France, c'est *P. Bezier* qui est à l'origine de ce genre d'idées avec le logiciel *UNISURF* qu'il développa pour le compte de la régie Renault à partir de 1962 (Cf. Math et C.A.O., vol. 4, Chap. 4).

La technique des courbes B-Splines est une généralisation des courbes de Bézier qui ont le désavantage d'avoir une nature globale, c'est-à-dire que la modification d'un point de contrôle  $M_i$  va modifier toute la courbe, contrairement aux B-Splines qui ont un caractère local.

### 3.2 Les bases de Bernstein. Polynômes de Bernstein

#### 3.2.1 Les bases de Bernstein

Pour tout  $n \geq 1$ , on note  $R_n[t]$  l'espace vectoriel des polynômes de degré au plus  $n$ .

*Définition* : On appelle *base de Bernstein d'ordre  $n$* , la base de  $R_n[t]$  formée des polynômes  $B_{n,k}$  définis par :

$$B_{n,k}(t) = C_n^k \cdot t^k \cdot (1-t)^{n-k} \quad (0 \leq k \leq n)$$

*Remarque* — En écrivant que, pour tout  $k = 0, 1, \dots, n$  :

$$t^k = (t + (1-t))^{n-k} \cdot t^k = \sum_{j=0}^{n-k} C_{n-k}^j \cdot t^j \cdot (1-t)^{n-k-j} \cdot t^k = \sum_{i=k}^n \lambda_{i,k} \cdot B_{n,i}(t)$$

on déduit que  $\{ B_{n,k} ; 0 \leq k \leq n \}$  engendrent  $R_n[t]$  et c'est donc bien une base.

Il est facile de vérifier les :

*Propriétés* :

(i)  $\forall t \in ]0,1[, B_{n,k}(t) > 0$  ;

(ii)  $B_{n,k}(0) = \begin{cases} 0 & \text{si } 1 \leq k \leq n \\ 1 & \text{si } k = 0 \text{ et } n \geq 1 \end{cases}$

(iii)  $B_{n,k}(t) = B_{n,n-k}(1-t)$  ( $0 \leq k \leq n$ )

(iv)  $\sum_{k=0}^n B_{n,k}(t) = 1$

(v) On a la récurrence qui permet de calculer les  $B_{n,k}(t)$  :

$$B_{n,n}(t) = t \cdot B_{n-1,n-1}(t)$$

$$B_{n,0}(t) = (1-t) \cdot B_{n-1,0}(t)$$

$$B_{n,k}(t) = (1-t) \cdot B_{n-1,k}(t) + t \cdot B_{n-1,k-1}(t) \quad (1 \leq k \leq n-1)$$

*Remarque* — Les fonctions  $B_{n,k}(t)$  sont aussi appelées *fonctions mélanges* ou *fonctions pondérantes*.

Sur la figure 5.4, on trouvera les graphes des fonctions  $B_{n,k}$  pour  $0 \leq k \leq n = 5$ .

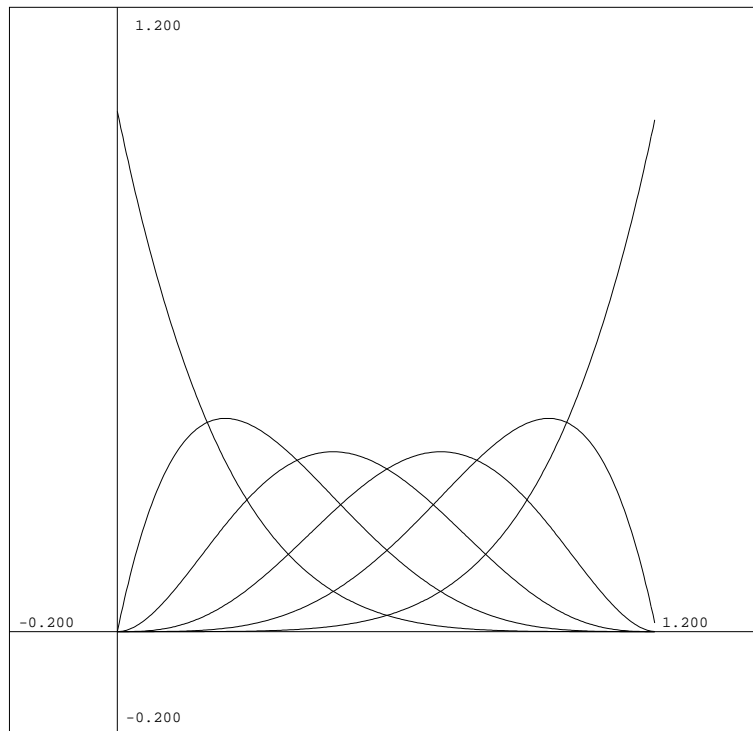


Figure 5.4

Les  $B_{n,k}(t)$  peuvent être calculées de façon récursive avec la procédure suivante.

*FONCTION BaseBernstein\_R(Entrée n, k : Entier ; t : Réel) : Réel ;*

*Début*

*Si n = 0*

*Alors Début*

*BaseBernstein\_R = 1 ;*

*Sortir ;*

*Fin ;*

*Si k = 0*

*Alors BaseBernstein\_R = (1 - t)<sup>n</sup>*

*Sinon Début*

*Aux1 = BaseBernstein\_R(n - 1, k, t) ;*

*Aux2 = BaseBernstein(n - 1, k - 1, t) ;*

*BaseBernstein\_R = (1 - t)·Aux1 + t·Aux2 ;*

*Fin ;*

*Fin ;*

Une autre façon de calculer est d'utiliser la récurrence (V) ci-dessus ce qui donnera, pour n grand, une procédure de calcul beaucoup plus rapide.

L'algorithme de calcul est le suivant :

$$\text{Etape } 0 \text{ — } B_{n,i}^{(0)} = \begin{cases} 1 & \text{si } i = k \\ 0 & \text{sinon} \end{cases} \quad (i = 0, 1, \dots, n)$$

$$\text{Etape } j \text{ — } B_{n,i}^{(j)} = (1 - t) \cdot B_{n,i}^{(j-1)} + t \cdot B_{n,i+1}^{(j-1)} \quad (i = 0, 1, \dots, n - j) \text{ pour } j = 1, 2, \dots, n.$$

Et alors :

$$B_{n,k}(t) = B_{n,0}^{(n)}(t)$$

Cet algorithme est connu sous le nom d'algorithme de *De Casteljeau*. On le retrouvera pour calculer les polynômes de Bernstein et de Bézier.

### 3.2.2 Polynômes de Bernstein associés à une fonction continue

Dans ce paragraphe, on se donne une fonction :  $f : [0,1] \rightarrow \mathbb{R}$ .

*Définition* : Pour tout  $n \leq 1$ , le  $n^{\text{ème}}$  polynôme de Bernstein associé à  $f$  est défini par :

$$B_n(f, t) = \sum_{k=0}^n f\left(\frac{k}{n}\right) \cdot B_{n,k}(t)$$

L'intérêt de cette suite de polynômes est lié au :

*Théorème* : La suite  $(B_n)_{n \in \mathbb{N}}$  converge uniformément vers  $f$  sur  $[0,1]$ .

*Démonstration* — Cf. Davis p. 109.

*Remarque* — Malheureusement la convergence n'est pas très rapide, on peut montrer que l'erreur est un  $O(1/n)$ .

De plus, du point de vue numérique il ne sera pas très pratique de manipuler des polynômes de degré trop élevé.

Un algorithme de calcul est celui de De Casteljeau qui est décrit dans le paragraphe suivant dans le cas des courbes de Bézier.

## 3.3 Les courbes de Bézier

### 3.3.1 Données du problème

On suppose que l'on dispose d'un nuage de points expérimentaux (ou définis par l'utilisateur) :

$$S = \{ P_i = (x_i, y_i) ; 0 \leq i \leq n \}$$

et on voudrait trouver une courbe d'équations paramétriques :

$$\begin{cases} x = x(t) \\ y = y(t) \end{cases} \quad t \in [0,1]$$

qui vérifie :

$$\begin{cases} x\left(\frac{i}{n}\right) = x_i \\ y\left(\frac{i}{n}\right) = y_i \end{cases} \quad (0 \leq i \leq n)$$

En approximant les fonctions  $x$  et  $y$  par leurs approximants de Bernstein respectifs, on aboutit à la définition des courbes de Bézier.

*Remarque* — Les points  $P_i$  sont appelés *points de contrôles* ou *pôles* et le polygone de  $\mathbb{R}^2$  qu'ils définissent est le *polygone descripteur*.

Les courbes de Bézier font partie d'une famille de courbes appelées *courbes à pôles*.

### 3.3.2 Les courbes de Bézier

*Définition* : La *courbe de Bézier* associée aux  $n$  pôles  $P_i$  est la courbe polynomiale de degré  $n$  définie par les équations paramétriques :

$$\begin{cases} P_x(t) = B_n(x, t) \\ P_y(t) = B_n(y, t) \end{cases} \quad t \in [0, 1]$$

ou encore par :

$$\forall t \in [0, 1], P(t) = \sum_{k=0}^n B_{n,k}(t) \cdot P_k$$

Les propriétés de ces courbes sont résumées dans le :

*Théorème* : (i)  $P(0) = P_0$  ;  $P(1) = P_n$  ( $P_0$  et  $P_n$  sont donc les extrémités de la courbe) ;

(ii) la courbe de Bézier est contenue dans l'enveloppe convexe des points de contrôles ;

(iii) la tangente en  $P_0$  à la courbe de Bézier est dirigée par  $\vec{P_0P_1}$  et la tangente en  $P_n$  est dirigée par  $\vec{P_{n-1}P_n}$  ;

(iv) chaque point de contrôle exerce sur la courbe une attraction qui est proportionnelle à  $B_{n,k}(t)$ .

*Démonstration* — Dony p. 111.

On a donc l'allure de courbe suivante :

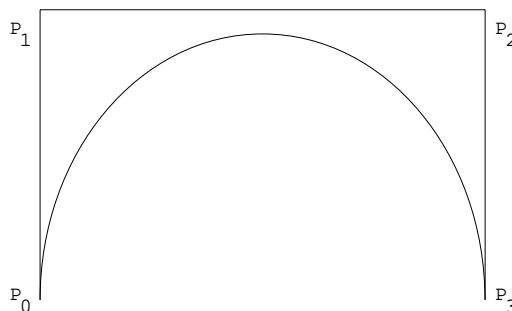


Figure 5.5



*Remarque* — Avec la récurrence (v) du paragraphe (3.2.1) sur les polynômes de Bernstein, en notant  $P^{(u,v)}$  la courbe de Bézier associée au polygone  $\{P_i; u \leq i \leq v\}$ , on déduit que :

$$P^{(0,n)}(t) = (1-t) \cdot P^{(0,n-1)}(t) + t \cdot P^{(1,n)}(t)$$

c'est-à-dire que la courbe de Bézier associée au polygone  $\{P_i; 0 \leq i \leq n\}$  se déduit des courbes associées aux polygones  $\{P_i; 0 \leq i \leq n-1\}$  et  $\{P_i; 1 \leq i \leq n\}$ , ce qui permet de donner un algorithme de construction de telles courbes.

### 3.3.3 Algorithme de De Casteljeau

De ce qui précède, on déduit que si on pose, pour  $t \in [0,1]$  :

$$P_k^{(0)}(t) = P_k \quad (0 \leq k \leq n)$$

$$P_k^{(j)}(t) = (1-t) \cdot P_k^{(j-1)}(t) + t \cdot P_{k+1}^{(j-1)}(t) \quad (0 \leq k \leq n-j; 1 \leq j \leq n)$$

on a alors :

$$P^{(0,n)}(t) = \sum_{k=0}^{n-j} B_{n-j,k}(t) \cdot P_k^{(j)}(t) \quad (j = 0, 1, \dots, n)$$

Ce qui se démontre par récurrence sur  $j$  en utilisant la remarque précédente.

Et pour  $j = n$ , on a alors :

$$P^{(0,n)}(t) = P_0^{(n)}$$

Ce qui nous donne la programmation structurée :

*PROCEDURE Bézier*(Entrée  $n$  : Entier ;  $P$  : TableauDePoints ;  $t$  : Réel ; Sortie  $B$  : Point) ;

Début

$Q = P ; \{ Q_k = P_k^{(0)} \}$

Pour  $j$  allant de 1 à  $n$  faire

Début

Pour  $k$  allant de 0 à  $n-j$  faire

Début

$P_k = (1-t) \cdot Q_k + t \cdot Q_{k+1} ; \{ P_k = P_k^{(j)}, Q_k = P_k^{(j-1)} \}$

Fin ;

Pour  $k$  allant de 0 à  $n-j+1$  faire  $Q_k = P_k$  ;

Fin ;

$B = P_0$  ;

Fin ;

*Remarque 1* — Cet algorithme peut aussi être utilisé pour calculer l'approximant de Bernstein, les calculs se faisant seulement sur la deuxième composante  $y$ .

*Remarque 2* — Cet algorithme est très performant car il évite le calcul récursif des fonctions de base, qui peut être très lent pour  $n$  grand.

*Remarque 3* — Si on désire que la courbe de Bézier passe plus près de l'un des points de contrôle, il suffit de le compter plusieurs fois, on dit que l'on augmente sa multiplicité.

*Remarque 4* — Dans le cas des courbes fermées, on aura en général la situation de la figure 5.6 :

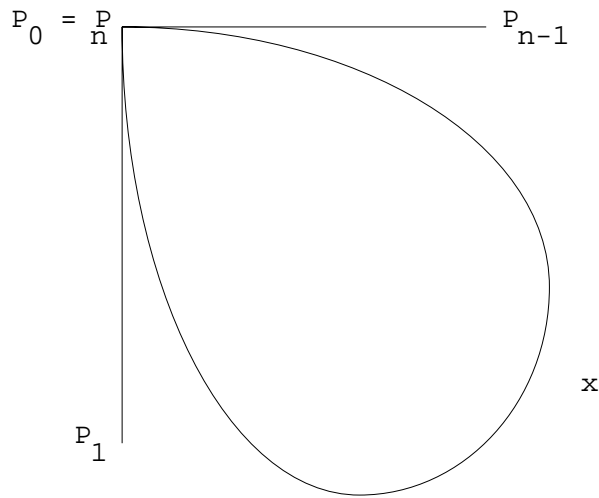


Figure 5.6

Des exemples de tracés sont donnés sur la figure 5.7 où on travaille avec le même fichier de points, en affectant successivement à tous les points les multiplicités 1, 2, 3 et 4.

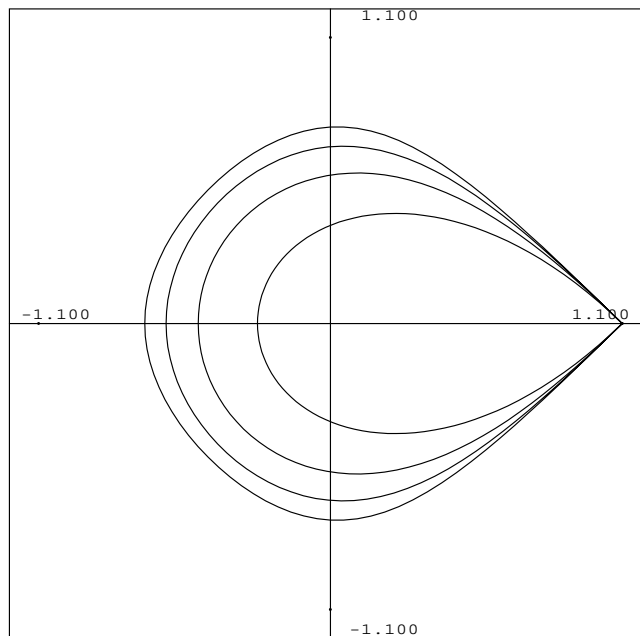


Figure 5.7

### 3.4 Les surfaces de Bézier

#### 3.4.1 Définition

Les définitions précédentes s'étendent au cas des surfaces de  $\mathbb{R}^3$ .

Les données sont un nuage de points de contrôle de  $\mathbb{R}^3$ , formant le *polyèdre descripteur* :

$$S = \{ P_{i,j} ; 0 \leq i \leq n ; 0 \leq j \leq m \}$$

et le *carreau de Bézier* correspondant est la surface d'équations paramétriques :

$$B: [0,1]^2 \rightarrow \mathbb{R}^3$$

$$(u, v) \mapsto \sum_{\substack{0 \leq i \leq n \\ 0 \leq j \leq m}} B_{n,i}(u) \cdot B_{m,j}(v) \cdot P_{i,j}$$

Une carrosserie automobile peut être définie avec une vingtaine de carreaux de Bézier correctement raccordés.

*Remarque* — On peut écrire que  $B(u, v) = \sum_{i=0}^n B_{n,i}(u) \cdot P_i(v)$ , où on a noté

$$P_i(v) = \sum B_{m,j}(v) \cdot P_{i,j}. \text{ C'est-à-dire que pour tout } v \text{ fixé, } u \rightarrow B(u, v) \text{ est la courbe}$$

de Bézier associée au polygone descripteur  $\{ P_i(v) ; 1 \leq i \leq n \}$ . On a un résultat analogue à  $u$  fixé.

### 3.4.2 Propriétés

Les propriétés de ces surfaces sont résumées dans le théorème ci-dessous.

*Théorème* : (i) Le bord de la surface de Bézier définie par le polyèdre descripteur  $S$  est formé des courbes de Bézier définies par les polygones  $\Gamma_0 = \{P_{0,0}, P_{0,1}, \dots, P_{0,m}\}$ ,  $\Gamma_n = \{P_{n,0}, P_{n,1}, \dots, P_{n,m}\}$ ,  $\Gamma'_0 = \{P_{0,0}, P_{1,0}, \dots, P_{n,0}\}$  et  $\Gamma'_m = \{P_{0,m}, P_{1,m}, \dots, P_{n,m}\}$ . Soit :

$$B(0, v) = \sum_{j=0}^m B_{m,j}(v) \cdot P_{0,j}$$

$$B(1, v) = \sum_{j=0}^m B_{m,j}(v) \cdot P_{n,j}$$

$$B(u, 0) = \sum_{i=0}^n B_{n,i}(u) \cdot P_{i,0}$$

$$B(u, 1) = \sum_{i=0}^n B_{n,i}(u) \cdot P_{i,m}$$

(ii) la surface de Bézier est contenue dans l'enveloppe convexe des points de contrôle ;

(iii) le plan tangent à la surface de Bézier en  $P_{0,0}$  est défini par  $\left( \begin{matrix} \overrightarrow{P_{0,0}P_{0,1}} & \overrightarrow{P_{0,0}P_{1,0}} \end{matrix} \right)$ ; en  $P_{n,0}$  par  $\left( \begin{matrix} \overrightarrow{P_{n,0}P_{n,1}} & \overrightarrow{P_{n,0}P_{n-1,0}} \end{matrix} \right)$ ; en  $P_{n,m}$  par  $\left( \begin{matrix} \overrightarrow{P_{n,m}P_{n,m-1}} & \overrightarrow{P_{n,m}P_{n-1,m}} \end{matrix} \right)$ ; et en  $P_{0,m}$  par  $\left( \begin{matrix} \overrightarrow{P_{0,m}P_{0,m-1}} & \overrightarrow{P_{0,m}P_{1,m}} \end{matrix} \right)$ ;

### 3.4.3 Algorithme de De-Casteljeau

En utilisant la remarque du paragraphe (3.4.1), il suffit d'utiliser l'algorithme de De-Casteljeau  $n + 1$  fois pour calculer les  $P_i(v)$ , puis de nouveau cet algorithme pour les  $B(u,v)$ .

## 3.5 Les courbes B-Splines

### 3.5.1 Introduction

Les courbes de Bézier présentent les inconvénients suivants :

- ce sont des courbes polynomiales de degré  $n$  pour un polygone à  $n + 1$  points, ce qui entraînera une forte instabilité numérique pour  $n$  grand ;
- la modification d'un seul point de contrôle va se répercuter sur toute la courbe (caractère global des courbes de Bézier) ;
- comme on peut le voir sur les exemples de tracés du paragraphe précédent, la courbe peut passer assez loin des points de contrôle ; pour s'en rapprocher on devra augmenter la multiplicité des points ce qui entraînera une augmentation du degré des fonctions polynomiales à traiter.

C'est en relation avec ces problèmes que dans les années 1970, *M.G. Cox* et *C. De Boor* ont introduit la notion de fonction B-Spline comme généralisation des courbes de Bézier.

Les courbes B-Splines sont définies localement par des expressions polynomiales de degré  $j$  assez faible (usuellement 3 ou 4), ce dernier étant indépendant du nombre de points et la modification d'un point de contrôle ne se fait ressentir qu'au voisinage de ce dernier. De plus elles ont les mêmes propriétés que celles de Bézier.

### 3.5.2 Fonctions de base B-Splines

*Définition* : Soient  $n \geq 1$  et  $j$  dans  $\mathbb{N}$ , avec  $2 \leq j \leq n + 1$ . On appelle alors *vecteur nodal* associé à  $(n,j)$ , le vecteur  $v$  de  $\mathbb{N}^{n+j+1}$  défini par :

$$\begin{cases} v_i = 0; & 0 \leq i \leq j - 1 \\ v_i = i - j + 1; & j \leq i \leq n \\ v_i = n - j + 2; & n + 1 \leq i \leq n + j \end{cases}$$

*Définition* : Soient  $n \geq 1$  et  $j$  dans  $\mathbb{N}$ , avec  $2 \leq j \leq n + 1$ . En notant  $v$  le vecteur nodal associé à  $(n,j)$ , on appelle *base B-Spline d'ordre  $j$* , la suite de fonctions  $\{ N_{i,j} ; 0 \leq i \leq n \}$ , définie sur  $[0, t_{\max}]$ , avec  $t_{\max} = n - j + 2$ , par la récurrence :

$$N_{i,j}(t) = \begin{cases} 1 & \text{si } v_i \leq t < v_{i+1} \\ 0 & \text{sinon} \end{cases} \quad (0 \leq i \leq n)$$

$$N_{i,k}(t) = \alpha_{i,k} \cdot (t - v_i) \cdot N_{i,k-1}(t) + \beta_{i,k} \cdot (v_{i+k} - t) \cdot N_{i+1,k-1}(t)$$

pour  $2 \leq k \leq j$ ,  $0 \leq i \leq n$ , où :

$$\alpha_{i,k} = \begin{cases} \frac{1}{v_{i+k-1} - v_i} & \text{si } v_{i+k-1} \neq v_i \\ 0 & \text{sinon} \end{cases}$$

$$\beta_{i,k} = \begin{cases} \frac{1}{v_{i+k} - v_{i+1}} & \text{si } v_{i+k} \neq v_{i+1} \\ 0 & \text{sinon} \end{cases}$$

*Remarque 1* — L'algorithme de construction des  $N_{i,j}$  défini ci-dessus est dû à Cox et De Boor.

*Remarque 2* — Les fonctions  $N_{i,j}$  sont aussi appelées *fonctions mélanges* ou *fonctions pondérantes*.

*Exemple* — Pour  $n = 3$ ,  $j = 2$ , le vecteur nodal est  $v = (0,0,1,2,3,3)$  et :

$$N_{0,2}(t) = \begin{cases} 1-t & \text{si } t \in [0,1] \\ 0 & \text{sinon} \end{cases}$$

$$N_{1,2}(t) = \begin{cases} t & \text{si } t \in [0,1] \\ 2-t & \text{si } t \in [1,2] \\ 0 & \text{si } t \in [2,3] \end{cases}$$

$$N_{2,2}(t) = \begin{cases} 0 & \text{si } t \in [0,1] \\ t-1 & \text{si } t \in [1,2] \\ 3-t & \text{si } t \in [2,3] \end{cases}$$

$$N_{3,2}(t) = \begin{cases} t-2 & \text{si } t \in [2,3] \\ 0 & \text{sinon} \end{cases}$$

Les propriétés de ces fonctions sont résumées dans le :

*Théorème* : (i)  $\forall t \in [0, t_{\max}]$ ,  $N_{i,j}(t) \geq 0$ , avec, pour  $j \geq 2$  :

$$\begin{cases} N_{i,j}(t) > 0 & \text{si } t \in ]v_i, v_{i+j}[ \\ 0 & \text{sinon} \end{cases}$$

et sur  $]v_i, v_{i+j}[$ ,  $N_{i,j}$  est polynomiale de degré  $j - 1$  ;

(ii)  $\forall t \in ]0, t_{\max}[$ ,  $\sum_{i=0}^n N_{i,j}(t) = 1$  ;

(iii)

$$\left\{ \begin{array}{l} N_{i,j}(0) = 0 \quad (1 \leq i \leq n) \\ N_{0,j}(0) = 1 \\ N_{i,j}(t_{\max}) = 0 \quad (0 \leq i \leq n - 1) \\ N_{n,j}(t_{\max}) = 1 \end{array} \right.$$

(iv)  $\{ N_{i,n+1} ; 0 \leq i \leq n \}$  est la base de Bernstein de degré  $n$  sur  $[0,1]$  ( $t_{\max} = 1$ ).

*Tracé des fonctions de base B-Spline* — Tout d'abord on définit une procédure qui calcule les composantes du vecteur nodal  $v$ , soit :

*PROCEDURE VecteurNodal(Entrée  $n, j$  : Entier ; Sortie  $v$  : VecteurEntier) ;*

*Début*

*Pour  $i$  allant de  $0$  à  $j - 1$  faire  $v_i = 0$  ;*

*Pour  $i$  allant de  $j$  à  $n$  faire  $v_i = i - j + 1$  ;*

*Pour  $i$  allant de  $n + 1$  à  $n + j$  faire  $v_i = n - j + 2$  ;*

*Fin ;*

Ensuite, on peut définir les fonctions de base de manière récursive de la manière suivante :

*FONCTION BaseBSpline(Entrée  $n, i, j$  : Entier ;  $v$  : VecteurEntier ;  $t$  : Réel) : Réel ;*

*Début*

*Si  $(t < v_j)$  ou  $(t \geq v_{i+j})$*

*Alors Début*

*BaseBSpline = 0 ;*

*Sortir ;*

*Fin*

*Sinon Début*

*Si  $j = 1$*

*Alors Début*

*BaseBSpline = 1 ; { puisque  $t \in [v_i, v_{i+1}]$  }*

*Sortir ;*

*Fin*

*Sinon Début*

*Aux1 = BaseBSpline( $n, i, j - 1, v, t$ ) ;*

*Si (Aux1 = 0) ou  $(v_{i+j-1} = v_i)$*

*Alors S1 = 0*

*Sinon S1 =  $\frac{t - v_i}{v_{i+j-1} - v_i} \cdot \text{Aux1}$  ;*

*Aux2 = BaseBSpline( $n, i+1, j - 1, v, t$ ) ;*

*Si (Aux2 = 0) ou  $(v_{i+j} = v_{i+1})$*

*Alors S2 = 0*

*Sinon S2 =  $\frac{v_{i+j} - t}{v_{i+j} - v_{i+1}} \cdot \text{Aux2}$  ;*

*BaseBSpline = S1 + S2 ;*

*Fin ;*

*Fin ;*

*Fin ;*

En fait, comme on le verra dans le paragraphe suivant, on peut se passer de la récursivité qui peut donner un calcul très long pour de grandes valeurs de  $j$ .

Sur la figure 5.8, on donne les graphes des fonctions de base pour  $n = 5$  et  $j = 3$ .

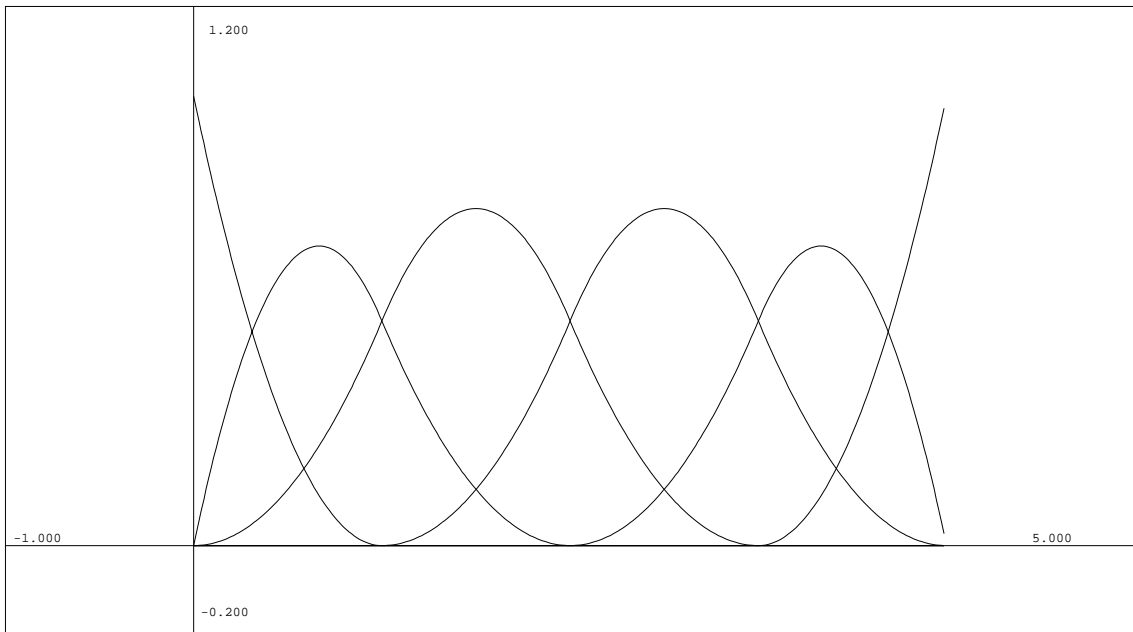


Figure 5.8

### 3.5.3 Courbes B-Splines

On reprend les notations du paragraphe (3.3.1) et on se donne  $n \geq 1$  et  $2 \leq j \leq n + 1$  dans  $\mathbb{N}$ .

*Définition :* On appelle *courbe B-Spline d'ordre  $j$*  (ou de degré  $j - 1$ ) associée au polygone descripteur  $S$ , la courbe d'équation paramétrique :

$$B(t) = \sum_{i=0}^n N_{i,j}(t) \cdot P_i \quad \text{si } t \in [0, t_{\max}] = [0, n - j + 2]$$

où  $\{N_{i,j}(t) ; 0 \leq i \leq n\}$  est la base B-Spline associée à  $(n, j)$ .

*Remarque* — On a  $[0, t_{\max}] = [v_{j-1}, v_{n+1}]$  et pour  $k = j - 1, \dots, n$ , si  $t$  est dans  $[v_k, v_{k+1}]$ , on aura  $t \notin ]v_i, v_{i+j}[$  pour  $i \leq k - j$  ou  $i \geq k + 1$ , ce qui donnera  $N_{i,j}(t) = 0$  pour de tels  $i$ . En définitive on a donc :

$$B(t) = \sum_{i=k-j+1}^k N_{i,j}(t) \cdot P_i \quad (k = j - 1, \dots, n; t \in [v_k, v_{k+1}])$$

On a donc seulement  $j$  termes à calculer dans la somme et non pas  $n + 1$  comme dans le cas des courbes de Bézier. De plus les polynômes que l'on manipule sont de degré au plus  $j$ . Tout cela étant indépendant du nombre de points  $n + 1$ .

*Exemple* — Pour  $n = 3, j = 2$  :  $B(t) = \sum_{i=0}^3 N_{i,3}(t) \cdot P_i \quad t \in [0,3]$ .

Sur  $[0,1] = [v_1, v_2]$  ( $k = j - 1 = 1$ ), on a :

$$B(t) = N_{0,2}(t) \cdot P_0 + N_{1,2}(t) \cdot P_1 = (1 - t) \cdot P_0 + t \cdot P_1 ;$$

sur  $[1,2] = [v_2, v_3]$  ( $k = j = 2$ ), on a :

$$B(t) = N_{1,2}(t) \cdot P_1 + N_{2,2}(t) \cdot P_2 = (2 - t) \cdot P_1 + (t - 1) \cdot P_2 ;$$

Sur  $[2,3] = [v_3, v_4]$  ( $k = n$ ), on a :

$$B(t) = N_{2,2}(t) \cdot P_2 + N_{3,2}(t) \cdot P_3 = (3 - t) \cdot P_2 + (t - 2) \cdot P_3 ;$$

En utilisant la procédure de calcul des fonctions de base B-Spline, on en déduit immédiatement un algorithme de calcul des courbes B-Splines.

L'indice  $k \in \{j - 1, \dots, n\}$  tel que  $t \in [v_k, v_{k+1}[$ , pour  $t$  dans  $[0, t_{\max}[$ , est simplement donné par :  $k = \text{PartieEntière}(t) + j - 1$

De plus, pour  $t = 0$ , on a :

$$B(0) = \sum_{i=0}^n N_{i,j}(0) \cdot P_i = 1 \cdot P_0 = P_0$$

et pour  $t = t_{\max}$ , on a :

$$B(t_{\max}) = \sum_{i=0}^n N_{i,j}(t_{\max}) \cdot P_i = 1 \cdot P_n = P_n$$

Si on ne veut pas utiliser la récursivité (qui donnera un calcul assez lent pour de grandes valeurs de  $j$  et même très lent pour  $j = n + 1$ ), on peut s'inspirer de l'algorithme de De Casteljeau. L'algorithme obtenu est connu sous le nom d'algorithme d'Oslo. (Cf. Math. et C. A. O., Vol. 6, p. 316 et 317).

La procédure de calcul est alors la suivante :

*PROCEDURE B-Spline*(Entrée  $n, j$  : Entier ;  $v$  : VecteurEntier ;  $P$  : TableauDePoints ;  $t$  : Réel ;  
Sortie  $B$  : Point) ;

*Début*

    Si  $t \geq n - j + 2$

        Alors *Début*

$B = P_n$  ;

*Sortir* ;

*Fin* ;

$k = \text{PartieEntière}(t) + j - 1$  ;

    Pour  $i$  allant de  $k - j + 1$  à  $k$  faire  $Q_i = P_i$  ;

    Pour  $p$  allant de  $j$  à 2 faire

*Début*

            Pour  $i$  allant de  $k - p + 1$  à  $k$  faire

*Début*

                    Si  $v_{i+p-1} = v_i$

                        Alors *Début*

$S1 = 0$  ;  $S2 = 0$  ;

*Fin*

                    Sinon *Début*

$$S1 = \frac{t - v_i}{v_{i+p-1} - v_i} \cdot Q_i ;$$

$$S2 = \frac{v_{i+p-1} - t}{v_{i+p-1} - v_i} \cdot Q_{i-1} ;$$

*Fin*

*Fin*

*Fin*



*Fin ;*  
 $P = S1 + S2 ;$   
*Fin ;*  
*Pour i Allant de  $k - p + 1$  à  $k$  faire  $Q_i = P_i ;$*   
*Fin ;*  
 $B = P_k ;$   
*Fin ;*

Sur la figure 5.9, on donne un exemple de tracé, avec  $n = 17$ ,  $j = 3$  et le fichier de points suivant : 1 -2 ; 1 -1 ; 1.4 0.6 ; 1 1.8 ; -0.6 1.8 ; -1 1.6 ; -1 1 ; -0.6 0.8 ; -1.6 0 ; -1.4 -0.2 ; -1 0 ; -1 -1 ; -0.6 -1.2 ; -1 -1.2 ; -1 -1.8 ; -0.4 -1.6 ; -0.4 -2 ; 1 -2

Les temps de calcul respectifs, sur un micro-ordinateur équipé d'un coprocesseur 80387, sont : en utilisant la récursivité : 82 centièmes de secondes ; en utilisant l'algorithme d'Oslo : 33 centièmes de secondes.

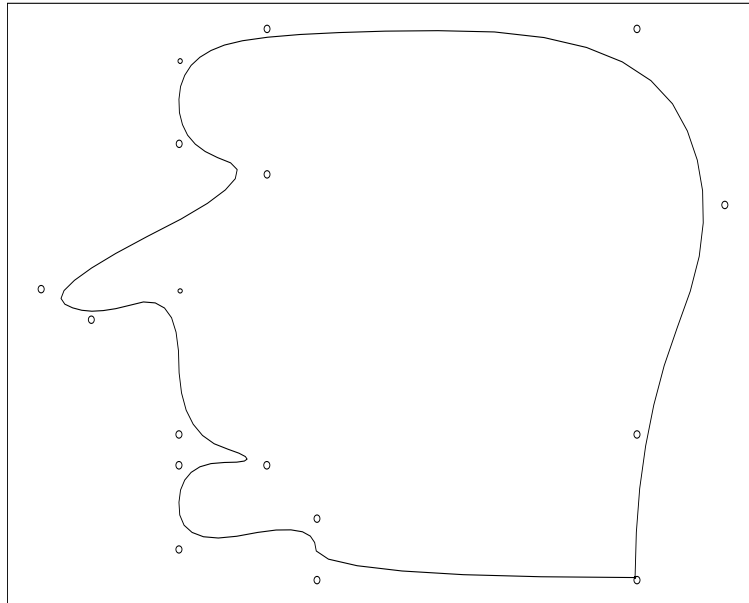


Figure 5.9

## 4. Problèmes d'interpolation

### 4.1 Introduction. Position des problèmes

On se place toujours dans l'hypothèse où une fonction continue  $f$ , définie sur un intervalle  $[a,b]$  est connue en un nombre fini de points  $x_i$  avec  $i = 0, 1, \dots, n$ . Mais dans ce paragraphe, on cherche une fonction  $\phi$  qui coïncide avec  $f$  en tous les points  $x_i$ . Cette fonction  $\phi$  sera cherchée aussi simple que possible (polynomiale ou polynomiale par morceaux). On dira alors que  $\phi$  *interpole*  $f$ .

Historiquement ce sont les problèmes de *tabulations numériques* des fonctions transcendentes élémentaires, puis à partir du 17<sup>ème</sup> siècle le calcul approché des dérivées et des intégrales qui sont à l'origine des méthodes d'interpolation.

Un problème important qui se posera pour toute méthode d'interpolation sera d'évaluer *l'erreur d'interpolation*  $R(x) = |f(x) - \Phi(x)|$  en tout point  $x$  de  $[a,b]$ .

On peut aussi se demander s'il est possible de reconstituer  $f$  comme limite des fonctions  $\phi$  quand le module de la subdivision :

$$\Delta(S) = \text{Sup} \{ |x_{j+1} - x_j|; 0 \leq j \leq n-1 \}$$

tend vers 0.

Par exemple dans le cas de l'interpolation polynomiale (Cf. (4.2)), en prenant une subdivision à pas constant (i.e.  $x_{j+1} - x_j = \Delta = \frac{b-a}{n}$ ) et en notant  $P_n$  le polynôme de degré  $n$  qui coïncide avec  $f$  en tous les  $x_i$ , peut-on espérer que la suite  $(P_n)_{n \in \mathbb{N}}$  converge uniformément vers  $f$  sur  $[a,b]$ ?

Depuis Lagrange en 1795 jusqu'à Runge en 1901 on pensait que la réponse était positive. Mais en fait il n'en est rien comme le montre le contre-exemple de Runge  $f(x) = \frac{1}{1+8 \cdot x^2}$  sur  $[0,1]$  : en traçant les graphes de  $f$  et des polynômes interpolateurs successifs, on constate des oscillations aux bords de l'intervalle (effets de bord ou phénomène de Runge : Cf. paragraphe (4.2)).

Dans le Calcul Infinitésimal de Dieudonné (p. 319), on montre que  $\lim_{n \rightarrow +\infty} |f(1) - P_n(1)| = +\infty$ . Ce qui peut se constater sur la figure 5.10.

Pour résoudre ce genre de problème on aura recours à un procédé plus élaboré qui consiste à approcher  $f$  par des fonctions polynomiales (usuellement de degré 3) sur chacun des intervalles  $[x_i, x_{i+1}]$  en imposant des conditions de régularités aux points communs à deux intervalles. C'est le principe de l'interpolation spline cubique (Cf. paragraphe (4.3)).

## 4.2 L'interpolation polynomiale de Lagrange

### 4.2.1 L'interpolation linéaire

L'interpolation linéaire permet de calculer approximativement une valeur intermédiaire dans une table numérique. Par exemple, pour trouver une valeur approchée de la fonction  $f$  en  $x = 1.5$  sachant que  $f(1) = 2.434$  et  $f(2) = 5.842$ , le principe consiste à remplacer le graphe de  $f$  sur  $[1,2]$  par la droite passant par les deux points  $A = (1, f(1))$  et  $B = (2, f(2))$  (on dit qu'on a linéarisé  $f$ ), ce qui donne  $y = 4.138$  comme valeur approchée de  $f(x)$ .

Cette valeur approchée s'écrit sous la forme :

$$y = \frac{x - x_1}{x_0 - x_1} \cdot y_0 + \frac{x - x_0}{x_1 - x_0} \cdot y_1$$

On peut penser que cette approximation sera meilleure si au lieu d'utiliser des polynômes de degré 1 on utilise des polynômes de degré  $p > 1$ . C'est le principe des polynômes d'interpolation de Lagrange.

### 4.2.2 Le théorème d'interpolation de Lagrange

*Théorème :* Etant donnés  $n + 1$  réels distincts  $x_0 < x_1 < \dots < x_n$  et  $n + 1$  réels  $y_i$  avec  $0 \leq i \leq n$ , il existe un unique polynôme  $P$  de degré au plus  $n$  tel que :

$$P(x_i) = y_i \quad (0 \leq i \leq n)$$

Ce polynôme est donné par :

$$P(x) = \sum_{i=0}^n y_i \cdot L_i(x)$$

où :

$$L_i(x) = \frac{\prod_{\substack{j=0 \\ j \neq i}}^n (x - x_j)}{\prod_{\substack{j=0 \\ j \neq i}}^n (x_i - x_j)} \quad (0 \leq i \leq n)$$

*Démonstration* — En notant  $E = \mathbb{R}[x]$  l'espace vectoriel des polynômes à coefficients réels et  $u : E \rightarrow \mathbb{R}^{n+1}$  définie par :

$$\forall P \in E, u(P) = (P(x_0), P(x_1), \dots, P(x_n))$$

tout revient à montrer que  $u$  est surjective et que sa restriction à  $E_n = \mathbb{R}_n[x]$  est bijective.

Le noyau de  $u$  est formé des multiples du polynôme  $\omega(x) = \prod_{i=0}^n (x - x_i)$  et avec la division euclidienne, on peut écrire tout polynôme  $P$  sous la forme :

$$P = Q \cdot \omega + R \quad (R = 0 \text{ ou } d^\circ(R) \leq n)$$

On a donc :

$$E = E_n \oplus \text{Ker}(u)$$

ce qui suffit à prouver que  $u$  est un isomorphisme de  $E_n$  sur  $\mathbb{R}^{n+1}$ .

*Remarque* — En prenant  $\{1, x, \dots, x^n\}$  comme base de  $E_n$  et la base canonique sur  $\mathbb{R}^{n+1}$ , la matrice de la restriction de  $u$  à  $E_n$  est une matrice de *Van Der Monde* dont le déterminant est  $\prod_{j < i} (x_i - x_j) \neq 0$ . On a donc ainsi une deuxième preuve de

théorème.

Dans le cas où  $y_i = f(x_i)$ , avec  $f$  assez régulière, on peut donner une majoration de l'erreur d'interpolation :

*Théorème* : Soit  $f : [x_0, x_n] \rightarrow \mathbb{R}$  de classe  $C^{n+1}$  et  $P$  le polynôme d'interpolation de Lagrange défini par  $P(x_i) = f(x_i)$ , pour  $i = 0, 1, \dots, n$ , avec  $x_0 < x_1 < \dots < x_n$ . L'erreur d'interpolation vérifie alors :

$$R(x) = |f(x) - P(x)| \leq M_{n+1} \cdot \frac{|\omega(x)|}{(n+1)!}$$

où :  $\omega(x) = \prod_{i=0}^n (x - x_i)$  et  $M_{n+1} = \text{Sup} \left\{ |f^{(n+1)}(x)| ; x \in [x_0, x_n] \right\}$

*Démonstration* — Théodor p. 104.

### 4.2.3 Algorithme de Neville pour calculer le polynôme de Lagrange

Une programmation directe du polynôme de Lagrange est élémentaire, mais peu performante. Des algorithmes plus performants ont été donnés par Aitken (Cf. Théodor p. 94) puis par Neville. Nous utiliserons celui de Neville qui est plus performant. Le principe en est résumé dans le :

*Théorème* : Soit  $(P_{ik})$  la suite de polynômes définie par la récurrence :

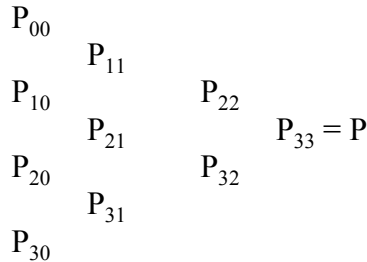
$$P_{i0}(x) = y_i \quad (0 \leq i \leq n)$$

$$P_{ik}(x) = \frac{(x - x_{i-k}) \cdot P_{i,k-1}(x) - (x - x_i) \cdot P_{i-1,k-1}(x)}{x_i - x_{i-k}} \quad (1 \leq k \leq i \leq n)$$

Alors  $P_{n,n}$  est le polynôme de Lagrange.

*Démonstration* — Cf. Stoer et Burlisch, par. (2.1.2).

Une construction peut être visualisée par la disposition triangulaire ci-dessous :



L'algorithme de Neville peut aussi s'écrire sous la forme :

$$P_{i,k} = P_{i,k-1} + (x - x_i) \cdot \frac{P_{i,k-1} - P_{i-1,k-1}}{x_i - x_{i-k}} \quad (1 \leq k \leq i \leq n)$$

Ce qui nous donne la programmation structurée :

*FONCTION* LagrangeNeville(*Entrée*  $n$  : Entier ;  $x, y$  : Vecteur ;  $t$  : Réel) : Réel ;

*Début*

*Pour*  $i$  allant de 0 à  $n$  faire

*Début*

$P_i = y_i$  ;

*Pour*  $j$  allant de  $i - 1$  à 0 faire

*Début*

$$P_j = P_{j+1} + \frac{P_{j+1} - P_j}{x_i - x_j} \cdot (t - x_j) ;$$

*Fin* ;

*Fin* ;

    LagrangeNeville =  $P_0$  ;

*Fin* ;

#### 4.2.4 Forme de Newton du polynôme de Lagrange

L'algorithme de Neville est pratique pour calculer quelques valeurs du polynôme de Lagrange, mais les calculs sont relativement lents. Si on doit calculer un grand nombre de valeurs (par exemple en vue du tracé graphique), on lui préférera un algorithme du même type déduit de la forme de Newton du polynôme de Lagrange.

L'idée de Newton est d'écrire le polynôme d'interpolation sous la forme :

$$P(x) = P_0 + P_1 \cdot (x - x_0) + \dots + P_n \cdot (x - x_0) \cdot \dots \cdot (x - x_{n-1})$$

et les conditions  $P(x_i) = y_i$  nous conduisent à la résolution d'un système triangulaire inférieur aux inconnues  $P_i$ .

La encore la résolution directe est élémentaire, mais en s'inspirant de l'algorithme de Neville on aboutit à l'algorithme plus performant décrit ci-dessous (Cf. Stoer et Burlisch : § (2.1.3)).

*PROCEDURE PolynômeLagrange(Entrée  $n$  : Entier ;  $x, y$  : Vecteur ; Sortie  $P$  : Vecteur) ;*

*Début*

*Pour  $i$  allant de 0 à  $n$  faire*

*Début*

$Q_i = y_i$  ;

*Pour  $j$  allant de  $i - 1$  à 0 faire  $Q_j = \frac{Q_{j+1} - Q_j}{x_i - x_j}$  ;*

$P_i = Q_0$  ;

*Fin ;*

*Fin ;*

L'évaluation du polynôme  $P$  se fait alors en utilisant l'algorithme de Horner, soit :

*FONCTION LagrangeNewton(Entrée  $n$  : Entier ;  $x, P$  : Vecteur ;  $t$  : Réel) : Réel ;*

*Début*

$Aux = P$  ;

*Pour  $i$  allant de  $n - 1$  à 0 faire  $Aux = P_i + Aux \cdot (t - x_i)$  ;*

$LagrangeNewton = Aux$  ;

*Fin ;*

*Remarque* — Les coefficients  $P_k$  qui interviennent dans la forme de Newton sont aussi appelés les *différences divisées* associées aux  $x_i$ .

La figure 5.10 donne un exemple de tracé du polynôme de Lagrange. Il s'agit de l'interpolée de la fonction  $f(x) = \frac{1}{1+x^2}$  sur  $[-5,5]$  en prenant 11 points équidistants ( $n = 10$ ). On peut constater le phénomène de Runge aux bords de l'intervalle.

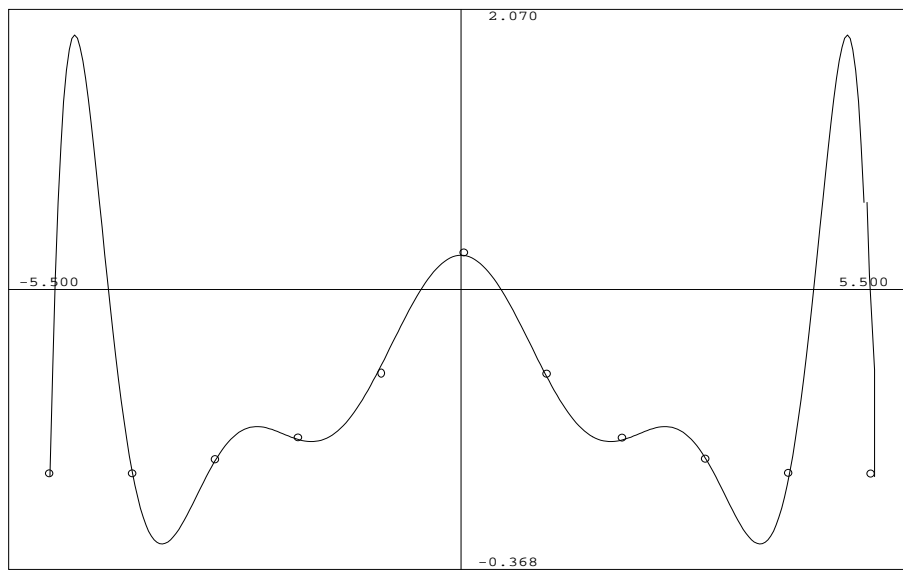


Figure 5.10

#### 4.2.5 Sur le choix des abscisses d'interpolation. Les polynômes de Tchébychev

Le problème est de savoir s'il existe un choix optimal des abscisses  $x_i$ , c'est-à-dire un choix qui minimisera l'erreur d'interpolation.

Une réponse positive à ce problème est donnée par les zéros des *polynômes de Tchébychev* sur  $[-1,1]$ .

Les polynômes de Tchébychev, que nous étudierons plus en détail dans le chapitre sur le calcul numérique des intégrales, peuvent être définis sur  $[-1,1]$  par :

$$\forall x \in [-1,1], T_n(x) = \text{Cos}(n \cdot \text{ArcCos}(x)) \quad (n \geq 0)$$

ou encore :

$$\forall x \in [-1,1], T_n(x) = \sum_{k=0}^{\lfloor \frac{n}{2} \rfloor} (-1)^k \cdot C_n^{2k} \cdot x^{n-2k} \cdot (1-x^2)^k$$

Pour  $n \geq 1$ ,  $T_n$  admet  $n$  racines distinctes dans  $[-1,1]$  données par :

$$x_k = \text{Cos}\left(\frac{2 \cdot k - 1}{2 \cdot n} \cdot \pi\right) \quad (1 \leq k \leq n)$$

On a alors le :

**Théorème :** Sur  $[-1,1]$  le choix optimal des abscisses d'interpolation est donné par les racines du polynôme de Tchébychev de degré  $n$ .

*Démonstration* — Théodor par. (3.5).

*Remarque* — Dans le cas d'un intervalle  $[a,b]$  quelconque, on se ramène facilement à  $[-1,1]$  en faisant le changement de variable :

$$t = \frac{2}{b-a} \cdot x - \frac{b+a}{b-a} \quad (x \in [a,b])$$

Ci-dessous on reprend l'exemple de tracé du paragraphe précédent mais avec les abscisses de Tchébychev. Et on constate que les effets de bords sont moins importants.

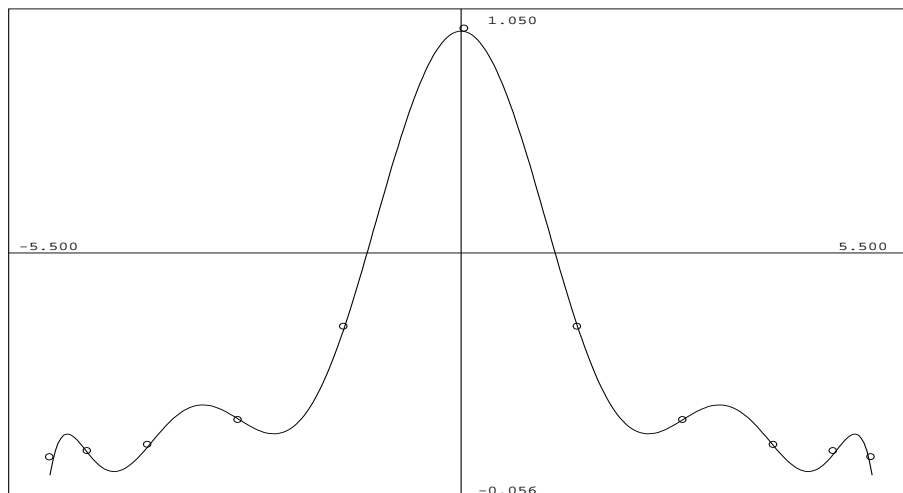


Figure 5.11

### 4.3 Interpolation spline cubique

#### 4.3.1 Introduction

L'interpolation polynomiale présente deux gros inconvénients :

- le coût des calculs est trop élevé dès que le nombre de points est important ;
- pour un grand intervalle d'interpolation, on maîtrise mal les effets de bords (phénomène de Runge).

L'idée des fonctions splines est de travailler localement, c'est-à-dire qu'on cherchera plutôt à approcher le nuage de points expérimentaux par des fonctions polynomiales par morceaux, de degré assez faible. L'exemple le plus simple étant l'interpolation par une fonction affine par morceaux (c'est cette interpolation que l'on utilise pour tracer le graphe d'une fonction sur l'écran d'un ordinateur ou sur une table traçante) :

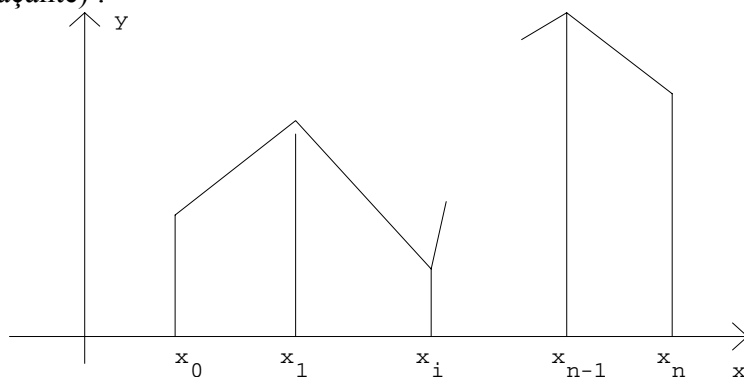


Figure 5.12

Du point de vue du calcul approché des intégrales, cette méthode aboutit à la méthode des trapèzes (Cf. le chapitre « Calcul Numérique des Intégrales »).



Dans ce paragraphe on étudiera le cas des fonctions *splines cubiques*, c'est-à-dire celui où les morceaux de polynômes sont de degré au plus 3.

Le mot anglais « spline » désigne une latte de bois flexible qui était utilisée pour tracer des courbes passant par des points donnés (les noeuds).

Vers 1946 *I. J. Schoenberg* a appelé fonctions splines des fonctions formées de morceaux de polynômes se raccordant correctement entre eux. Ces fonctions sont solutions du problème d'élasticité correspondant à la tige flexible du dessinateur. On a ensuite appelé fonctions splines des solutions approchées des extensions et généralisations du problème précédent. Le résultat n'ayant plus rien à voir avec les polynômes.

### 4.3.2 Position du problème. Notations

Etant donné le nuage de points :

$$S = \{ P_i = (x_i, y_i) ; 1 \leq i \leq n \}$$

on cherche une fonction polynomiale par morceaux, de degré  $\leq 3$ , et de classe  $C^2$ , qui passe exactement par ces points.

De manière plus précise une telle courbe sera définie par une paramétrisation :

$$\begin{aligned} \gamma: [a, b] &\longrightarrow \mathbb{R}^2 \\ t &\longrightarrow \gamma(t) = (x(t), y(t)) \end{aligned}$$

avec une subdivision  $a = t_1 < t_2 < \dots < t_n = b$ , telle que la restriction de  $\gamma$  à  $[t_i, t_{i+1}]$  est polynomiale de degré au plus 3.

En raisonnant séparément sur  $x(t)$  et  $y(t)$ , on est donc ramené à résoudre le problème suivant :

Etant donné un nuage de points :

$$S = \{ P_i = (t_i, x_i) ; 1 \leq i \leq n \}$$

avec :

$$t_1 < t_2 < \dots < t_n$$

on cherche une fonction  $s$  de classe  $C^2$ , polynomiale de degré au plus 3 sur chaque  $[t_i, t_{i+1}]$  telle que  $s(t_i) = x_i$ , pour tout  $i = 1, \dots, n$ .

*Remarque 1* — Le choix de l'espacement des  $t_i$  va influencer la forme de la courbe.

*Remarque 2* — Dans le cas des courbes paramétriques, on pourra prendre  $[1, n]$  comme espace des paramètres avec  $t_i = i$ , si on dispose seulement des  $P_i = (x_i, y_i)$  sans autre indication. Un autre choix judicieux sera fourni par les angles polaires.

On note  $h_i = t_{i+1} - t_i$  pour  $1 \leq i \leq n$ .

On cherche  $s : [t_1, t_n] \longrightarrow \mathbb{R}$  telle que :

$$\begin{cases} s \text{ est de classe } C^2 \text{ sur } [t_1, t_n] \\ s(t) = s_i(t) = a_i \cdot (t - t_i)^3 + b_i \cdot (t - t_i)^2 + c_i \cdot (t - t_i) + d_i \text{ sur } [t_i, t_{i+1}] \\ s(t_i) = x_i, \text{ pour tout } i = 1, \dots, n \end{cases}$$

Le fait que l'on cherche  $s$  de classe  $C^2$  impose les conditions de recollement :

$$s_i^{(k)}(t_{i+1}) = s_{i+1}^{(k)}(t_{i+1}) \text{ pour } k = 0, 1, 2 \text{ et } i = 1, \dots, n-2$$

ce qui revient à imposer que la pente et la courbure de deux cubiques voisines soient identiques.

On va d'abord calculer les coefficients  $a_i$ ,  $b_i$ ,  $c_i$  et  $d_i$  en fonction des  $x_i$  et des  $x''_i = x''(t_i)$ , puis les  $x''_i$  apparaîtront comme solutions d'un système linéaire tridiagonal.

#### 4.3.3 Calcul des coefficients $d_i$

On a  $s(t_i) = d_i = x_i$  pour  $i = 1, \dots, n-1$ . Soit :

$$d_i = x_i \quad (i = 1, \dots, n-1)$$

#### 4.3.4 Calcul des $b_i$ et des $a_i$

On a, sur  $[t_i, t_{i+1}]$ ,  $s''(t) = 6 \cdot a_i(t - t_i) + 2 \cdot b_i$ , ce qui donne :

$$\begin{cases} x''_i = s''(t_i) = 2 \cdot b_i \\ x''_{i+1} = s''(t_{i+1}) = 6 \cdot a_i \cdot h_i + 2 \cdot b_i \end{cases}$$

d'où :

$$\begin{cases} b_i = \frac{x''_i}{2} \\ a_i = \frac{x''_{i+1} - x''_i}{6 \cdot h_i} \end{cases} \quad (i = 1, \dots, n-1)$$

#### 4.3.5 Calcul des $c_i$

On a  $x_{i+1} = a_i \cdot h_i^3 + b_i \cdot h_i^2 + c_i \cdot h_i + d_i$ , ce qui donne avec (4.3.3) et (4.3.4) :

$$c_i = \frac{x_{i+1} - x_i}{h_i} - \frac{x''_i}{2} \cdot h_i - \frac{x''_{i+1} - x''_i}{6} \cdot h_i$$

soit :

$$c_i = \frac{x_{i+1} - x_i}{h_i} - \frac{h_i}{2} \cdot \frac{x''_{i+1} + 2 \cdot x''_i}{3}$$

On calcule donc les coefficients des polynômes  $s_i$  grâce aux formules ci-dessus en fonction des  $x''_i$ . Il reste donc à déterminer ces derniers.

#### 4.3.6 Calcul des $s'_i = x''_i$

Avec les conditions  $y'_i = s'(t_i) = s'_{i-1}(t_i) = s'_i(t_i)$ , on déduit que les  $x''_i$  sont solutions du système :

$$h_{i-1} \cdot x''_{i-1} + 2 \cdot (h_{i-1} + h_i) \cdot x''_i + h_i \cdot x''_{i+1} = 6 \cdot \left( \frac{x_{i+1} - x_i}{h_i} - \frac{x_i - x_{i-1}}{h_{i-1}} \right)$$

pour  $i = 2, \dots, n-1$ .

Soit un système de  $n-2$  équations à  $n$  inconnues. Il nous manque donc deux équations. On complète alors en posant des *conditions aux limites* qui devront fournir deux équations indépendantes des  $n-2$  précédentes. Classiquement, on a les conditions suivantes :

*Splines naturels* — On pose  $x''_1 = x''_n = 0$  ;

*Splines périodiques* — Si la fonction considérée est périodique, on aura  $x''(t_1) = x''(t_n)$  et il nous manque seulement une équation que l'on peut obtenir en disant que dans le système ci-dessus on peut faire  $i = n$  (Cf. § (4.3.10)) ;

*Conditions de De Boor* — On impose une condition de continuité  $C^3$  en  $t_2$  et  $t_{n-1}$ , ce qui se traduit par  $a_1 = a_2$  et  $a_{n-1} = a_n$  et fournit les deux équations manquantes sur les  $s_i$ .

*Remarque* — Ces conditions peuvent aussi porter sur des points intérieurs, mais dans tous les cas des conditions différentes donneront des courbes d'allures différentes. On choisira ces conditions en fonction de la nature du problème.

### 4.3.7 L'interpolation spline naturelle

On pose donc, dans ce paragraphe :  $x''_1 = x''_n = 0$  et les inconnues sont donc  $x''_2, \dots, x''_{n-1}$ .

Pour se ramener à des notations plus conventionnelles (inconnues numérotées de 1 à p), on pose :

$$u_i = x''_{i+1} \quad (i = 1, 2, \dots, p = n - 2)$$

avec :

$$u_0 = u_{p+1} = 0$$

Alors le vecteur  $u$  de composantes  $u_i$  dans  $R^p$  est solution du système tridiagonal :

$$T \cdot u = v$$

avec :

$$T = \begin{pmatrix} T_{12} & T_{13} & 0 & \dots & 0 & 0 \\ T_{21} & T_{22} & T_{23} & 0 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot & 0 & T_{p1} & T_{p2} \end{pmatrix} \quad v = \begin{pmatrix} v_1 \\ \cdot \\ \cdot \\ \cdot \\ v_p \end{pmatrix}$$

où :

$$\begin{cases} T_{i1} = h_i & (i = 2, \dots, p) \\ T_{i2} = 2 \cdot (h_i + h_{i+1}) & (i = 1, \dots, p) \\ T_{i3} = h_{i+1} & (i = 1, \dots, p - 1) \end{cases}$$

$$v_i = 6 \cdot \left( \frac{x_{i+2} - x_{i+1}}{h_{i+1}} - \frac{x_{i+1} - x_i}{h_i} \right) \quad (i = 1, \dots, p)$$

En remarquant que, pour tout  $i = 1, \dots, p$  (où on a posé :  $T_{12} = T_{p3} = 0$ ) :

$$|T_{i2}| > |T_{i1}| + |T_{i3}|$$

on déduit que la matrice  $T$  est à diagonale strictement dominante, ce qui entraîne que le système admet une unique solution. De plus si on utilise la méthode de résolution de Gauss, il n'est pas nécessaire de faire des permutations de lignes de sorte que l'on pourra utiliser la méthode de double balayage de Cholesky pour résoudre ce système (Cf. le chapitre « Analyse Numérique Linéaire », § (3.7.3)).

### 4.3.8 Programmation structurée pour l'interpolation spline naturelle

Tout d'abord, on écrit une procédure de formation du système donnant les  $x''_i$ .

*PROCEDURE SystemeSplineNaturel(Entrée n : Entier ; x, y : Vecteur ;  
Sortie T : MatriceTridiagonale ; v : Vecteur) ;*

*Début*

*Pour i allant de 1 à n - 2 Faire*

*Début*

$$T_{i1} = x_{i+1} - x_i ;$$

$$T_{i3} = x_{i+2} - x_{i+1} ;$$

$$T_{i2} = 2 \cdot (T_{i1} + T_{i3}) ;$$

$$v_i = 6 \cdot \left( \frac{y_{i+2} - y_{i+1}}{T_{i3}} - \frac{y_{i+1} - y_i}{T_{i1}} \right) ;$$

*Fin ;*

*Fin ;*

Enfin on écrit une procédure qui donne les coefficients des polynômes  $s_i$  en fonction des  $x''_i$ .

*PROCEDURE CalculCoeff(Entrée n : Entier ; x, y, s : Vecteur ; Sortie a, b, c, d : Vecteur) ;*

*Début*

*Pour i allant de 1 à n - 1 Faire*

*Début*

$$h = x_{i+1} - x_i ;$$

$$a_i = \frac{s_{i+1} - s_i}{6 \cdot h} ;$$

$$b_i = \frac{s_i}{2} ;$$

$$c_i = \frac{y_{i+1} - y_i}{h} - h \cdot \frac{2 \cdot s_i + s_{i+1}}{6} ;$$

$$d_i = y_i ;$$

*Fin ;*

*Fin ;*

Le calcul des coefficients, dans le cas des splines naturels, se fait alors simplement avec la :

*PROCEDURE CalculCoeffSplineNaturel(Entrée n : Entier ; x, y : Vecteur ;  
Sortie a, b, c, d : Vecteur) ;*

*Début*

*SystemeSplineNaturel(n, x, y, T, v) ;*

*SystTridiag(n-2, T, v, s) ;*

*Pour i allant de n - 1 à 2 faire  $s_i = s_{i-1}$  ;*

*$s_1 = s_n = 0$  ;*

*CalculCoeff(n, x, y, s, a, b, c, d) ;*

*Fin ;*

En reprenant l'exemple du paragraphe (4.2.4), on obtient le résultat de la figure 5.13.

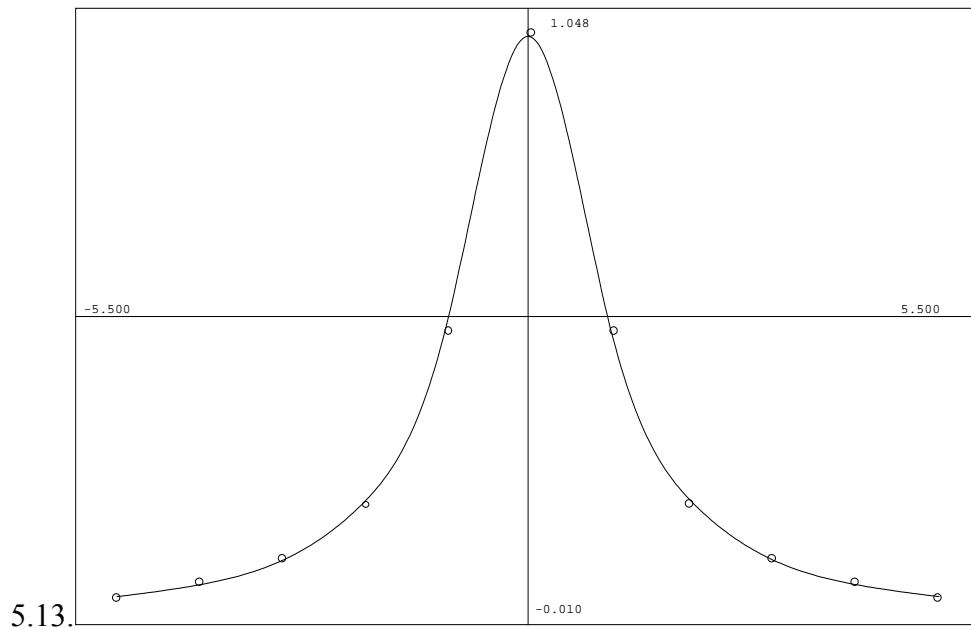


Figure 5.13

### 4.3.9 Application à un calcul d'intégrale

En notant  $S_i$  une primitive de  $s_i$ , une approximation de l'intégrale de la fonction interpolée est donnée par :

$$I = \sum_{i=1}^{n-1} (S_i(t_{i+1}) - S_i(t_i))$$

avec :

$$S_i(t) = \frac{a_i}{4} \cdot (t - t_i)^4 + \frac{b_i}{3} \cdot (t - t_i)^3 + \frac{c_i}{2} \cdot (t - t_i)^2 + d_i \cdot (t - t_i)$$

Ce qui donne, en tenant compte des calculs précédents, la formule d'intégration approchée :

$$I = \sum_{i=1}^{n-1} \frac{x_{i+1} + x_i}{2} \cdot h_i - \frac{1}{24} \cdot \sum (x_i'' + 2 \cdot x_{i+1}'') \cdot h_i^3$$

On peut alors compléter la procédure CalculCoeff en ajoutant une variable de sortie Int, qui donnera la valeur de l'intégrale du spline. Le calcul se faisant comme suit :

Int1 = 0 ; Int2 = 0 ;

Dans la boucle sur i, on ajoute les instructions :

h3 = h<sup>3</sup> ;

Int1 = Int1 + (y<sub>i</sub> + y<sub>i+1</sub>) · h ;

Int2 = (s<sub>i</sub> + 2 · s<sub>i+1</sub>) · h3 ;

Ce qui donne, après la fin de la boucle sur i : Int = Int1/2 – Int2/24 ;

### 4.3.10 Cas des fonctions périodiques

On suppose donc ici que le nuage de points provient d'une fonction périodique de période  $P = t_n - t_1$  et de classe  $C^2$ .

Le système de  $n - 2$  équations à  $n$  inconnues du paragraphe (4.3.6) :

$$h_{i-1} \cdot x''_{i-1} + 2 \cdot (h_{i-1} + h_i) \cdot x''_i + h_i \cdot x''_{i+1} = 6 \cdot \left( \frac{x_{i+1} - x_i}{h_i} - \frac{x_i - x_{i-1}}{h_{i-1}} \right)$$

est encore valable pour  $i = 2, \dots, n-1$ .

Comme  $x''_1 = x''_n$ , il nous manque ici seulement une équation. Pour obtenir cette dernière, on demande que (4.3.6) soit encore valable pour  $i = n$ , en posant, du fait de la périodicité :

$$x_n = x_1, x''_1 = x''_n, h_n = h_1, x''_{n+1} = x''_2, h_{n+1} = h_2 \text{ et } x_{n+1} = x_2$$

L'équation N°  $(n-1)$  (correspondante à  $i = n$ ) va donc s'écrire :

$$h_1 \cdot x''_2 + h_{n-1} \cdot x''_{n-1} + 2 \cdot (h_{n-1} + h_1) \cdot x''_n = 6 \cdot \left( \frac{x_2 - x_1}{h_1} - \frac{x_1 - x_{n-1}}{h_{n-1}} \right)$$

et l'équation N° 1 s'écrit :

$$2 \cdot (h_1 + h_2) \cdot x''_2 + h_2 \cdot x''_3 + h_1 \cdot x''_n = 6 \cdot \left( \frac{x_3 - x_2}{h_2} - \frac{x_2 - x_1}{h_1} \right)$$

Là encore, en posant  $u_i = x''_{i+1}$  pour  $i = 1, \dots, p = n-1$ , le vecteur  $u$  est solution d'un système de  $p$  équations à  $p$  inconnues du type :

$$T \cdot u = v \quad (*)$$

avec :

$$T = \begin{pmatrix} T_{12} & T_{13} & 0 & \dots & 0 & T_{11} \\ T_{21} & T_{22} & T_{23} & 0 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ T_{p3} & 0 & \cdot & 0 & T_{p1} & T_{p2} \end{pmatrix} \quad v = \begin{pmatrix} v_1 \\ \cdot \\ \cdot \\ \cdot \\ v_p \end{pmatrix}$$

où :

$$\begin{cases} T_{i1} = h_i \\ T_{i2} = 2 \cdot (h_i + h_{i+1}) \quad (i = 1, \dots, p) \\ T_{i3} = h_{i+1} \end{cases}$$

$$v_i = 6 \cdot \left( \frac{x_{i+2} - x_{i+1}}{h_{i+1}} - \frac{x_{i+1} - x_i}{h_i} \right) \quad (i = 1, \dots, p)$$

Là encore la matrice  $A$  est symétrique définie positive et le système  $(*)$  admet une unique solution.

Avec ces notations, la construction de la matrice  $T$  est identique à celle utilisée dans le cas du spline naturel, à la différence que les coefficients  $T_{11}$  et  $T_{p3}$  interviendront dans la résolution du système linéaire.

Pour résoudre un tel système, on considère l'inconnue  $u_p$  comme un paramètre et on retire un instant la dernière ligne pour se ramener au cas du paragraphe (3.7).

La résolution de ce système se fait en utilisant la méthode des pivots de Gauss.

En notant  $\gamma$  la dernière colonne de  $T$ , soit :

$\gamma_1 = T_{11}, \gamma_0 = 0$  ( $i = 2, \dots, p-2$ ),  $\gamma_{p-1} = T_{p-1,3}, \gamma_p = T_{p2}$   
 on a alors les formules de résolution suivantes :

$$m = \frac{T_{i,1}}{T_{i-1,2}} ; T_{i,1} = 0 ;$$

$$T_{i,2} = T_{i,2} - m \cdot T_{i-1,3} ;$$

$T_{i,3}$  inchangé

$$\gamma_i = \gamma_i - m \cdot \gamma_{i-1} ;$$

$$v_i = v_i - m \cdot v_{i-1}$$

$$T_{p,2} = T_{p,2} - m \cdot \gamma_{i-1}$$

pour  $i = 2, \dots, p-1$ .

Ce qui transforme le système en système triangulaire supérieur qui se résout « en remontée », ce qui donne :

$$\begin{cases} u_p = \frac{v_p}{T_{p,2}} \\ u_i = \frac{v_i - T_{i,3} \cdot u_{i+1} - \gamma_i \cdot u_p}{T_{i,2}} \quad (i = p-1, \dots, 1) \end{cases}$$

D'où la programmation structurée :

*PROCEDURE SystemeSplinePeriodique*(Entrée  $n$  : Entier ;  $P$  : Réel ;  $x, y$  : Vecteur ;  
 Sortie  $T$  : MatriceTridiagonale ;  $v$  : Vecteur) ;  
 {  $T$  n'est pas vraiment tridiagonale, c'est la donnée de trois vecteurs }  
 Début  
 $x_{n+1} = x_2 + P$  ;  
 $y_{n+1} = y_2$  ;  
 SystemeSplineNaturel( $n+1, x, y, T, v$ ) ;  
 Fin ;

*PROCEDURE ResolutionSystemeSplinePeriodique*(  
 Entrée  $p$  : Entier ;  $T$  : MatriceTridiagonale ;  $v$  : Vecteur ;  
 Sortie  $u$  : Vecteur) ;

Début  
 $\gamma_1 = T_{11}, \gamma_{p-1} = T_{p-1,3}, \gamma_p = T_{p2}$  ;  
 Pour  $i$  allant de 2 à  $p-2$  faire  $\gamma_i = 0$  ;  
 Pour  $i$  allant de 2 à  $p-1$  faire

Début  
 $m = \frac{T_{i,1}}{T_{i-1,2}} ;$   
 $T_{i,2} = T_{i,2} - m \cdot T_{i-1,3} ;$   
 $\gamma_i = \gamma_i - m \cdot \gamma_{i-1} ;$   
 $v_i = v_i - m \cdot v_{i-1} ;$   
 $T_{p,2} = T_{p,2} - m \cdot \gamma_{i-1}$   
 Fin ;

Pour  $i$  allant de  $p-1$  à 1 faire  $u_i = \frac{v_i - T_{i,3} \cdot u_{i+1} - \gamma_i \cdot u_p}{T_{i,2}}$

Fin ;

*PROCEDURE* CalculCoeffSplinePeriodique(Entrée  $n$  : Entier ;  $P$  : Réel ;  $x, y$  : Vecteur ;  
Sortie  $a, b, c, d$  : Vecteur) ;

*Début*

SystemeSplinePeriodique( $n, P, x, y, T, v$ ) ;  
ResolutionSystemeSplinePeriodique( $n - 1, T, v, s$ ) ;  
Pour  $i$  allant de  $n$  à 2 faire  $s_i = s_{i-1}$  ;  
 $s_1 = s_n$  ;  
CalculCoeff( $n, x, y, s, a, b, c, d$ ) ;

*Fin* ;

#### 4.3.11 Interpolation spline pour les courbes fermées

Si le nuage de points correspond à une courbe fermée, soit  $(x_1, y_1) = (x_n, y_n)$ , une idée pour réaliser l'interpolation spline est alors de passer en coordonnées paramétriques. Si la courbe est paramétrée par  $x = x(t)$  et  $y = y(t)$  où  $t$  est l'angle polaire avec l'origine ramenée au centre de gravité du nuage de points, alors il suffit d'utiliser l'interpolation spline pour chacune des fonctions périodiques de période  $2 \cdot \pi$   $x(t)$  et  $y(t)$ , si on suppose que la courbe est de classe  $C^2$  et sans points doubles.

On calcule donc d'abord les coordonnées du centre de gravité :

$$(x_g, y_g) = \frac{1}{n-1} \cdot \left( \sum_{i=1}^{n-1} x_i, \sum_{i=1}^{n-1} y_i \right)$$

Puis on écrit une procédure de passage en coordonnées polaires, en prenant soin de ranger ensuite les points par angles croissants.

Ensuite on utilise la procédure d'interpolation spline pour  $x$  et  $y$  en fonction de l'angle polaire  $t$ , ce qui donne deux fonctions polynomiales par morceaux  $s_x$  et  $s_y$ .

Et enfin on trace la courbe d'équation polaire :

$$\begin{cases} x = x_g + s_x(t) \\ y = y_g + s_y(t) \end{cases}$$

Une autre façon est de prendre  $t \in [0, n]$  comme paramètre en posant  $x(i) = x_i$  et  $y(i) = y_i$ , puis on utilise la procédure d'interpolation spline périodique pour  $x$  et  $y$  séparément, ce qui donne une autre représentation paramétrique de l'interpolée spline (et une courbe d'allure différente en général). Cela permet de traiter le cas des courbes avec points doubles.

La figure 5.14 donne un exemple de tracé correspondant au nuage de points du paragraphe (3.4.3).



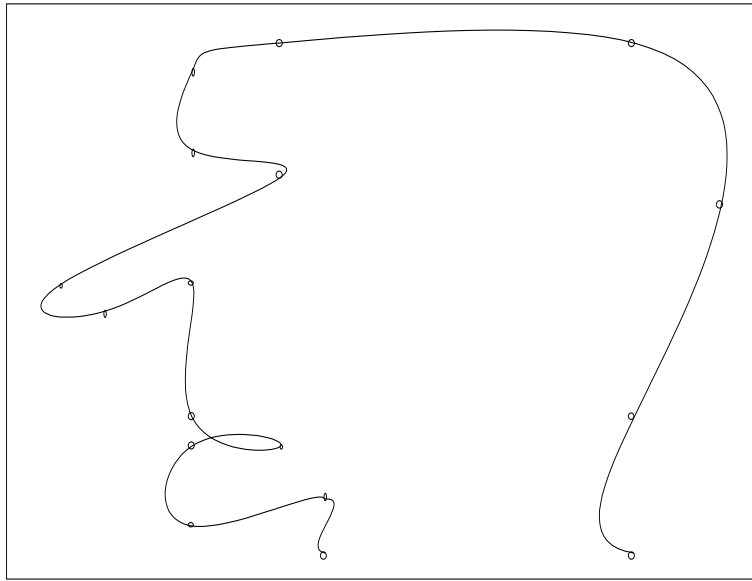


Figure 5.14

#### 4.3.12 Majoration de l'erreur d'interpolation

Si le nuage de points  $S$  provient d'une fonction  $f$  de classe  $C^4$  sur  $[0,1]$ , avec les  $x_i$  dans  $[0,1]$  (par changement de variable, on peut toujours se ramener à ce cas), on peut alors montrer que, si  $s$  est la fonction spline obtenue, on a les majorations suivantes, où  $\|\cdot\|$  désigne la norme du sup. pour les fonctions bornées sur  $[0,1]$  :

$$\begin{aligned} \|f - s\| &\leq C_1 \cdot h^4 \\ \|f' - s'\| &\leq C_2 \cdot h^3 \\ \|f'' - s''\| &\leq C_3 \cdot h^2 \end{aligned}$$

avec :

$$C_1 = \frac{13}{64} \cdot \|f^{(4)}\|, \quad C_2 = C_3 = \frac{13}{8} \cdot \|f^{(4)}\|$$

(Cf. Alhberg, Nilson et Walsh, § (2.3)).

## 5. Exercices

### 5.1 Ellipse, plan et sphère des moindres carrés

Ecrire les procédures correspondantes aux problèmes suivants :

- (i) Recherche d'une ellipse des moindres carrés ;
- (ii) Recherche d'un plan des moindres carrés ;
- (iii) Recherche d'une sphère des moindres carrés.

### 5.2 Interpolation par fractions continues

Première partie

Si  $(a_j)_{j \geq 1}$  et  $(b_j)_{j \geq 1}$  sont deux suites de fonctions, les fractions continues associées sont définies par :

$$f_1(x) = a_1(x)$$

$$f_n(x) = a_1(x) + \frac{b_1(x)}{a_2(x) + \frac{b_2(x)}{a_3(x) + \frac{b_3(x)}{\dots \frac{b_{n-1}(x)}{a_n(x)}}}}$$

Cette fraction peut aussi s'écrire :

$$f_n(x) = \frac{P_n(x)}{Q_n(x)} \quad (n \geq 1)$$

1°) Donner  $P_n$  et  $Q_n$  pour : (a)  $n = 1$  ; (b)  $n = 2$  ; (c)  $n = 3$ .

2°) Montrer par récurrence que les suites  $(P_n)$  et  $(Q_n)$  peuvent être définies par :

$$\begin{cases} P_0(x) = 1 \\ P_1(x) = a_1(x) \end{cases}, \begin{cases} Q_0(x) = 0 \\ Q_1(x) = 1 \end{cases}$$

$$\begin{cases} P_n(x) = a_n(x) \cdot P_{n-1}(x) + b_{n-1}(x) \cdot P_{n-2}(x) \\ Q_n(x) = a_n(x) \cdot Q_{n-1}(x) + b_{n-1}(x) \cdot Q_{n-2}(x) \end{cases} \quad (n \geq 2)$$

3°) En probabilités, la fonction de répartition de la loi Gamma de paramètre  $\alpha > 0$  est définie par :

$$\forall x > 0, P(\alpha, x) = 1 - Q(\alpha, x)$$

où :

$$Q(\alpha, x) = \frac{1}{\Gamma(\alpha)} \cdot \int_x^{+\infty} e^{-t} \cdot t^{\alpha-1} dt, \text{ avec : } \Gamma(\alpha) = \int_0^{+\infty} e^{-t} \cdot t^{\alpha-1} dt$$

On peut montrer que la fonction  $Q$  admet le développement :

$$Q(\alpha, x) = \frac{e^{-x} \cdot x^\alpha}{\Gamma(\alpha)} \cdot F(x) \text{ où : } F(x) = \lim_{n \rightarrow +\infty} (f_n(x))$$

les  $f_n$  étant les fractions continues définies par :

$$\left\{ \begin{array}{l} a_1 = 0 \\ a_{2k} = x \quad (k \geq 1) \\ a_{2k+1} = 1 \quad (k \geq 1) \end{array} \right. ; \left\{ \begin{array}{l} b_1 = 1 \\ b_{2k} = k - \alpha \quad (k \geq 1) \\ b_{2k+1} = k \quad (k \geq 1) \end{array} \right.$$

Ecrire une fonction qui permet de calculer une valeur approchée de  $P(\alpha, x)$  à une précision  $\varepsilon > 0$  donnée.

### Deuxième partie

Soit  $S = \{ P_i = (x_i, y_i) ; 1 \leq i \leq n \}$  un nuage de points expérimentaux.

On cherche à interpoler  $S$  par une fonction rationnelle définie par la fraction continue :

$$P(x) = a_1 + \frac{x - x_1}{a_2 + \frac{x - x_2}{a_3 + \frac{x - x_3}{\dots + \frac{x - x_{n-1}}{a_n}}}}$$

En vue de calculer les  $a_i$ , on définit la suite  $(b_{ij})$  par :

$$b_{i1} = y_i \quad (1 \leq i \leq n)$$

$$b_{i,j+1} = \frac{x_i - x_j}{b_{i,j} - b_{j,j}} \quad (1 \leq j < i \leq n) \quad (1 < j < i < n)$$

1°) (a) Calculer les  $b_{ij}$  pour  $n = 3$ .

(b) Quel lien a-t-on entre les  $a_i$  et les  $b_{ij}$ , pour  $n = 3$ ?

2°) Ecrire une procédure de calcul des  $b_{ij}$ .

3°) Montrer par récurrence sur  $i$  que  $a_i = b_{ii}$  ( $i = 1, \dots, n$ ).

4°) Ecrire une fonction qui permet de calculer  $P(x)$  pour tout  $x$ .

5°) Ecrire un programme permettant de tracer le graphe d'une fonction  $f$  et de son interpolée en prenant des abscisses  $x_i$  équidistantes.

### 5.3 Interpolation spline d'ordre 5

Ecrire les procédures correspondantes à l'interpolation spline de degré au plus 5.

## 6. Programmation Ada

### 6.1 Spécification du paquetage *REGRESSIONS*

Dans ce paquetage on trouve la programmation des procédures d'approximation du paragraphe 2.

```
with COMMON_MATRIX, COMMON_POLY ;
use COMMON_MATRIX, COMMON_POLY ;
package REGRESSIONS is
procedure REGRESSION_AFFINE(x,y : in VECTEUR ; a,b,Rho : in out FLOAT);
procedure REGRESSION_POLYNOMIALE(x, y : in VECTEUR ; U : out POLYNOE);
procedure REGRESSION_CIRCULAIRE(x,y : in VECTEUR ;
      a,b,r : in out FLOAT) ;
end REGRESSIONS ;
```

### 6.2 Le paquetage *DONNEES\_POINTS*

Dans ce paquetage on décrit des procédures permettant l'entrée de nuages de points.

```
with COMMON_MATRIX ;
use COMMON_MATRIX ;
package DONNEES_POINTS is
procedure GET(x, y : in out VECTEUR ; ORDRE : in BOOLEAN := FALSE) ;
procedure POINTS_DU_FICHER(x, y : in out VECTEUR ;
      ORDRE : in BOOLEAN := FALSE;
      NOM_DE_REPERTOIRE : in STRING := "A :\DONNEES") ;
procedure ENTRER_POINTS(x, y : in out VECTEUR ;
      ORDRE : in BOOLEAN := FALSE) ;
procedure EXTREMA_POINTS(x, y : in VECTEUR ;
      xMin, xMax, yMin, yMax : in out FLOAT) ;
end DONNEES_POINTS ;
```

### 6.3 Démonstration de la régression affine

```
with TEXT_IO, FIO, COMMON_MATH0, MATH0, MATH3, CURVE,
      COMMON_MATRIX, REGRESSIONS, DONNEES_POINTS ;
use TEXT_IO, FIO, COMMON_MATH0, MATH0, MATH3, CURVE,
      COMMON_MATRIX, REGRESSIONS, DONNEES_POINTS ;
procedure DEM_AFF is
n : NATURAL ;
procedure TRACE_REGRESSION_AFFINE(x, y : in VECTEUR ;
      a, b : in FLOAT) is
xMin,yMin,xMax,yMax,xOrigine,yOrigine,xUnite,yUnite : FLOAT ;
xFenMin, xFenMax, yFenMin, yFenMax : FLOAT ;
begin
      EXTREMA_POINTS(x,y,xMin,xMax,yMin,yMax) ;
      AXES_PAR_DEFAUT(xMin,yMin,xMax,yMax,xOrigine,yOrigine,
            xUnite,yUnite,xFenMin, xFenMax, yFenMin, yFenMax) ;
      INIT_GRAPHIQUE ;
      FENETRE_GRAPHIQUE(xFenMin,xFenMax,yFenMin,yFenMax) ;
      XY_AXES(xOrigine,yOrigine,xUnite,yUnite,FALSE) ;
```

```

    for j in x'range loop
        CROIX(x(j),y(j)) ;
    end loop ;
    DEPLACE(xMin,a*xMin + b) ; TRACE(xMax,a*xMax + b) ;
    SORTIE_GRAPHIQUE ;
end TRACE_REGRESSION_AFFINE ;
begin
    MODE_AFFICHAGE ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Nombre de points du nuage : ") ;
    declare
        x, y : VECTEUR(1..n) ;
        a, b, Rho : FLOAT ;
    begin
        ENTRER_POINTS(x,y) ;
        REGRESSION_AFFINE(x,y,a,b,Rho) ;
        PUT_LINE(IMP,"DROITE DES MOINDRES CARREES") ;
        PUT(IMP," y = ") ; FIO.PUT(IMP,a,6,4,0) ;
        PUT(IMP,"*x + ") ; FIO.PUT(IMP,b,6,4,0) ; NEW_LINE(IMP) ;
        PUT(IMP,"Coefficient de regression : Rho = ") ;
        FIO.PUT(IMP,Rho,6,4,0) ; NEW_LINE(IMP) ; PAUSE ;
        TRACE_REGRESSION_AFFINE(x,y,a,b) ;
    end ;
    CLOSE(IMP) ;
end DEM_AFF ;

```

#### 6.4 Démonstration de la régression polynomiale

```

with TEXT_IO, COMMON_MATH0, MATH0, MATH3, CURVE, COMMON_MATRIX,
    COMMON_POLY, POLY, REGRESSIONS, DONNEES_POINTS ;
use TEXT_IO, COMMON_MATH0, MATH0, MATH3, CURVE, COMON_MATRIX,
    COMMON_POLY, POLY, REGRESSIONS, DONNEES_POINTS ;
procedure DEM_RPOL is
    n : NATURAL ;
    p : DEGRE_MAX ;
begin
    MODE_AFFICHAGE ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Nombre de points du nuage : ") ;
    ENTRER_ENTIER_BORNE(1,n - 1,p,
        "Degre du polynome des moindres carres : ") ;
    declare
        x, y : VECTEUR(1..n) ;
        U : POLYNOME(0..p) ;
        xMin, xMax, yMin, yMax : FLOAT ;
    begin
        ENTRER_POINTS(x,y) ;
        REGRESSION_POLYNOMIALE(x,y,U) ;
        TEXT_IO.PUT_LINE(IMP,"POLYNOME DES MOINDRES CARREES") ;
        TEXT_IO.PUT(IMP,"P(x) = ") ; PUT_LINE(U) ; NEW_LINE(IMP) ;
        PAUSE ;
        EXTREMA_POINTS(x,y,xMin,xMax,yMin,yMax) ;
        TRACE_POLYNOME(U,xMin,xMax) ;
        for j in x'range loop
            CROIX(x(j),y(j)) ;
        end loop ;
    end ;
end DEM_RPOL ;

```

```

        end loop ;
        SORTIE_GRAPHIQUE ;
    end ;
    CLOSE(IMP) ;
end DEM_RPOL ;

```

## 6.5 Démonstration de la régression circulaire

```

with TEXT_IO, FIO, COMMON_MATH0, MATH0, MATH1, MATH3, CURVE,
     COMMON_MATRIX, REGRESSIONS, DONNEES_POINTS ;
use TEXT_IO, FIO, COMMON_MATH0, MATH0, MATH1, MATH3, CURVE,
     COMMON_MATRIX, REGRESSIONS, DONNEES_POINTS ;
procedure DEM_CIRC is
n : NATURAL ;
procedure TRACE_REGRESSION_CIRCULAIRE(x, y : in VECTEUR ;
                                       a, b, r : in FLOAT) is
xMin,yMin,xMax,yMax : FLOAT ;
begin
    EXTREMA_POINTS(x,y,xMin,xMax,yMin,yMax) ;
    xMin := MIN(xMin,a - r) - 0.1*r ; xMax := MAX(xMax,a + r) + 0.1*r ;
    yMin := MIN(yMin,b - r) - 0.1*r ; yMax := MAX(yMax,b + r) + 0.1*r ;
    INIT_GRAPHIQUE ;
    FENETRE(xMin,xMax,yMin,yMax) ;
    CADRE_GRAPHIQUE(0,Y_MAXIMUM - 2*MARGE,MARGE, Y_MAXIMUM - MARGE) ;
    XY_AXES(a,b,0.1*r,0.1*r,FALSE) ;
    for j in x'range loop
        CROIX(x(j),y(j)) ;
    end loop ;
    CERCLE(a,b,r) ;
    SORTIE_GRAPHIQUE ;
end TRACE_REGRESSION_CIRCULAIRE ;
begin
    MODE_AFFICHAGE ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Nombre de points du nuage : ") ;
    declare
        x, y : VECTEUR(1..n) ;
        a, b, r : FLOAT ;
    begin
        ENTRER_POINTS(x,y) ;
        REGRESSION_CIRCULAIRE(x,y,a,b,r) ;
        PUT_LINE(IMP,"CERCLE DES MOINDRES CARREES") ;
        PUT(IMP," Centre, (a,b) = (" ; FIO.PUT(IMP,a,6,4,0) ;
        PUT(IMP,",") ; FIO.PUT(IMP,b,6,4,0) ; PUT_LINE(IMP,")") ;
        PUT(IMP,"Rayon, r = ") ; FIO.PUT(IMP,r,6,4,0) ; NEW_LINE(IMP) ;
        PAUSE ;
        TRACE_REGRESSION_CIRCULAIRE(x,y,a,b,r) ;
    end ;
    CLOSE(IMP) ;
end DEM_CIRC ;

```

## 6.6 Spécification du paquetage BEZIER

```

with CURVE, COMMON_MATRIX ;

```

```

use CURVE, COMMON_MATRIX ;
package BEZIER is
function POLYNOME_BERNSTEIN(P : in VECTEUR ; t : in FLOAT)
    return FLOAT ; -- § 3.2
procedure CALCUL_BEZIER(x,y : in VECTEUR ; t : in FLOAT ;
    Bx,By : out FLOAT) ; -- § 3.2
procedure TRACE_BASE_BERNSTEIN(n : in NATURAL) ; -- § 3.2
procedure COURBE_BEZIER(x, y : in VECTEUR ;
    C : in out COURBE) ; -- § 3.3
end BEZIER ;

```

## 6.7 Démonstration de l'approximation de Bézier

```

with TEXT_IO, MATH0, MATH3, CURVE, COMMON_MATRIX,
    DONNEES_POINTS, BEZIER ;
use TEXT_IO, MATH0, MATH3, CURVE, COMMON_MATRIX,
    DONNEES_POINTS, BEZIER ;
procedure DEM_BEZI is
CHOIX : NATURAL range 0..2 ;
n, j : NATURAL ;
procedure MENU_BEZIER(CHOIX : out NATURAL) is
begin
    PUT_LINE("-0- FIN") ;
    PUT_LINE("-1- Trace des fonctions de base de Bernstein") ;
    PUT_LINE(
        "-2- Tracer une courbe de Bezier a partir d'un nuage de points") ;
    NEW_LINE ;
    ENTRER_ENTIER_BORNE(0,2,CHOIX,"Votre choix : ") ;
end MENU_BEZIER ;
begin
    loop
        MENU_BEZIER(CHOIX) ;
        case CHOIX is
            when 0 => exit ;
            when 1 =>
                ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Nombre de fonctions a tracer : ") ;
                n := n - 1 ;
                TRACE_BASE_BERNSTEIN(n) ;
            when 2 =>
                ENTRER_ENTIER_BORNE(3,DIM_MAX,n,"Nombre de points du nuage (> 2) : ") ;
                n := n - 1 ;
                declare
                    x, y : VECTEUR(0..n) ;
                    C : COURBE(300) ;
                begin
                    ENTRER_POINTS(x,y,FALSE) ;
                    COURBE_BEZIER(x, y, C) ;
                    TRACE_UNE_COURBE(C) ;
                    SORTIE_GRAPHIQUE ;
                end ;
            end case ;
        end loop ;
end DEM_BEZI ;

```

## 6.8 Spécification du paquetage *B\_SPLINE*

```
with CURVE, COMMON_MATRIX ;
use CURVE, COMMON_MATRIX ;
package B_SPLINE is
type VECTEUR_NODAL is array(NATURAL range <>) of NATURAL ;
function BASE_BSPLINE(i, j : NATURAL ; Ne : VECTEUR_NODAL ; t : FLOAT)
return FLOAT ; -- § 3.4.2
procedure CALCUL_BSPLINE(j : in NATURAL ; Ne : in VECTEUR_NODAL ;
x, y : in VECTEUR ; t : in FLOAT ; Bx, By : out FLOAT) ; -- § 3.4.3
procedure COURBE_BSPLINE(j : in NATURAL ; x, y : in VECTEUR ;
C : in out COURBE) ; -- § 3.4.3
procedure TRACE_BASE_BSPLINE(n, j : in NATURAL) ; -- § 3.4.2
end B_SPLINE ;
```

## 6.9 Démonstration de l'approximation *B\_Spline*

```
with TEXT_IO, MATH0, MATH3, CURVE, COMMON_MATRIX,
DONNEES_POINTS, B_SPLINE ;
use TEXT_IO, MATH0, MATH3, CURVE, COMMON_MATRIX,
DONNEES_POINTS, B_SPLINE ;
procedure DEM_BSPL is
CHOIX : NATURAL range 0..2 ;
n, j : NATURAL ;
procedure MENU_BSPLINE(CHOIX : out NATURAL) is
begin
PUT_LINE("-0- FIN") ;
PUT_LINE("-1- Trace des fonctions de base B_Splines") ;
PUT_LINE(
"-2- Tracer une courbe B_Spline a partir d'un nuage de points") ;
NEW_LINE ;
ENTRER_ENTIER_BORNE(0,2,CHOIX,"Votre choix : ") ;
end MENU_BSPLINE ;
begin
loop
MENU_BSPLINE(CHOIX) ;
case CHOIX is
when 0 => exit ;
when 1 =>
ENTRER_ENTIER_BORNE(2,DIM_MAX,n,"Nombre de fonctions a tracer : ") ;
n := n - 1 ;
ENTRER_ENTIER_BORNE(2,n + 1,j,"Ordre des B_Splines : ") ;
TRACE_BASE_BSPLINE(n,j) ;
when 2 =>
ENTRER_ENTIER_BORNE(3,DIM_MAX,n,"Nombre de points du nuage (> 2) : ") ;
n := n - 1 ;
ENTRER_ENTIER_BORNE(2,n + 1,j,"Ordre de la B_Spline : ") ;
declare
x, y : VECTEUR(0..n) ;
C : COURBE(300) ;
begin
ENTRER_POINTS(x,y,FALSE) ;
```



```

        COURBE_BSPLINE(j, x, y, C) ;
        TRACE_UNE_COURBE(C) ;
        SORTIE_GRAPHIQUE ;
    end ;
end case ;
end loop ;
end DEM_BSPL ;

```

## 6.10 Spécification du paquetage LAGRANGE

```

with CURVE, COMMON_MATRIX ;
use CURVE, COMMON_MATRIX ;
package LAGRANGE is
function LAGRANGE_NEVILLE(x, y : VECTEUR ; t : FLOAT)
    return FLOAT ; -- § 4.2.3
procedure POLYNOме_LAGRANGE(x, y : in VECTEUR ;
    P : out VECTEUR) ; -- § 4.2.3
function LAGRANGE_NEWTON(x, P : in VECTEUR ; t : in FLOAT)
    return FLOAT ; -- § 4.2.4
procedure COURBE_INTERPOLATION_LAGRANGE(x,y : in VECTEUR ;
    xMin, xMax : in FLOAT ; C : in out COURBE) ;
end LAGRANGE ;

```

## 6.11 Démonstration de l'interpolation de Lagrange

```

with TEXT_IO, MATH0, COMMON_MATH1, MATH1, MATH2, MATH3, CURVE,
    COMMON_MATRIX, LAGRANGE, DONNEES_POINTS ;
use TEXT_IO, MATH0, COMMON_MATH1, MATH1, MATH2, MATH3, CURVE,
    COMMON_MATRIX, LAGRANGE, DONNEES_POINTS ;
procedure DEM_LAGR is
CHOIX : NATURAL range 0..2 ;
n : NATURAL ;
f : FONCTION ;
procedure MENU_LAGRANGE(CHOIX : out NATURAL) is
begin
    PUT_LINE("-0- FIN") ;
    PUT_LINE(
        "-1- Interpolation d'une fonction connue en quelques points") ;
    PUT_LINE("-2- Interpolation d'une fonction donnee analytiquement");
    NEW_LINE ;
    ENTRER_ENTIER_BORNE(0,2,CHOIX,"Votre choix : ") ;
end MENU_LAGRANGE ;
begin
    loop
        MENU_LAGRANGE(CHOIX) ;
        case CHOIX is
            when 0 => exit ;
            when 1 =>
                ENTRER_ENTIER_BORNE(3,DIM_MAX,n,"Nombre de points du nuage (> 2) : ") ;
                declare
                    x, y : VECTEUR(0..n) ;
                    C : COURBE(300) ;
                begin

```

```

        ENTRER_POINTS(x,y,TRUE) ;
        COURBE_INTERPOLATION_LAGRANGE(x, y, x(1), x(n), C) ;
        TRACE_UNE_COURBE(C) ;
        SORTIE_GRAPHIQUE ;
    end ;
when 2 =>
    GET_LINE(f,"Fonction a interpoler, f(x) = ") ;
ENTRER_ENTIER_BORNE(3,DIM_MAX,n,"Nombre de points du nuage (> 2) : ") ;
n := n - 1 ;
declare
    x, y : VECTEUR(0..n) ;
    C : COURBE(300) ;
    Pas, a, b : FLOAT ;
begin
    PUT_LINE("L'intervalle de definition de f est note [a,b]") ;
    ENTRER_REEL(a,"a = ") ;
    ENTRER_REEL_BORNE(a,1.0E+10,b,"b = ") ;
    if LIRE_REPONSE(
        "Interpolation avec les abscisses de Tchebychev? ") then
        for i in x'range loop
            x(i) := 0.5*((b-a)*Cos(FLOAT(2*i+1)/FLOAT(2*n+2)*Pi) + b + a) ;
            y(i) := EVALUE(f,x(i)) ;
        end loop ;
    else
        Pas := (b - a)/FLOAT(n) ;
        x(0) := a ; y(0) := EVALUE(f,a) ;
        for i in 1..n loop
            x(i) := x(i-1) + Pas ;
            y(i) := EVALUE(f,x(i)) ;
        end loop ;
    end if ;
    COURBE_INTERPOLATION_LAGRANGE(x, y, a, b, C) ;
    TRACE_UNE_COURBE(C) ;
    Pas := (b - a)/FLOAT(C.n) ;
    for i in 0..C.n loop
        CROIX(a,EVALUE(f,a)) ;
        a := a + Pas ;
    end loop ;
    SORTIE_GRAPHIQUE ;
end ;
end case ;
end loop ;
end DEM_LAGR ;

```

## 6.12 Spécification du paquetage *INTERPOLATION\_SPLINE*

```

with COMMON_MATRIX ;
use COMMON_MATRIX ;
package INTERPOLATION_SPLINE is
procedure CALCUL_COEFFICIENTS_SPLINE_NATUREL(x, y : in VECTEUR ;
    a,b,c,d : out VECTEUR) ; -- § 4.3.7
function SPLINE(x,y,a,b,c,d : in VECTEUR ; t : in FLOAT)
return FLOAT ; -- § 4.3.2

```

```

procedure TRACE_SPLINE_NATUREL(x, y : in VECTEUR) ;
procedure CALCUL_COEFFICIENTS_SPLINE_PERIODIQUE(PERIODE : in FLOAT ;
          x, y : in VECTEUR ; a,b,c,d : out VECTEUR) ; -- § 4.3.10
procedure TRACE_SPLINE_PERIODIQUE(PERIODE : in FLOAT ;
          x, y : in VECTEUR) ;
procedure TRACE_SPLINE_PARAMETRIQUE(x, y : in VECTEUR) ; -- § 4.3.11
end SPLINE ;

```

### 6.13 Démonstration de l'interpolation spline

```

with TEXT_IO, MATH0, COMMON_MATRIX, DONNEES_POINTS,
     INTERPOLATION_SPLINE ;
use TEXT_IO, MATH0, COMMON_MATRIX, DONNEES_POINTS,
     INTERPOLATION_SPLINE ;
procedure DEM_SPL is
CHOIX : NATURAL range 0..2 ;
procedure DEMO_SPLINE_CARTESIENNE is
n : NATURAL ;
begin
  ENTRER_ENTIER_BORNE(4,DIM_MAX,n,"Nombre de points du nuage (> 3) : ");
  declare
    x, y : VECTEUR(1..n) ;
  begin
    ENTRER_POINTS(x,y,TRUE) ;
    TRACE_SPLINE_NATUREL(x,y) ;
  end ;
end DEMO_SPLINE_CARTESIENNE ;
procedure DEMO_SPLINE_PARAMETRIQUE is
n : NATURAL ;
begin
  ENTRER_ENTIER_BORNE(4,DIM_MAX,n,"Nombre de points du nuage (> 3) : ");
  declare
    x, y : VECTEUR(1..n) ;
  begin
    ENTRER_POINTS(x,y,FALSE) ;
    TRACE_SPLINE_PARAMETRIQUE(x,y) ;
  end ;
end DEMO_SPLINE_PARAMETRIQUE ;
procedure MENU_SPLINE(CHOIX : out NATURAL) is
begin
  PUT_LINE("-1- Interpolation spline d'une courbe cartesienne") ;
  PUT_LINE("-2- Interpolation spline d'une courbe parametrique") ;
  NEW_LINE ;
  ENTRER_ENTIER_BORNE(0,2,CHOIX,"Votre choix : ") ;
end MENU_SPLINE ;

begin
  loop
    MENU_SPLINE(CHOIX) ;
    case CHOIX is
      when 0 => exit ;
      when 1 => DEMO_SPLINE_CARTESIENNE ;
      when 2 => DEMO_SPLINE_PARAMETRIQUE ;
    end case ;
  end loop ;
end ;

```

```
        end case ;  
    end loop ;  
end DEM_SPL ;
```

## CHAPITRE 6

# Calcul numérique des intégrales

### 1. Introduction. Position des problèmes

#### 1.1 Remarques préliminaires sur le calcul des primitives

##### 1.1.1 Calcul direct des intégrales

On sait que, si  $f : [a,b] \rightarrow \mathbb{R}$  est continue, alors elle admet des primitives et si  $F$  est l'une d'elles, l'intégrale de  $f$  sur  $[a,b]$  est :

$$\int_a^b f(x)dx = F(b) - F(a)$$

Puis, pour  $f : [a,b[ \rightarrow \mathbb{R}$  continue avec  $a < b \leq +\infty$ , en cas de convergence, l'intégrale impropre de  $f$  sur  $[a,b[$  est définie par :

$$\int_a^b f(x)dx = \lim_{x \rightarrow b} \int_a^x f(t)dt$$

Enfin, le théorème de Fubini va ramener le calcul des intégrales multiples à celui des intégrales simples.

Tout semble donc résolu. Mais en réalité, les cas où l'on sait calculer explicitement les primitives sont rares. De plus on a souvent affaire à des fonctions  $f$  qui ne sont connues explicitement qu'en un nombre fini de points. Il faudra donc développer des méthodes de calcul adaptées à ces cas.

##### 1.1.2 A propos des primitives élémentaires

On dira qu'une fonction  $f$  admet des primitives élémentaires si elle admet des primitives qui peuvent s'exprimer à l'aide des fonctions polynomiales, des fonctions transcendantes usuelles (Log, trigonométriques et leurs inverses) et des quatre opérations usuelles sur  $\mathbb{R}$ .

En ce sens, la fonction  $f: x \mapsto f(x) = \frac{1}{\text{Ln}(x)}$  n'admet pas de primitives élémentaires sur  $]1, +\infty[$ . Cette dernière étant continue, on peut poser pour  $x > 1$ ,  $L_i(x) = \int_2^x \frac{dt}{\text{Ln}(t)}$ , ce qui définit une fonction appelée « *logarithme intégral* ».

De manière plus générale, Liouville a montré le résultat suivant :

**Théorème** : Si  $f$  et  $g$  sont deux fonctions rationnelles, alors  $f \cdot e^g$  admet des primitives élémentaires si, et seulement si, il existe une fonction rationnelle  $h$  telle que  $f = h' + h \cdot g'$ .

*Démonstration* — B. Randé : Primitives élémentaires. Revue des Mathématiques spéciales, Mai 1984.

*Application* — En prenant  $f(x) = \frac{1}{x}$  et  $g(x) = x$ , on déduit que  $\frac{e^x}{x}$  n'admet pas de primitives élémentaires et en faisant des changements de variables, on déduit qu'il en est de même pour les fonctions  $\frac{1}{\text{Ln}(x)}$ ,  $\frac{e^x}{x^2}$ ,  $e^{\frac{1}{x}}$  et  $\text{Ln}(\text{Ln}(x))$ .

## 1.2 Elaboration de méthodes numériques

Il est donc nécessaire de développer des méthodes numériques de calcul approché des intégrales définies, impropres et multiples. On parle aussi de méthodes de « *quadrature numérique* », pour les intégrales simples et de « *cubature numérique* » pour les intégrales doubles.

Une façon de résoudre le problème qui nous intéresse est de se ramener à une équation différentielle avec conditions initiales, soit à calculer  $y(b)$ , où  $y$  est solution de :

$$\begin{cases} y'(x) = f(x) & (a \leq x \leq b) \\ y(a) = 0 \end{cases}$$

et utiliser les méthodes de résolution des équations différentielles (Cf. le chapitre sur la résolution numérique des équations différentielles).

On pourra aborder le problème de ce point de vue, par exemple, pour des fonctions ayant de nombreux pics, mais en général, vue la forme particulière de l'équation différentielle, on aura intérêt à utiliser une méthode spécifique.

Comme dans le cas des équations différentielles, à la base on a une idée de « *discrétisation* » de l'intervalle  $[a, b]$ , à partir de laquelle on calculera des sommes qui vont approximer l'intégrale.

Dans ce chapitre, on s'intéressera d'abord au calcul numérique des intégrales définies en utilisant des interpolations polynomiales, puis on étudiera des méthodes utilisant les polynômes orthogonaux, qui permettent également de calculer certaines intégrales impropres et enfin on dira quelques mots sur le calcul numérique des intégrales doubles.

Comme toujours, à chaque fois on essaiera de donner une évaluation des erreurs de méthode.

## 2. Méthodes de calcul par interpolation polynomiale. Schémas d'intégration classiques

### 2.1 Idée des méthodes par interpolation

On se donne, pour tout le paragraphe 2, une fonction  $f : [a,b] \rightarrow \mathbb{R}$  supposée continue et connue en un nombre fini de points  $x_0 < \dots < x_n$  dans  $[a,b]$ .

L'idée est alors d'utiliser les valeurs connues  $(f(x_i))_{0 \leq i \leq n}$  pour calculer une valeur approchée de  $I(f) = \int_a^b f(x)dx$ .

Dans un premier temps, on cherchera des formules d'intégration du type :

$$(1) \quad I(f) \cong \hat{I}(f) = \sum_{i=0}^n a_i \cdot f(x_i)$$

où les constantes  $a_i$  ne dépendent que des  $x_i$  et pas de la fonction  $f$ .

On parle alors de « formule d'intégration élémentaire ».

Dans les méthodes par interpolation polynomiale, on demandera que la formule (1) soit exacte pour les polynômes de degré au plus  $n$ , soit en notant  $R_n[x]$  l'espace des polynômes de degré au plus  $n$  :

$$(2) \quad \forall P \in R_n[x], \int_a^b P(x)dx = \sum_{i=0}^n a_i \cdot P(x_i)$$

Cette idée est justifiée par le théorème de Stone-Weirstrass qui nous dit que toute fonction continue sur  $[a,b]$  est limite uniforme d'une suite de polynômes.

Donc si pour une précision  $\varepsilon > 0$  donnée,  $P_n$  est un polynôme de degré  $n$  tel que  $\text{Sup}\{|f(x) - P_n(x)|; a \leq x \leq b\} < \varepsilon$  et si  $(a_i)_{0 \leq i \leq n}$  est une suite de coefficients vérifiant (2), on écrira :

$$\int_a^b f(x)dx \cong \int_a^b P_n(x)dx = \sum_{i=0}^n a_i \cdot P_n(x_i) \cong \sum_{i=0}^n a_i \cdot f(x_i)$$

soit la formule (1).

En fait, numériquement il n'est pas pratique de raisonner comme indiqué ci-dessus.

L'idée que l'on retiendra est de remplacer  $f$  par son polynôme d'interpolation en  $(x_i)_{0 \leq i \leq n}$ , ce qui théoriquement est moins satisfaisant, car quand  $n$  tend vers  $+\infty$  une telle suite ne converge pas vers  $f$  et donc la convergence de la méthode n'est pas assurée dans tous les cas. Mais d'un point de vue pratique cela est en général satisfaisant.

Il sera nécessaire d'évaluer l'erreur de méthode :

$$(3) \quad R(f) = |I - \hat{I}|$$

## 2.2 Résolution du problème (2)

### 2.2.1 Existence et unicité d'une solution

L'existence et l'unicité des coefficients  $a_i$  est assurée par le :

**Théorème :** Il existe des constantes uniques  $(a_i)_{0 \leq i \leq n}$  telles que :

$$(2) \quad \forall P \in R_n[x], \int_a^b P(x) dx = \sum_{i=0}^n a_i \cdot P(x_i)$$

*Démonstration* — En utilisant la base canonique  $\{1, x, \dots, x^n\}$  de  $R_n[x]$ , on voit que (2) équivaut à trouver une solution au système linéaire :

$$(4) \quad \sum_{j=0}^n x_j^i \cdot a_j = \frac{b^{i+1} - a^{i+1}}{b - a} \quad (i = 0, \dots, n)$$

La matrice de ce système étant du type Vandermonde de déterminant  $\Delta_i = \prod_{0 \leq i < j \leq n} (x_j - x_i) \neq 0$ , on déduit que (4) et donc (2) admet une unique solution.

### 2.2.2 Détermination pratique des coefficients $a_j$

*Exemple* — Prenons  $a = 0$ ,  $b = 2$  et  $x_0 = 0$ ,  $x_1 = 1$ ,  $x_2 = 2$ .

Le système à résoudre est alors :

$$\begin{cases} a_0 + a_1 + a_2 = 2 \\ a_1 + 2 \cdot a_2 = 2 \\ a_1 + 4 \cdot a_2 = \frac{8}{3} \end{cases}$$

de solution  $a_0 = \frac{1}{3}$ ,  $a_1 = \frac{4}{3}$ ,  $a_2 = \frac{1}{3}$ .

On a donc pour tout polynôme  $P = a \cdot x^2 + b \cdot x + c$  :

$$\int_0^2 P(x) dx = \frac{1}{3} \cdot (P(0) + 4 \cdot P(1) + P(2))$$

et la formule d'intégration associée sera alors :

$$\int_0^2 f(x) dx \cong \frac{1}{3} \cdot (f(0) + 4 \cdot f(1) + f(2))$$

Il s'agit en fait de la règle de Simpson que nous décrirons plus loin.

Pour  $f(x) = \text{Ch}(x)$ , on a  $I = \text{Sh}(2) \cong 3,6269$  et  $\hat{I} \cong 3,6449$ .

De manière générale, le système (4) étant mal conditionné, sa résolution posera des problèmes pour  $n$  grand.

La matrice de (4) étant symétrique définie positive, on peut penser à utiliser la méthode de Cholesky, mais pour  $n$  grand les résultats ne sont pas satisfaisants.

En fait, dans la pratique, on ne considère pas ces formules pour de grandes valeurs de  $n$  comme on le verra plus loin.

D'autre part, on peut calculer les  $a_i$  d'une autre façon.

Soit  $\{L_0, \dots, L_n\}$  la base de Lagrange de  $R_n[x]$  définie par :



$$L_i(x) = \frac{\prod_{\substack{j=0 \\ j \neq i}}^n (x - x_j)}{\prod_{\substack{j=0 \\ j \neq i}}^n (x_i - x_j)} \quad (i = 0, \dots, n)$$

on a alors :

$$\forall P \in R_n[x], P(x) = \sum_{i=0}^n P(x_i) \cdot L_i(x)$$

et :

$$\int_a^b P(x) dx = \sum_{i=0}^n \left( \int_a^b L_i(x) dx \right) \cdot P(x_i)$$

et de l'unicité des coefficients  $a_i$ , on déduit que :

$$a_i = \int_a^b L_i(x) dx, \quad i = 0, \dots, n$$

ce qui détermine de façon « explicite » les  $a_i$ .

### 2.2.3 Majoration de l'erreur

Pour  $f : [a,b] \rightarrow R$ , continue, on pose :

$$E(f) = \int_a^b f(x) dx - \sum_{i=0}^n a_i \cdot f(x_i)$$

où les  $a_i$  sont définis par  $E(P) = 0$ , pour tout  $P$  dans  $R_n[x]$ .

Pour  $t$  fixé dans  $R$ , on note :

$$\varphi_t : t \mapsto \varphi_t(x) = (x - t)_+^n = \begin{cases} (x - t)^n & \text{si } x \geq t \\ 0 & \text{si } x < t \end{cases}$$

et :

$$K_n(t) = E(\varphi_t)$$

La fonction  $K_n$  ainsi définie s'appelle le « noyau de Péano » de la méthode.

On peut alors montrer le :

*Théorème (Péano) :* Si  $f$  est de classe  $C^{n+1}$ , alors :

$$E(f) = \frac{1}{n!} \cdot \int_a^b K_n(t) \cdot f^{(n+1)}(t) dt$$

et donc :

$$R(f) = |E(f)| \leq C_n \cdot M_{n+1} \cdot \frac{1}{n!}$$

où  $C_n = \int_a^b |K_n(t)| dt$  et  $M_{n+1} = \text{Sup} \left\{ |f^{(n+1)}(x)| ; a \leq x \leq b \right\}$ .

Si de plus,  $K_n$  garde un signe constant sur  $[a,b]$ , alors il existe  $\xi$  dans  $[a,b]$  tel que :

$$E(f) = \frac{E(x^{n+1})}{(n+1)!} \cdot f^{(n+1)}(\xi)$$

où :

$$E(x^{n+1}) = \frac{b^{n+2} - a^{n+2}}{b - a} - \sum_{i=0}^n a_i \cdot x_i^{n+1}$$

Démonstration — Voir Crouzeix et Mignot p.39.

**2.3 Cas particulier où les  $x_j$  sont équidistants.  
Méthodes de Newton et Cotes**

**2.3.1 Calcul des coefficients  $a_i$**

On suppose ici que les  $x_i$  sont équidistants dans  $[a,b]$ , soit : en posant  $h = \frac{b - a}{n}$

:

$$x_i = a + i \cdot h \quad (i = 0, 1, \dots, n)$$

La méthode obtenue est appelée « méthode fermée de Newton-Cotes à  $n + 1$  points ». On parle de méthode fermée car les extrémités  $a$  et  $b$  font partie des points utilisés. On peut aussi considérer des méthodes qui n'utilisent que les points  $x_1, \dots, x_{n-1}$  et on parle de « méthode ouverte de Newton-Cotes à  $n - 1$  points », mais ces méthodes sont peu utilisées.

En prenant pour base de  $R_n[x]$ ,  $\{1, (x - a), \dots, (x - a)^n\}$ , on déduit que les coefficients  $a_i$  sont solutions du système linéaire :

$$(5) \quad \begin{cases} a_0 + a_1 + \dots + a_n = n \cdot h \\ a_1 + 2 \cdot a_2 + \dots + n \cdot a_n = \frac{n^2}{2} \cdot h \\ \dots \\ a_1 + 2^n \cdot a_2 + \dots + n^n \cdot a_n = \frac{n^{n+1}}{n+1} \cdot h \end{cases}$$

En écrivant les coefficients  $a_i$  sous la forme :

$$a_i = \int_a^b L_i(x) dx, \quad i = 0, \dots, n$$

avec les notations de (2.2.2), puis en faisant le changement de variable  $x = a + t \cdot h$ , on déduit que :

$$a_i = \frac{b - a}{n} \cdot b_i, \quad \text{avec } b_i = \int_0^n \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - j}{i - j} dt \quad (i = 0, \dots, n)$$

la suite  $(b_i)_{0 \leq i \leq n}$  ne dépendant que de  $n$  et pas de l'intervalle  $[a,b]$ , on peut en calculer ces valeurs une fois pour toutes. Les  $b_i$  sont aussi solutions de (5) avec  $h = 1$ .

En utilisant un programme de résolution des systèmes linéaires à coefficients entiers (Cf. (3.4.6.2) du chapitre sur l'analyse numérique linéaire), on obtient les résultats du tableau 6.1 :

n	b <sub>0</sub>	b <sub>1</sub>	b <sub>2</sub>	b <sub>3</sub>	b <sub>4</sub>	b <sub>5</sub>
1	1/2	1/2				
2	1/3	4/3	1/3			
3	3/8	9/8	9/8	3/8		
4	14/45	64/45	8/15	64/45	14/45	
5	95/288	125/96	125/144	125/144	125/96	95/28

Tableau 6.1

### 2.3.2 Majoration de l'erreur dans les méthodes de Newton\_Cotes

On peut montrer les résultats suivants, qui dépendent de la parité de n :

*Théorème :* (1) Si n est pair et f de classe  $C^{n+2}$ , alors il existe x dans [a,b] tel que :

$$E(f) = K_n \cdot f^{(n+2)}(\xi) \cdot \frac{h^{n+3}}{(n+2)!}, \text{ où } K_n = \int_0^n t^2 \cdot (t-1) \dots (t-n) dt$$

(2) Si n est impair et f de classe  $C^{n+2}$ , alors il existe x dans [a,b] tel que :

$$E(f) = K_n \cdot f^{(n+1)}(\xi) \cdot \frac{h^{n+2}}{(n+1)!}, \text{ où } K_n = \int_0^n t \cdot (t-1) \dots (t-n) dt$$

*Démonstration* — Cf Stoer et Burlisch p. 123.

En fait, comme la suite des polynômes d'interpolation de f ne tend pas vers f quand n tend vers  $+\infty$ , il y a peu de chances que les méthodes de Newton-Cotes convergent quand n tend vers  $+\infty$ . Le bon point de vue sera, non pas de travailler sur tout l'intervalle [a,b], mais de découper cet intervalle en un grand nombre de petits intervalles et d'utiliser une méthode de Newton-Cotes sur chacun de ces intervalles, ce qui donne des résultats très satisfaisants en général.

Nous allons maintenant regarder ce qui se passe pour quelques valeurs classiques de n.

## 2.4 Méthodes classiques d'intégration

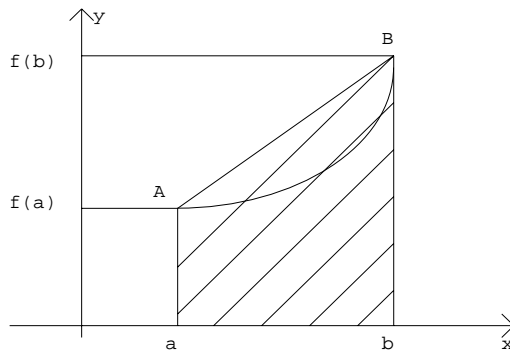
### 2.4.1 Cas n = 1. Méthode des trapèzes

Dans ce cas, on a  $h = b - a$ ,  $x_0 = a$  et  $x_1 = b$ , ce qui donne :

$$a_0 = a_1 = \frac{b-a}{2}$$

et la formule d'intégration approchée :

$$\int_a^b f(x) dx \cong \frac{b-a}{2} \cdot (f(a) + f(b))$$



Cette valeur approchée  
représente l'aire du trapèze  
hachuré.  
Ce qui consiste aussi à  
remplacer l'arc AB par sa  
corde.

Figure 6.1

*Méthode composite des trapèzes* — L'erreur de méthode est alors donnée par :

$$R(f) \leq \frac{M_2}{12} \cdot (b-a)^3$$

où  $M_2 = \text{Sup}\{|f''(x)| ; a \leq x \leq b\}$ .

Pour un intervalle  $[a, b]$  trop grand, la précision est donc médiocre.

La première idée est donc de découper  $[a, b]$  en un grand nombre de sous-intervalles et d'appliquer la méthode des trapèzes sur chacun de ces sous-intervalles.

On pose donc, pour  $n > 0$  dans  $\mathbb{N}$  :

$$h = \frac{b-a}{n}, \quad x_i = a + i \cdot h \quad (i = 0, \dots, n)$$

ces notations étant totalement indépendantes des précédentes.

Et en écrivant :

$$I(f) = \int_a^b f(x) dx = \sum_{i=0}^{n-1} \int_{x_i}^{x_{i+1}} f(x) dx$$

on aboutit à la formule d'intégration approchée :

$$I(f) \cong T(h) = h \cdot \left( \frac{f(a) + f(b)}{2} + \sum_{i=0}^{n-1} f(x_i) \right)$$

Une majoration de l'erreur étant :

$$R(f) \leq \frac{M_2 \cdot (b-a)}{12} \cdot h^2$$

On a donc une méthode d'ordre 2.

*Remarque* — A l'aide de la « formule d'Euler-Maclaurin », on peut préciser le reste dans la méthode des trapèzes. Cette formule s'écrit, pour une fonction  $f$  de classe  $C^{2n+2}$  :

$$\int_a^b f(x) dx =$$

$$T(h) - \sum_{k=1}^n \frac{B_k}{(2 \cdot k)!} \left( f^{(2k-1)}(b) - f^{(2k-1)}(a) \right) \cdot h^{2k} - \frac{B_{2n+2}}{(2 \cdot n + 1)!} \cdot (b-a) \cdot h^{2n+2} \cdot f^{(2n+2)}(\xi)$$

où  $\xi$  est dans  $[a, b]$  et les  $B_{2k}$  sont les nombres de Bernoulli définis par le développement en série entière :

$$\frac{t}{e^t - 1} = \sum_{k=0}^{+\infty} B_k \cdot \frac{t^k}{k!}$$

soit, pour les premières valeurs :

$$B_0 = 1, B_1 = -\frac{1}{2}, B_2 = \frac{1}{6}, B_4 = -\frac{1}{30}, B_6 = \frac{1}{42}, B_8 = -\frac{1}{30}, B_{10} = \frac{5}{66}, B_{12} = -\frac{691}{2730}, \dots$$

et pour  $k \geq 1, B_{2k+1} = 0$ .

*Méthode des trapèzes dichotomique* — On peut remarquer, dans la méthode composite des trapèzes, que si l'on double le nombre de points, soit avec  $n' = 2 \cdot n$ ,  $h' = \frac{h}{2}$ , alors dans le calcul de  $T(h')$ , on peut exploiter les calculs utilisés pour  $T(h)$ , soit :

$$T(h') = \frac{h}{2} \cdot \left( \frac{f(a) + f(b)}{2} + \sum_{i=0}^{n'-1} f(x'_i) \right)$$

et dans  $\sum_{i=1}^{n'-1} f(x'_i)$ , tous les  $x'_i$  avec  $i$  pair sont les termes précédents, il suffit donc de calculer les  $n$  termes  $f(x'_{2k-1})$ , pour  $k = 1, \dots, n$ . Et à ce stade l'erreur est un  $O\left(\left(\frac{h}{2}\right)^2\right)$ .

D'où l'algorithme de calcul :

*Etape 1* —  $n = 1, h = (b - a)$  et  $T_1 = \frac{b - a}{2} \cdot (f(a) + f(b))$

*Etape 2* —  $n = 2, h = \frac{b - a}{2}$  et  $T_2 = \frac{1}{2} \cdot (T_1 + (b - a) \cdot f(x_1))$

*Etape 3* —  $n = 4, h = \frac{b - a}{4}$  et  $T_4 = \frac{1}{2} \cdot (T_2 + \frac{b - a}{2} \cdot (f(x_1) + f(x_3)))$

Dans les programmations structurées qui suivent, la procédure Trapèze Dichotomique décrit une itération de la méthode, puis la procédure Integration Trapèze décrit le calcul de l'intégrale, un test d'arrêt étant :

$$\text{Iter} = \text{MaxIteration} \text{ ou } |T_{\text{Iter}} - T_{\text{Iter}-1}| < \text{Eps}$$

où le nombre maximum d'itérations MaxIter et la précision Eps sont donnés en constantes. MaxIter ne devra pas être trop grand car  $n = 2^p$  avec  $p \leq \text{MaxIter}$ . MaxIter = 10 est plus que suffisant.

*PROCEDURE* Trapèze Dichotomique(Entrée  $f$  : Fonction,  $a, b$  : Réel ; Iter : Entier ;  
Entrée\_Sortie  $n$  : Entier ;  $T$  : Réel) ;

Début

Si Iter = 1

Alors Début

$$T = \frac{b - a}{2} \cdot (f(a) + f(b)) ;$$

```

    n = 1 ;
Fin
Sinon Début
    h =  $\frac{b-a}{n}$  ; x = a +  $\frac{h}{2}$  ;
    Somme = 0 ;
    Pour j allant de 1 à n Faire
    Début
        Somme = Somme + f(x) ;
        x = x + h ;
    Fin ;
    T =  $\frac{1}{2} \cdot (T + h \cdot \text{Somme})$  ;
    n = 2 · n ;
Fin ;
Fin ;

PROCEDURE Integration_Trapèze(Entrée f : Fonction ; a, b : Réel ;
                               Entrée_Sortie T : Réel) ;

Début
    Iter = 0 ;
    T1 = -1038 ;
    Répéter
        Iter = Iter + 1 ;
        Trapèze_Dichotomique(f,a,b,Iter,n,T) ;
        Correct = (Iter = MaxIter) ou ( $|T_{\text{Iter}} - T_{\text{Iter-1}}| < \text{Eps}$ ) ;
        T1 = T ;
    Jusqu'à Correct ;
Fin ;

```

### 2.4.2 Cas n = 2. Méthode de Simpson

Dans ce cas, on a  $h = \frac{b-a}{2}$ ,  $x_0 = a$ ,  $x_1 = \frac{a+b}{2}$  et  $x_2 = b$ , ce qui donne :

$$a_0 = \frac{1}{3} \cdot h, \quad a_1 = \frac{4}{3} \cdot h \quad \text{et} \quad a_2 = \frac{1}{3} \cdot h$$

et la formule d'intégration approchée :

$$\int_a^b f(x) dx \cong \frac{b-a}{6} \cdot \left( f(a) + 4 \cdot f\left(\frac{a+b}{2}\right) + f(b) \right)$$

L'erreur de méthode est alors donnée par :

$$R(f) \leq \frac{M_4}{90} \cdot h^5 = \frac{M_4}{2880} \cdot (b-a)^5$$

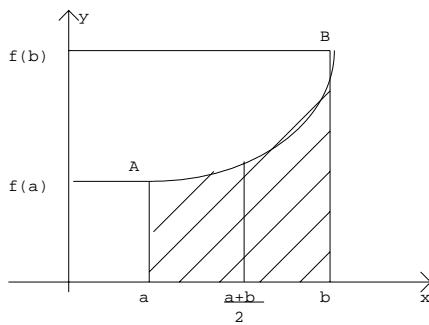


Figure 6.2

Cette valeur approchée représente l'aire du trapèze parabolique hachuré.  
Ce qui consiste aussi à remplacer l'arc AB par l'arc de parabole passant par  $a$ ,  $\frac{a+b}{2}$  et  $b$ .

*Méthode de Simpson composite* — Avec les notations utilisées pour la méthode composite des trapèzes, en découpant  $[a,b]$  en  $n$  sous-intervalles de même longueur et en utilisant la formule de Simpson sur chacun de ces sous-intervalles, on aboutit à la formule approchée :

$$\int_a^b f(x)dx \cong S(h) = \frac{h}{6} \cdot \left( f(a) + f(b) + 2 \cdot \sum f(x_i) + 2 \cdot f(m_i) \right)$$

où  $m_i = \frac{x_i + x_{i+1}}{2}$  ] est le point milieu et  $h = \frac{b-a}{n}$ .

Une majoration de l'erreur étant :

$$R(f) \frac{M_4 \cdot (b-a)}{2880} \cdot h^4$$

Pour la programmation, on peut exploiter les considérations précédentes en remarquant que :

$$S(h) = \frac{4}{3} \cdot T\left(\frac{h}{2}\right) - \frac{1}{3} \cdot T(h)$$

Comme pour la méthode des trapèzes, on peut utiliser une méthode dichotomique, ce qui donne la procédure :

*PROCEDURE* Integration\_Simpson(Entrée  $f$  : Fonction ;  $a, b$  : Réel ;  
Entrée\_Sortie  $S$  : Réel) ;

Début

Iter = 1 ;

Trapèze\_Dichotomique( $f,a,b,Iter,n,T$ ) ;

T1 := T ;

S1 =  $-10^{38}$  ; { Valeur précédente de  $S(h)$  }

Répéter

Iter = Iter + 1 ;

Trapèze\_Dichotomique( $f,a,b,Iter,n,T$ ) ; {  $T = T(h/2)$  }

$S = (4 \cdot T - T1)/3$  ; {  $S = S(h)$  }

Correct = (Iter = MaxIter) ou  $(|S - S_1| < Eps)$  ;

T1 = T ;

S1 = S ;

Jusqu'à Correct ;

Fin ;

### 2.4.3 Méthode de Romberg

En utilisant un procédé, due à *Richardson* en 1927, d'accélération de la convergence de certaines suites, on peut améliorer la méthode dichotomique des trapèzes. La méthode obtenue est due à Romberg (1955) et le principe consiste à construire la suite double  $(T(k,m))_{(k,m) \in \mathbb{N}^2}$  définie par la récurrence :

$$(1) \quad \begin{cases} T(0, m) = T\left(\frac{b-a}{2^m}\right), m \geq 0 \\ T(k, m) = \frac{1}{4^k - 1} \cdot (4^k \cdot T(k-1, m+1) - T(k-1, m)) \end{cases}$$

où  $T\left(\frac{b-a}{2^m}\right)$  est l'approximation de l'intégrale obtenue par la méthode des trapèzes avec le pas  $\frac{b-a}{2^m}$ .

Les propriétés de la suite  $(T(k,m))$  sont résumées dans le :

*Théorème* : (1) Il existe une suite double  $(a(k,m))_{(k,m) \in \mathbb{N}^2}$  telle que :

$$T(k,m) = \sum_{p=0}^k a(k,p) \cdot T(0, p+m)$$

où :

$$a(k, m) = (-1)^{k+m} \cdot \frac{b(k, m)}{d_k}$$

avec :

$$d_k = 1 \cdot 3 \cdot 15 \cdot \dots \cdot (4k-1) \cong \text{Cste} \cdot 2^{k(k+1)}$$

et les  $b(k,m)$  définis par :

$$\begin{cases} b(k, 0) = 1 \\ b(k, m) = 4^k \cdot b(k-1, m-1) + b(k-1, m) \end{cases}$$

en particulier :

$$\begin{cases} b(k, m) > 0, \text{ pour } m = 0, 1, \dots, k \\ b(k, m) = 0, \text{ pour } m > k \end{cases}$$

(2)  $\lim_{m \rightarrow +\infty} T(0, m) = I(f) = \int_a^b f(x) dx$

et avec (1), on déduit que :

$$\forall k \in \mathbb{N}, \quad \lim_{m \rightarrow +\infty} T(k, m) = I(f)$$

(3)  $T(k,m) - I(f) = O(4^{-(k+1)m})$

(4)  $\forall m \in \mathbb{N}, \quad \lim_{k \rightarrow +\infty} T(k, m) = I(f)$

*Démonstration* — Cf Warusfel, R. M. S. Janvier 1986.

*Remarque 1* — Avec (3), on déduit que la convergence de  $T(k,m)$  vers  $I(f)$ , à  $k$  fixé, est d'autant plus rapide que  $k$  est grand.



Pour  $k = 0$ , on retrouve la méthode des trapèzes dichotomiques, pour  $k = 1$ , c'est la méthode de Simpson, pour  $k = 2$ , c'est la méthode de Boole-Villarceau et pour  $k \geq 3$  il ne s'agit plus d'une méthode de Newton-Cotes.

*Remarque 2* — La suite qui est généralement utilisée pour approcher  $I(f)$  est celle définie par :  $R_k = T(k,0)$ .

La convergence de  $(R_k)$  vers  $I(f)$  est, en général très rapide, mais on ne sait rien sur la vitesse de convergence (dixit Warusfel).

La relation (3) ne permet pas de montrer cette convergence, la démonstration se faisant par une étude directe (Cf. Warusfel).

L'algorithme de calcul des  $R_k$  est alors le suivant :

$k = 0$  — Calcul de  $T(0,0) = R_0$  par Trapèze-Dichotomique ;

$k = 1$  — Calcul de  $T(0,1)$  par Trapèze-Dichotomique,  
Calcul de  $T(1,0) = R_1$  en fonction de  $T(0,1)$  et  $T(0,0)$  ;

$k = 2$  — Calcul de  $T(0,2)$  par Trapèze-Dichotomique,  
Calcul de  $T(1,1)$  en fonction de  $T(0,2)$  et  $T(0,1)$ ,  
Calcul de  $T(2,0) = R_2$  en fonction de  $T(1,1)$  et  $T(1,0)$  ;

$k$  quelconque — Calcul de  $T(0,k)$  par Trapèze-Dichotomique,  
Calcul de  $T(1,k - 1)$  en fonction de  $T(0,k)$  et  $T(0,k-1)$ ,

.....  
Calcul de  $T(k-j,j)$  en fonction de  $T(k-j-1,j+1)$  et  $T(k-j-1,j)$ ,

.....  
Calcul de  $T(k,0) = R_k$  en fonction de  $T(k - 1,1)$  et  $T(k - 1,0)$ .

En fait, on peut remarquer qu'il est inutile de travailler avec un tableau à deux dimensions, ce qui revient à conserver en mémoire des valeurs qui ne sont pas utilisées par la suite.

Il est plus astucieux d'utiliser un seul tableau :

$$T = (T(0), \dots, T(kMax))$$

et à l'étape  $k$  du calcul de stocker  $T(k - j, j)$  dans  $T(j)$ , pour  $j$  allant de  $k$  à  $0$ . Ce qui donne, pour l'étape  $k$  :

Calcul de  $T(k)$  par Trapèze-Dichotomique ;

Pour  $j$  allant de  $k - 1$  à  $0$   $T(j)$  est fonction de  $T(j + 1) = T(k - (j + 1), j+1)$ , qui vient d'être calculé, et de  $T(j) = T((k - 1) - j, j)$ , qui a été calculé à l'étape  $k - 1$ .

Ce qui donne en définitive, la procédure ci-dessous, où  $kMax$  est une constante (par ex.  $kMax = 15$ ) ;  $T$  est un tableau de réels numérotés de  $0$  à  $kMax$  ; pour  $k = 0$ ,  $T(0)$  contient  $T(k,0)$ , c'est-à-dire la valeur  $R$  de l'intégrale obtenue par la méthode de Romberg, et  $T(k) = T(0,k)$  est la valeur calculée par la méthode des trapèzes.

Comme test d'arrêt, on prendra :  $|T(k,0) - T(k - 1,0)| < Eps$ , où  $Eps$  est une précision donnée en constante.

*PROCEDURE Intégration\_Romberg*(Entrée  $f$  : Fonction ;  $a, b$  : Réel ;  
Entrée\_Sortie  $R$  : Réel) ;

*Début*

*Trapèze\_Dichotomique(f,a,b,l,n,Tl) ;*

```

T(0) = T1 ; { T1 est réinjecté à l'étape suivante d'où la nécessité de cette variable. }
R1 = -MaxRéel ;
k = 0 ;
Répéter
  k = k + 1 ;
  Trapèze_Dichotomique(f,a,b,k,n,T1) ;
  T(k) = T1 ;
  p = 1 ;
  Pour j allant de k - 1 à 0 Faire
  Début
    p = 4p ;
    T(j) = (p·T(j + 1) - T(j))/(p - 1) ;
  Fin ;
  R = T(0) ;
  Correct = (k = kMax) ou (|R - R1| < Eps) ;
  R1 = R ;
Jusqu'à Correct ;
Fin ;

```

### 3. Utilisation des polynômes orthogonaux. Quadratures de Gauss

#### 3.1 Introduction. Idées des méthodes

Pour calculer, de façon approchée,  $I(f) = \int_a^b f(x)dx$ , l'idée des méthodes de Newton-Cotes est, étant données des abscisses  $x_1, \dots, x_n$  dans  $[a,b]$ , de trouver des coefficients  $a_1, \dots, a_n$  tels que

$$(1) \quad \int_a^b P(x)dx = \sum_{i=1}^n a_i \cdot P(x_i), \text{ pour tout } P \text{ dans } R_{n-1}[x]$$

et de prendre comme formule de quadrature :

$$(2) \quad I(f) = \int_a^b f(x)dx \cong \hat{I}(f) = \sum_{i=1}^n a_i \cdot f(x_i)$$

La question que l'on peut se poser est alors : peut-on choisir les  $x_i$  de façon optimale ?

Pour répondre à la question ci-dessus, l'idée des méthodes de Gauss est de demander que (1) soit vraie pour des polynômes de degré maximum. Pour ce faire, on prendra aussi comme inconnues les abscisses  $x_i$ , ce qui va doubler le nombre de degrés de liberté et, de ce fait va aussi augmenter le temps de calcul.

Une fois les  $x_i$  et  $a_i$  déterminés, pour le même nombre d'évaluations que dans une méthode de Newton-Cotes correspondante au même nombre de points, on aura une formule d'intégration exacte pour des polynômes de degré élevé, ce qui peut faire espérer à une meilleure approximation de  $I(f)$ , du moins pour les fonctions « lisses », c'est-à-dire qui peuvent être bien approchées par des polynômes. Mais pour d'autres cas il se pourra que l'approximation ne soit pas meilleure.

Un autre intérêt des méthodes de Gauss est qu'elles permettent aussi de calculer certaines intégrales impropres du type

$$I(f) = \int_a^b f(x) \cdot \pi(x) dx$$

où  $\pi : ]a, b[ \rightarrow \mathbb{R}^+$  est une fonction continue positive appelée « *fonction poids* » et  $-\infty \leq a < b \leq +\infty$ . La fonction poids contiendra, en général, les singularités de l'intégrale à calculer.

On pourra, par exemple, considérer des intégrales impropres du type :

$$\int_{-1}^{+1} \frac{f(x)}{\sqrt{1-x^2}} dx, \int_0^{+\infty} f(x) \cdot e^{-x} dx \text{ ou } \int_{-\infty}^{+\infty} f(x) \cdot e^{-x^2} dx$$

Dans la mesure où les abscisses  $x_i$  seront imposées ces méthodes s'appliqueront aux fonctions définies analytiquement, mais pas aux fonctions connues en quelques points si ces derniers ne sont pas les points privilégiés trouvés.

## 3.2 Notations. Position du problème

### 3.2.1 Notations

On note  $\mathbb{R}[x]$  l'espace vectoriel des polynômes à coefficients réels et pour  $n \geq 0$ ,  $\mathbb{R}_n[x]$  est le sous-espace vectoriel formé des polynômes de degré au plus  $n$ .

On se donne une fonction poids continue :

$$\pi : ]a, b[ \rightarrow \mathbb{R}^+$$

où  $-\infty \leq a < b \leq +\infty$  et telle que :

$$\forall k \in \mathbb{N}, \int_a^b |\pi(x) \cdot x^k| dx < +\infty$$

Enfin, on pose :

$$E = \{ f : ]a, b[ \rightarrow \mathbb{R} ; \text{continue} ; \int_a^b |f(x)|^2 \cdot \pi(x) dx < +\infty \}$$

alors,  $E$  est un espace vectoriel sur  $\mathbb{R}$ , contenant  $\mathbb{R}[x]$  et qui peut être muni d'un produit scalaire en posant :

$$\forall f, g \in E, \langle f | g \rangle = \int_a^b f(x) \cdot g(x) \cdot \pi(x) dx$$

### 3.2.2 Position du problème

Pour  $n \geq 1$  donné, on cherche des coefficients  $x_1, \dots, x_n$  dans  $]a, b[$ ,  $a_1, \dots, a_n$  dans  $\mathbb{R}$  et  $m$  dans  $\mathbb{N}$  aussi grand que possible tels que :

$$\forall P \in \mathbb{R}_m, \int_a^b P(x) \cdot \pi(x) dx = \sum_{i=1}^n a_i \cdot P(x_i) \quad (3)$$

Ce qui équivaut encore au système de  $(m+1)$  équations à  $2 \cdot n$  inconnues :

$$\sum_{i=1}^n a_i \cdot x_i^k = \int_a^b \pi(x) \cdot x^k dx \quad (k = 0, 1, \dots, m) \quad (4)$$

Si on veut avoir une chance de résoudre ce système, il est raisonnable de prendre  $m = 2 \cdot n - 1$ .

Ensuite, pour  $f$  dans  $E$ , on utilisera la formule de quadrature :

$$I(f) = \int_a^b f(x) \cdot \pi(x) dx \cong \hat{I}(f) = \sum_{i=1}^n a_i \cdot f(x_i) \quad (5)$$

Enfin, comme toujours, il faudra évaluer l'erreur :

$$E(f) = |I(f) - \hat{I}(f)| \quad (6)$$

### 3.3 Calcul des abscisses $x_i$ et des coefficients de Gauss $a_i$

#### 3.3.1 Remarque

Calculer les  $x_i$  équivaut à trouver un polynôme  $P_n$  de degré  $n$  qui s'annule en tous les  $x_i$ . Et il s'agit alors de traduire (4) sur  $P_n$ .

On pose donc :

$$P_n(x) = \alpha_n \cdot \prod_{i=1}^n (x - x_i) \quad (7)$$

où  $\alpha_n \in \mathbb{R}$ ,  $x_1 < \dots < x_n$  dans  $]a, b[$  sont les inconnues.

#### 3.3.2 Condition nécessaire et suffisante sur $P$

*Lemme* : Les abscisses  $x_i$  vérifient (4) si, et seulement si :

$$\forall Q \in R_{n-1}, \int_a^b P_n(x) \cdot Q(x) \cdot \pi(x) dx = 0 \quad (8)$$

Ce qui peut se traduire en disant que  $P_n$  est orthogonal à  $R_{n-1}[x]$ .

*Démonstration* — C. N. Si les  $x_i$  vérifient (4), comme  $P_n \cdot Q$  est dans  $R_{2n-1}[x]$ , pour tout  $Q$  dans  $R_{n-1}[x]$ , on a nécessairement :

$$\int_a^b P_n(x) \cdot Q(x) \cdot \pi(x) dx = \sum_{i=1}^n a_i \cdot P_n(x_i) \cdot Q(x_i) = 0$$

C. S. Supposons que  $P_n$  vérifie (8). Soit  $P$  dans  $R_{2n-1}[x]$  ; par division euclidienne, on peut écrire :

$$P = Q \cdot P_n + R_n \quad (d^\circ(Q) \leq n-1, d^\circ(R_n) \leq n-1)$$

et avec (8) :

$$\int_a^b P(x) \cdot \pi(x) dx = \int_a^b \pi(x) \cdot R_n(x) dx$$

Il suffit donc de vérifier (4) pour  $k = 0, 1, \dots, n-1$ , car  $R_n(x_i) = P(x_i)$  pour tout  $i = 1, \dots, n$ .

Ce qui donne le système linéaire, aux inconnues  $a_1, \dots, a_n$  :

$$\sum_{i=1}^n x_i^k \cdot a_i = \int_a^b \pi(x) \cdot x^k dx \quad (k = 0, 1, \dots, n-1)$$

Le déterminant de ce dernier étant du type Vandermonde, il est non nul et il existe une unique solution  $a_1, \dots, a_n$ .

#### 3.3.3 Détermination explicite de $P_n$

Il s'agit de trouver un polynôme  $P_n$  orthogonal à  $R_{n-1}[x]$  et n'ayant que des racines simples. Puis la recherche des racines de  $P_n$  nous donnera les  $x_i$  et enfin la résolution d'un système de Vandermonde nous donnera les  $a_i$ .

En fait, on verra que pour les méthodes de Gauss classiques, les  $a_i$  seront donnés explicitement et il est inutile de résoudre le système linéaire.

**Lemme :** Si  $P_n$  est orthogonal à  $R_{n-1}[x]$ , alors il n'a que des racines simples.

*Démonstration* — On a :

$$\forall Q \in R_{n-1}[x], \int_a^b P_n(x) \cdot Q(x) \cdot \pi(x) dx = 0$$

Pour  $Q = 1$ , on a donc  $\int_a^b P_n(x) \cdot \pi(x) dx = 0$ , avec  $\pi$  continue strictement positive et  $P_n$  continue, donc  $P_n$  ne peut garder un signe constant c'est-à-dire que  $P_n$  admet au moins une racine d'ordre impair.

En notant  $x_1 < \dots < x_k$  toutes les racines d'ordre impair de  $P_n$ , avec  $k \geq 1$ , on a :

$$P_n(x) = (x - x_1) \cdot \dots \cdot (x - x_k) \cdot Q(x)$$

où  $Q$  est dans  $R_{n-1}[x]$  et garde un signe constant sur  $]a, b[$ .

Si  $k \leq n - 1$ , alors  $R(x) = (x - x_1) \cdot \dots \cdot (x - x_k)$  est dans  $R_{n-1}[x]$  et  $\int_a^b P_n(x) \cdot R(x) \cdot \pi(x) dx = 0$  avec  $R(x) \cdot P_n(x)$  de signe constant, ce qui est impossible, donc  $k = n$ .

On a donc :  $P_n(x) = \alpha_n \cdot \prod_{i=1}^n (x - x_i)$ .

*Remarque* — Localisation des zéros de  $P_n$ .

Pour localiser les racines de  $P_n$ , on dispose du :

**Théorème :** En notant  $x_1, \dots, x_n$  les zéros de  $P_n$  et  $y_1, \dots, y_{n-1}$  ceux de  $P_{n-1}$ , on a :

$$0 < x_1 < y_1 < x_2 < y_2 < \dots < x_{n-1} < y_{n-1} < x_n$$

Ce qui permet de les localiser de proche en proche.

*Conclusion* — Pour construire  $P_n$ , il suffit de se donner une base orthogonale de  $R_n[x]$ ,  $\{P_0, P_1, \dots, P_n\}$ .

Le procédé d'orthogonalisation de Gram-Schmidt nous permet toujours de déterminer une telle base à partir de la base canonique de  $R_n[x]$ .

En pratique, on verra que dans les cas classiques, une telle base est connue. Il suffira alors de trouver les  $n$  racines de  $P_n$ , en utilisant la méthode de Newton (par exemple), puis de trouver les coefficients  $a_n$ .

### 3.4 Propriétés des polynômes orthogonaux. Calcul des coefficients de Gauss

Dans ce paragraphe, on se donne une famille orthonormée de polynôme  $(P_n)_{n \geq 0}$ , relativement au produit scalaire défini par une fonction poids  $\pi$  sur  $E$ .

On a donc :

$$\langle P_k | P_j \rangle = \int_a^b P_k(x) \cdot P_j(x) dx = \begin{cases} 1 & \text{si } k = j \\ 0 & \text{si } k \neq j \end{cases} \quad (0 \leq j, k \leq n)$$

et  $d^\circ(P_k) = k$ .

En particulier, on en déduit que :

$$\langle P_k | x^j \rangle = 0, \text{ pour } j = 0, 1, \dots, k - 1$$

et si  $P_k(x) = \alpha_k \cdot x^k + \beta_k \cdot x^{k-1} + \dots + \gamma_0$  alors :  $\langle P_k | x^k \rangle = \frac{1}{\alpha_k}$

### 3.4.1 Récurrence vérifiée par les polynômes orthogonaux

Un algorithme de calcul des polynômes orthogonaux est donné par le :

*Théorème* : Si  $(P_n)$  est une famille orthonormée de polynômes dans  $E$ , on a alors la relation de récurrence :

$$x \cdot P_n(x) = a_n \cdot P_{n+1}(x) + b_n \cdot P_n(x) + c_n \cdot P_{n-1}(x) \quad (n \geq 1)$$

avec :

$$a_n = \frac{\alpha_n}{\alpha_{n+1}}, \quad b_n = \frac{\beta_n}{\alpha_n} - \frac{\beta_{n+1}}{\alpha_{n+1}} \quad \text{et} \quad c_n = \frac{\alpha_{n-1}}{\alpha_n} = a_{n-1}$$

où  $P_n(x) = \alpha_n \cdot x^n + \beta_n \cdot x^{n-1} + \dots + \gamma_0$ .

*Démonstration* — Cf Ayant et Borg p. 117.

### 3.4.2 Formule de Darboux-Christoffel

*Théorème* : Si  $(P_n)$  est une famille orthonormée de polynômes, avec les notations ci-dessus, pour  $x, y$  dans  $R$ , on a :

$$(x - y) \cdot \sum_{k=0}^n P_k(x) \cdot P_k(y) = \frac{\alpha_n}{\alpha_{n+1}} \cdot (P_{n+1}(x) \cdot P_n(y) - P_n(x) \cdot P_{n+1}(y))$$

*Démonstration* — Cf. Ayant et Borg p. 119.

### 3.4.3 Calcul des coefficients de Gauss $a_k$

En (2.2.2), on a vu que les coefficients  $a_k$  sont donnés par :

$$a_k = \int_a^b L_k(x) \cdot \pi(x) dx \quad (k = 1, \dots, n)$$

où :

$$L_k(x) = \prod_{\substack{j=1 \\ j \neq k}}^n \frac{x - x_j}{x_k - x_j} = \frac{P_n(x)}{x - x_k} \cdot \frac{1}{P_n'(x_k)}$$

donc,

$$a_k = \frac{1}{P_n'(x_k)} \cdot \int_a^b \frac{P_n(x)}{x - x_k} \cdot \pi(x) dx$$

Avec l'identité de Darboux-Christoffel, on a :

$$(x - x_k) \cdot \sum_{j=0}^n P_j(x) \cdot P_j(x_k) = -\frac{\alpha_n}{\alpha_{n+1}} \cdot P_{n+1}(x_k) \cdot P_n(x)$$

et avec  $\int_a^b P_j(x) \cdot \pi(x) dx = \langle P_j | P_0 \rangle = 0$ , pour  $j \geq 1$ , on a :

$$P_0(x_k) \cdot \int_a^b P_0(x) \cdot \pi(x) dx = -\frac{\alpha_n}{\alpha_{n+1}} \cdot P_{n+1}(x_k) \cdot \int_a^b \frac{P_n(x)}{x - x_k} \cdot \pi(x) dx$$

puis, avec  $P_0(x_k) = P_0(x) = P_0$ , on déduit que :

$$P_0(x_k) \cdot \int_a^b P_0(x) \cdot \pi(x) dx = \langle P_0 | P_0 \rangle = 1$$

d'où :

$$a_k = -\frac{\alpha_{n+1}}{\alpha_n} \cdot \frac{1}{P_n'(x_k) \cdot P_{n+1}(x_k)}$$

Enfin, avec la récurrence vérifiée par les  $P_n$ , on a :

$$0 = a_n \cdot P_{n+1}(x_k) + c_n \cdot P_{n-1}(x_k)$$

soit,

$$a_k = \frac{\alpha_n}{\alpha_{n-1}} \cdot \frac{1}{P_n'(x_k) \cdot P_{n-1}(x_k)}$$

*Remarque* — Les polynômes orthogonaux classiques sont, en général, plus faciles à manipuler non normalisés.

Si  $(P_n)$  est une famille de polynômes orthogonaux, on les normalise en posant :

$$P_n^* = \frac{1}{\lambda_n} \cdot P_n, \text{ où } \lambda_n = \sqrt{\langle P_n | P_n \rangle}$$

et la formule ci-dessus devient :

$$a_k = (\lambda_{n-1})^2 \cdot \frac{\alpha_n}{\alpha_{n-1}} \cdot \frac{1}{P_n'(x_k) \cdot P_{n-1}(x_k)}$$

### 3.5 Majoration de l'erreur de quadrature

On peut montrer le résultat suivant :

*Théorème (Markov)* : Si  $f$  est de classe  $C^{2n}$  sur  $]a, b[$ , alors l'erreur de quadrature dans la méthode de Gauss est donnée par :

$$E(f) = \frac{f^{(2n)}(\xi)}{(2n)!} \cdot \frac{\langle P_n | P_n \rangle}{\alpha_n^2}$$

où  $\xi$  est dans  $]a, b[$ .

*Démonstration* — (Cf. Stoer et Burlisch p. 150, ou Deheuvelds p. 217)

### 3.6 Exemples classiques de polynômes orthogonaux et formules de quadrature correspondantes

Pour les propriétés des polynômes orthogonaux, on pourra se reporter à Ayant et Borg.

La programmation des méthodes correspondantes se trouve dans le paquetage INTEGRATION\_GENERIQUE sur la disquette.



### 3.6.1 Les polynômes de Legendre

Ce sont les polynômes définis, à une constante multiplicative près, par :

$$[a, b] = [-1, 1] \text{ et } \pi(x) = 1$$

Ils sont classiquement définis par la formule de Rodrigues :

$$P_n(x) = \frac{1}{2^n \cdot n!} \cdot \left( \frac{d}{dx} \right)^n \left( (x^2 - 1)^n \right)$$

soit :

$$P_n(x) = \frac{1}{2^n} \cdot \sum_{k=0}^{\left[ \frac{n}{2} \right]} (-1)^k \cdot C_n^k \cdot C_{2(n-k)}^n \cdot x^{n-2k}$$

où  $[ ]$  désigne la fonction partie entière.

On peut aussi les définir par la récurrence :

$$\begin{cases} P_0(x) = 1, P_1(x) = x \\ (n+1) \cdot P_{n+1}(x) = (2 \cdot n + 1) \cdot x \cdot P_n(x) - n \cdot P_{n-1}(x) \quad (n \geq 1) \end{cases}$$

On peut encore définir  $P_n$  comme la solution polynomiale de l'équation différentielle :

$$\begin{cases} (1-x^2) \cdot y''(x) - 2 \cdot x \cdot y'(x) + n \cdot (n+1) \cdot y(x) = 0 \\ y(1) = 1 \end{cases}$$

ou par :

$$\varphi(x, z) = (1 - 2 \cdot x \cdot z + z^2)^{-\frac{1}{2}} = \sum_{n=0}^{+\infty} P_n(x) \cdot z^n \quad (|z| < 1)$$

On dit que  $\varphi$  est la « fonction génératrice » des  $P_n$ .

On a la majoration :

$$\forall x \in [-1, 1], |P_n(x)| \leq 1$$

Enfin, on peut localiser les zéros de  $P_n$  grâce aux *inégalités de Brun* (Cf. Ayant et Borg p. 138) :

$$\cos\left(\frac{2k}{2n+1} \cdot \pi\right) \leq x_k \leq \cos\left(\frac{2k-1}{2n+1} \cdot \pi\right)$$

valables pour tous les zéros strictement positifs de  $P_n$ , et :

$$x_k = -x_{n+1-k} \quad (k = 1, \dots, n)$$

ce qui résulte de  $P_n(-x) = (-1)^n \cdot P_n(x)$  (donc les racines vont former un ensemble symétrique par rapport à 0).

On déduit alors de ce qui précède que les  $a_k$  sont donnés par :

$$a_k = \frac{2 \cdot (1 - x_k^2)}{(n \cdot P_{n-1}(x_k))^2} \quad (n = 1, \dots, n)$$

avec aussi :  $a_k = a_{n+1-k}$  ( $k = 1, \dots, n$ ).

La programmation du calcul des  $P_n$ , de leurs racines  $x_k$  et des coefficients  $a_k$  est alors élémentaire. L'exécution d'un tel programme nous donne, par exemple, les résultats suivants :

*Polynôme  $L_{24}$*

$$L_{24}(x) = 0.16118 - 48.35408 \cdot x^2 + 2393.52683 \cdot x^4$$

$$\begin{aligned}
& -46274.85201 \cdot x^6 + 461095.84685 \cdot x^8 \\
& -2705095.63485 \cdot x^{10} + 10041642.88695 \cdot x^{12} - 24497194.73519 \cdot x^{14} \\
& + 39807941.44469 \cdot x^{16} - 42669950.30673 \cdot x^{18} \\
& + 28970650.47141 \cdot x^{20} - 11287266.41743 \cdot x^{22} + 1922106.96239 \cdot x^{24}
\end{aligned}$$

Zéros du polynôme  $L_{24}$

$$\begin{array}{ll}
x_1 = -x_{24} = -0.995187220000022 & x_2 = -x_{23} = -0.974728555978706 \\
x_3 = -x_{22} = -0.938274551976997 & x_4 = -x_{21} = -0.886415527020260 \\
x_5 = -x_{20} = -0.820001985964673 & x_6 = -x_{19} = -0.740124191575882 \\
x_7 = -x_{18} = -0.648093651937057 & x_8 = -x_{17} = -0.545421471388688 \\
x_9 = -x_{16} = -0.433793507626036 & x_{10} = -x_{15} = -0.315042679696163 \\
x_{11} = -x_{14} = -0.191118867473616 & x_{12} = -x_{13} = -0.064056892862606
\end{array}$$

Coefficients de Gauss du polynôme  $L_{24}$

$$\begin{array}{ll}
a_1 = a_{24} = 0.012341229289022 & a_2 = a_{23} = 0.028531388532063 \\
a_3 = a_{22} = 0.044277438545067 & a_4 = a_{21} = 0.059298584856864 \\
a_5 = a_{20} = 0.073346481465724 & a_6 = a_{19} = 0.086190161546685 \\
a_7 = a_{18} = 0.097618652106859 & a_8 = a_{17} = 0.107444270116811 \\
a_9 = a_{16} = 0.115505668053805 & a_{10} = a_{15} = 0.121670472927805 \\
a_{11} = a_{14} = 0.125837456346828 & a_{12} = a_{13} = 0.127938195346752
\end{array}$$

Les polynômes de Legendre peuvent être utilisés pour calculer une intégrale définie sur un intervalle  $[a,b]$ , le changement de variable :

$$x = \frac{b+a}{2} + \frac{b-a}{2} \cdot t$$

nous ramenant à  $[-1,1]$ .

Ce qui donne la formule de quadrature de Legendre :

$$I(f) \cong \hat{I}(f) = \frac{b-a}{2} \cdot \sum_{k=1}^n a_k \cdot f(\xi_k)$$

où :

$$\xi_k = \frac{b+a}{2} + \frac{b-a}{2} \cdot x_k \quad (k = 1, \dots, n)$$

Là encore la programmation est évidente.

Pour l'évaluation de l'erreur, avec :

$$\langle P_n | P_n \rangle = \int_{-1}^1 P_n(x)^2 dx = \frac{2}{2 \cdot n + 1}$$

et :

$$\alpha_n = \frac{(2n)!}{2^n \cdot (n!)^2}$$

on déduit que :

$$E(f) = \frac{2^{2n+1} \cdot (n!)^2 \cdot f^{(2n)}(\xi)}{(2 \cdot n + 1) \cdot ((2n)!)^3}$$

où  $\xi$  est dans  $]-1,1[$ .

### 3.6.2 Polynômes de Tchébychev

Ils sont définis, à une constante multiplicative près par :

$$[a,b] = [-1,1] \text{ et } \pi(x) = \frac{1}{\sqrt{1-x^2}}$$

Une définition classique étant :

$$T_n(x) = \text{Cos}(n \cdot \text{ArcCos}(x))$$

soit

$$T_n(x) = \frac{n}{2} \cdot \sum_{k=0}^{\lfloor \frac{n}{2} \rfloor} (-1)^k \cdot \frac{(n-k-1)!}{k! \cdot (n-2k)!} \cdot (2x)^{n-2k}$$

On peut aussi les définir par la récurrence :

$$\begin{cases} T_0(x) = 1, T_1(x) = x \\ T_{n+1}(x) = 2 \cdot x \cdot T_n(x) + T_{n-1}(x) \quad (n \geq 1) \end{cases}$$

On peut encore définir  $T_n$  comme la solution polynomiale de l'équation différentielle :

$$\begin{cases} (1-x^2) \cdot y''(x) - x \cdot y'(x) + n^2 \cdot y(x) = 0 \\ y(1) = 1 \end{cases}$$

ou par :

$$(1-t^2) \cdot (1-2 \cdot x \cdot z + z^2)^{-1} = T_0(x) + 2 \cdot \sum_{n=1}^{+\infty} T_n(x) \cdot z^n \quad (|z| < 1)$$

Les racines de  $T_n$  sont données par :

$$x_k = \text{Cos}\left(\frac{2 \cdot k - 1}{2 \cdot n} \cdot \pi\right) \quad (k = 1, \dots, n)$$

Là encore, on a :

$$|T_n(x)| \leq 1$$

Le calcul des  $a_k$  est ici très simple car ils sont tous égaux, soit :

$$a_k = \frac{\pi}{n} \quad (k = 1, \dots, n)$$

La formule d'intégration approchée est alors :

$$I(f) = \int_{-1}^1 \frac{f(x)}{\sqrt{1-x^2}} dx \cong \frac{\pi}{n} \cdot \sum_{i=1}^n f(x_i)$$

Avec  $\langle T_n | T_n \rangle = \frac{\pi}{2}$ , et  $\alpha_n = 2^{n-1}$ , on déduit que l'erreur de quadrature est :

$$E(f) = \frac{f^{(2n)}(\xi)}{(2 \cdot n)! \cdot 2^{2n-2}} \cdot \frac{\pi}{2}$$

### 3.6.3 Polynômes de Laguerre

Ils sont définis, à une constante multiplicative près par :

$$[a,b] = [0, +\infty[ \text{ et } \pi(x) = e^{-x}$$

Soit classiquement par :

$$L_n(x) = \frac{e^x}{n!} \cdot \left( \frac{d}{dx} \right)^n (x^n \cdot e^{-x})$$

donc :

$$L_n(x) = \sum_{k=0}^n (-1)^k \cdot C_n^{n-k} \cdot \frac{x^k}{k!}$$

Ils sont aussi définis par la récurrence :

$$\begin{cases} L_0(x) = 1, L_1(x) = 1 - x \\ (n+1) \cdot L_{n+1}(x) = (2 \cdot n + 1 - x) \cdot L_n(x) - n \cdot L_{n-1}(x) \quad (n \geq 1) \end{cases}$$

En utilisant le fait qu'entre deux racines de  $L_{n-1}$  il y a une seule racine de  $L_n$ , on peut localiser, de proche en proche les racines de  $L_n$  et les calculer par la méthode de Newton.

Les coefficients  $a_k$  sont donnés par :

$$a_k = \frac{x_k}{(n \cdot L_{n-1}(x_k))^2} \quad (k = 1, \dots, n)$$

Ce qui donne la formule de quadrature de Laguerre :

$$I(f) = \int_0^{+\infty} f(x) \cdot e^{-x} dx \cong \hat{I}(f) = \sum_{i=1}^n a_i \cdot f(x_i)$$

L'erreur de quadrature étant donnée par :

$$E(f) = \frac{(n!)^2}{(2 \cdot n)!} \cdot f^{(2n)}(\xi)$$

où  $\xi > 0$ .

### 3.6.4 Polynômes d'Hermite

Ils sont définis, à une constante multiplicative près par :

$$]a, b[ = \mathbb{R} \text{ et } \pi(x) = e^{-x^2}$$

Soit classiquement par :

$$H_n(x) = (-1)^n \cdot e^{x^2} \cdot \left( \frac{d}{dx} \right)^n (e^{-x^2})$$

C'est-à-dire :

$$H_n(x) = n! \cdot \sum_{k=0}^{\lfloor \frac{n}{2} \rfloor} (-1)^k \cdot \frac{(2 \cdot x)^{n-2k}}{k! \cdot (n-2 \cdot k)!}$$

Ils sont aussi définis par la récurrence :

$$\begin{cases} H_0(x) = 1, H_1(x) = 2 \cdot x \\ H_{n+1}(x) = 2 \cdot x \cdot H_n(x) - 2 \cdot n \cdot H_{n-1}(x) \quad (n \geq 1) \end{cases}$$

Comme  $H_n(-x) = (-1)^n \cdot H_n(x)$ , les zéros forment un ensemble symétrique par rapport à 0, soit :

$$x_i = -x_{n+1-i} \quad (i = 1, \dots, n)$$

En utilisant le fait qu'entre deux racines de  $H_{n-1}$  il y a une seule racine de  $H_n$ , on peut localiser, de proche en proche les racines de  $H_n$  et les calculer par la méthode de Newton.

Les coefficients  $a_k$  étant donnés par :

$$a_k = \frac{\sqrt{\pi} \cdot 2^{n-1} \cdot (n-1)!}{(H_{n-1}(x_k))^2 \cdot n}$$

L'erreur de quadrature étant :

$$E(f) = \frac{f^{(2n)}(\xi) \cdot \sqrt{\pi} \cdot n!}{(2 \cdot n)! \cdot 2^n}$$

## 4. La méthode probabiliste de Monte-Carlo

### 4.1 Nombres pseudo-aléatoires

Un ordinateur ne peut pas véritablement générer des nombres aléatoires, il ne peut générer que des nombres « *pseudo-aléatoires* » en ce sens que la connaissance d'un ou plusieurs premiers d'entre eux va entraîner la connaissance de tous les suivants. Les nombres obtenus auront une bonne qualité aléatoire si leur comportement, dans des simulations d'expériences aléatoires, est proche de véritables nombres aléatoires. On peut, par exemple, simuler une partie de « pile ou face » avec des nombres aléatoires entiers entre 0 et 1, et pour un grand nombre de tirages le nombre de 1 doit être sensiblement équivalent au nombre de 0.

Une méthode classique pour engendrer une suite de tels nombres est d'utiliser l'algorithme suivant :

$$\begin{cases} x_0 \text{ est donné impair} \\ x_n = a \cdot x_{n-1} \pmod{2^p} \end{cases}$$

avec  $p$  donné en général entre 30 et 40 ( $p$  dépend de la machine, on prend souvent  $p = 32$ ), le réel  $a$  donné comme puissance impaire élevée de 5 et la valeur initiale  $x_0$  (en anglais « *initial seed* ») également donnée.

L'avantage de cette congruence linéaire est d'être assez rapide.

$x_n$  étant le reste dans la division euclidienne de  $a \cdot x_{n-1}$  par  $2^p$ , on a :  $x_n < 2^p$ , donc la quantité  $r = \frac{x_n}{2^p - 1}$  va donner un nombre aléatoire entre 0 et 1.

On peut donc obtenir par cette méthode un moyen de faire un tirage aléatoire sur  $[0,1[$  et ce tirage peut être considéré comme uniforme.

On peut également simuler un tirage aléatoire suivant une loi de Gauss de moyenne et d'écart type donné ou toute autre loi classique (Cf Nougier, p. 305).

### 4.2 La méthode de Monte-Carlo.

#### *Première version : tirage par « noir ou blanc »*

#### 4.2.1 Equiprobabilité sur un segment, ou un pavé de $\mathbb{R}^n$

On appelle tirage uniforme sur un intervalle  $[a,b]$  un tirage aléatoire de réels  $x$  dans  $[a,b]$  tel que la probabilité que  $x$  soit dans l'intervalle  $I$  est proportionnelle à la longueur de  $I$ .

Pour un tel tirage on a donc :

$$\text{Prob}(x \in I) = \frac{|I|}{b - a}$$

De la même façon pour un tirage uniforme sur un pavé  $P = \prod_{i=1}^n [a_i, b_i]$  de  $R^n$ , on aura :

$$\text{Prob}(x \in D) = \frac{m(D)}{m(P)}$$

### 4.2.2 Calcul d'une intégrale simple

*Premières hypothèses* — On se donne une fonction continue  $f$  de  $[a,b]$  dans  $[0,1]$ .

La mesure du domaine hachuré est alors :  $m(D) = \int_a^b f(x)dx$

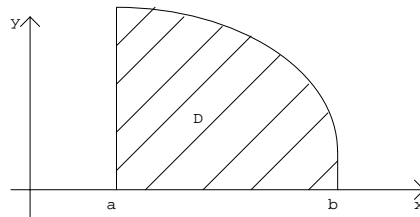


Figure 6.3

Et pour un tirage uniforme sur le rectangle  $P$ , on aura :

$$\text{Prob}((x,y) \in D) = \frac{\int_a^b f(x)dx}{b-a}$$

Donc, si on effectue  $n$  tirages dans  $P$  et si  $k$  est le nombre de fois où le point est tiré dans  $D$ , pour  $n$  grand on aura :

$$\text{Prob}((x,y) \in D) = \frac{k}{n} \cong \frac{\int_a^b f(x)dx}{b-a}$$

*Cas général* — Si la fonction  $f$  n'est pas à valeurs dans  $[0,1]$ , on s'y ramène facilement en faisant le changement de fonction :

$$g(x) = \frac{f(x) - y_{\text{Min}}}{y_{\text{Max}} - y_{\text{Min}}}$$

où  $y_{\text{Min}}$  [Resp.  $y_{\text{Max}}$ ] est la borne inférieure [Resp. Supérieure] de  $f$  sur  $[a,b]$ .

### 4.2.3 Calcul d'une intégrale multiple

Soit à calculer  $\int_{\Delta} f(x)dx$ , où  $\Delta = \prod_{i=1}^n [a_i, b_i]$  est un pavé de  $R^p$  et  $f$  est continue de  $\Delta$  dans  $[0,1]$ .

En reprenant les idées précédentes, une méthode de calcul approché consistera donc à faire  $n$  tirages aléatoires d'un  $p$ -uplet  $x = (x_1, \dots, x_n)$  et d'un réel  $y$  dans  $[0,1]$  puis à compter le nombre  $k$  de fois où  $y \leq f(x)$ . Alors  $\frac{k}{n} \cdot m(\Delta)$  donnera une valeur

approchée de l'intégrale où  $m(\Delta) = \prod_{i=1}^n (b_i - a_i)$ .

Comme dans le cas d'une variable on peut aussi traiter le cas où  $f$  n'est pas à valeurs dans  $[0,1]$ .

#### 4.2.4 Remarques

(i) — Cette méthode est compétitive pour les intégrales multiples car pour chaque tirage aléatoire on ne fait qu'une évaluation de  $f$  (opération coûteuse en temps de calcul). Le nombre d'opérations étant proportionnel à la dimension de l'intégrale.

(ii) — Comme il s'agit de tirages aléatoires, deux tirages différents donneront des résultats différents. Il est donc conseillé de faire plusieurs expériences et de faire la moyenne des résultats obtenus. Si  $I_k$  est la valeur obtenue à l'expérience  $k$  et  $m$  le nombre d'expériences, alors on prendra pour valeur approchée :

$$I = \frac{\sum_{i=1}^m I_k}{m}$$

Une mesure de la dispersion des  $I_k$  autour de  $I$  est alors donnée par l'écart type  $\sigma$  avec :

$$\sigma^2 = \frac{\sum_{i=1}^m (I_k - I)^2}{m}$$

La méthode est d'autant plus efficace que  $\sigma$  est petit. En augmentant le nombre d'expériences, on diminuera  $\sigma$  mais on augmentera le temps de calcul  $T$ . Une mesure de l'efficacité de la méthode est donnée par :

$$E = \frac{1}{T \cdot \sigma^2}$$

### 4.3 Méthode de Monte Carlo, deuxième version

#### 4.3.1 Remarques préliminaires

Si  $f : [a,b] \rightarrow \mathbb{R}$  est continue, sa valeur moyenne est :

$$m(f) = \frac{1}{b-a} \cdot \int_a^b f(x) dx = f(c)$$

pour un point  $c$  de  $[a,b]$ .

L'idée est alors de trouver cette valeur  $c$  « au hasard ».

Pour ce faire, on effectue  $n$  tirages, suivant une même loi uniforme, sur  $[a,b]$  et si au tirage  $n^\circ k$ ,  $x_k$  est le point obtenu, on prend comme valeur approchée de  $m(f)$  la quantité :

$$m_{k,n} = f(x_k) \quad (k = 1, \dots, n)$$

Une valeur approchée de  $m(f)$  sera alors donnée par la valeur moyenne :

$$m_n = \frac{1}{n} \cdot \sum_{k=1}^n m_{k,n}$$

L'idée de cette approximation est justifiée par la loi des grands nombres.

En effet,  $(x_k)_{k \geq 1}$  peut être considérée comme une suite de variables aléatoires suivant une même loi uniforme sur  $[a,b]$  et  $(f(x_k))_{k \geq 1}$  est alors une suite de



variables aléatoires suivant une même loi de moyenne  $m(f)$  et d'écart-type  $\sigma(f)$  définie par :

$$\sigma(f)^2 = \int_a^b (f(x) - m(f))^2 \cdot \frac{dx}{b-a}$$

la loi des grands nombres nous dit alors que la suite  $(m_n)$  va converger en probabilité vers la moyenne  $m(f)$  et la différence  $m_n - m(f)$  a pour écart-type  $\frac{\sigma(f)}{\sqrt{n}}$ , c'est-à-dire que pour tout

$\lambda > 0$ , on a :

$$\lim_{n \rightarrow +\infty} \text{Pr ob}(|m_n - m(f)| > \lambda) = 0$$

l'inégalité de Bienaymé-Tchebychev nous donnant la majoration :

$$\text{Pr ob}(|m_n - m(f)| > \lambda) < \frac{\sigma(f)^2}{n \cdot \lambda^2}$$

### 4.3.2 Méthode de Monte Carlo avec échantillonnage simple

Au tirage  $n^\circ k$ , on dispose de  $x_k$  dans  $[a,b]$ , ce qui donne  $y_k = f(x_k)$ , puis on fait la moyenne des  $y_k$ , ce qui donne  $m_n$  et enfin l'intégrale de  $f$  est approximée par  $(b-a) \cdot m_n$ , soit :

$$I(f) = \int_a^b f(x) dx \cong \frac{b-a}{n} \cdot \sum_{k=1}^n f(x_k)$$

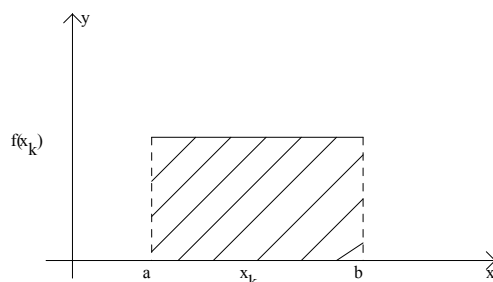


Figure 6.4

*Ce qui revient à remplacer le domaine délimité par les droites  $x = a$ ,  $x = b$ , l'axe des  $x$  et le graphe de  $f$  sur  $[a,b]$  par le rectangle hachuré.*

**Remarque : Echantillonnage stratifié** — Pour améliorer la précision on peut découper l'intervalle  $[a,b]$  en sous-intervalles plus petits et appliquer la méthode sur chacun de ces sous-intervalles. On parle alors d'échantillonnage stratifié. Pour un découpage très fin, on retrouve la définition d'une intégrale à l'aide des sommes de Riemann.

La programmation structurée est alors la suivante :

*FONCTION MonteCarlo(Entrée  $a, b$  : Réel ;  $n$  : Entier ;  $f$  : Fonction) ;*

*Début*

*$S = 0$*

*Pour  $k$  allant de 1 à  $n$  faire*

*Début*

*$r$  = Nombre aléatoire dans  $[0,1]$  ;*

*$x = a + (b - a) \cdot r$  ;*

*$S = S + f(x)$  ;*

Fin ;  
 $S = (b - a) \cdot S/n$  ;  
 MonteCarlo = S ;  
 Fin ;

### 4.3.3 Utilisation de transformations antithétiques

Au tirage n° k, on dispose d'un nombre aléatoire  $r_k$  dans  $[0,1[$  auquel on va associer deux points de  $[a,b[$  :

$$x'_k = a + (b - a) \cdot r_k \text{ et } x''_k = a + (b - a) \cdot (1 - r_k)$$

ce qui donne deux approximations de la valeur moyenne :

$$y'_k = f(x'_k) \text{ et } y''_k = f(x''_k)$$

on prend alors comme nouvelle valeur approchée de la valeur moyenne :

$$y_k = \frac{1}{2} \cdot (y'_k + y''_k)$$

La méthode obtenue se trouve alors être plus efficace que la précédente.

On peut encore augmenter l'efficacité de cette méthode en généralisant ce qui précède de la façon suivante : au tirage n°k, on construit à partir de  $r_k$  dans  $[0,1[$  les réels  $x_{k,j}$  de  $[a,b[$  en posant :

$$x_{k,j} = a + (b - a) \cdot \frac{r_k + j}{m} \quad (j = 0, \dots, m - 1)$$

m étant donné. Ce qui va donner les approximations de  $m(f)$  :

$$y_{k,j} = f(x_{k,j}) \quad (j = 0, \dots, m - 1)$$

on prend alors, à l'étape k, comme nouvelle approximation de  $m(f)$  la moyenne des  $y_{k,j}$ , soit :

$$y_k = \frac{1}{m} \cdot \sum_{j=0}^{m-1} y_{k,j} \quad (k = 1, \dots, n)$$

En combinant cette méthode avec un échantillonnage stratifié, on gagne encore en efficacité.

La méthode Monte Carlo que nous retiendrons, pour calculer des intégrales simples, est finalement la suivante :

(a) on coupe  $[a,b]$  en deux sous-intervalles  $[a,c]$  et  $[c,b]$  avec  $c = \frac{a+b}{2}$  ;

(b) sur chacun de ces deux sous-intervalles, on applique ce qui précède, c'est-à-dire qu'à partir de  $r_k$  dans  $[0,1[$ , on pose  $r_{k,j} = \frac{r_k + j}{m}$  pour  $j = 0, 1, \dots, m - 1$  et

$$y'_{k,j} = f(a + (c - a) \cdot r_{k,j}), \quad y''_{k,j} = f(c + (b - c) \cdot (1 - r_{k,j}))$$

puis :

$$y'_k = \frac{1}{m} \cdot \sum_{j=0}^{m-1} y'_{k,j}, \quad y''_k = \frac{1}{m} \cdot \sum_{j=0}^{m-1} y''_{k,j}$$

$$y_k = \frac{1}{2} \cdot (y'_k + y''_k)$$

enfin, une valeur approchée de l'intégrale est :

$$S = \frac{b-a}{n} \cdot \sum_{k=1}^n y_k$$

D'où la programmation structurée :

FONCTION MonteCarlo1(Entrée  $a, b$  : Réel ;  $m, n$  : Entier ;  $f$  : Fonction) ;

Début

$S = 0$  ;  $L = (b - a)/2$  ;

Pour  $k$  allant de 1 à  $n$  Faire

Début

$r =$  Nombre aléatoire dans  $[0,1]$  ;

$S_k = 0$  ;

Pour  $j$  Allant de 0 à  $m - 1$  Faire

Début

$S_k = S_k + f(a + r \cdot L) + f(c + (1 - r) \cdot L)$  ;

Fin ;

$S_k = S_k / (2 \cdot m)$  ;

$S = S + S_k$  ;

Fin ;

$S = (b - a) \cdot S / n$  ;

Fin ;

#### 4.3.4 Calcul d'intégrales multiples

Les méthodes utilisant les variables antithétiques et l'échantillonnage stratifié sont difficiles à mettre en oeuvre pour les intégrales multiples. On utilisera donc la méthode avec échantillonnage simple.

Tout d'abord, on se ramène à l'intégration d'une fonction sur un pavé  $P = \prod_{i=1}^n [a_i, b_i]$  de  $\mathbb{R}^p$  aussi petit que possible et contenant le domaine d'intégration

$\Omega$  et on prolonge la fonction à intégrer en posant  $f = 0$  sur  $P - \Omega$ .

Au tirage  $n^{\circ} k$ , on extrait au hasard  $p$  valeurs  $x_{i,k}$  de  $[a_i, b_i]$  ( $i = 1, \dots, p$ ), ce qui donne :

$$y_k = f(x_{1,k}, \dots, x_{p,k}) \quad (k = 1, \dots, n)$$

et l'intégrale à calculer est approchée par :

$$\frac{m(P)}{n} \cdot \sum_{k=1}^n y_k$$

où :

$$m(P) = \prod_{i=1}^p (b_i - a_i)$$

#### 4.3.5 Exemple : calcul des coordonnées du centre de gravité d'un corps

On considère le volume  $\Omega$  défini comme l'intersection du tore de centre 0 et de rayons 4 et 2 et du pavé  $[1,4] \times [-3,4] \times [-1,1]$ . Il est donc défini par les relations :

$$\begin{cases} (\sqrt{x^2 + y^2} - 3)^2 + z^2 \leq 1 \\ 2 \leq x \leq 4, \quad -3 \leq y \leq 4, \quad -1 \leq z \leq 1 \end{cases}$$

Les coordonnées du centre de gravité sont alors données par :

$$x_G = \frac{I_x}{I}, \quad y_G = \frac{I_y}{I}, \quad z_G = \frac{I_z}{I}$$

où :

$$I = \int_{\Omega} \rho(x, y, z) dx dy dz; \quad I_x = \int_{\Omega} x \cdot \rho(x, y, z) dx dy dz$$
$$I_y = \int_{\Omega} y \cdot \rho(x, y, z) dx dy dz; \quad I_z = \int_{\Omega} z \cdot \rho(x, y, z) dx dy dz$$

$\rho$  étant la densité volumique définie, pour notre exemple par :

$$\rho(x,y,z) = e^{5 \cdot z}$$

En posant  $e^{5 \cdot z} dz = ds$ , soit  $s = \frac{1}{5} \cdot e^{5 \cdot z}$ , on se ramène à :

$$\rho(x,y,z) dx dy dz = dx dy ds$$

ce qui va faciliter les calculs numériques.

Le nouveau domaine d'intégration est alors :

$$P = [1,4] \times [-3,4] \times [s_0, s_1]$$

où  $s_0 = \frac{1}{5} \cdot e^{-5} = 0.00135$  et  $s_1 = \frac{1}{5} \cdot e^5 = 29.682$ .

Le volume  $V$  de ce domaine étant  $V = 3 \times 7 \times 0.2 \times (e^5 - e^{-5})$ .

La fonction à intégrer valant 1 dans  $\Omega$  et 0 ailleurs.

*Remarque* — Dans la pratique, il faudra toujours essayer par changement de variable ou autre de simplifier l'intégrale multiple à calculer.

## 5. Transformation de Fourier rapide

### 5.1 Position du problème. Notations

Soit  $f : \mathbb{R} \rightarrow \mathbb{C}$  un signal dépendant du temps vérifiant les propriétés suivantes :

$$\left\{ \begin{array}{l} f(t) = 0 \text{ pour } t < 0 \\ f \text{ continue sur } ]0, +\infty[ \\ f \in L^1(\mathbb{R}) \\ \lim_{t \rightarrow +\infty} f(t) = 0 \end{array} \right.$$

En général ce signal n'est pas connu de façon explicite mais seulement en quelques points équidistants  $t_j$ .

Pour  $\varepsilon > 0$  assez petit, on peut trouver  $t_{\max} > 0$  assez grand tel que :

$$|f(t)| < \varepsilon \text{ pour } t > t_{\max}$$

Pour un tel  $t_{\max}$ , on pose alors :

$$\left\{ \begin{array}{l} \delta = \frac{t_{\max}}{n} \\ t_j = j \cdot \delta \\ f_j = f(t_j) \end{array} \right. \quad (j = 0, 1, \dots, n)$$

On définit la *fréquence de Nyquist* par :

$$x_c = \frac{1}{2 \cdot \delta}$$

Comme la transformée de Fourier  $\hat{f}$  est nulle à l'infini, pour  $t_{\max}$  donné on pourra trouver  $n$  assez grand pour que :

$$|\hat{f}(x)| < \varepsilon \text{ pour } |x| > x_c$$

C'est-à-dire que  $f$  est négligeable en dehors de  $[0, t_{\max}]$  et  $\hat{f}$  est négligeable en dehors de  $[-x_c, x_c]$  :

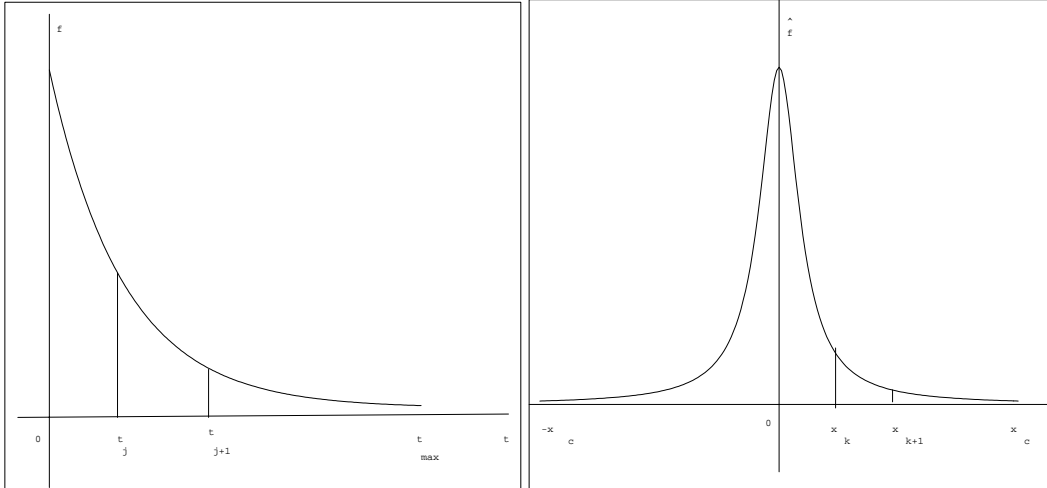


Figure 6.5

Le problème sera donc de calculer de façon approchée les valeurs de la transformée de Fourier de  $f$  en  $n$  points équidistants  $x_k$  dans l'intervalle  $[-x_c, x_c]$ .

On note donc :

$$x_k = k \cdot \frac{2 \cdot x_c}{n} = \frac{k}{n \cdot \delta} \left( k = -\frac{n}{2}, \dots, \frac{n}{2} \right)$$

## 5.2 Approximation des $\hat{f}(x_k)$

Le calcul approché des intégrales de Fourier se fera en utilisant la méthode des rectangles à gauche, soit :

$$\hat{f}(x_k) = \int_{-\infty}^{+\infty} f(t) \cdot e^{-2i\pi x_k t} dt \cong \int_0^{t_{\max}} f(t) \cdot e^{-2i\pi x_k t} dt$$

et :

$$\hat{f}(x_k) \cong \frac{t_{\max}}{n} \cdot \sum_{j=0}^{n-1} f(t_j) \cdot e^{-2i\pi x_k t_j}$$

en tenant compte de :

$$x_k \cdot t_j = \frac{k}{n \cdot \delta} \cdot j \cdot \delta = \frac{k \cdot j}{n}$$

et en posant :

$$\omega = e^{-\frac{2i\pi}{n}} \quad (\text{racine } n^{\text{ème}} \text{ de l'unité})$$

il s'agit de calculer les quantités :

$$\hat{y}_k = \delta \cdot \sum_{j=0}^{n-1} \omega^{jk} \cdot f_j = \delta \cdot y_k \quad \left( k = -\frac{n}{2}, \dots, \frac{n}{2} \right)$$

*Remarque* — Avec la  $2\pi$ -périodicité de la fonction  $t \rightarrow e^{it}$ , on voit que les  $y_k$  sont définis pour tout  $k$  dans  $Z$  et qu'ils forment une suite  $n$ -périodique, c'est-à-dire que  $y_{k+n} = y_k$  ( $k \in Z$ )

Il nous suffit donc, pour déterminer complètement cette suite de calculer les  $y_k$  pour  $k = 0, 1, \dots, n-1$ .

Du point de vue pratique, les valeurs approchées des  $\hat{f}(x_k)$  sont données par :

$$\left\{ \begin{array}{ll} k = 0 & \text{pour } x_k = 0 \\ k = 1, 2, \dots, \frac{n}{2} - 1 & \text{pour } x_k \in ]0, x_c[ \\ k = \frac{n}{2} + 1, \dots, n - 1 & \text{pour } x_k \in ]-x_c, 0[ \\ k = \frac{n}{2} & \text{pour } x_k = x_c \text{ et } x_k = -x_c \end{array} \right.$$

### 5.3 La transformation de Fourier discrète

#### 5.3.1 Définition, propriétés

*Définition* : Soit  $n$  un entier positif non nul et  $\omega = e^{-\frac{2i\pi}{n}}$  une racine  $n^{\text{ème}}$  de l'unité. La *transformation de Fourier discrète* est l'application linéaire :

$$F: \begin{array}{l} C^n \rightarrow C^n \\ x \mapsto y = F(x) \end{array}$$

qui associe au vecteur  $x$  de  $C^n$  de composantes  $x_0, x_1, \dots, x_{n-1}$  le vecteur  $y$  de composantes  $y_0, y_1, \dots, y_{n-1}$  définies par :

$$y_k = \sum_{j=0}^{n-1} \omega^{jk} \cdot x_j \quad (k = 0, 1, \dots, n-1)$$

On peut aussi l'écrire sous forme matricielle :  $y = F \cdot x$  où  $F$  est la matrice à coefficients dans  $C$  :

$$F = \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & \omega & \dots & \omega^{n-1} \\ \cdot & \cdot & \dots & \cdot \\ 1 & \omega^{n-1} & \dots & \omega^{(n-1)^2} \end{pmatrix}$$

Comme pour la transformation intégrale de Fourier, on a une formule d'inversion et un théorème de Plancherel, ce qui est résumé dans le :

*Théorème* : (i) *Formule d'inversion* — La transformation de Fourier discrète est un isomorphisme de  $C^n$ , d'inverse  $F^{-1}$  définie par :

$$(x = F^{-1}(y)) \Leftrightarrow \left( x_j = \frac{1}{n} \cdot \sum_{k=0}^{n-1} \omega^{-jk} \cdot y_k \right) \quad (j = 0, 1, \dots, n-1)$$

(ii) *Formule de Plancherel* — Si  $y = F(x)$ , alors :

$$\sum_{j=0}^{n-1} |x_j|^2 = \frac{1}{n} \cdot \sum_{k=0}^{n-1} |y_k|^2$$

*Démonstration* — Théodor, p. 561.

### 5.3.2 Calcul direct de la transformée de Fourier discrète

Tout d'abord, on écrit une procédure de calcul des  $\omega^r$ , pour  $r = 0, 1, \dots, n - 1$ , les valeurs de  $\omega^r$  étant stockées dans un vecteur  $W$  de  $\mathbb{C}^n$ .

Ce qui donne, en disposant d'opérations sur les nombres complexes :

*PROCEDURE PuissancesOmega(Entrée  $n$  : Entier ; Sortie  $W$  : VecteurComplexe) ;*

*Début*

*$W_0 = 1$  ;  $W_1 = \text{Cos}(2\pi/n) - i \cdot \text{Sin}(2\pi/n)$  ;*

*Pour  $r$  allant de 2 à  $n - 1$  Faire  $W_r = W_{r-1} \cdot W_1$  ;*

*Fin ;*

Ce calcul nécessitant  $n - 2$  multiplications complexes.

En remarquant que pour  $0 \leq j, k \leq n - 1$ , on a :  $\omega^{k \cdot j} = \omega^r$ , où  $r$  est le reste dans la division euclidienne de  $j \cdot k$  par  $n$ , on en déduit la procédure de calcul du vecteur  $y$  :

*PROCEDURE FourierDiscret\_1(Entrée  $n$  : Entier ;  $W, x$  : VecteurComplexe ;*

*Sortie  $y$  : VecteurComplexe) ;*

*Début*

*Pour  $k$  Allant de 0 à  $n - 1$  Faire  $y_k = x_0$  ;*

*Pour  $j$  Allant de 1 à  $n - 1$  Faire  $y_0 = y_0 + x_j$  ;*

*Pour  $k$  Allant de 1 à  $n - 1$  Faire*

*Début*

*Pour  $j$  Allant de 1 à  $n - 1$  Faire*

*Début*

*$r = j \cdot k \text{ Modulo } n$  ;*

*$y_k = y_k + \omega_r \cdot x_j$  ;*

*Fin ;*

*Fin ;*

*Fin ;*

Cette procédure de calcul très simple à le gros inconvénient d'être lente. En effet, le calcul de  $y_0$  nécessite  $n - 1$  additions complexes et pour  $k = 1, \dots, n - 1$ , le calcul de chaque  $y_k$  nécessite  $n - 1$  multiplications complexes et  $n - 1$  additions complexes, ce qui donne en ajoutant le nombre d'opérations nécessaires au calcul de  $W$  un total de  $(n - 2) + (n - 1) + 2 \cdot (n - 1)^2$  opérations complexes. Soit un nombre d'opérations de l'ordre de  $n^2$ , ce qui sera peu performant pour  $n$  grand (par exemple  $n = 10^6$  demandera plusieurs jours de calcul).

Un algorithme très performant est celui de Cooley et Tukey.



## 5.4 L'algorithme FFT de Cooley et Tukey

### 5.4.1 Introduction

Cet algorithme, aussi connu sous le nom de *FFT (Fast Fourier Transform)*, a été mis au point dans les années 1960 par Cooley et Tukey sur une idée de Danielson et Lanczos en 1942. Il demandera, pour le calcul d'une transformée de Fourier à  $n$  points, un nombre d'opérations de l'ordre de  $n \cdot \ln(n)$ , ce qui est beaucoup plus performant que l'algorithme direct, le prix à payer étant une programmation plus délicate.

L'idée de base est d'écrire une transformée à  $n$  points comme somme de deux transformées à  $n/2$  points et ainsi de suite. Pour ce faire on supposera que  $n$  est une puissance de 2.

On notera donc pour tout ce qui suit :

$$n = 2^p \text{ (} p \geq 1 \text{) et } \omega_p = e^{-\frac{2i\pi}{n}}$$

### 5.4.2 Cas particulier $n = 2$

Dans ce cas, on a  $p = 1$  et  $\omega_1 = -1$ , avec les formules :

$$\begin{cases} y_0 = x_0 + x_1 \\ y_1 = x_0 - x_1 \end{cases}$$

ce qui s'écrit globalement :

$$(1) y_{k_0} = \sum_{j_0=0}^1 (-1)^{k_0 \cdot j_0} \cdot x_{j_0} \quad (k_0 = 0, 1)$$

### 5.4.3 Cas particulier $n = 4$

Dans ce cas, on a  $p = 2$  et  $\omega_2 = -i$ , avec les formules :

$$\begin{cases} y_0 = x_0 + x_1 + x_2 + x_3 \\ y_1 = x_0 + \omega_2 \cdot x_1 + \omega_2^2 \cdot x_2 + \omega_2^3 \cdot x_3 \\ y_2 = x_0 + \omega_2^2 \cdot x_1 + \omega_2^4 \cdot x_2 + \omega_2^6 \cdot x_3 \\ y_3 = x_0 + \omega_2^3 \cdot x_1 + \omega_2^6 \cdot x_2 + \omega_2^9 \cdot x_3 \end{cases}$$

en séparant dans chaque somme les termes d'indices pairs et impairs et en tenant compte de la périodicité, on a :

$$\begin{cases} y_0 = (x_0 + x_2) + (x_1 + x_3) \\ y_1 = (x_0 - x_2) + \omega_2 \cdot (x_1 - x_3) \\ y_2 = (x_0 + x_2) + \omega_2^2 \cdot (x_1 + x_3) \\ y_3 = (x_0 - x_2) + \omega_2^3 \cdot (x_1 - x_3) \end{cases}$$

En décomposant tout nombre entier  $r \in \{0, 1, \dots, n-1\}$  en base deux sous la forme :

$$r = r_0 + r_1 \cdot 2 + \dots + r_{p-2} \cdot 2^{p-2} + r_{p-1} \cdot 2^{p-1} \quad (r_j = 0 \text{ ou } 1)$$

et en notant, pour  $j = 0, 1, \dots, p-1$  :

$$X_j = X_{j_0, j_1, \dots, j_{p-1}}$$

les formules vont s'écrire :

$$\begin{cases} y_{0,0} = (x_{0,0} + x_{0,1}) + (x_{1,0} + x_{1,1}) \\ y_{1,0} = (x_{0,0} - x_{0,1}) + \omega_2 \cdot (x_{1,0} - x_{1,1}) \\ y_{0,1} = (x_{0,0} + x_{0,1}) + \omega_2^2 \cdot (x_{1,0} + x_{1,1}) \\ y_{1,1} = (x_{0,0} - x_{0,1}) + \omega_2^3 \cdot (x_{1,0} - x_{1,1}) \end{cases}$$

ce qui peut s'écrire globalement :

$$y_k = y_{k_0, k_1} = \sum_{j_0=0}^1 \omega^{2 \cdot k_1 \cdot j_0} \cdot \sum_{j_1=0}^1 \omega^{k_0 \cdot (j_0 + 2 \cdot j_1)} \cdot x_{j_0, j_1} \quad (k_0, k_1 = 0, 1)$$

#### 5.4.4 Cas général $n = 2^p$

Pour tout  $k = 0, 1, \dots, n-1$ ,  $y_k$  est une somme du type :

$$S = \sum_{j=0}^{n-1} S_j = \sum_{j_0, \dots, j_{p-1}} S_{j_0, \dots, j_{p-1}}$$

qui peut s'écrire sous la forme :

$$S = \sum_{j_0=0}^1 \sum_{j_1=0}^1 \dots \sum_{j_{p-1}=0}^1 S_{j_0, \dots, j_{p-1}}$$

ce qui donne :

$$y_k = \sum_{j_0=0}^1 \sum_{j_1=0}^1 \dots \sum_{j_{p-1}=0}^1 \omega_p^{k \cdot j} \cdot x_{j_0, \dots, j_{p-1}}$$

En remarquant que pour  $j$ ,  $k = 0, 1, \dots, n-1$ , on a :

$$j \cdot k \equiv k_0 \cdot (j_0 + \dots + j_{p-1} \cdot 2^{p-1}) + 2 \cdot k_1 \cdot (j_0 + \dots + j_{p-2} \cdot 2^{p-2}) + \dots + 2^{p-1} \cdot k_{p-1} \cdot j_0 \quad (\text{Mod. } n)$$

on déduit que :

$$\omega_p^{k \cdot j} = \omega_p^{2^{p-1} \cdot k_{p-1} \cdot j_0} \cdot \omega_p^{2^{p-2} \cdot k_{p-2} \cdot (j_0 + 2 \cdot j_1)} \cdot \dots \cdot \omega_p^{k_0 \cdot (j_0 + \dots + 2^{p-1} \cdot j_{p-1})}$$

et donc que la transformation de Fourier discrète peut s'écrire :

$$(p) \quad y_k = \sum_{j_0=0}^1 \omega_p^{2^{p-1} \cdot k_{p-1} \cdot j_0} \cdot \sum_{j_1=0}^1 \omega_p^{2^{p-2} \cdot k_{p-2} \cdot (j_0 + 2 \cdot j_1)} \cdot \dots \cdot \sum_{j_{p-1}=0}^1 \omega_p^{k_0 \cdot (j_0 + \dots + 2^{p-1} \cdot j_{p-1})} \cdot x_{j_0, \dots, j_{p-1}}$$

#### 5.4.5 Utilisation de (p) pour calculer la transformation de Fourier discrète

On décrit tout d'abord l'algorithme pour  $n = 4$ .

*Etape 0* — On pose  $x^{(0)} = x$ .

Les composantes de  $x^{(0)}$  sont donc les  $x_{j_0, j_1}$  pour  $j_0, j_1 = 0, 1$ .

*Etape 1* — On construit  $x^{(1)}$  à partir de  $x^{(0)}$  en s'inspirant de la récurrence (p), soit en notant  $x_{j_0, k_0}^{(1)}$  les coordonnées de  $x^{(1)}$  :

$$x_{j_0, k_0}^{(1)} = \sum_{j_1=0}^1 \omega_2^{k_0 \cdot (j_0 + 2 \cdot j_1)} \cdot x_{j_0, j_1}^{(0)}$$

*Etape 2* — On construit de la même façon  $x^{(2)}$  en fonction de  $x^{(1)}$  en notant  $x_{k_1, k_0}^{(2)}$  ses composantes, soit :

$$x_{k_0, k_0}^{(2)} = \sum_{j_0=0}^1 \omega_2^{2 \cdot k_1 \cdot j_0} \cdot x_{j_0, k_0}^{(1)}$$

Le vecteur  $x^{(2)}$  donne donc les composantes de  $y$ , mais elles sont dans l'ordre binaire inverse :

$$x_{k_1, k_0}^{(2)} = y_{k_0, k_1}$$

Dans le cas général, on procède comme suit :

*Etape 0* — On pose  $x^{(0)} = x$ .

*Etape  $r + 1$*  ( $r = 0, 1, \dots, p - 1$ ) — L'étape précédente a donné le vecteur  $x^{(r)}$  de composantes  $x_{j_0, \dots, j_{p-r-2}, k_{r-1}, k_0}^{(r)}$  et on veut calculer celles de  $x^{(r+1)}$  notées  $x_{j_0, \dots, j_{p-r-2}, k_{r-1}, k_0}^{(r+1)}$

Avec la récurrence (p), on déduit alors que :

$$x_{j_0, \dots, j_{p-r-2}, k_{r-1}, k_0}^{(r+1)} = \sum_{j_{p-r-1}=0}^1 \omega_p^{2^r \cdot k_r \cdot (j_0 + \dots + 2^{p-r-1} \cdot j_{p-r-1})} \cdot x_{j_0, \dots, j_{p-r-1}, k_{r-1}, \dots, k_0}^{(r)}$$

On pose alors :

$$m = \frac{n}{2^{r+1}} = 2^{p-r-1}$$

$$L = j_0 + \dots + j_{p-r-2} \cdot 2^{p-r-2} \in \{0, \dots, m-1\}$$

$$h = k_{r-1} + \dots + k_0 \cdot 2^{r-1} \in \{0, \dots, 2^r - 1 = \frac{n}{2 \cdot m} - 1\}$$

$$i = j_0 + \dots + j_{p-r-2} \cdot 2^{p-r-2} + 0 \cdot 2^{p-r-1} + k_{r-1} \cdot 2^{p-r} + \dots + k_0 \cdot 2^{p-1} = L + 2 \cdot m \cdot h$$

$$j = j_0 + \dots + j_{p-r-2} \cdot 2^{p-r-2} + 1 \cdot 2^{p-r-1} + k_{r-1} \cdot 2^{p-r} + \dots + k_0 \cdot 2^{p-1} = i + m$$

et les formules précédentes vont s'écrire :

$$\begin{cases} x_i^{(r+1)} = x_i^{(r)} + x_j^{(r)} & (k_r = 0) \\ x_j^{(r+1)} = \omega_p^{2^r \cdot L} \cdot x_i^{(r)} + \omega_p^{2^r \cdot (L + 2^{p-r-1})} \cdot x_j^{(r)} & (k_r = 1) \end{cases}$$

Avec :

$$\omega_p^{2^{r+p-r-1}} = -1 \text{ et } \omega_p^{2^r} = \omega_{p-r}$$

on déduit que :

$$\begin{cases} x_i^{(r+1)} = x_i^{(r)} + x_j^{(r)} \\ x_j^{(r+1)} = (\omega_{p-r})^L \cdot (x_i^{(r)} - x_j^{(r)}) \end{cases}$$

Et au bout de  $p$  étapes, on a :

$$x_{k_{p-1}, \dots, k_0}^{(p)} = y_{\sigma(k)}$$

où  $\sigma$  est l'inversion binaire définie par :

$$\sigma(k_0 + \dots + k_{p-1} \cdot 2^{p-1}) = k_{p-1} + \dots + k_0 \cdot 2^{p-1}$$

Pour retrouver les composantes de  $y$  dans le bon ordre on devra donc faire une inversion binaire sur celles de  $x^{(p)}$ .

#### 5.4.6 Nombre d'opérations élémentaires dans l'algorithme de Cooley et Tukey

En notant, pour  $p \geq 1$  :

$$\begin{cases} A_p = \text{Nombre d'additions complexes;} \\ M_p = \text{Nombre de multiplications complexes;} \end{cases}$$

on a :

$$\begin{cases} A_1 = 2 \\ M_1 = 0 \end{cases}$$

et la récurrence :

$$\begin{cases} A_p = 2 \cdot A_{p-1} + 2^p \\ M_p = 2 \cdot M_{p-1} + 2^{p-1} \end{cases}$$

d'où :

$$\begin{cases} M_p = A_{p-1} \\ A_p = p \cdot 2^p \end{cases}$$

Ce qui donne un nombre d'opérations élémentaires de l'ordre de  $n \cdot \text{Log}_2(n)$ . (Cf. Théodor, p. 563)

### 5.4.7 Programmation structurée

*PROCEDURE FFT(Entrée p : Entier ; x : VecteurComplexe ; Sortie y : VecteurComplexe) ;*

*Début*

$$n = 2^p ;$$

$$m = \frac{n}{2} ; \{ \text{Etape } r + 1 = 1 \}$$

*Pour r allant de 0 à p - 1 Faire*

*Début*

$$w1 = 1 ; \{ w1 = \omega_{p-r}^L, \text{ avec } L = 0 \}$$

$$w2 = \text{Cos}(\pi/m) - i \cdot \text{Sin}(\pi/m) ; \{ w2 = \omega_{p-r} \}$$

*Pour L allant de 0 à m - 1 Faire*

*Début*

$$\text{Pour } h \text{ Allant de } 0 \text{ à } \frac{n}{2 \cdot m} - 1 \text{ Faire}$$

*Début*

$$i = L + 2 \cdot m \cdot h ;$$

$$j = i + m ;$$

$$\text{Aux} = (x_i - x_j) \cdot w1 ;$$

$$x_i = x_i + x_j ;$$

$$x_j = \text{Aux} ;$$

*Fin ;*

$$w1 = w1 \cdot w2 ; \{ w1 = \omega_{p-r}^L \}$$

*Fin ;*

$$m = \frac{m}{2} ;$$

*Fin ;*

*InversionBinaire(p,x,y) ;*

*Fin ;*

Pour retrouver les  $\hat{y}_k$ , il ne faudra pas oublier de multiplier les  $y_k$  par le coefficient  $\delta$ .

En notant :

$$k = k_0 + k_1 \cdot 2 + \dots + k_{p-1} \cdot 2^{p-1}$$

le calcul de l'inverse binaire  $h = \sigma(k)$ , se fait en utilisant les remarques suivantes :

$$k_r = (k \text{ Divise } 2^r) \text{ Modulo } 2 ; h_{p-r-1} = k_r \quad (r = 0, 1, \dots, p-1)$$

Ce qui donne la procédure :

*PROCEDURE InversionBinaire(Entrée  $p$  : Entier ;  $x$  : VecteurComplexe ;*

*Sortie  $y$  : VecteurComplexe) ;*

*Début*

*$n = 2^p ;$*

*Pour  $k$  Allant de 0 à  $n-1$  Faire*

*Début*

*$h = 0 ;$*

*$Aux = n/2 ; \{Aux = 2^{p-1}\}$*

*Pour  $r$  Allant de 0 à  $p-1$  Faire*

*Début*

*$c = (k \text{ Divise } 2^r) \text{ Modulo } 2 ; \{c = k_r\}$*

*$h = h + c \cdot Aux ; \{c \cdot Aux = k_r \cdot 2^{p-r-1}\}$*

*$Aux = Aux/2 ; \{Aux = 2^{p-1-(r+1)}\}$*

*Fin ;*

*$y_h = x_k ;$*

*Fin ;*

*Fin ;*

### 5.5 Application au calcul des coefficients de Fourier d'une fonction périodique

Soit  $f : \mathbb{R} \rightarrow \mathbb{C}$ , continue par morceaux et  $t_{\max}$ -périodique. Ces coefficients de Fourier complexes sont définis par :

$$c_k = \frac{1}{t_{\max}} \cdot \int_0^{t_{\max}} f(t) \cdot e^{-\frac{2ik\pi}{t_{\max}} t} dt = \frac{1}{t_{\max}} \cdot \int_0^{t_{\max}} f(t) \cdot e^{-2i\pi x_k \cdot t} dt$$

Soit en notant encore  $f$  la fonction coïncidant avec la précédente sur  $[0, t_{\max}]$  et nulle ailleurs :

$$c_k = \frac{1}{t_{\max}} \cdot \hat{f}(x_k) \quad (k = -\frac{n}{2}, \dots, 0, \dots, \frac{n}{2})$$

C'est-à-dire que si  $y = \left( y_{-\frac{n}{2}}, \dots, y_0, \dots, y_{\frac{n}{2}} \right)$  est la transformée de Fourier

discrète de  $f = (f_0, \dots, f_{n-1})$ , alors :

$$c_k = \frac{\delta}{t_{\max}} \cdot y_k = \frac{1}{n} \cdot y_k \quad (k = -\frac{n}{2}, \dots, 0, \dots, \frac{n}{2})$$

Soit avec la  $n$ -périodicité de la suite  $(y_k)$  :

$$\begin{cases} c_k = \frac{1}{n} \cdot y_k & \left( k = 0, \dots, \frac{n}{2} - 1 \right) \\ c_k = \frac{1}{n} \cdot y_{k+n} & \left( k = -\frac{n}{2}, \dots, -1 \right) \end{cases}$$

On retrouve les coefficients de Fourier trigonométriques avec :

$$a_0 = c_0 = \operatorname{Re}(c_0)$$

$$a_k = c_k + c_{-k} = c_k + \bar{c}_k = 2 \cdot \operatorname{Re}(c_k) \quad \left( k = 1, \dots, \frac{n}{2} - 1 \right)$$

$$b_k = i \cdot (c_k - c_{-k}) = i \cdot (c_k - \bar{c}_k) = -2 \cdot \operatorname{Im}(c_k) \quad \left( k = 1, \dots, \frac{n}{2} - 1 \right)$$

En cas de convergence de la série de Fourier, on a alors :

$$f(t) \cong a_0 + \sum_{k=1}^{\frac{n}{2}-1} \left\{ a_k \cdot \operatorname{Cos} \left( k \cdot \frac{2 \cdot \pi}{t_{\max}} \cdot t \right) + b_k \cdot \operatorname{Sin} \left( k \cdot \frac{2 \cdot \pi}{t_{\max}} \cdot t \right) \right\}$$

ou encore :

$$f(t) \cong a_0 + \sum_{k=1}^{\frac{n}{2}-1} \left\{ a_k \cdot \operatorname{Cos}(2 \cdot \pi \cdot x_k \cdot t) + b_k \cdot \operatorname{Sin}(2 \cdot \pi \cdot x_k \cdot t) \right\}$$

En posant :

$$\begin{aligned} A_0 &= a_0 \\ A_k &= 2 \cdot |c_k| = |a_k + i \cdot b_k| \\ \varphi_k &= \operatorname{Arg}(a_k + i \cdot b_k) = -\operatorname{Arctg} \left( \frac{\operatorname{Im}(c_k)}{\operatorname{Re}(c_k)} \right) \end{aligned}$$

on a aussi :

$$f(t) \cong A_0 + \sum_{k=1}^{\frac{n}{2}-1} A_k \cdot \operatorname{Cos}(2 \cdot \pi \cdot x_k \cdot t - \varphi_k)$$

## 6. Exercices

### 6.1 Aire intérieure à une courbe fermée

On considère une courbe définie par une équation polaire  $\rho = f(\theta)$ , avec  $\theta \in [0, 2 \cdot \pi[$ .

On veut calculer  $S = \operatorname{Aire} \{ (\rho, \theta) \in \operatorname{Rx}[0, 2 \cdot \pi[ ; \rho < f(\theta) \}$

1°) En s'inspirant de la méthode des trapèzes, écrire un algorithme de calcul de S.

2°) Application : aire de la cardiode  $\rho = 1 + \cos(\theta)$ .

### 6.2 Séries de Fourier

1°) Ecrire une procédure Fourier, qui permet de calculer les coefficients du développement en série de Fourier d'une fonction périodique f (Les intégrales donnant les coefficients de Fourier pourront être calculées successivement par les

méthodes des trapèzes, de Simpson et de Romberg, puis on comparera les temps de calcul).

2°) Ecrire un programme qui permet de comparer la fonction  $T$ -périodique  $f$  avec son développement en série de Fourier. On pourra tracer, sur  $[-T, T]$ , les courbes  $y = f(x)$  et

$$y = a_0 + \sum_{i=1}^n \left( a_i \cdot \cos\left(i \cdot \frac{2 \cdot \pi}{T} \cdot x\right) + b_i \cdot \sin\left(i \cdot \frac{2 \cdot \pi}{T} \cdot x\right) \right)$$

les  $a_i$  et  $b_i$  étant les coefficients de Fourier de  $f$ .

3°) *Application* — Appliquer ce qui précède à la fonction  $f$ ,  $2 \cdot \pi$ -périodique définie par :

$$f_p(x) = \frac{\sin(p \cdot x)}{x}, \text{ pour } x \in [-\pi, \pi]$$

### 6.3 Circuit R-C

On considère le circuit R-C suivant, soumis à une entrée  $e(t)$  :

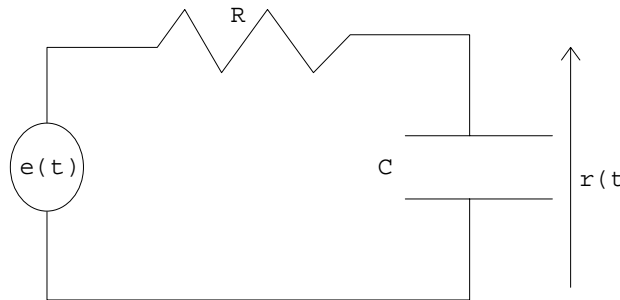


Figure 6.6

On pose  $R \cdot C = \tau$ .

On a alors :  $r = h * e$ , avec  $h(t) = 0$  pour  $t < 0$  et  $h(t) = \frac{1}{\tau} \cdot e^{-\frac{t}{\tau}}$ , pour  $t > 0$ , le signe

\* désignant le *produit de convolution* défini dans le cas où  $f(x) = g(x) = 0$  pour  $x < 0$  par :

$$f * g(x) = \int_0^x f(t) \cdot g(x-t) dt$$

1°) Ecrire une procédure de calcul de  $r(t)$ , pour tout  $t$  dans  $\mathbb{R}$ .

*Applications*

$$(a) e(t) = \begin{cases} 0 & \text{si } t < 0 \\ 1 & \text{si } 0 \leq t \leq 1. \\ 0 & \text{si } t > 1 \end{cases}$$

(b)  $e(t) = \sin(t)$ .

(c)  $e$  est  $T$ -périodique, affine par morceaux sur  $[0, T]$ , avec  $e(0) = e(T) = 0$  et  $e\left(\frac{T}{2}\right) = 1$ .

2°) Ecrire une procédure de tracé des graphes de  $e$  et de  $r$ .

3°) Ecrire l'équation différentielle vérifiée par  $r(t)$  sous la forme :  $r'(t) = f(t, r(t))$ .

Résoudre cette équation de façon approchée en utilisant la méthode de Runge-Kutta d'ordre 4 (Cf. le chapitre : « Résolution numérique des équations différentielles »).

Tracer les graphes de  $e$  et  $r$ .

Comparer à la méthode précédente.

### 6.4 Intégrales elliptiques

Une intégrale elliptique de deuxième espèce est définie par :

$$e_2(y, k_c, a, b) = \int_0^{+\infty} \frac{a + b \cdot x^2}{(1 + x^2) \cdot \sqrt{(1 + x^2) \cdot (1 + k_c^2 \cdot x^2)}} dx$$

où  $y > 0$ ,  $a, b, c, k_c \in \mathbb{R}$ .

1°) Ecrire  $e_2$  en utilisant le changement de variable  $x = \operatorname{tg}(\phi)$  et en notant  $k^2 = 1 - k_c^2$ .

2°) L'intégrale elliptique complète est obtenue en faisant tendre  $y$  vers  $+\infty$  (ou  $\phi$  vers  $\frac{\pi}{2}$ ). De manière plus générale, l'intégrale elliptique complète est définie par :

$$e(k_c, p, a, b) = \int_0^{+\infty} \frac{a + b \cdot x^2}{(1 + p \cdot x^2) \cdot \sqrt{(1 + x^2) \cdot (1 + k_c^2 \cdot x^2)}} dx$$

Ecrire  $e$  en utilisant le changement de variable  $x = \operatorname{tg}(\varpi)$  et en notant  $k^2 = 1 - k_c^2$ .

3°) Ecrire un programme permettant de calculer chacune de ces intégrales.

### 6.5 Calcul d'intégrales doubles par la méthode des éléments finies

Soit  $\Omega$  un domaine polygonal à cotés parallèles aux axes dans  $\mathbb{R}^2$  :

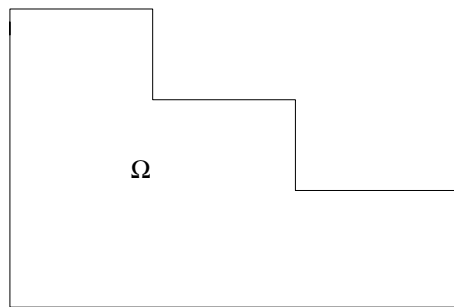


Figure 6.7

et  $f : \Omega \rightarrow \mathbb{R}$ , une fonction continue.

Le but est de calculer, de façon approchée :  $I(f) = \iint_{\Omega} f(x, y) dx dy$ .

I. On suppose dans cette partie que  $\Omega$  est le carré unité, noté  $\Omega_1$  :

$$\Omega_1 = \{ (x, y) \in \mathbb{R}^2 ; 0 \leq x \leq 1, 0 \leq y \leq 1 \}$$

Ce carré est appelé élément de référence.

On note  $A_1 = (0, 0)$ ,  $A_2 = (1, 0)$ ,  $A_3 = (0, 1)$ ,  $A_4 = (1, 1)$ .

$E_1$  désigne l'ensemble des fonctions polynomiales du type :



$$g(x,y) = a + b \cdot x + c \cdot y + d \cdot x \cdot y$$

$E_1$  est appelé l'ensemble des fonctions de base.

1°) Déterminer des fonctions  $g_1, g_2, g_3$  et  $g_4$  dans  $E_1$  telles que :

$$g_i(A_j) = \begin{cases} 1 & \text{si } j = i \\ 0 & \text{si } j \neq i \end{cases}$$

2°) Calculer  $\iint_{\Omega} g_j(x,y) dx dy$ , pour  $j = 1, 2, 3, 4$ .

3°) Déterminer  $g$  dans  $E_1$  qui interpole  $f$  aux sommets de  $\Omega_1$ , c'est-à-dire telle que :

$$g(A_j) = f(A_j), \text{ pour } j = 1, 2, 3, 4$$

4°) En déduire une approximation  $I_1(f)$  de  $I(f)$  sur  $\Omega_1$ .

II. On suppose dans cette partie que  $\Omega$  est un rectangle à cotés parallèles aux axes, noté  $\Omega_2$  :

$$\Omega_2 = \{ (x,y) \in \mathbb{R}^2 ; a \leq x \leq b, c \leq y \leq d \}$$

1°) Transformer  $I(f)$  en une intégrale sur  $\Omega_1$ .

2°) En déduire une formule d'approximation de  $I(f)$ , faisant intervenir les valeurs de  $f$  aux sommets de  $\Omega_2$ .

III. Dans cette partie  $\Omega$  est un domaine polygonal à cotés parallèles aux axes.

On se donne un maillage de  $\Omega$  par des rectangles à cotés parallèles aux axes. Si  $n$  est le nombre de noeuds de ce maillage, on note  $S_i$  chacun des sommets (on pourra les numéroter de haut en bas et de gauche à droite).

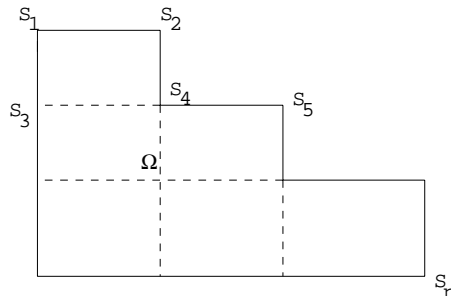


Figure 6.8

Donner une formule d'approximation de  $I(f)$  du type :

$$I(f) \cong \sum_{i=1}^n \mu_i \cdot f(S_i)$$

où les  $\mu_i$  sont à préciser.

IV. Reprendre les questions I, II et III avec pour élément de référence le carré unité à 9 noeuds : les 4 sommets, les 4 milieux de chaque coté et le centre du carré, et pour espace de fonctions de base l'espace  $E$  des fonctions polynomiales engendré par  $1, x, y, x \cdot y, x^2, x^2 \cdot y, x^2 \cdot y^2, y^2$  et  $x \cdot y^2$ .

## 7. Programmation Ada

### 7.1 Spécification du paquetage *INTEGRATION\_GENERIQUE*

```

generic
  with function f(x : in FLOAT) return FLOAT ;
package INTEGRATION_GENERIQUE is
function INTEGRATION_TRAPEZES(a, b : in FLOAT) return FLOAT; -- § 2.4.1
function INTEGRATION_SIMPSON(a, b : in FLOAT) return FLOAT ; -- § 2.4.2
function INTEGRATION_ROMBERG(a,b : in FLOAT) return FLOAT ; -- § 2.4.3
function INTEGRATION_LEGENDRE(a, b : in FLOAT) return FLOAT; -- § 3.6.1
function INTEGRATION_LAGUERRE return FLOAT ; -- § 3.6.3
function INTEGRATION_HERMITE return FLOAT ; -- § 3.6.4
function INTEGRATION_TCHEBYCHEV return FLOAT ; -- § 3.6.2
end INTEGRATION_GENERIQUE ;

```

### 7.2 Démonstration du paquetage *INTEGRATION\_GENERIQUE*

```

with TEXT_IO, FIO, CRT, COMMON_MATH0, MATH0, MATH2,
  INTEGRATION_GENERIQUE ;
use TEXT_IO, FIO, CRT, COMMON_MATH0, MATH0, MATH2 ;

procedure DEM_INT is
CHOIX : NATURAL range 0..7 ;
ff : FONCTION ;

function f(x : in FLOAT) return FLOAT is
begin
  return EVALUE(ff,x) ;
end ;

package INTEGRATION_F is new INTEGRATION_GENERIQUE(f) ;
use INTEGRATION_F ;

procedure DONNEES_INTEGRATION(a, b : in out FLOAT) is
begin
  PUT_LINE("          Calculs d'integrales definies") ;
  NEW_LINE ;
  GET_LINE(ff,"Donner la fonction a integrer, f(x) = ") ;
  ENTRER_REEL(a,"Donner la valeur minimale de l'abscisse : ") ;
  ENTRER_REEL_BORNE(a,1.0E+12,b,
    "Donner la valeur maximale de l'abscisse : ") ;
end DONNEES_INTEGRATION ;

procedure DEMO_TRAPEZES is
a, b, T : FLOAT ;
begin
  MODE_AFFICHAGE ;
  DONNEES_INTEGRATION(a,b) ;
  T := INTEGRATION_TRAPEZES(a,b,) ;
  PUT(IMP,"  Integrale de f(x) = ") ;
  PUT(TEXTE_DE_FONCTION(ff)) ;

```

```

    PUT(IMP," sur [" ) ;
    PUT(IMP,a,4,2,0) ;
    PUT(IMP," ," ) ;
    PUT(IMP,b,4,2,0) ;
    PUT(IMP," ] = " ) ;
    PUT(IMP,T,6,4,0) ; NEW_LINE(IMP) ;
    CLOSE(IMP) ;
    PAUSE ;
end DEMO_TRAPEZES ;

procedure DEMO_SIMPSON is
a, b, S : FLOAT ;
begin
    MODE_AFFICHAGE ;
    DONNEES_INTEGRATION(a,b) ;
    S := INTEGRATION_SIMPSON(a,b) ;
    PUT(IMP," Integrale de f(x) = " ) ;
    PUT(TEXTE_DE_FONCTION(ff)) ;
    PUT(IMP," sur [" ) ;
    PUT(IMP,a,4,2,0) ;
    PUT(IMP," ," ) ;
    PUT(IMP,b,4,2,0) ;
    PUT(IMP," ] = " ) ;
    PUT(IMP,S,6,4,0) ;
    NEW_LINE(IMP) ;
    CLOSE(IMP) ;
    PAUSE ;
end DEMO_SIMPSON ;

procedure DEMO_ROMBERG is
a, b, R : FLOAT ;
begin
    MODE_AFFICHAGE ;
    DONNEES_INTEGRATION(a,b) ;
    R := INTEGRATION_ROMBERG(a,b) ;
    PUT(IMP," Integrale de f(x) = " ) ;
    PUT(TEXTE_DE_FONCTION(ff)) ;
    PUT(IMP," sur [" ) ;
    PUT(IMP,a,4,2,0) ;
    PUT(IMP," ," ) ;
    PUT(IMP,b,4,2,0) ;
    PUT(IMP," ] = " ) ;
    PUT(IMP,R,6,4,0) ;
    NEW_LINE(IMP) ;
    CLOSE(IMP) ;
    PAUSE ;
end DEMO_ROMBERG ;

procedure DEMO_LEGENDRE is
a, b, R : FLOAT ;
begin
    MODE_AFFICHAGE ;
    DONNEES_INTEGRATION(a,b) ;

```

```

R := INTEGRATION_LEGENDRE(a,b) ;
PUT(IMP," Integrale de f(x) = ") ;
PUT(TEXTE_DE_FONCTION(ff)) ;
PUT(IMP," sur [") ;
PUT(IMP,a,4,2,0) ;
PUT(IMP,",") ;
PUT(IMP,b,4,2,0) ;
PUT(IMP,"] = ") ;
PUT(IMP,R,6,4,0) ;
NEW_LINE(IMP) ;
CLOSE(IMP) ;
PAUSE ;
end DEMO_LEGENDRE ;

procedure DEMO_LAGUERRE is
S : FLOAT ;
begin
MODE_AFFICHAGE ;
PUT_LINE(" CALCUL DE L'INTEGRALE IMPROPRE DE EXP(-x).f(x) SUR R+") ;
NEW_LINE ;
PUT_LINE(" Utilisation du polynome de Laguerre Numero 12") ;
NEW_LINE ;
GET_LINE(ff," Entrer la fonction a integrer f(x) = ") ;
NEW_LINE ;
S := INTEGRATION_LAGUERRE ;
PUT(IMP," Valeur approchee de l'integrale, I = ") ;
PUT(IMP,S,6,4,0) ;
NEW_LINE(IMP) ;
CLOSE(IMP) ;
PAUSE ;
end DEMO_LAGUERRE ;

procedure DEMO_HERMITE is
S : FLOAT ;
begin
MODE_AFFICHAGE ;
PUT_LINE(" CALCUL DE L'INTEGRALE IMPROPRE DE EXP(-x**2).f(x) sur R");
NEW_LINE ;
PUT_LINE(" Utilisation du polynome d'Hermite Numero 10") ;
NEW_LINE ;
GET_LINE(ff," Entrer la fonction a integrer f(x) = ") ;
NEW_LINE ;
S := INTEGRATION_HERMITE ;
PUT(IMP," Valeur approchee de l'integrale, I = ") ;
PUT(IMP,S,6,4,0) ;
NEW_LINE(IMP) ;
CLOSE(IMP) ;
PAUSE ;
end DEMO_HERMITE ;

procedure DEMO_TCHEBYCHEV is
S : FLOAT ;
begin

```

```

MODE_AFFICHAGE ;
PUT_LINE(
" CALCUL DE L'INTEGRALE IMPROPRE de f(x)/SQRT(1 - x*x) sur [-1,1]" ) ;
NEW_LINE ;
PUT_LINE(" Utilisation du polynome de Tchebychev Numero 30" ) ;
NEW_LINE ;
GET_LINE(ff," Entrer la fonction a integrer f(x) = " ) ;
NEW_LINE ;
S := INTEGRATION_TCHEBYCHEV ;
PUT(IMP," Valeur approchee de l'integrale, I = " ) ;
PUT(IMP,S,6,4,0) ;
NEW_LINE(IMP) ;
CLOSE(IMP) ;
PAUSE ;
end DEMO_TCHEBYCHEV ;

procedure MENU_INTEGRATION(CHOIX : out NATURAL) is
begin
  CLRSCR ;
  PUT_LINE("Calcul d'integrales" ) ;
  PUT_LINE("-0- FIN" ) ;
  PUT_LINE("-1- Methode des trapezes" ) ;
  PUT_LINE("-2- Methode de Simpson" ) ;
  PUT_LINE("-3- Methode de Romberg" ) ;
  PUT_LINE("-4- Methode de Legendre" ) ;
  PUT_LINE("-5- Methode de Tchebychev (sur [-1,1])" ) ;
  PUT_LINE("-6- Methode de Laguerre ( sur [0,+8])" ) ;
  PUT_LINE("-7- Methode d'Hermite (Sur R )" ) ;
  ENTRER_ENTIER_BORNE(0,7,CHOIX,"Votre choix : " ) ;
end MENU_INTEGRATION ;

begin
  loop
    MENU_INTEGRATION(CHOIX) ;
    case CHOIX is
      when 0 => exit ;
      when 1 => DEMO_TRAPEZES ;
      when 2 => DEMO_SIMPSON ;
      when 3 => DEMO_ROMBERG ;
      when 4 => DEMO_LEGENDRE ;
      when 5 => DEMO_TCHEBYCHEV ;
      when 6 => DEMO_LAGUERRE ;
      when 7 => DEMO_HERMITE ;
    end case ;
  end loop ;
end DEM_INT ;

```

### 7.3 Spécifications des paquetages *COMMON\_FOURIER* et *FOURIER\_GENERIQUE*

```

with NOMBRES_COMPLEXES ;
use NOMBRES_COMPLEXES ;
package COMMON_FOURIER is

```

```

P_MAX : constant INTEGER := 8 ;
N_MAX : constant INTEGER := 2**P_MAX ;
subtype INDICE is INTEGER range 0..N_MAX - 1 ;
type VECTEUR_COMPLEXE is array(INDICE range <>) of COMPLEXE ;
end COMMON_FOURIER ;

with COMMON_MATRIX, COMMON_FOURIER ;
use COMMON_MATRIX, COMMON_FOURIER ;
generic
  with function f(x : in FLOAT) return FLOAT ;
package FOURIER_GENERIQUE is
procedure ECHANTILLONNAGE(tMax : in FLOAT ;
  x : out VECTEUR_COMPLEXE ; D : out FLOAT) ; -- § 5.1
-- x = (x(0),...,x(n-1)), avec x(j) = f(j*D) et D = tMax/n
function SERIE_DE_FOURIER(a, b : in VECTEUR ; x : in FLOAT)
  return FLOAT ; -- § 5.5
procedure TRACE_UNE_SERIE_DE_FOURIER(tMax : in FLOAT ;
  a, b : in VECTEUR ; NB_POINTS : in INTEGER := 300) ; -- § 5.5
procedure FFT(x : in VECTEUR_COMPLEXE ;
  y : out VECTEUR_COMPLEXE) ; -- § 5.4
procedure CALCUL_COEFFICIENTS_DE_FOURIER(y : in VECTEUR_COMPLEXE ;
  a, b : out VECTEUR) ; -- § 5.5
function ENERGIE(x : in VECTEUR_COMPLEXE) return FLOAT ; -- § 5.3
end FOURIER_GENERIQUE ;

```

## 7.4 Démonstration du paquetage *FOURIER\_GENERIQUE*

```

with TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH1, GRAPH, MATH2,
  MATH3, CURVE, NOMBRES_COMPLEXES, COMMON_FOURIER,
  FOURIER_GENERIQUE, COMMON_MATRIX, MATRIX ;
use TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH1, GRAPH, MATH2,
  MATH3, CURVE, NOMBRES_COMPLEXES, COMMON_FOURIER,
  COMMON_MATRIX, MATRIX ;

procedure DEM_FFT is

MAX_FONCTIONS : constant INTEGER := 10 ;
subtype INDICE_FONCTION is INTEGER range 0..MAX_FONCTIONS ;
type TABLEAU_DE_FONCTIONS is array(INDICE_FONCTION range <>)
  of FONCTION ;

NBR_INT, p, n, m : INTEGER ;
procedure AFFICHER_FFT(D : in FLOAT ; y : in VECTEUR_COMPLEXE) is
--      ^      ^
-- Affichage des  $y(k) = f(k/tMax) = D*y(k)$ , ou des  $c(k) = y(k)/n$ , ou :
--  $n = 2^p$ ,  $k = -n/2, \dots, 0, \dots, n/2 - 1$ .

m : INTEGER ;
begin
  m := y'LENGTH/2 ;
  for i in m..y'LENGTH - 1
  loop
    PUT(IMP,"x(") ; PUT(IMP,i - n,3) ; PUT(IMP,") = ") ;

```

```

        PUT(D*y(i)) ; NEW_LINE ;
        if (i + 1) mod 20 = 0 then
            PAUSE ;
        end if ;
    end loop ;
    for i in 0..m - 1
    loop
        PUT(IMP,"x(") ; PUT(i,3) ; PUT(") = " ) ;
        PUT(D*y(i)) ;
        NEW_LINE ;
        if (i + 1) mod 20 = 0 then
            PAUSE ;
        end if ;
    end loop ;
end AFFICHER_FFT ;

begin
-- DONNEES POUR LE CALCUL D'UNE F.F.T. ET DES COEFFICIENTS DE FOURIER
-- La fonction f(t) est defini comme une fonction continue par morceaux
-- sur un intervalle [0,tn], avec : 0 < t1 < ... < tn et :
-- f(t) = f(k-1)(t), si t dans [t(k-1),t(k)] (k = 1, ..., n),
-- les fonctions f(k) etant continues pour k = 0, ..., n - 1.
    CLRSCR ;
    PUT_LINE(
"Calculs de la FFT et des coefficients de Fourier d'une fonction f" ) ;
    NEW_LINE ;
    PUT_LINE("La fonction est definie sur [0,tMax], de periode tMax,");
    PUT_LINE("par des fonctions continues fk sur [tk,t(k+1)], avec" ) ;
    PUT_LINE("          0 < t1 < ... < tn = tMax" ) ;
    NEW_LINE ;
    ENTRER_ENTIER_BORNE(1,MAX_FONCTIONS,NBR_INT,
        "Nombre d'intervalles : " ) ;
    NEW_LINE ;
    PUT_LINE("Le nombre de points de l'echantillonnage de f est n = 2^p.");
    ENTRER_ENTIER_BORNE(1,p_Max,p,"p = " ) ;
    n := 2**p ;
    m := n/2 ;
    declare
        ff : TABLEAU_DE_FONCTIONS(0..NBR_INT - 1) ;
        t : VECTEUR(0..NBR_INT) ;
        x, y : VECTEUR_COMPLEXE(0..n - 1) ;
        a : VECTEUR(0..m- 1) ;
        b : VECTEUR(1..m - 1) ;
        C : COURBE(n) ;
        PAS, D, xOrigine, yOrigine, xUnite, yUnite,
        xFenMin, xFenMax, yFenMin, yFenMax : FLOAT ;

    function f(x : in FLOAT) return FLOAT is
        k : INTEGER := 0 ;
    begin
        loop
            k := k + 1 ;
            exit when (t(k) >= x) or (k = NBR_INT) ;

```

```

        end loop ; -- x est dans [t(k-1),t(k)]
        return EVALUE(ff(k - 1),x) ;
    end f ;

package FOURIER_F is new FOURIER_GENERIQUE(f) ;
use FOURIER_F ;

begin
    t(0) := 0.0 ;
    PUT_LINE("Entree des bornes des intervalles : t1, t2, ..., tn");
    for k in 1..NBR_INT
    loop
        ENTRER_REEL_BORNE(t(k-1),1.0E+32,t(k),"t = ") ;
    end loop ;
    NEW_LINE ;
    PUT_LINE("Definition de la fonction f = (f0,...,f(n-1))") ;
    for k in 0..NBR_INT - 1
    loop
        GET_LINE(ff(k),"f(x) = ") ;
    end loop ;
    NEW_LINE ;
    MODE_AFFICHAGE ;
    PUT_LINE(IMP,"Calcul de l'echantillonnage : Vecteur x :") ;
    ECHANTILLONNAGE(t(NBR_INT),x, D) ;
    for i in x'range
    loop
        PUT(IMP,"x") ; PUT(IMP,i,3) ; PUT(IMP," = ") ;
        PUT(IMP,x(i)) ;
        NEW_LINE(IMP) ;
        if (i + 1) mod 20 = 0 then
            PAUSE ;
        end if ;
    end loop ;
    NEW_LINE(IMP) ;
    PUT_LINE(IMP,"Calcul de la transformee de Fourier discrete") ;
    NEW_LINE(IMP) ;
    FFT(x,y) ;
    PUT_LINE(IMP,"VERIFICATION : FORMULE DE PLANCHEREL") ;
    PUT(IMP,"Ex = ") ; PUT(IMP,ENERGIE(x),6,6) ;
    PUT(IMP," ; Ey/n = ") ; PUT(IMP,ENERGIE(y)/FLOAT(n),6,6) ;
    PAUSE ;
    CLRSCR ;
    PUT_LINE(IMP,"Calcul des coefficients de Fourier trigonometriques") ;
    CALCUL_COEFFICIENTS_DE_FOURIER(y,a,b) ;
    CLRSCR ;
    PUT_LINE(IMP,"Transformee de Fourier ") ;
    AFFICHER_FFT(D,y) ;
    PAUSE ;
    CLRSCR ;
    PUT_LINE(IMP,"Coefficients de Fourier complexes") ;
    AFFICHER_FFT(1.0/FLOAT(n),y) ;
    PAUSE ;
    CLRSCR ;

```



```

PUT_LINE(IMP,"Coefficients de Fourier trigonometriques") ;
PUT_LINE("Coefficients a(k), pour k = 0, ..., n/2 - 1") ;
PUT_LINE(a) ;
NEW_LINE ;
PUT_LINE("Coefficients b(k), pour k = 1, ..., n/2 - 1") ;
PUT_LINE(b) ;
NEW_LINE ;
PAUSE ;
CLRSCR ;
C.uMin := 0.0 ;
C.uMax := t(NBR_INT) ;
PAS := (C.uMax - C.uMin)/FLOAT(C.n) ;
C.vMIN := 1.0E+37 ;
C.vMAX := -C.vMIN ;
C.u(0) := C.uMin ;
C.v(0) := f(C.u(0)) ;
for i in 1..C.n
loop
  C.u(i) := C.u(i-1) + PAS ;
  C.v(i) := f(C.u(i)) ;
  C.vMIN := MIN(C.vMIN,C.v(i-1)) ;
  C.vMAX := MAX(C.vMAX,C.v(i-1)) ;
end loop ;
TRACE_UNE_COURBE(C) ;
TITRE("Signal initial") ;
PAUSE_GRAPHIQUE ;
CLEAR_SCREEN ;
TITRE("Reconstitution par les series de Fourier") ;
TRACE_UNE_SERIE_DE_FOURIER(t(NBR_INT),a,b,n) ;
XY_AXES(xOrigine,yOrigine,xUnite,yUnite,TRUE) ;
PAUSE_GRAPHIQUE ;
end ;
MODE_TEXTE ;
end DEM_FFT ;

```

## CHAPITRE 7

# Résolution numérique des équations différentielles

### 1. Introduction.

#### Origines des problèmes d'équations différentielles

Les équations différentielles interviennent naturellement dans l'étude des *problèmes d'évolution*, que l'on peut formuler de la façon suivante :

- l'état d'un système (mécanique, électrique, économique, ...) est complètement décrit par  $p$  degrés de liberté  $y_1, \dots, y_p$  qui sont fonctions du temps  $t$  ;
- les  $y_i$  sont tous connus à un instant donné  $t$  ;
- la loi d'évolution du système est connue, c'est-à-dire qu'à chaque instant  $t > t_0$ , on connaît une relation entre les dérivées  $y_i^{(k)}(t)$  et les valeurs des paramètres  $y_j(t)$ .

Dans le cas des systèmes d'ordre 1, les  $y_i$  sont alors solutions d'un *système différentiel avec conditions initiales* :

$$\begin{cases} f_i(t, y_1(t), \dots, y_p(t), y_1'(t), \dots, y_p'(t)) = 0 & (t \in I) \\ y_i(t_0) \text{ sont donnés pour } i = 1, \dots, p \end{cases}$$

où  $I$  est un intervalle réel contenant  $t_0$ .

*Exemples* — En *mécanique*, l'exemple classique est donné par la loi de Newton :  $f(t, x(t), y(t), z(t)) = m \cdot \gamma(t)$ , où  $\gamma(t) = (x''(t), y''(t), z''(t))$  et le point de masse  $m$  et de coordonnées  $(x, y, z)$  est soumis à un instant  $t$  à une force  $f(t, x, y, z)$ .

En *électricité*, les équations différentielles interviennent, entre autres, dans l'étude des phénomènes transitoires.

Considérons, par exemple, une résistance métallique de coefficient de température  $\alpha$ , de capacité calorifique  $C$  et telle que la résistance soit reliée à la température par une loi du type :  $R = R_0 \cdot (1 + \alpha \cdot (T - T_0))$ , où  $R_0$  et  $T_0$  sont les

résistance et température initiales. On suppose, de plus, qu'elle transmet au milieu extérieur, à la température  $T_0$ , la puissance thermique

$$P_T = G \cdot (T - T_0).$$

Alors  $\theta = T - T_0$  est solution de l'équation différentielle :

$$\begin{cases} \theta'(t) = -\frac{G}{C} \cdot \theta(t) + \frac{u^2}{C \cdot R_0 \cdot (1 + \alpha \cdot \theta(t))} \\ \theta(0) = 0 \end{cases}$$

En *théorie mathématique des populations* (Cf. A. Hillion), on peut montrer que l'évolution de la taille  $x(t)$  d'une population simple, en fonction du temps, suit une loi du type :

$$\begin{cases} x'(t) = r(t, x(t)) \cdot x(t) \\ x(0) \text{ donné} \end{cases}$$

où  $r(t, x)$  représente un taux d'accroissement.

Quand le taux d'accroissement est constant, la solution est  $x(t) = x(0) \cdot e^{rt}$ , c'est-à-dire qu'on a une croissance exponentielle de la population. Mais dans la réalité une telle situation ne se produit pas et l'équation différentielle obtenue ne peut se résoudre explicitement.

Dans le cas de populations en interaction, on aboutit à un système d'équations différentielles. Un exemple classique étant le modèle proie-prédateur qui peut se décrire de la façon suivante : supposons qu'à un instant initial  $t_0 = 0$ , il y ait sur un territoire  $x(0)$  prédateurs et  $y(0)$  proies. Le problème est alors de connaître l'évolution de ces populations (si, par exemple, il y a trop de prédateurs alors les proies disparaîtront, ce qui entraînera aussi la disparition des prédateurs faute de nourriture).

On peut alors montrer que  $x$  et  $y$  sont solutions d'un système différentiel du type :

$$\begin{cases} x'(t) = r_1(t, x) \cdot x(t) - b_1(t, x) \cdot x(t) \cdot y(t) \\ y'(t) = -r_2(t, x) \cdot y(t) + b_2(t, x) \cdot x(t) \cdot y(t) \end{cases}$$

La résolution de ce système permet alors de prévoir la situation à chaque instant.

## 2. Problème de Cauchy

### 2.1 Position du problème

Soient  $f_j: [a, b] \times \mathbb{R}^p \rightarrow \mathbb{R}$   
 $(t, y_1, \dots, y_p) \mapsto f_j(t, y_1, \dots, y_p)$  des fonctions continues, avec

$j = 1, \dots, p$ .

Le *problème de Cauchy* associé aux  $f_j$ , consiste à trouver des fonctions  $y_1, y_2, \dots, y_p$  de classe  $C^1$  (i. e. continuellement dérivable) de  $[a, b]$  dans  $\mathbb{R}$  telles que :



où :

$$f(t, z) = \begin{pmatrix} z_2 \\ \cdot \\ \cdot \\ z_q \\ g_1(t, z_1, \dots, z_{p,q}) \\ z_{q+2} \\ \cdot \\ \cdot \\ z_{p,q} \\ g_p(t, z_1, \dots, z_{p,q}) \end{pmatrix}$$

*Remarque 2* — On peut aussi considérer les *problèmes aux limites*, où les conditions initiales ne sont pas nécessairement données en  $a$ . Les exemples classiques sont :

- le *problème de Dirichlet* : 
$$\begin{cases} y''(t) = f(t, y(t), y'(t)) \text{ sur } [a, b] \\ y(a) \text{ et } y(b) \text{ donnés} \end{cases}$$
- le *problème de Neumann* : 
$$\begin{cases} y''(t) = f(t, y(t), y'(t)) \text{ sur } [a, b] \\ y'(a) \text{ et } y'(b) \text{ donnés} \end{cases}$$
- le *problème de Dirichlet-Neumann* : 
$$\begin{cases} y''(t) = f(t, y(t), y'(t)) \text{ sur } [a, b] \\ y(a) \text{ et } y'(b) \text{ donnés} \end{cases}$$

*Remarque 3* — On peut aussi considérer des équations différentielles sous forme « non résolue » :  $\varphi(t, y(t), y'(t), \dots, y^{(q)}(t)) = 0$ . Mais dans ce cours on s'intéressera surtout au cas résolu, le cas général étant beaucoup plus difficile.

## 2.2 Problème de l'existence et l'unicité de solutions

*Définition* : On dit que  $f$  est *Lipschitzienne* en  $y$ , s'il existe une constante  $L \geq 0$ , telle que :

$$\forall t \in [a, b], \forall y, z \in \mathbb{R}^p, \|f(t, y) - f(t, z)\| \leq L \cdot \|y - z\|$$

*Remarque 1* — Si  $L = 0$ , dans la définition ci-dessus, alors  $f$  est constante. On supposera donc  $L > 0$ .

*Remarque 2* — *Critère pour vérifier qu'une application est Lipschitzienne en  $y$*

Avec le théorème des accroissements finis, on voit que si  $f$  est continûment dérivable par rapport à  $y$  et s'il existe  $L \geq 0$  tel que :

$$\left\| \frac{\partial f}{\partial y_j}(t, y) \right\| \leq L \quad (t \in [a, b], y \in \mathbb{R}^p, j = 1, \dots, p)$$

alors  $f$  est lipschitzienne en  $y$ .

*Théorème (Cauchy-Lipschitz) :* Si  $f$  est continue et  $L$ -Lipschitzienne en  $y$ , alors le problème de Cauchy (1) admet une unique solution  $y$  de classe  $C^1$ .

*Démonstration* — L'idée est de se ramener à un problème de point fixe que l'on résout par la méthode des approximations successives. Le problème (1) équivaut à trouver  $y$  dans  $E = C^0([a, b], \mathbb{R}^p)$  (espace des fonctions continues de  $[a, b]$  dans  $\mathbb{R}^p$ ) solution de :

$$(2) \quad y(t) = \int_a^t f(x, y(x)) dx + y_a$$

En notant  $\Phi$  l'application de  $E$  dans  $E$  définie par :

$$(\forall y \in E, z = \Phi(y)) \Leftrightarrow \left( z(t) = \int_a^t f(x, y(x)) dx + y_a \right)$$

alors (2) équivaut à :  $y = \Phi(y)$ . C'est-à-dire que  $y$  est un point fixe de  $\Phi$ .

On veut alors écrire  $y$  comme limite de la suite d'approximations successives  $(y_k)$  définie par :

$$\begin{cases} y_0 \in E \\ y_{k+1} = \Phi(y_k) \quad (k \geq 0) \end{cases}$$

Comme  $E$  est un espace de Banach, on sait qu'une telle suite converge vers l'unique point fixe de  $\Phi$ ,  $y$ , si  $\Phi$  ou l'une de ses itérées  $\Phi^n$  est contractante (théorème du point fixe de Banach).

Du fait que  $\Phi$  est  $L$ -lipschitzienne en  $y$ , on déduit que pour tout  $n > 0$  et tous  $y, z$  dans  $E$ , on a :

$$\|\Phi^n(y) - \Phi^n(z)\| \leq \frac{(L \cdot (b-a))^n}{n!} \cdot \|y - z\|$$

et tenant compte de  $\lim_{n \rightarrow +\infty} \frac{(L \cdot (b-a))^n}{n!} = 0$ , on a le résultat.

*Exemple 1* — Considérons le problème, sur  $\mathbb{R}$  :

$$\begin{cases} y'(t) = y(t) \quad (t \in \mathbb{R}) \\ y(0) = 1 \end{cases}$$

La méthode du point fixe consiste à étudier la suite de fonctions définie par :

$$\begin{cases} y_0(t) = 1 \\ y_{n+1}(t) = \int_0^t y_n(x) dx + 1 \quad (n \geq 0) \end{cases}$$

Et, par récurrence, on a :  $y_n(t) = \sum_{k=0}^n \frac{t^k}{k!}$ . On retrouve ainsi une définition de l'exponentielle.

*Exemple 2* — Considérons le problème de Cauchy :

$$\begin{cases} y'(t) = y(t)^\alpha & (t \in [0,1]) \\ y(0) = 0 \end{cases}$$

où  $\alpha \in \mathbb{R}$ .

Pour  $\alpha = 1$ , il admet une unique solution donnée par  $y = 0$ , mais pour  $\alpha \in ]0,1[$ , il admet une infinité de solutions données par :

$$y(t) = \begin{cases} 0 & \text{si } t \in [0, a] \\ ((1-\alpha) \cdot (t-a))^{1/(1-\alpha)} & \text{si } t \in [a, 1] \end{cases}$$

où  $0 < a < 1$ .

On déduit donc que l'application  $f : (t,y) \rightarrow y^\alpha$  ne vérifie pas la condition de Lipschitz.

*Corollaire (Equations différentielles linéaires)* : Si  $f(t,y) = A(t) \cdot y + b(t)$ , avec  $A : [a,b] \rightarrow \text{End}(\mathbb{R}^p)$  et  $b : [a,b] \rightarrow \mathbb{R}^p$  continues, alors le problème de Cauchy linéaire :

$$\begin{cases} y'(t) = A(t) \cdot y(t) + b(t) & t \in [a, b] \\ y(a) = y_a \end{cases}$$

admet une unique solution.

*Remarque* — Le corollaire est encore valable en remplaçant  $[a,b]$  par un intervalle quelconque.

### 2.3 Approximation de la solution d'un problème de Cauchy par discrétisation

On suppose que  $f$  est L-Lipschitzienne par rapport à  $y$ , de sorte que le problème de Cauchy associé admet une unique solution.

Pour trouver une approximation de la solution  $y$ , on utilisera une *méthode de discrétisation* dont le principe est le suivant : on se donne une subdivision de l'intervalle  $[a,b]$   $a = t_0 < t_1 < \dots < t_n = b$  et pour tout  $k = 1, \dots, n$ , on cherche une valeur approchée  $y_k$  de la valeur exacte  $y(t_k)$ .

La quantité  $h = \text{Max} \{ t_{i+1} - t_i ; 0 \leq i \leq n-1 \}$  est le *pas d'intégration*.

Quand  $h$  sera très petit, on espère que la fonction affine par morceaux définie par  $\varphi(t_k) = y_k$  pour  $k = 0, 1, \dots, n$  nous donnera une bonne approximation de  $y$ .

En général, on utilisera une subdivision à pas constant, soit :

$$t_k = a + k \cdot \frac{b-a}{n}, \text{ pour } k = 0, 1, \dots, n.$$

L'idée sera alors de calculer les  $y_k$  par des formules de récurrence.

On distingue deux types de schémas d'intégration :

- les *schémas à un pas* (ou à pas séparés) :  $y_k$  est fonction de  $y_{k-1}$  ;
- les *schémas à plusieurs pas* (ou à pas liés) :  $y_k$  est fonction de  $y_{k-i}$  où  $i = 1, 2, \dots, p$  et  $p$  est le nombre de pas.



### 3. Généralités sur les méthodes d'intégration à un pas

#### 3.1 Définitions

Une méthode à un pas est définie par les formules de récurrence :

$$(1) \quad \begin{cases} y_0 \text{ donné} \\ y_{k+1} = y_k + h \cdot \Phi(t_k, y_k, h), \text{ pour } k = 0, 1, \dots, n-1 \end{cases}$$

où  $\Phi$  est une fonction continue qui ne dépend que de  $f$ .

*Définition :* L'erreur de consistance en  $t_k$  de la méthode est la quantité :

$$\tau_k = y(t_{k+1}) - y(t_k) - h \cdot \Phi(t_k, y(t_k), h)$$

Et on dira que la méthode est consistante avec l'équation différentielle si  $\Phi(t, y, 0) = f(t, y)$  pour tous  $t$  et  $y$ .

*Remarque* — L'erreur de consistance permet de mesurer de combien la solution exacte est éloignée du schéma d'intégration.

*Définition :* L'erreur de discrétisation en  $t_k$  de la méthode est la quantité :

$$e_k = \|y(t_k) - y_k\|$$

#### 3.2 Lien entre l'erreur de consistance et l'erreur de discrétisation

##### 3.2.1 Hypothèses

On fait les hypothèses suivantes :

- $p = 1$  (pour simplifier) ;
- $f$  est lipschitzienne en  $y$  ;
- $\Phi$  est lipschitzienne en  $y$  ;
- la méthode est consistante avec l'équation différentielle ;
- $f$  et  $\Phi$  sont de classe  $C^1$ .

##### 3.2.2 Première évaluation de l'erreur de consistance

L'erreur de consistance en  $t_k$  est :

$$\tau_k(h) = y(t_{k+1}) - y(t_k) - h \cdot \Phi(t_k, y(t_k), h)$$

En faisant un développement limité au voisinage de 0, on a :

$$\Phi(t_k, y(t_k), h) = \Phi(t_k, y(t_k), 0) + h \cdot \frac{\partial \Phi}{\partial h}(t_k, y(t_k), \xi_k) \quad (\xi_k \text{ entre } 0 \text{ et } h)$$

Avec (iv), on déduit alors que :

$$\tau_k(h) = y(t_{k+1}) - y(t_k) - h \cdot f(t_k, y(t_k)) - h^2 \cdot \frac{\partial \Phi}{\partial h}(t_k, y(t_k), \xi_k)$$

avec  $f(t, y(t)) = y'(t)$  et  $f$  de classe  $C^1$ , on déduit que  $y$  est de classe  $C^2$ , puis avec la formule de Taylor, on a :

$$\tau_k(h) = h^2 \cdot \left\{ \frac{1}{2} \cdot y''(v_k) - \frac{\partial \Phi}{\partial h}(t_k, y(t_k), \xi_k) \right\} = O(h^2)$$

### 3.2.3 Hypothèses supplémentaires

Si on suppose que  $f$  et  $\Phi$  sont de classe  $C^p$ , on peut faire des développements limités à l'ordre  $p$ .

### 3.2.4 Deuxième évaluation de l'erreur de consistance

Un développement limité en  $h = 0$ , nous donne :

$$\Phi(t_k, y(t_k), h) = \Phi(t_k, y(t_k), 0) + h \cdot \frac{\partial \Phi}{\partial h}(t_k, y(t_k), 0) + \dots + \frac{h^{p-1}}{(p-1)!} \cdot \frac{\partial^{p-1} \Phi}{\partial h^{p-1}}(t_k, y(t_k), 0) + \frac{h^p}{p!} \cdot \frac{\partial^p \Phi}{\partial h^p}(t_k, y(t_k), \xi_k)$$

Pour calculer un développement limité de  $y$  en  $t_k$ , il nous faut d'abord calculer les dérivées successives de  $y$  en fonction de  $f$ .

Au départ, on a :  $y'(t) = f(t, y(t))$ . Donc :

$$y''(t) = \frac{\partial f}{\partial t}(t, y(t)) + \frac{\partial f}{\partial y}(t, y(t)) \cdot y'(t) = \frac{\partial f}{\partial t}(t, y(t)) + \frac{\partial f}{\partial y}(t, y(t)) \cdot f(t, y(t))$$

Et par récurrence, on a :

$$y^{(k)}(t) = f^{(k-1)}(t, y(t)) \quad (k = 1, \dots, p+1)$$

où les  $f^{(i)}$  sont définis par :

$$\begin{cases} f^{(0)}(t, y) = f(t, y) \\ f^{(k)}(t, y) = \frac{\partial f^{(k-1)}}{\partial t}(t, y) + \frac{\partial f^{(k-1)}}{\partial y}(t, y) \cdot f(t, y) \quad (k = 1, \dots, p+1) \end{cases}$$

Un développement de Taylor, au voisinage de  $t_k$ , va alors nous donner :

$$y(t_{k+1}) = y(t_k) + h \cdot f^{(0)}(t_k, y(t_k)) + \dots + \frac{h^p}{p!} \cdot f^{(p-1)}(t_k, y(t_k)) + O(h^{p+1})$$

Ce qui donne alors, pour  $\tau_k(h)$  :

$$\begin{aligned} \tau_k(h) = h \cdot \left\{ f^{(0)}(t_k, y(t_k)) - \Phi(t_k, y(t_k), 0) \right\} \\ + \dots + \frac{h^p}{p!} \cdot \left\{ f^{(p-1)}(t_k, y(t_k)) - p \cdot \frac{\partial^{p-1} \Phi}{\partial h^{p-1}}(t_k, y(t_k), 0) \right\} + O(h^{p+1}) \end{aligned}$$

On déduit alors le :

*Théorème* : Avec les hypothèses ci-dessus, si de plus, on a :

$$\begin{cases} 2 \cdot \frac{\partial \Phi}{\partial h}(t, y, 0) = f^{(1)}(t, y) \\ p \cdot \frac{\partial^{p-1} \Phi}{\partial h^{p-1}}(t, y, 0) = f^{(p-1)}(t, y) \end{cases}$$

alors,  $\tau_k(h) = O(h^{p+1})$ .

*Remarque* — La réciproque est aussi vraie. (Cf. Crouzeix et Mignot, p. 95 et 96, Théorème (5.6))

### 3.2.5 Majoration de l'erreur de discrétisation

On note :

$$E_k = y(t_k) - y_k \text{ et } e_k = \|E_k\|$$

Avec :

$$\begin{cases} y_{k+1} = y_k + h \cdot \Phi(t_k, y_k, h) \\ y(t_{k+1}) = y(t_k) + h \cdot \Phi(t_k, y(t_k), h) + \tau_k(h) \end{cases}$$

on déduit que :

$$E_{k+1} = E_k + h \cdot \{\Phi(t_k, y_k, h) - \Phi(t_k, y(t_k), h)\} - \tau_k(h)$$

En supposant les hypothèses du théorème (2.2.4) vérifiées, on a :

$$|\tau_k(h)| \leq C_1 \cdot h^{p+1}$$

et avec  $\Phi$   $C_2$ -Lipschitzienne, on déduit que :

$$e_{k+1} < (1 + h \cdot C_2) \cdot e_k + C_1 \cdot h^{p+1}$$

Mais on a le :

*Lemme (Gronwall) :* Si  $(e_k)$  est une suite de réels positifs telle que :

$$e_{k+1} < (1 + \alpha) \cdot e_k + \beta \quad (\alpha, \beta > 0)$$

Alors :  $e_k \leq (e_0 + \frac{\beta}{\alpha}) \cdot e^{k \cdot \alpha}$ , pour tout  $k > 0$ .

*Démonstration* — Par récurrence, on déduit que :

$$e_k \leq (1 + \alpha)^k \cdot e_0 + \beta \cdot (1 + (1 + \alpha) + \dots + (1 + \alpha)^{k-1})$$

soit :

$$e_k \leq (1 + \alpha)^k \cdot e_0 + \beta \cdot \frac{(1 + \alpha)^k - 1}{\alpha}$$

et :

$$e_k \leq (e_0 + \frac{\beta}{\alpha}) \cdot (1 + \alpha)^k$$

Comme  $\alpha = h \cdot C_2$  est petit, on majore  $(1 + \alpha)$  par  $e^\alpha$ , ce qui donne le résultat.

*Théorème :* Si l'erreur de consistance est un  $O(h^{p+1})$ , pour tout  $k$ , alors l'erreur de discrétisation est un  $O(h^p)$ .

*Démonstration* — Avec ce qui précède et la condition initiale  $e_0 = 0$ , on déduit que

$$e_k \leq C \cdot e^{k \cdot h \cdot C_2} \cdot h^p$$

et avec  $k \cdot h \leq (b - a)$ , on déduit que :

$$e_k \leq C \cdot e^{(b-a) \cdot C_2} \cdot h^p$$

C'est-à-dire que  $e_k = O(h^p)$ .

*Remarques* — (i) En général, on a seulement  $e_0 \cong 0$ , car  $y_a$  est seulement connu de façon approchée. Et cela peut parfois poser des problèmes.

(ii) Si on fixe le pas  $h$  et si la longueur de l'intervalle tend vers l'infini, alors le majorant de l'erreur ci-dessus va croître de façon exponentielle avec  $(b - a)$ .

*Définition* : Avec les notations et conditions du théorème ci-dessus, on dit que la méthode est d'ordre  $p$ .

*Exemple* — La méthode d'Euler est définie par  $\Phi(x,y) = y$  et si on considère l'équation :

$$\begin{cases} y' = y \text{ sur } [0, b] \\ y(0) = 1 \end{cases}$$

le schéma d'Euler s'écrit :

$$\begin{cases} y_0 = 1 \\ y_{k+1} = (1 + h) \cdot y_k \end{cases}$$

donc  $y_k = (1 + h)^k = \left(1 + \frac{b}{n}\right)^k$ .

Et pour  $k = n$ , on a :

$$y_n = \left(1 + \frac{b}{n}\right)^n \xrightarrow{n \rightarrow +\infty} e^b = y(b)$$

L'erreur de discrétisation en  $b$  est alors :

$$e_n = e^b - \left(1 + \frac{b}{n}\right)^n = e^{n \cdot h} \cdot \left\{1 - \left(\frac{1+h}{e^h}\right)^n\right\} \xrightarrow{n \rightarrow +\infty} +\infty$$

pour un pas  $h$  fixé.

### 3.3 Les méthodes de Runge-Kutta

On suppose ici que  $f$  est de classe  $C^2$  et  $L$ -lipschitzienne en  $y$ .

#### 3.3.1 Principe des méthodes de Runge-Kutta

A partir de l'équation différentielle  $y'(t) = f(t, y(t))$  sur  $[a, b]$ , on peut écrire que :

$$y(t_{k+1}) = y(t_k) + \int_{t_k}^{t_{k+1}} f(x, y(x)) dx \quad (0 \leq k \leq n-1)$$

Et l'idée est alors de calculer une valeur approchée de  $y_{k+1}$  en fonction de  $y_k$  à partir d'une formule d'intégration approchée.

#### 3.3.2 Exemples classiques

(i) *Méthode des rectangles à gauche* —  $\int_a^b g(t) dt \cong (b - a) \cdot g(a)$ .

On obtient alors le schéma :

$$y_{k+1} = y_k + h \cdot f(t_k, y_k) \text{ pour } k = 0, 1, \dots, n-1$$

C'est la *méthode d'Euler* qui est d'ordre 1.

(ii) *Méthode du point milieu* —  $\int_a^b g(t) dt \cong (b - a) \cdot g\left(\frac{a+b}{2}\right)$ .

On obtient alors le schéma :

$$y_{k+1} = y_k + h \cdot f\left(t_k + \frac{h}{2}, y_k + \frac{h}{2} \cdot f(t_k, y_k)\right)$$

C'est la *méthode d'Euler-Cauchy* qui est d'ordre 2.

$$(iii) \text{ Méthode des trapèzes — } \int_a^b g(t) dt \cong (b-a) \cdot \frac{g(a)+g(b)}{2}.$$

On obtient alors le schéma ;

$$y_{k+1} = y_k + \frac{h}{2} \cdot \{f(t_k, y_k) + f(t_k + h, y_k + h \cdot f(t_k, y_k))\}$$

C'est la *méthode de Heun* qui est d'ordre 2.

$$(iv) \text{ Méthode de Simpson — } \int_a^b g(t) dt \cong \frac{b-a}{6} \cdot \left\{g(a) + 4 \cdot g\left(\frac{a+b}{2}\right) + g(b)\right\}.$$

On peut alors écrire :

$$y(t_{k+1}) \cong y(t_k) + \frac{h}{6} \cdot \left\{f(t_k, y(t_k)) + 2 \cdot f\left(t_k + \frac{h}{2}, y\left(t_k + \frac{h}{2}\right)\right) + 2 \cdot f\left(t_k + \frac{h}{2}, y\left(t_k + \frac{h}{2}\right)\right) + f(t_k + h, y(t_k + h))\right\}$$

*Calcul approché de  $f(t_k, y(t_k))$*  — On a :

$$f(t_k, y(t_k)) \cong K_1 = f(t_k, y_k)$$

*Premier calcul approché de  $f\left(t_k + \frac{h}{2}, y\left(t_k + \frac{h}{2}\right)\right)$ .*

On a :  $y(t_k + \frac{h}{2}) = y(t_k) + \int_{t_k}^{t_k + \frac{h}{2}} f(x, y(x)) dx$  et la première idée est de calculer cette intégrale avec la méthode des rectangles à gauche, soit :

$$\int_{t_k}^{t_k + \frac{h}{2}} f(x, y(x)) dx \cong \frac{h}{2} \cdot f(t_k, y(t_k)) \cong \frac{h}{2} \cdot f(t_k, y_k)$$

ce qui donne :

$$(1) \quad y\left(t_k + \frac{h}{2}\right) \cong y_k + \frac{h}{2} \cdot K_1$$

et

$$f\left(t_k + \frac{h}{2}, y\left(t_k + \frac{h}{2}\right)\right) \cong K_2 = f\left(t_k + \frac{h}{2}, y_k + \frac{h}{2} \cdot K_1\right)$$

*Deuxième calcul approché de  $f\left(t_k + \frac{h}{2}, y\left(t_k + \frac{h}{2}\right)\right)$ .*

On utilise la méthode des rectangles à droite pour le calcul d'une intégrale, ce qui donne :

$$y(t_k + \frac{h}{2}) = y(t_k) + \int_{t_k}^{t_k + \frac{h}{2}} f(x, y(x)) dx \cong y_k + \frac{h}{2} \cdot f\left(t_k + \frac{h}{2}, y\left(t_k + \frac{h}{2}\right)\right)$$

et en utilisant (1), il vient :

$$y(t_k + \frac{h}{2}) \cong y_k + \frac{h}{2} \cdot f\left(t_k + \frac{h}{2}, y_k + \frac{h}{2} \cdot K_1\right)$$

c'est-à-dire :

$$(2) \quad y\left(t_k + \frac{h}{2}\right) \cong y_k + \frac{h}{2} \cdot K_2$$

et :

$$f\left(t_k + \frac{h}{2}, y\left(t_k + \frac{h}{2}\right)\right) \cong K_3 = f\left(t_k + \frac{h}{2}, y_k + \frac{h}{2} \cdot K_2\right)$$

Calcul approché de  $f(t_k + h, y(t_k + h))$ .

On a :

$$y(t_k + h) = y(t_k) + \int_{t_k}^{t_k+h} f(x, y(x)) dx$$

l'intégrale se calculant avec la méthode du point milieu, soit :

$$y(t_k + h) \cong y_k + h \cdot f\left(t_k + \frac{h}{2}, y\left(t_k + \frac{h}{2}\right)\right)$$

ce qui donne avec (2) :

$$y(t_k + h) \cong y_k + h \cdot f\left(t_k + \frac{h}{2}, y_k + \frac{h}{2} \cdot K_2\right)$$

c'est-à-dire :

$$(3) \quad y(t_k + h) \cong y_k + h \cdot K_3$$

et :

$$f(t_k + h, y(t_k + h)) \cong K_4 = f(t_k + h, y_k + h \cdot K_3)$$

La méthode obtenue est d'ordre 4. On l'appelle méthode *RK4* (pour méthode de Runge-Kutta d'ordre 4). C'est cette méthode que l'on utilise le plus souvent. Elle est donc définie par :

$$\Phi(t, y, h) = \frac{1}{6} \cdot \{K_1 + 2 \cdot K_2 + 2 \cdot K_3 + K_4\}$$

où :

$$\begin{cases} K_1 = f(t, y) \\ K_2 = f\left(t + \frac{h}{2}, y + \frac{h}{2} \cdot K_1\right) \\ K_3 = f\left(t + \frac{h}{2}, y + \frac{h}{2} \cdot K_2\right) \\ K_4 = f(t + h, y + h \cdot K_3) \end{cases}$$

Les calculs des  $K_j$  devant se faire dans l'ordre indiqué.

### 3.3.3 Remarque

Les méthodes de Runge-Kutta plus générales sont basées sur le même principe. On calcule tout d'abord  $\int_{t_k}^{t_k+h} f(x, y(x)) dx$  avec une méthode du type :

$$\int_{t_k}^{t_k+h} f(x, y(x)) dx \cong \sum_{i=1}^q b_i \cdot f(t_{i,k}, y(t_{i,k}))$$

avec  $t_k \leq t_{1,k} \leq \dots \leq t_{q,k} \leq t_{k+1}$ .

Le calcul de chaque  $f(t_{i,k}, y(t_{i,k}))$  se faisant avec :  $y(t_{i,k}) = y(t_k) + \int_{t_k}^{t_{i,k}} f(x, y(x)) dx$  en utilisant une formule d'intégration du type :

$$\int_{t_k}^{t_{i,k}} f(x, y(x)) dx \cong \sum_{j=1}^i a_{i,j} \cdot f(t_{j,k}, y(t_{j,k}))$$

et en exploitant les calculs précédents.

### 3.3.4 Les schémas RK2

L'idée est de généraliser la méthode d'Euler-Cauchy en cherchant tous les quadruplés  $a_1, a_2, a_3, a_4$  tels que le schéma d'intégration défini par :

$$\Phi(t, y, h) = a_1 \cdot f(t, y) + a_2 \cdot f(t + a_3 \cdot h, y + a_4 \cdot h \cdot f(t, y))$$

soit d'ordre 2, ce qui équivaut à dire que l'erreur de consistance est d'ordre 3.

On a :

$$\begin{cases} \Phi(t, y, 0) = (a_1 + a_2) \cdot f(t, y) \\ \frac{\partial \Phi}{\partial h}(t, y, 0) = a_2 \cdot a_3 \cdot \frac{\partial f}{\partial t}(t, y) + a_2 \cdot a_4 \cdot \frac{\partial f}{\partial y}(t, y) \cdot f(t, y) \end{cases}$$

Donc la méthode sera d'ordre 2 si, et seulement si :

$$\begin{cases} a_1 + a_2 = 1 \\ a_2 \cdot a_3 = a_2 \cdot a_4 = \frac{1}{2} \end{cases}$$

En posant  $\lambda = a_2$ , on déduit alors que les méthodes de Runge-Kutta d'ordre 2 sont définies par :

$$\Phi(t, y, h) = (1 - \lambda) \cdot f(t, y) + \lambda \cdot f\left(t + \frac{h}{2 \cdot \lambda}, y + \frac{h}{2 \cdot \lambda} \cdot f(t, y)\right)$$

Pour  $\lambda = 1$ , on a le schéma d'Euler Modifié, pour  $\lambda = 1/2$ , on a le schéma d'Euler-Cauchy et pour  $\lambda = 3/4$ , on a le schéma de Heun.

### 3.4 Programmation structurée de la méthode RK4 pour les équations différentielles d'ordre 1

Dans ce paragraphe, on décrit une itération de la méthode de Runge-Kutta d'ordre 4 pour résoudre une équation différentielle d'ordre 1 :

$$\begin{cases} y'(t) = f(t, y(t)) \text{ sur } [a, b] \\ y(a) = y_a \end{cases}$$

où  $f : [a, b] \rightarrow \mathbb{R}$  est continue et lipschitzienne en  $y$ .

A l'étape  $i$  du calcul, on dispose d'une abscisse initiale  $t\_Initial$  et d'une ordonnée  $y\_Initial$ , valeur approchée de  $y(t\_Initial)$ , on déduit alors les valeurs  $t\_Final = t\_Initial + h$ , où  $h$  est un pas d'intégration, et la valeur approchée de  $y(t\_Final)$ ,  $y\_Final$ , obtenue par la méthode de Runge\_Kutta d'ordre 4.

Il suffit ensuite d'exécuter cette procédure sur chaque intervalle  $[t_i, t_{i+1}]$ , ( $0 \leq i \leq n - 1$ , où  $h = \frac{b-a}{n}$ ) pour avoir une solution approchée sur  $[a, b]$ .

La variable  $f$ , de type Fonction, représente une fonction de deux variables réelles  $x$  et  $y$  définissant l'équation différentielle.

*PROCEDURE RK4\_11(Entrée  $f$ : Fonction ;  $t\_Initial, y\_Initial, h$ : Réel ;  
Sortie  $t\_Final, y\_Final$ : Réel) ;*

*Début*

$$k1 = h \cdot f(t\_Initial, y\_Initial) ;$$

$$t\_Final = t\_Initial + \frac{h}{2} ;$$

$$k2 = h \cdot f(t\_Final, y\_Initial + \frac{k1}{2}) ;$$

$$k3 = h \cdot f(t\_Final, y\_Initial + \frac{k2}{2}) ;$$

$$t\_Final = t\_Initial + h ;$$

$$k4 = h \cdot f(t\_Final, y\_Initial + k3) ;$$

$$y\_Final = y\_Initial + \frac{k1 + 2 \cdot (k2 + k3) + k4}{6} ;$$

*Fin ;*

### 3.5 Programmation structurée de la méthode RK4 pour les systèmes de $p$ équations différentielles d'ordre 1

Dans ce paragraphe, on s'intéresse à un système de  $p$  équations différentielles d'ordre 1 :

$$\begin{cases} y'_1 = f_1(t, y_1, \dots, y_p) \\ \dots \dots \dots t \in [a, b] \\ y'_p = f_p(t, y_1, \dots, y_p) \\ y_k(a) = y_{k,a} \text{ donnés dans } \mathbb{R} \text{ pour } k = 1, \dots, p \end{cases}$$

Il suffit de reprendre la procédure précédente de « façon vectorielle ».

Pour ce faire on définit un type SystemeDiff comme un tableau de fonctions des variables  $t$  de type réel et  $y$  de type Variable\_2 lui même défini comme un tableau de réels. Ainsi  $f(j)(t, y)$  représentera  $f_j(t, y_1, \dots, y_p)$ , pour tout  $j = 1, \dots, p$ . Ce qui donne la procédure ci-dessous décrivant une itération de la méthode :

*PROCEDURE RK4\_p1(Entrée  $f$ : SystemeDiff ;  $t\_Initial, h$ : Réel ;  
 $y\_Initial$ : Variable\_2 ;  $p$ : Entier ;  
Sortie  $t\_Final$ : Réel ;  $y\_Final$ : Variable\_2) ;*

*Début*

*Pour  $j$  allant de 1 à  $p$  faire  $k1(j) = h \cdot f(j)(t\_Initial, y\_Initial)$  ;*

$$t\_Final = t\_Initial + \frac{h}{2} ;$$

$$\text{Pour } j \text{ allant de 1 à } p \text{ faire } y\_Final(j) = y\_Initial(j) + \frac{k1(j)}{2} ;$$

*Pour  $j$  allant de 1 à  $p$  faire  $k2(j) = h \cdot f(j)(t\_Final, y\_Final)$  ;*

$$\text{Pour } j \text{ allant de 1 à } p \text{ faire } y\_Final(j) = y\_Initial(j) + \frac{k2(j)}{2} ;$$



Pour  $j$  allant de 1 à  $p$  faire  $k3(j) = h \cdot f(j)(t\_Final, y\_Final)$  ;

$t\_Final = t\_Initial + h$  ;

Pour  $j$  allant de 1 à  $p$  faire  $y\_Final(j) = y\_Initial(j) + k3(j)$  ;

Pour  $j$  allant de 1 à  $p$  faire  $k4(j) = h \cdot f(j)(t\_Final, y\_Final)$  ;

Pour  $j$  allant de 1 à  $p$  faire

$$y\_Final(j) = y\_Initial(j) + \frac{k1(j) + 2 \cdot (k2(j) + k3(j)) + k4(j)}{6} ;$$

Fin ;

### 3.6 Programmation structurée pour les équations différentielles d'ordre $q$

Dans ce paragraphe, on s'intéresse aux équations différentielles d'ordre  $q \geq 1$ .

$$\begin{cases} y^{(q)}(t) = f(t, y(t), y'(t), \dots, y^{(q-1)}(t)) \text{ sur } [a, b] \\ y(a) = y_a, \dots, y^{(q-1)}(a) = y_a^{q-1} \end{cases}$$

Les conditions initiales sont stockées dans un vecteur  $y\_Initial$  de composantes  $y^{(j)}(a)$  pour  $j = 0, \dots, q - 1$ .

Comme indiqué au paragraphe (2.1), on se ramène, par changement de fonction inconnue, à un système de  $q$  équations différentielles d'ordre 1 en posant  $z_i = y^{(i)}$  pour  $i = 0, \dots, q - 1$ , ce qui donne le système :

$$\begin{cases} z'_0 = z_1 \\ \dots\dots\dots \\ z'_{q-2} = z_{q-1} \\ z'_{q-1} = f(t, z_0, \dots, z_{q-1}) \end{cases}$$

noté  $u' = g(t, u)$  que l'on résout par la méthode de Runge-Kutta pour les systèmes.

La variable  $f$ , de type Fonction, représente une fonction des variables  $t$ , de type réel et  $y$ , de type Variable\_2 définissant l'équation différentielle. De sorte que  $f(t, y)$ , avec  $y$  de composantes  $y_0, \dots, y_{q-1}$  représente  $f(t, y, y', \dots, y^{(q-1)})$ .

PROCEDURE RK4\_1q(Entrée  $f$  : Fonction ;  $t\_Initial, h$  : Réels ;

$y\_Initial$  : Variable\_2 ;  $q$  : Entier ;

Sortie  $t\_Final$  : Réel ;  $y\_Final$  : Variable\_2 ;

Début

Pour  $j$  allant de 0 à  $q - 2$  Faire  $k1(j) = h \cdot y\_Initial(j+1)$  ;

$k1(q-1) = h \cdot f(t\_Initial, y\_Initial)$  ;

$$t\_Final = t\_Initial + \frac{h}{2} ;$$

$$\text{Pour } j \text{ allant de } 0 \text{ à } q - 1 \text{ faire } y\_Final(j) = y\_Initial(j) + \frac{k1(j)}{2} ;$$

Pour  $j$  allant de 0 à  $q - 2$  Faire  $k2(j) = h \cdot y\_Final(j+1)$  ;

$k2(q-1) = h \cdot f(t\_Final, y\_Final)$  ;

$$\text{Pour } j \text{ allant de } 0 \text{ à } q - 1 \text{ faire } y\_Final(j) = y\_Initial(j) + \frac{k2(j)}{2} ;$$

Pour  $j$  allant de 0 à  $q - 2$  Faire  $k3(j) = h \cdot y\_Final(j+1)$  ;

$k3(q-1) = h f(t\_Final, y\_Final) ;$   
 Pour  $j$  allant de 0 à  $q-1$  faire  $y\_Final(j) = y\_Initial(j) + k3(j) ;$   
 $t\_Final = t\_Initial + h ;$   
 Pour  $j$  allant de 0 à  $q-2$  Faire  $k4(j) = h y\_Final(j+1) ;$   
 $k4(q-1) = h f(t\_Final, y\_Final) ;$   
 Pour  $j$  allant de 0 à  $q-1$  Faire  

$$y\_Final(j) = y\_Initial(j) + \frac{k1(j) + 2 \cdot (k2(j) + k3(j)) + k4(j)}{6} ;$$

Fin ;

### 3.7 Programmation structurée pour les systèmes de $p$ équations différentielles d'ordre $q$

On s'intéresse ici à la résolution d'un système de  $p$  équations différentielles d'ordre  $q$  du type :

$$y_j^{(q)} = f_j(t, y_1, \dots, y_1^{(q-1)}, \dots, y_p, \dots, y_p^{(q-1)}) \quad (j = 1, \dots, p)$$

avec les conditions initiales :  $y_j^{(k)}(a)$  donnés pour  $j = 1, \dots, p$  et  $k = 0, \dots, q-1$ .

Un tel système sera représenté par un tableau de fonctions  $f$ . Pour tous  $j = 1, \dots, p$  et  $i = 0, \dots, q-1$ , la dérivée  $y_j^{(i)}$  sera représentée par la composante  $y_{j,i}$  d'une matrice  $y$  à  $p$  lignes et  $q$  colonnes. La fonction  $f_j$  a pour paramètre un réel  $x$  et un vecteur  $y$  de composantes  $(y_1, \dots, y_{p,q})$  où  $y_{q,(j-1)+i+1}$  représente la dérivée numéro  $i$  de la fonction numéro  $j$ . La représentation matricielle de ces dérivées est faite pour faciliter l'entrée des données et la représentation vectorielle pour faciliter la programmation.

Ce qui donne la procédure suivante, toujours pour une itération de la méthode :

PROCEDURE RK4\_pq(Entrée  $f$  : SystemeDiff ;  $t\_Initial$ ,  $h$  : Réel ;  
 $y\_Initial$  : Matrice ;  $p$ ,  $q$  : Entier ;  
 Sortie  $t\_Final$  : Réel ;  $y\_Final$  : Matrice) ;

Début

Pour  $j$  allant de 1 à  $p$  faire

Début

Pour  $i$  allant de 0 à  $q-1$  faire

Début

$y_{q,(j-1)+i+1} = y\_Initial_{j,i} ;$

Fin ;

Fin ;

Pour  $j$  allant de 1 à  $p$  faire

Début

Pour  $i$  allant de 0 à  $q-2$  faire

Début

$k1_{q,(j-1)+i+1} = h y_{q,(j-1)+i+2} ;$

Fin ;

$k1_{q,j} = h f(j)(t\_Initial, y_i) ;$

Fin ;

$t\_Final = t\_Initial + \frac{h}{2} ;$

Pour  $j$  allant de 1 à  $p$  faire

*Début*

$$y_j = y_i + 0.5 \cdot k_1 j;$$

*Fin ;*

*Pour j allant de 1 à p faire*

```

Début
  Pour i allant de 0 à q - 2 faire
    Début
       $k2_{q \cdot (j-1) + i + 1} = h yf_{q \cdot (j-1) + i + 2};$ 
    Fin ;
     $k2_{qj} = hf(j)(t\_Initial, yf);$ 
  Fin ;
  Pour j allant de 1 à p · q faire
    Début
       $yf_j = y_i + 0.5 \cdot k2_j;$ 
    Fin ;
    Pour j allant de 1 à p faire
      Début
        Pour i allant de 0 à q - 2 faire
          Début
             $k3_{q \cdot (j-1) + i + 1} = h yf_{q \cdot (j-1) + i + 2};$ 
          Fin ;
           $k3_{qj} = hf(j)(t\_Initial, yf);$ 
        Fin ;
         $t\_Final = t\_Initial + h;$ 
        Pour j allant de 1 à p · q faire
          Début
             $yf_j = y_i + k3_j;$ 
          Fin ;
          Pour j allant de 1 à p faire
            Début
              Pour i allant de 0 à q - 2 faire
                Début
                   $k4_{q \cdot (j-1) + i + 1} = h yf_{q \cdot (j-1) + i + 2};$ 
                Fin ;
                 $k4_{qj} = hf(j)(t\_Initial, yf);$ 
              Fin ;
              Pour j allant de 1 à p faire
                Début
                  Pour i allant de 0 à q - 1 faire
                    Début
                       $y\_Final(j, i) = y_{i \cdot q \cdot (j-1) + i + 1}$ 
                       $+ \frac{k1_{q \cdot (j-1) + i + 1} + 2 \cdot (k2_{q \cdot (j-1) + i + 1} + k3_{q \cdot (j-1) + i + 1}) + k4_{q \cdot (j-1) + i + 1}}{6}$ 
                    Fin ;
                  Fin ;
                Fin ;
              Fin ;
            Fin ;
          Fin ;
        Fin ;
      Fin ;
    Fin ;
  Fin ;

```

## 4. Contrôle du pas d'intégration

### 4.1 Position du problème

A l'étape  $i$  du calcul, dans une méthode à un pas de résolution d'une équation différentielle on voudrait modifier le pas d'intégration  $h_i$  de façon à atteindre une précision déterminée  $\epsilon > 0$ . L'idée de prendre un pas  $h$  très petit présente

l'inconvénient de nécessiter beaucoup de calculs et de ce fait va entraîner un cumul trop important des erreurs d'arrondis.

On prendra donc, au départ, un pas  $h_0$  « raisonnable » et ce dernier sera modifié en cours de route (diminué ou augmenté) de façon à avoir la précision désirée.

Pour ce faire, il est nécessaire de connaître une estimation de l'erreur de discrétisation en fonction du pas  $h_i$  à chaque étape  $i$  du calcul. Il est clair que cette estimation va entraîner une augmentation du temps de calcul.

## 4.2 Choix du pas d'intégration à l'étape $i$ du calcul

Considérons une méthode d'intégration à un pas définie par :

$$\begin{cases} y_0 = y(t_0) \text{ donné} \\ y_{i+1} = y_i + h \cdot \Phi(t_i, y_i, h_i) \quad (i = 0, \dots, n-1) \end{cases}$$

où  $h_i = t_{i+1} - t_i$  est un pas non constant.

En supposant que, à l'étape  $i$  du calcul,  $y_i$  est une approximation de  $y(t_i)$  avec une précision inférieure ou égale à  $\epsilon_0$ , on voudrait modifier le pas  $h_i$  de sorte que  $y_{i+1}$  donne une approximation de  $y(t_{i+1})$  avec une précision  $\leq \epsilon_0$ .

Notons  $h_{\text{Initial}}$  le pas utilisé à l'étape  $i$  du calcul (au départ, on prendra un pas  $h_0$  arbitraire).

En reprenant le raisonnement du paragraphe (2.2.4), on voit que pour une méthode d'ordre  $p$  et un pas  $h$  donné, l'erreur de consistance est donnée par :

$$\tau(t, h) = y(t+h) - y(t) - h \cdot \Phi(t, y(t), h) = C_x \cdot h^{p+1} + O(h^{p+2})$$

où  $C_x$  est une constante qui ne dépend que de  $y^{(p+1)}(t)$  et non de  $h$  (on suppose évidemment qu'il existe une seule solution assez régulière).

En appliquant la méthode une fois avec le pas  $h = h_{\text{Initial}}$  et deux fois avec le pas  $h/2$  on obtient :

$$\begin{cases} y(t_i + h) \cong y_2 + C_1 \cdot h^{p+1} + O(h^{p+2}) \\ y(t_i + h) \cong y_2 + \left( C_1 + C_{i+\frac{1}{2}} \right) \cdot \left( \frac{h}{2} \right)^{p+1} + O(h^{p+2}) \end{cases}$$

où  $y_1$  et  $y_2$  sont les valeurs approchées de  $y(t_i + h)$  obtenues avec chacune des deux méthodes.

En supposant  $h$  assez petit, on peut écrire que  $C_{i+\frac{1}{2}} \cong C_1$ , de sorte que l'on a :

$$\begin{cases} y(t_i + h) \cong y_2 + C_1 \cdot h^{p+1} + O(h^{p+2}) \\ y(t_i + h) \cong y_2 + 2 \cdot C_1 \cdot \left( \frac{h}{2} \right)^{p+1} + O(h^{p+2}) \end{cases}$$

Si on pose :  $\epsilon = y_2 - y_1$  on a :

$$\epsilon \cong 2 \cdot C_1 \cdot (2^p - 1) \cdot \left( \frac{h}{2} \right)^{p+1} + O(h^{p+2})$$

de sorte que

$$\frac{\varepsilon}{2^p - 1} \cong 2 \cdot C_1 \cdot \left(\frac{h}{2}\right)^{p+1} + O(h^{p+2})$$

est une estimation de l'erreur que l'on commet en remplaçant  $y(t_i + h)$  par  $y_2$ .

De plus, on a :

$$y(t_i + h) \cong \left(y_2 + \frac{\varepsilon}{2^p - 1}\right) + O(h^{p+2})$$

c'est-à-dire que  $\left(y_2 + \frac{\varepsilon}{2^p - 1}\right)$  est une meilleure approximation de  $y(t_i + h)$ .

D'autre part, comme  $|\varepsilon|$  est proportionnel à  $h^{p+1}$ , en notant  $\varepsilon_1$  l'approximation de l'erreur correspondante à un autre pas  $h_1$ , on aura :

$$\frac{|\varepsilon_1|}{|\varepsilon|} = \left(\frac{h_1}{h}\right)^{p+1}$$

Si pour  $h = h\_Initial$ , on a  $|\varepsilon| = |y_2 - y_1| > \varepsilon_0$ , cela signifie que le pas  $h$  ne convient pas et il faudra alors le remplacer par le pas  $h\_Final$  qui va donner une erreur de l'ordre de  $\varepsilon_0$ . En définitive, on prendra :

$$h\_Final = h\_Initial \cdot \left(\frac{\varepsilon_0}{|\varepsilon|}\right)^{\frac{1}{p+1}}$$

c'est-à-dire qu'on diminue le pas.

Dans le cas contraire, on augmente le pas en prenant :  $h\_Final = 2 \cdot h\_Initial$

### 4.3 Programmation structurée de la méthode de Runge-Kutta d'ordre 4 avec contrôle du pas

On note  $h_0$  la valeur maximale autorisée du pas et la précision minimale  $\varepsilon_0$  est supposée donnée en constante.

On décrit la procédure dans le cas d'une équation différentielle d'ordre 1 à une variable, les trois autres cas s'en déduisant facilement.

*PROCEDURE RK4\_11\_Contrôle(Entrée f : Fonction ; t\_initial, y\_Initial, h0 : Réel ;  
Sortie t\_Final, y\_Final, Entrée\_Sortie h\_Initial : Réel) ;*

*Début*

*Répéter*

*RK4\_11(f,t\_Initial,y\_Initial,h\_Initial/2,t\_Final,y1) ;  
RK4\_11(f,t\_Initial + h\_Initial/2,y1,h\_Initial/2,t\_Final,y2) ;  
RK4\_11(f,t\_Initial,y\_Initial,h\_Initial,t\_Final,y1) ;*

$$\varepsilon = \frac{|y_2 - y_1|}{\varepsilon_0} ;$$

*Si  $\varepsilon > 1$*

$$\text{Alors } h\_Initial = \frac{h\_Initial}{\varepsilon^{\frac{1}{p+1}}}$$

*Sinon Si  $h\_Initial < h_0/2$  Alors  $h\_Initial = h\_Initial \cdot 2$  ;*

*Jusqu'à ( $\varepsilon \leq 1$ ) ;*

$$y\_Final = \frac{16 \cdot y_2 - y_1}{15};$$

Fin ;

## 5. Problèmes avec conditions aux limites. Méthode du tir

### 5.1 Introduction : position du problème

Dans les problèmes que nous avons considérés jusqu'à présent, les conditions aux limites étaient des « conditions initiales », c'est-à-dire que pour un système différentiel  $y' = f(t, y)$  on connaît toutes les valeurs initiales  $y_j(a)$ , pour  $j = 1, 2, \dots, p$ .

Mais on peut aussi considérer des problèmes différentiels où les *conditions aux limites* portent sur les deux bornes de l'intervalle, ou sur ces deux bornes et des points intérieurs.

Par exemple, la distribution de température à l'intérieur d'un cylindre de rayon 1, en mode stationnaire, est décrite par l'équation différentielle, avec conditions aux limites :

$$\begin{cases} y''(t) = -\frac{y'(t)}{t} - \alpha \cdot e^{y(t)} & (t \in ]0, 1[) \\ y'(1) = y(1) = 0 \end{cases}$$

où  $\alpha \in ]0, 0.8]$  est une constante proportionnelle à l'inverse de la conductivité.

Dans ce paragraphe, on s'occupera de problèmes aux limites du type :

$$(1) \quad \begin{cases} y'(t) = f(t, y(t)) & (t \in ]a, b[) \\ y_j(a) = y_{j,a} & (j \in I_1) \\ y_j(b) = y_{j,b} & (j \in I_2) \end{cases}$$

où  $I_1 \cup I_2 = \{1, \dots, p\}$  et  $I_1, I_2$  disjoints.

On notera  $p_1$  le nombre d'éléments de  $I_1$  et  $p_2$  celui de  $I_2$ . De plus, les indices dans  $I_2$  seront notés :  $j_1 < j_2 < \dots < j_{p_2}$ .

Une équation différentielle d'ordre  $q \geq 2$ , où les conditions aux limites portent sur les deux bornes de l'intervalle peut toujours se ramener à (1) par le procédé classique.

L'idée de la méthode du tir, pour résoudre (1) est de démarrer en posant les conditions initiales manquantes de façon arbitraire, c'est-à-dire qu'on se donne un vecteur  $\alpha = (\alpha_1, \dots, \alpha_{p_2})$ , et on pose :

$$y_{j_k}(a) = \alpha_k \quad (k = 1, \dots, p_2)$$

On résout alors le système différentiel par une méthode à un pas, ce qui va donner une valeur en  $b$  de la solution que l'on note  $y(b, \alpha)$ .

Le problème est alors de trouver  $\alpha \in \mathbb{R}^{p_2}$  solution du système de  $p_2$  équations à  $p_2$  inconnues :

$$e_k(\alpha) = y_{j_k}(b, \alpha) - y_{j_k,b} = 0 \quad (k = 1, \dots, p_2)$$

*Remarque : Cas des systèmes linéaires* — Dans le cas où le système est linéaire, la méthode peut être simplifiée. En prenant  $p$  valeurs distinctes  $\alpha^1, \dots, \alpha^p$  de sorte que les valeurs initiales  $y^1(a), \dots, y^p(a)$  correspondantes soient linéairement indépendantes, on construit  $p$  solutions linéairement indépendantes  $y^1, \dots, y^p$  et la solution cherchée va s'écrire :

$$y(t) = \sum_{i=1}^p \lambda_i \cdot y^i(t)$$

les constantes  $\lambda_i$  étant uniquement déterminées par les conditions initiales  $y_{j,a}$  et  $y_{k,b}$ , pour  $j$  dans  $I_1$  et  $k$  dans  $I_2$ .

## 5.2 La méthode du tir

### 5.2.1 Rappels

Il s'agit donc de résoudre le problème (1).

On démarre avec  $\alpha$  choisi arbitrairement dans  $\mathbb{R}^{p^2}$ , ce qui va donner une erreur sur  $b$ ,  $e(\alpha) = \left( (e_k(\alpha)) \right)_{1 \leq k \leq p_2}$ , définie par :

$$e_k(\alpha) = y_{j_k}(b, \alpha) - y_{j_k,b} \quad (k = 1, \dots, p)$$

En utilisant la méthode de Newton-Raphson, on construit donc la suite de vecteurs  $(\alpha^k)$ , de la façon suivante :

$$\begin{cases} \alpha^0 = \alpha \\ \alpha^{k+1} = \alpha^k - \delta^k \quad (k \geq 0) \end{cases}$$

où  $\delta^k$  est solution de :

$$A_k \cdot \delta^k = e(\alpha^k)$$

avec  $A_k = \left( (a_{i,j}) \right)_{1 \leq i,j \leq p_2}$ ,  $a_{i,j} = \frac{e_i(\beta^{k,j}) - e_i(\alpha^k)}{h}$  et

$\beta^{k,j} = (\alpha_1^k, \dots, \alpha_{j-1}^k, \alpha_j^k + h, \alpha_{j+1}^k, \dots, \alpha_j^n)$ ,  $h$  étant un pas assez petit.

### 5.2.2 Programmation structurée de la méthode du tir

Pour reconnaître les indices, on introduit un tableau de booléens  $B$  défini par :

$$\begin{cases} B_j = \text{Faux} \text{ si } j \in I_1 \text{ (i. e. } y_j(a) \text{ connue)} \\ B_j = \text{Vrai} \text{ si } j \in I_2 \text{ (i. e. } y_j(a) \text{ inconnue)} \end{cases}$$

Les  $y_j(a)$ , pour  $j = 1, \dots, p$ , c'est-à-dire ceux qui sont connus et ceux qui sont donnés arbitrairement, seront stockés dans la variable  $y\_Initial$  de type `Variable_2` (tableau de réels numérotés de 1 à  $p$ ). Les  $y_j(b)$  connus seront stockés dans la variable  $y\_Final$  de type `Variable_2`.

On suppose que l'on dispose d'une procédure de résolution d'un système de  $p$  équations différentielles définie par :

*PROCEDURE RésolSystDiff(Entrée  $f$  : SystemeDiff ;  $a, b$  : Réel ;*



$y\_Initial : Variable\_2 ; p : Entier ; Sortie y : Matrice) ;$

où  $y_{j,k}$  est une approximation de  $y_j(t_k)$ , pour  $j$  allant de 1 à  $p$  et  $k$  de 0 à  $n$ , avec

$$t_k = a + k \cdot \frac{b-a}{n}.$$

En particulier,  $y_{j,n}$  est l'approximation trouvée de  $y_j(b)$ .

On suppose également que l'on dispose d'une procédure de résolution d'un système linéaire définie par :

*PROCEDURE SystLin(Entrée  $n : Entier ; A : Matrice ; b : Variable\_2 ;$   
Sortie  $x : Variable\_2$ );*

On se reportera au chapitre sur l'analyse numérique linéaire pour l'écriture d'une telle procédure.

On définit tout d'abord une procédure de calcul de l'erreur commise à l'issue d'un tir. Soit :

*PROCEDURE Calcul\_Erreur(Entrée  $y\_Final : Variable\_2 ;$   
 $B : TableauBooléens ; p : Entier ; y : Matrice ;$   
Sortie  $e : Variable\_2$ );*

*Début*

*$i = 0 ;$*

*Pour  $j$  allant de 1 à  $p$  Faire*

*Début*

*Si  $B_j$*

*Alors Début*

*$i = i + 1 ;$*

*$e_i = y_{j,n} - y\_Final_j ;$*

*Fin ;*

*Fin ;*

*Fin ;*

On écrit ensuite une procédure de correction du tir, qui réalise une itération dans l'algorithme décrit ci-dessus, c'est-à-dire qu'on calcule  $\alpha^{k+1}$  en fonction de  $\alpha^k$ . Ce qui revient à remplacer  $y\_Initial$  par sa nouvelle valeur définie par :

$$y\_Initial_j = y\_Initial_j - \delta_j \quad (j \in I_2)$$

où  $\delta$  est solution du système :  $M \cdot \delta = e$ ,  $e$  étant l'erreur précédemment calculée.

*PROCEDURE Correction\_Tir(Entrée  $f : SystemeDiff ; a, b : Réel ; p : Entier ;$   
 $y\_Final : Variable\_2 ; B : TableauBooléens ;$   
Sortie  $y\_Initial : Variable\_2 ; y : Matrice ; Arret : Booléen$ );*

*Début*

*$p_2 = 0 ;$*

*Pour  $j$  allant de 1 à  $p$  Faire*

*Début*

*Si  $B_j$  Alors  $p_2 = p_2 + 1 ;$*

*Fin ;*

*ResolSystDiff( $f, a, b, y\_Initial, p, y$ );*

*Calcul\_Erreur( $y\_Final, B, p, y, e$ );*

*Arret = ( $\|e\| < \epsilon$ );*

*Si Arret*

*Alors Sortir de la procédure ;*

```

{ Formation de la matrice jacobienne }
k = 0 ;
Pour j Allant de 1 à p Faire
Début
  Si Bj
  Alors Début
    k = k + 1 ;
    Aux = y_Initialj ;
    y_Initialj = y_Initialj + ε ;
    ResolSystDiff(f,a,b,y_Initial,p,y) ;
    Calcul_Erreur(y_Final,B,p,y,e1) ;

    Pour i Allant de 1 à p2 Faire  $M_{i,k} = \frac{eI_i - e_i}{\epsilon}$  ;

    y_Initialj = Aux ;
  Fin ;
Fin ;
SystLin(p2,M,e,δ) ;
{ Calcul de la nouvelle valeur de y_Initial }
k = 0 ;
Pour j Allant de 1 à p Faire
Début
  Si Bj
  Alors Début
    k = k + 1 ;
    y_Initialj = y_Initialj - δj ;
  Fin ;
Fin ;
Fin ;

```

La procédure finale de résolution d'un système de p équations différentielles d'ordre 1 va alors s'écrire :

```

PROCEDURE Résolution_Tir(Entrée f : SystemeDiff ; a, b : Réel ; p : Entier ;
                          y_Initial, y_Final : Variable_2 ; B : TableauBooléens ;
                          Sortie y : Matrice) ;
Début
  Répéter
    Correction_Tir(f,a,b,p,y_Final,B,y_Initial,y,Arret) ;
  Jusqu'à Arret ;
  ResolSystDiff(f,a,b,y_Initial,p,y) ;
Fin ;

```

## 6. Un exemple d'application. Mouvement de translation d'un corps sphérique pesant dans un fluide au repos

### 6.1 Position du problème et notations

On considère de nouveau le problème étudié au paragraphe (2.7) du chapitre sur l'approximation et l'interpolation.

On a un corps sphérique que l'on laisse tomber à l'instant  $t = 0$ , sans vitesse initiale, dans un fluide au repos et on veut calculer la distance à l'origine  $z(t)$  et la vitesse  $v(t) = z'(t)$  du corps à chaque instant  $t$ .

On utilisera les notations suivantes :

- $m$  est la *masse* du corps sphérique ;
- $g$  est l'*accélération due à la pesanteur*, dans la direction de l'axe  $Oz$  ;
- à l'instant  $t = 0$ , le centre de la sphère est à  $z = 0$  et la vitesse est  $v = 0$  ;
- $d$  est le *diamètre* de la sphère ;
- $\rho$  est la *densité volumique de la sphère* ;
- $\rho_f$  est la *densité volumique du fluide* ;
- $\nu$  est la *viscosité cinématique du fluide* ;
- $C_d$  est le *coefficient de traînée* de la sphère en régime stationnaire ;

*Remarque* — Le coefficient  $C_d$  est lié au *nombre de Reynolds* de la sphère lui-même relié à la vitesse de cette dernière.

Dans le paragraphe (2.7) du chapitre 3, on a donné un programme qui permet de calculer  $C_d$  en fonction de  $Re$ .

## 6.2 Equations du mouvement de la sphère

On sait que dans le vide la loi du mouvement de la sphère est :

$$z(t) = \frac{1}{2} \cdot g \cdot t^2$$

Dans un fluide plus général, on doit prendre en compte les forces suivantes :

- la *poussée d'Archimède*, égale au poids du fluide déplacé par la particule, soit :

$$-m_f g = -\frac{1}{6} \cdot \pi \cdot d^3 \cdot \rho_f g \quad \uparrow -m_f g$$

- la *force induite par l'accélération* :

quand un corps plongé dans un liquide au repos est brusquement mis en mouvement avec une accélération  $\frac{dv}{dt}$ , il en résulte une force de réaction :

$$-\lambda \cdot \frac{dv}{dt} = -\frac{1}{2} \cdot m_f \cdot \frac{dv}{dt} \quad \uparrow -\lambda \cdot \frac{dv}{dt}$$

où  $\lambda$  est la « *masse virtuelle* » ou « *masse induite par accélération* ».

- la *force de réaction due à la viscosité du fluide* qui s'écrit sous la forme :

$$-C_p \cdot \frac{1}{2} \cdot \rho_f \cdot v \cdot |v| \cdot A_p$$

où

$$A_p = \frac{1}{4} \cdot \pi \cdot d^2$$

et  $C_d$  est le « *coefficient de traînée* ».

La loi de Newton nous permet d'écrire l'équation du mouvement de la particule :

$$m \cdot \frac{dv}{dt} = m \cdot g - m_f \cdot g - \frac{1}{2} \cdot m_f \cdot \frac{dv}{dt} - \frac{1}{2} \cdot \rho_f \cdot v \cdot |v| \cdot \frac{\pi}{4} \cdot d^2 \cdot C_d$$

soit l'équation différentielle du second ordre :

$$(2) \quad \left( m + \frac{1}{2} \cdot m_f \right) \cdot \frac{d^2 z}{dt^2} = (m - m_f) \cdot g - \frac{\pi}{8} \cdot \rho_f \cdot \frac{dz}{dt} \cdot \left| \frac{dz}{dt} \right| \cdot d^2 \cdot C_d \left( \frac{dz}{dt} \right)$$

encore équivalente au système de deux équations différentielles d'ordre 1 :

$$(3) \quad \begin{cases} v(t) = \frac{dz}{dt}(t) \\ \left( m + \frac{1}{2} \cdot m_f \right) \cdot \frac{dv}{dt} = (m - m_f) \cdot g - \frac{\pi}{8} \cdot \rho_f \cdot v \cdot |v| \cdot d^2 \cdot C_d(v) \end{cases}$$

*Remarque* — Le terme  $m + \frac{1}{2} \cdot m_f$  exprime que la particule se comporte comme si sa masse avait été augmentée (dans un mouvement accéléré ou décéléré).

En écrivant que :

$$m = \frac{1}{6} \cdot \pi \cdot d^3 \cdot \rho$$

et en posant :

$$\begin{cases} A = 1 + \frac{1}{2} \cdot \frac{\rho_f}{\rho} \\ B = \left( 1 - \frac{\rho_f}{\rho} \right) \cdot g \\ C = \frac{3}{4} \cdot \frac{\rho_f}{d \cdot \rho} \end{cases}$$

le système d'équations différentielles va s'écrire :

$$\begin{cases} v = \frac{dz}{dt} \\ \frac{dv}{dt} = \frac{1}{A} \cdot (B - C \cdot v \cdot |v| \cdot C_d(v)) \end{cases}$$

Le programme du paragraphe (2.7) du chapitre sur l'interpolation et l'approximation peut donc se compléter comme suit.

```
with TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH1,
COMMON_MATH3, MATH3, CURVE ;
use TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH1,
COMMON_MATH3, MATH3, CURVE ;
procedure BILLE is

NBR_INTERVALLES : constant NATURAL := 7 ;
  -- Nombre maximale d'intervalles dans [0,50000]
  -- Dans chaque intervalle on a des mesures experimentales de Re et Cd
NBR_MESURES : constant NATURAL := 6 ;
  -- Nombre maximale de mesure des Re et Cd dans chaque intervalle
type MATRICE_MESURES is array( POSITIVE range 1..NBR_INTERVALLES,
POSITIVE range 1..NBR_MESURES) of FLOAT ;
  -- Dans l'intervalle N° j, les mesures de Re et Cd sont notees
```

```

-- Re(j,i) et Cd(j,i) ou 1 <= i <= NBR_MESURES
type VECTEUR_COEFF is array(POSITIVE range 1..NBR_INTERVALLES,
                           NATURAL range 0..2) of FLOAT ;
-- Dans l'intervalle N° j, les coefficients de Morsi sont notes
-- K(j,i), ou i = 1, 2, 3.
type VECTEUR_ENTIER is array(POSITIVE range 1..NBR_INTERVALLES)
  of NATURAL ;
-- Dans l'intervalle N° j, le nombre de mesures est note n(j)
Re, Cd : MATRICE_MESURES ;
nn : VECTEUR_ENTIER ;
K : VECTEUR_COEFF ;
Rho, Rhof, g, Nu, d, t0, tMax, z0, v0, A, B, C : FLOAT ;

procedure DONNEES_DU_FICHIER(
  t0, tMax, z0, v0, Rho, Rhof, Nu, g, d, A, B, C : in out FLOAT ;
  nn : in out VECTEUR_ENTIER ;
  Re, Cd : out MATRICE_MESURES ; K : out VECTEUR_COEFF) is
Nom : STRING(1..50) ;
FichierDonnees : FILE_TYPE ;
RhoBarre : FLOAT ;
begin
  PUT_LINE("  Donnees du fichier") ;
  LIRE_FICHIER(FichierDonnees,Nom) ;
  GET(FichierDonnees,t0) ;
  GET(FichierDonnees,tMax) ;
  GET(FichierDonnees,z0) ;
  GET(FichierDonnees,v0) ;
  GET(FichierDonnees,Rho) ;
  GET(FichierDonnees,Rhof) ;
  GET(FichierDonnees,Nu) ;
  GET(FichierDonnees,g) ;
  GET(FichierDonnees,d) ;
  SKIP_LINE(FichierDonnees) ;
  for j in 1..NBR_INTERVALLES
  loop
    GET(FichierDonnees,nn(j)) ;
    SKIP_LINE(FichierDonnees) ;
    for i in 1..nn(j)
    loop
      GET(FichierDonnees,Re(j,i)) ;
    end loop ;
    SKIP_LINE(FichierDonnees) ;
    for i in 1..nn(j)
    loop
      GET(FichierDonnees,Cd(j,i)) ;
    end loop ;
    SKIP_LINE(FichierDonnees) ;
  end loop ;
  for i in 1..NBR_INTERVALLES
  loop
    for j in 0..2
    loop
      GET(FichierDonnees,K(i,j)) ;
    end loop ;
  end loop ;
end ;

```

```

        end loop ;
        SKIP_LINE(FichierDonnees) ;
    end loop ;
    CLOSE(FichierDonnees) ;
    RhoBarre := Rhof/Rho ;
    A := 1.0 + 0.5*RhoBarre ;
    B := g*(1.0 - RhoBarre) ;
    C := 0.75*RhoBarre/d ;
end DONNEES_DU_FICHER ;

procedure AFFICHAGE_DONNEES is
begin
    CLRSCR ;
    PUT_LINE(IMP,"                Donnees du probleme") ; NEW_LINE(IMP) ;
    PUT_LINE(IMP,"Conditions initiales :") ;
    PUT(IMP,"                t0 = ") ;
    PUT(IMP,t0,6,3,0) ; NEW_LINE(IMP) ;
    PUT(IMP,"                tMax = ") ;
    PUT(IMP,tMax,6,3,0) ; NEW_LINE(IMP) ;
    PUT(IMP,"                z0 = ") ;
    PUT(IMP,z0,6,3,0) ; NEW_LINE(IMP) ;
    PUT(IMP,"                v0 = ") ;
    PUT(IMP,v0,6,3,0) ; NEW_LINE(IMP,2) ;
    PUT_LINE(IMP,"Caracteristiques de la bille et du fluide :") ;
    NEW_LINE(IMP) ;
    PUT(IMP,"                Rho = ") ;
    PUT(IMP,Rho,6,3,0) ; NEW_LINE(IMP) ;
    PUT(IMP,"                Rhof = ") ;
    PUT(IMP,Rhof,6,3,0) ; NEW_LINE(IMP) ;
    PUT(IMP,"                Nu = ") ;
    PUT(IMP,Nu,6,6,0) ; NEW_LINE(IMP) ;
    PUT(IMP,"                g = ") ;
    PUT(IMP,g,6,3,0) ; NEW_LINE(IMP) ;
    PUT(IMP,"                d = ") ;
    PUT(IMP,d,6,3,0) ; NEW_LINE(IMP) ;
    PAUSE ; CLRSCR ;
end AFFICHAGE_DONNEES ;

function Cd_j(K : in VECTEUR_COEFF ; j : in INTEGER ; x : in FLOAT)
    return FLOAT is
-- Evaluation du polynome definissant la fonction Cd sur l'intervalle j
begin
    return K(j,0) + K(j,1)*x + K(j,2)*x*x ;
end Cd_j ;

function F(v : in FLOAT) return FLOAT is
r, C_d : FLOAT ;
j : INTEGER ;
begin
    r := v*d/Nu ;
    if r = 0.0 then
        return (B - C*K(1,2)*Nu/(d*d))/A ;
    else

```

```

    j := 0 ;
    loop
        j := j + 1 ;
        exit when (r >= Re(j,1)) and (r < Re(j,nn(j))) ;
    end loop ;
    C_d := Cd_j(K,j,1.0/r) ;
    return (B - C*v*ABS(v)*C_d)/A ;
end if ;
end F ;

procedure RK4(AFFICHER : in BOOLEAN ; Cz, Cv : in out COURBE) is
    Pas, z_k1, z_k2, z_k3, z_k4, v_k1, v_k2, v_k3, v_k4 : FLOAT ;
begin
    Cz.uMin := t0 ; Cz.uMax := tMax ;
    Cv.uMin := t0 ; Cv.uMax := tMax ;
    Cz.vMin := 1.0E+38 ; Cz.vMax := -1.0E+38 ;
    Cv.vMin := 1.0E+38 ; Cv.vMax := -1.0E+38 ;
    Cz.u(0) := t0 ;
    Cv.u(0) := t0 ;
    Cz.v(0) := z0 ;
    Cv.v(0) := v0 ;
    Pas := (tMax - t0)/FLOAT(Cz.n) ;
    for i in 1..Cz.n
    loop
        Cz.u(i) := Cz.u(i - 1) + Pas ;
        Cv.u(i) := Cz.u(i) ;
        z_k1 := Pas*Cv.v(i - 1) ;
        v_k1 := Pas*F(Cv.v(i - 1)) ;
        z_k2 := Pas*(Cv.v(i - 1) + 0.5*v_k1) ;
        v_k2 := Pas*F(Cv.v(i - 1) + 0.5*v_k1) ;
        z_k3 := Pas*(Cv.v(i - 1) + 0.5*v_k2) ;
        v_k3 := Pas*F(Cv.v(i - 1) + 0.5*v_k2) ;
        z_k4 := Pas*(Cv.v(i - 1) + v_k3) ;
        v_k4 := Pas*F(Cv.v(i - 1) + v_k3) ;
        Cz.v(i) := Cz.v(i - 1) + (z_k1 + 2.0*z_k2 + 2.0*z_k3 + z_k4)/6.0 ;
        Cz.vMin := MIN(Cz.vMin,Cz.v(i)) ;
        Cz.vMax := MAX(Cz.vMax,Cz.v(i)) ;
        Cv.v(i) := Cv.v(i - 1) + (v_k1 + 2.0*v_k2 + 2.0*v_k3 + v_k4)/6.0 ;
        Cv.vMin := MIN(Cv.vMin,Cv.v(i)) ;
        Cv.vMax := MAX(Cv.vMax,Cv.v(i)) ;
        if (AFFICHER) then
            PUT(IMP,"t = ") ;
            PUT(IMP,Cz.u(i),6,3,0) ;
            PUT(IMP," z = ") ;
            PUT(IMP,Cz.v(i),6,3,0) ;
            PUT(IMP," v = ") ;
            PUT(IMP,Cv.v(i),6,3,0) ;
            NEW_LINE(IMP) ;
            if (i mod 20 = 0) then
                PAUSE ;
            end if ;
        end if ;
    end loop ;
end loop ;

```



```

end RK4 ;

procedure TRACE_POSITION_VITESSE is
tFenMin, tFenMax, zFenMin, zFenMax, vFenMin, vFenMax, tUnite, zUnite,
vUnite : FLOAT ;
Afficher : BOOLEAN ;
Cz, Cv : COURBE(200) ;
begin
  Afficher := LIRE_REPONSE("Voulez vous l'affichage des resultats? ") ;
  RK4(AFFICHER,Cz,Cv) ;
  tUnite := (tMax - t0)/20.0 ;
  zUnite := (Cz.vMax - Cz.vMin)/20.0 ;
  vUnite := (Cv.vMax - Cv.vMin)/20.0 ;
  tFenMin := Cz.uMin - tUnite ;
  tFenMax := Cz.uMax + tUnite ;
  zFenMin := Cz.vMin - zUnite ;
  zFenMax := Cz.vMax + zUnite ;
  vFenMin := Cv.vMin - vUnite ;
  vFenMax := Cv.vMax + vUnite ;
  INIT_GRAPHIQUE ;
  FENETRE(tFenMin,tFenMax,zFenMin,zFenMax) ;
  CADRE_GRAPHIQUE(
    X_MINIMUM,X_MINIMUM + ((X_MAXIMUM - X_MINIMUM)/2) - 1,
    Y_MINIMUM + MARGE,Y_MAXIMUM - MARGE) ;
  DEPLACE(Cz.u(0),Cz.v(0)) ;
  for i in 1..Cz.n
  loop
    TRACE(Cz.u(i),Cz.v(i)) ;
  end loop ;
  XY_AXES(0.0,0.0,tUnite,zUnite,TRUE) ;
  FENETRE(tFenMin,tFenMax,vFenMin,vFenMax) ;
  CADRE_GRAPHIQUE(
    X_MINIMUM + ((X_MAXIMUM - X_MINIMUM)/2) + 1,
    X_MAXIMUM, Y_MINIMUM + MARGE, Y_MAXIMUM - MARGE) ;
  DEPLACE(Cv.u(0),Cv.v(0)) ;
  for i in 1..Cv.n
  loop
    TRACE(Cv.u(i),Cv.v(i)) ;
  end loop ;
  XY_AXES(0.0,0.0,tUnite,vUnite,TRUE) ;
  TITRE("Graphe de z(t)                               Graphe de v(t)") ;
  SORTIE_GRAPHIQUE ;
end TRACE_POSITION_VITESSE ;

begin
  MODE_AFFICHAGE ;
  DONNEES_DU_FICHER(t0,tMax,z0,v0,Rho,Rhof,Nu,g,d,A,B,C,nn,Re,Cd,K) ;
  AFFICHAGE_DONNEES ;
  PUT_LINE("Resolution du systeme differentiel donnant z(t) et v(t)") ;
  NEW_LINE ;
  TRACE_POSITION_VITESSE ;
  CLOSE(IMP) ;
end BILLE ;

```

## 7. Exercices

### 7.1 Le modèle « proies-prédateurs »

On considère le modèle « proies-prédateurs » du paragraphe 1, avec  $r_1(t) = 2$  ;  $r_2(t) = 1$  ;  $b_1(t) = 0.5$  ;  $b_2(t) = 0.25$  ; et les conditions initiales :  $x(0) = y(0) = 1$  (sous-entendu une unité).

On tracera les graphes d'évolution de  $x$  et  $y$  en fonction de  $t$  et de  $y$  en fonction de  $x$ .

### 7.2 Le problème des trois corps

On considère trois corps  $C_1$ ,  $C_2$  et  $C_3$  isolés dans l'univers et soumis à leurs attractions mutuelles. On veut déterminer le mouvement de deux de ces corps par rapport au troisième pris comme référence.

Dans un repère orthonormé d'origine  $C_1$ , on note  $(X_2, Y_2, Z_2)$  les coordonnées de  $C_2$  et  $(X_3, Y_3, Z_3)$  celles de  $C_3$ .

On note  $m_i$  la masse du corps  $C_i$  et  $d_{i,j}$  la distance de  $C_i$  à  $C_j$ . Enfin  $G$  désigne la constante de gravitation universelle.

1°) En utilisant la loi fondamentale de la dynamique et la loi de la gravitation universelle, montrer que les équations du mouvement de  $C_2$  et de  $C_3$  sont :

$$\left\{ \begin{array}{l} X_2''(t) = G \cdot m_3 \cdot \left\{ \frac{X_3(t) - X_2(t)}{(d_{2,3})^3} - \frac{X_3(t)}{(d_{1,3})^3} \right\} - G \cdot (m_1 + m_2) \cdot \frac{X_2(t)}{(d_{1,2})^3} \\ Y_2''(t) = G \cdot m_3 \cdot \left\{ \frac{Y_3(t) - Y_2(t)}{(d_{2,3})^3} - \frac{Y_3(t)}{(d_{1,3})^3} \right\} - G \cdot (m_1 + m_2) \cdot \frac{Y_2(t)}{(d_{1,2})^3} \\ Z_2''(t) = G \cdot m_3 \cdot \left\{ \frac{Z_3(t) - Z_2(t)}{(d_{2,3})^3} - \frac{Z_3(t)}{(d_{1,3})^3} \right\} - G \cdot (m_1 + m_2) \cdot \frac{Z_2(t)}{(d_{1,2})^3} \\ X_3''(t) = G \cdot m_2 \cdot \left\{ \frac{X_2(t) - X_3(t)}{(d_{2,3})^3} - \frac{X_2(t)}{(d_{1,2})^3} \right\} - G \cdot (m_1 + m_3) \cdot \frac{X_3(t)}{(d_{1,3})^3} \\ Y_3''(t) = G \cdot m_2 \cdot \left\{ \frac{Y_2(t) - Y_3(t)}{(d_{2,3})^3} - \frac{Y_2(t)}{(d_{1,2})^3} \right\} - G \cdot (m_1 + m_3) \cdot \frac{Y_3(t)}{(d_{1,3})^3} \\ Z_3''(t) = G \cdot m_2 \cdot \left\{ \frac{Z_2(t) - Z_3(t)}{(d_{2,3})^3} - \frac{Z_2(t)}{(d_{1,2})^3} \right\} - G \cdot (m_1 + m_3) \cdot \frac{Z_3(t)}{(d_{1,3})^3} \end{array} \right.$$

2°) On pose  $y_1 = X_2$ ,  $y_2 = X_2'$ , ...,  $y_{11} = Z_3$ ,  $y_{12} = Z_3'$ .

Ecrire le système différentiel que vérifient les fonctions  $y_1, \dots, y_{12}$ .

3°) Ecrire un programme de résolution de ce système en utilisant la méthode de Runge-Kutta d'ordre 4.

On pourra utiliser les valeurs numériques suivantes correspondantes au système {Terre, Lune, Vaisseau spatial} :  $m_1 = 5.98 \cdot 10^{24}$  Kg ;  $m_2 = 7.35 \cdot 10^{22}$  Kg ;  $m_3 = 10\ 000$  Kg ;  $X_2(0) = 383\ 000$  Km ;  $X'_2(0) = 0$  ;  $Y_2(0) = 0$  ;  $Y'_2(0) = 1.025$  ;  $Z_2(0) = 0$  ;  $Z'_2(0) = 0$  ;  $X_3(0) = 1197$  Km ;  $X'_3(0) = 8.49$  Km/s ;  $Y_3(0) = -9928$  ;  $Y'_3(0) = -2.455$  ;  $Z_3(0) = 0$  ;  $Z'_3(0) = 0.1$  ;  $G = 6.672 \cdot 10^{-11}$  m<sup>3</sup>·Kg<sup>-1</sup>·s<sup>-2</sup> (constante de gravitation).

### 7.3 Méthode de Stormer-Cowell

Dans ce problème, on étudie la méthode de Stormer-Cowell, qui est une méthode multipas explicite de résolution d'une équation différentielle du second ordre où  $y'(x)$  ne figure pas.

On considère l'équation différentielle, avec conditions initiales :

$$(1) \quad \begin{cases} y''(t) = f(x, y(x)) & (a < x < b) \\ y(a) = y_a \text{ et } y'(a) = y'_a \end{cases}$$

où  $y_a$  et  $y'_a$  sont donnés.

On pose  $h = \frac{b-a}{n+1}$  et pour  $i = 0, 1, \dots, n+1$ ,  $x_i = a + i \cdot h$ .

On veut alors calculer des valeurs approchées  $y_i$  des  $y(x_i)$  pour  $i = 1, \dots, n$ .

1°) Montrer que si  $y$  définie sur un intervalle  $I$  est assez régulière, alors pour tout  $x$  dans  $I$  et  $h$  assez petit, on a :

$$(2) \quad y''(x) = \frac{y(x+h) - 2 \cdot y(x) + y(x-h)}{h^2} + O(h^2)$$

2°) A l'étape  $i$  du calcul, on suppose que l'on connaît des valeurs approchées  $y_{i-1}$  de  $y(x_{i-1})$  et  $y_i$  de  $y(x_i)$ .

(a) Dédurre de (2) une valeur approchée  $y_{i+1}$  de  $y(x_{i+1})$ .

(b) En déduire un algorithme de résolution de (1).

3°) (a) Donner la valeur de  $y''(a)$ .

(b) En déduire une approximation d'ordre 2,  $y_1$  de  $y(x_1)$ .

4°) Appliquer la méthode en question au problème :

$$(3) \quad \begin{cases} y''(x) = y(x) & (0 < x < 1) \\ y(0) = 1, y'(0) = 1 \end{cases}$$

On prendra  $h = 0.1$ .

Pour tester le comportement de ce schéma numérique sur de grands intervalles de temps, on l'applique au problème :

$$(4) \quad \begin{cases} y''(x) = -y(x) & (x > 0) \\ y(0) = \alpha, y'(0) = \beta \end{cases}$$

On se donne un pas  $h > 0$  et on note  $x_n = n \cdot h$ , pour tout  $n$  dans  $\mathbb{N}$ .

5°) Ecrire la relation de récurrence vérifiée par la suite  $(y_n)_{n \in \mathbb{N}}$ .

6°) Pour quelles valeurs de  $q$  dans  $\mathbb{C}$ , la suite définie par  $u_n = q^n$ , vérifie-t-elle cette récurrence ?

7°) En notant  $q_1$  et  $q_2$  les racines de l'équation obtenue en 2°), montrer que la suite  $(y_n)$  est bornée si, et seulement si  $|q_i| \leq 1$ , pour  $i = 1, 2$ .

8°) En déduire que la solution numérique trouvée ne peut être bornée que si  $h$  est inférieur à une valeur critique que l'on précisera.

#### 7.4 La méthode du tir pour les équations différentielles d'ordre 2

On considère l'équation différentielle :

$$(1) \quad y'' = f(x, y, y') \quad (x \in [a, b])$$

avec les conditions aux limites :

$$y(a) = A, \quad y(b) = B$$

On suppose que cette équation admet une unique solution  $y$  deux fois continûment dérivable.

En posant, de façon arbitraire,  $y'(a) = \alpha$ , on note  $y(\alpha, x)$  la solution de (1) avec les conditions initiales :

$$y(a) = A, \quad y'(a) = \alpha$$

1°) Expliquer comment construire une suite  $(\alpha_n)_{n \geq 0}$  permettant « d'ajuster le tir », c'est-à-dire, sous de bonnes conditions, qui converge vers  $\alpha$  solution de  $y(\alpha, b) - B = 0$ .

Si cette suite converge vers  $\alpha$ , la solution  $y(\alpha, x)$  est alors solution de problème posé.

2°) Ecrire les 2 premières itérations de la méthode pour :

$$f(x, y, y') = \frac{2}{x^3}$$

sur  $[1, 2]$  avec  $y(1) = 1$ ,  $y(2) = 0.5$  et  $\alpha_0 = -0.5$ .

3°) Ecrire, une procédure de résolution d'un problème du type (1) par la méthode du tir.

## 8. Programmation Ada

Ci-dessous on décrit quatre paquetages génériques permettant de résoudre une équation différentielle d'ordre 1, un système de  $p$  équations différentielles d'ordre 1, une équation différentielle d'ordre  $q$  et un système de  $p$  équations différentielles d'ordre  $q$ .

En commentaire on indique comment utiliser les fonctions définissant les équations différentielles.

### 8.1 Spécification du paquetage EQUADIFF\_11\_GENERIQUE

Une équation différentielle d'ordre 1 s'écrit sous la forme :

$$y' = f(x, y) \quad x \text{ dans } (a, b)$$

$$y(a) \text{ donné dans } R$$

```
with COMMON_MATRIX ;
use COMMON_MATRIX ;
```

```

generic
  with function f11(x, y : in FLOAT) return FLOAT ;
package EQUADIFF_11_GENERIQUE is
procedure RESOLUTION_RK4_11(a, b, ya : in FLOAT ;
  AFFICHER : in BOOLEAN ; x, y : in out VECTEUR) ; -- § 3.4
procedure RESOLUTION_RK4_11_CONTROLE(a, b, ya : in FLOAT ;
  AFFICHER : in BOOLEAN ; x, y : in out VECTEUR) ; -- § 4.3
end EQUADIFF_11_GENERIQUE ;

```

## 8.2 Spécification du paquetage EQUADIFF\_P1\_GENERIQUE

Un système de p équations différentielles d'ordre 1 s'écrit sous la forme :

$$\begin{aligned}
 y(1)' &= f(1)(x, y(1), \dots, y(p)) \\
 &\dots\dots\dots x \text{ dans } (a, b) \\
 y(p)' &= f(p)(x, y(1), \dots, y(p)) \\
 y(1)(a), \dots, y(p)(a) &\text{ donnés dans } R
 \end{aligned}$$

```

with COMMON_MATRIX, MATH2 ;
use COMMON_MATRIX, MATH2 ;
generic
  with function fp1(j : in POSITIVE ; x : in FLOAT ; y : VARIABLE_2)
return FLOAT ;
package EQUADIFF_P1_GENERIQUE is
procedure RESOLUTION_RK4_P1(a, b : in FLOAT ;
  Y_INITIAL : in VARIABLE_2 ; AFFICHER : in BOOLEAN ;
  x : in out VECTEUR ; y : in out MATRICE) ; -- § 3.5
procedure RESOLUTION_RK4_P1_CONTROLE(a, b : in FLOAT ;
  Y_INITIAL : in VARIABLE_2 ; AFFICHER : in BOOLEAN ;
  x : in out VECTEUR ; y : in out MATRICE) ; -- § 4.3
end EQUADIFF_P1_GENERIQUE ;

```

## 8.3 Spécification du paquetage EQUADIFF\_1Q\_GENERIQUE

Une équation différentielle d'ordre q s'écrit sous la forme :

$$\begin{aligned}
 y^{(q)} &= f(x, y(0), y(1), \dots, y(q-1)) \quad x \text{ dans } (a, b) \\
 y(0)(a), \dots, y(q-1)(a) &\text{ donnés dans } R \\
 y(j) &\text{ représente la dérivée d'ordre } j \text{ de } y \text{ pour } j = 0, \dots, q - 1.
 \end{aligned}$$

```

with COMMON_MATRIX, MATH2 ;
use COMMON_MATRIX, MATH2 ;
generic
  with function f1Q(x : in FLOAT ; y : VARIABLE_2) return FLOAT ;
package EQUADIFF_1Q_GENERIQUE is

procedure RESOLUTION_RK4_1Q(a, b : in FLOAT ;
  Y_INITIAL : in VARIABLE_2 ; AFFICHER : in BOOLEAN ;
  x, y : in out VECTEUR) ; -- § 3.6
procedure RESOLUTION_RK4_1Q_CONTROLE(a, b : in FLOAT ;
  Y_INITIAL : in VARIABLE_2 ; AFFICHER : in BOOLEAN ;
  x, y : in out VECTEUR) ; -- § 4.3
end EQUADIFF_1Q_GENERIQUE ;

```

### 8.4 Spécification du paquetage *EQUADIFF\_PQ\_GENERIQUE*

Un système différentielle de  $p$  équations d'ordre  $q$  s'écrit sous la forme :

$$y(1)^{(q)} = f(1)(x, y(1,0), \dots, y(1, q-1), y(2,0), \dots, y(2, q-1), \dots, y(p,0), \dots, y(p, q-1))$$

$$\dots \dots \dots$$

$$y(p)^{(q)} = f(p)(x, y(1,0), \dots, y(1, q-1), y(2,0), \dots, y(2, q-1), \dots, y(p,0), \dots, y(p, q-1))$$

Où  $y(j,i)$  représente la dérivée d'ordre  $i$  de  $y(j)$ .

FPQ( $j,x,y$ ) représente équation N°  $j$  avec :  $y(q*(j-1) + i + 1)$  représentant la dérivée d'ordre  $i$  de  $y(j)$ .

```
with COMMON_MATRIX, MATH2 ;
use COMMON_MATRIX, MATH2 ;
generic
  with function FPQ(j : in POSITIVE ; x : in FLOAT ; y : VARIABLE_2)
    return FLOAT ;
package EQUADIFF_PQ_GENERIQUE is
procedure RESOLUTION_RK4_PQ(a, b : in FLOAT ; Y_INITIAL : in MATRICE ;
  AFFICHER : in BOOLEAN ; -- § 3.7
  x : in out VECTEUR ; y : in out MATRICE) ;
procedure RESOLUTION_RK4_PQ_CONTROLE(a, b : in FLOAT ;
  Y_INITIAL : in MATRICE ; AFFICHER : in BOOLEAN ;
  x : in out VECTEUR ; y : in out MATRICE) ; -- 4.3
end EQUADIFF_PQ_GENERIQUE ;
```

### 8.5 Démonstration du paquetage *EQUADIFF\_11\_GENERIQUE*

```
with TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH1, MATH2,
  MATH3, COMMON_MATRIX, EQUADIFF_11_GENERIQUE ;
use TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, MATH1, MATH2,
  MATH3, COMMON_MATRIX ;

procedure DEM_RK11 is
N_MAX : constant INTEGER := 500 ;
n : INTEGER ;
a, b, ya : FLOAT ;
AFFICHER, CONTROLE_DU_PAS : BOOLEAN ;
x, y : VECTEUR(0..N_MAX) ;
X_UNITE, Y_UNITE, X_FEN_MIN, X_FEN_MAX, Y_FEN_MIN, Y_FEN_MAX : FLOAT ;
FF : FONCTION ;

function f(x, y : in FLOAT) return FLOAT is
begin
  return EVALUE(FF,x,y) ;
end f ;

package EQUADIFF_11_F is new EQUADIFF_11_GENERIQUE(f) ;
use EQUADIFF_11_F ;
begin
  MODE_AFFICHAGE ;
  GET_LINE(FF,"Entrez l'equation : y' = f(x,y) = ") ;
  NEW_LINE ;
  ENTRER_REEL(a,"Valeur de a : ") ;
```

```

ENTRER_REEL_BORNE(a,1.0E+6,b,"Valeur de b : ") ;
ENTRER_REEL(ya,"Valeur initiale en a : ") ;
NEW_LINE ;
NEW_LINE ;
AFFICHER := LIRE_REPONSE("Voulez-vous l'affichage des resultats? ") ;
CONTROLE_DU_PAS := LIRE_REPONSE("Methode avec controle du pas? ") ;
CLRSCR ;
if (CONTROLE_DU_PAS) then
    RESOLUTION_RK4_11_CONTROLE(a,b,ya,AFFICHER,x,y) ;
else
    RESOLUTION_RK4_11(a,b,ya,AFFICHER,x,y) ;
end if ;
PAUSE ;
CLRSCR ;
if LIRE_REPONSE("Voulez vous le trace de la solution trouvee? ") then
    n := x'LAST ;
    X_UNITE := (b - a)/20.0 ;
    X_FEN_MIN := a - X_UNITE ;
    X_FEN_MAX := b + X_UNITE ;
    Y_FEN_MIN := 1.0E+38 ;
    Y_FEN_MAX := -1.0E+38 ;
    for i in y'range
    loop
        Y_FEN_MIN := MIN(Y_FEN_MIN,y(i)) ;
        Y_FEN_MAX := MAX(Y_FEN_MAX,y(i)) ;
        if (x(i) > b) then
            n := i - 1 ;
            exit ;
        end if ;
    end loop ;
    Y_UNITE := (Y_FEN_MAX - Y_FEN_MIN)/20.0 ;
    Y_FEN_MIN := Y_FEN_MIN - Y_UNITE ;
    Y_FEN_MAX := Y_FEN_MAX + Y_UNITE ;
    INIT_GRAPHIQUE ;
    FENETRE_GRAPHIQUE(X_FEN_MIN,X_FEN_MAX,Y_FEN_MIN,Y_FEN_MAX) ;
    XY_AXES(a,Y_FEN_MIN + Y_UNITE,X_UNITE,Y_UNITE,TRUE) ;
    DEPLACE(x(x'FIRST),y(y'FIRST)) ;
    for i in x'FIRST + 1..n
    loop
        TRACE(x(i),y(i)) ;
    end loop ;
    TITRE("Graphe de la solution approchee y = y(x)") ;
    SORTIE_GRAPHIQUE ;
end if ;
CLOSE(IMP) ;
end DEM_RK11 ;

```

## 8.6 Démonstration du paquetage EQUADIFF\_P1\_GENERIQUE

```

with TEXT_IO, IIO, CRT, COMMON_MATH0, MATH0, MATH1, GRAPH,
    MATH2, MATH3, COMMON_MATRIX, EQUADIFF_P1_GENERIQUE ;
use TEXT_IO, IIO, CRT, COMMON_MATH0, MATH0, MATH1, GRAPH,
    MATH2, MATH3, COMMON_MATRIX ;

```



```

procedure DEM_RKP1 is
N_MAX: constant INTEGER := 500 ;
MAX_FONCTIONS : constant INTEGER := 10 ;
subtype INDICE_FCT is INTEGER range 1..MAX_FONCTIONS ;
type SYSTEME_DIFF is array(INDICE_FCT range <>) of FONCTION ;
a, b : FLOAT ;
AFFICHER, CONTROLE_DU_PAS : BOOLEAN ;
p : INTEGER ;
X_UNITE, Y_UNITE, X_FEN_MIN, X_FEN_MAX, Y_FEN_MIN, Y_FEN_MAX : FLOAT ;
x : VECTEUR(0..N_MAX) ;
n : INTEGER ;
begin
  MODE_AFFICHAGE ;
  PUT_LINE(" Le systeme d'equations est ecrit sous la forme :") ;
  NEW_LINE ;
  PUT_LINE("  y' = f (x,y , ...,y  )") ;
  PUT_LINE("    j   j   1       p") ;
  NEW_LINE ;
  PUT(" pour j = 1, ..., p, ou 1 <= p <= ") ;
  PUT(MAX_FONCTIONS) ; NEW_LINE ;
  ENTRER_ENTIER_BORNE(1,MAX_FONCTIONS,p,"Nombre d'equations : ") ;
  NEW_LINE ;
  declare
    FF : SYSTEME_DIFF(1..p) ;
    ya : VARIABLE_2(1..p) ;
    y : MATRICE(0..N_MAX,1..p) ;
    function f(j : in POSITIVE ; x : in FLOAT ; y : in VARIABLE_2)
      return FLOAT is
    begin
      return EVALUE(FF(j),x,y) ;
    end f ;
    package EQUADIFF_P1_F is new EQUADIFF_P1_GENERIQUE(F) ;
    use EQUADIFF_P1_F ;
  begin
    for j in YA'range
      loop
        PUT("Equation Numero ") ; PUT(j,2) ; NEW_LINE ;
        PUT("f(x,") ;
        for k in 1..p - 1
          loop
            PUT("y") ; PUT(k,2) ; PUT(",") ;
          end loop ;
        PUT("y") ; PUT(p,2) ; PUT_LINE(") = ") ;
        GET_LINE(ff(j)," ") ;
      end loop ;
      NEW_LINE ;
      PUT_LINE(" Bornes de l'intervalle") ;
      NEW_LINE ;
      ENTRER_REEL(a,"Valeur de a : ") ;
      ENTRER_REEL_BORNE(a,1.0E+6,b,"Valeur de b : ") ;
      for j in ya'range
        loop
          PUT("Valeur initiale y") ; PUT(j,2) ; PUT_LINE("(a)") ;

```

```

        ENTRER_REEL(ya(j), " y = ") ;
    end loop ;
    NEW_LINE ;
    AFFICHER := LIRE_REPONSE("Voulez-vous l'affichage des resultats? ") ;
    CONTROLE_DU_PAS := LIRE_REPONSE("Methode avec controle du pas? ") ;
    CLRSCR ;
    if (CONTROLE_DU_PAS) then
        RESOLUTION_RK4_P1_CONTROLE(a,b,ya,AFFICHER,x,y) ;
    else
        RESOLUTION_RK4_P1(a,b,ya,AFFICHER,x,y) ;
    end if ;
    PAUSE ; CLRSCR ;
if LIRE_REPONSE("Voulez vous le trace de la solution trouvee? ") then
    n := x'LAST ;
    INIT_GRAPHIQUE ;
    X_UNITE := (b - a)/20.0 ;
    X_FEN_MIN := a - X_UNITE ;
    X_FEN_MAX := b + X_UNITE ;
    for j in YA'range
    loop
        Y_FEN_MIN := 1.0E+38 ;
        Y_FEN_MAX := -1.0E+38 ;
        for i in y'range(1)
        loop
            Y_FEN_MIN := MIN(Y_FEN_MIN,y(i,j)) ;
            Y_FEN_MAX := MAX(Y_FEN_MAX,y(i,j)) ;
            if (x(i) > b) then
                n := i - 1 ;
                exit ;
            end if ;
        end loop ;
        Y_UNITE := (Y_FEN_MAX - Y_FEN_MIN)/20.0 ;
        Y_FEN_MIN := Y_FEN_MIN - Y_UNITE ;
        Y_FEN_MAX := Y_FEN_MAX + Y_UNITE ;
        FENETRE_GRAPHIQUE(X_FEN_MIN,X_FEN_MAX,Y_FEN_MIN,Y_FEN_MAX) ;
        XY_AXES(a,Y_FEN_MIN + Y_UNITE,X_UNITE,Y_UNITE,TRUE) ;
        DEPLACE(x(x'FIRST),y(y'FIRST,j)) ;
        for i in x'FIRST + 1..n
        loop
            TRACE(x(i),y(i,j)) ;
        end loop ;
        if (j < p) then
            PAUSE_GRAPHIQUE ;
            CLEAR_SCREEN ;
        else
            SORTIE_GRAPHIQUE ;
        end if ;
    end loop ;
end if ;
end ;
CLOSE(IMP) ;
end DEM_RKp1 ;

```

## 8.7 Démonstration du paquetage EQUADIFF\_1Q\_GENERIQUE

```

with TEXT_IO, IIO, CRT, COMMON_MATH0, MATH0, MATH1, MATH2,
     MATH3, COMMON_MATRIX, EQUADIFF_1Q_GENERIQUE ;
use TEXT_IO, IIO, CRT, COMMON_MATH0, MATH0, MATH1, MATH2,
     MATH3, COMMON_MATRIX ;

procedure DEM_RK1Q is
N_MAX: constant INTEGER := 500 ;
MAX_ORDRE : constant INTEGER := 10 ;
a, b : FLOAT ;
q : INTEGER ;
AFFICHER, CONTROLE_DU_PAS : BOOLEAN ;
FF : FONCTION ;
x, y : VECTEUR(0..N_MAX) ;
n : INTEGER := x'LAST ;
X_UNITE, Y_UNITE, X_FEN_MIN, X_FEN_MAX, Y_FEN_MIN, Y_FEN_MAX : FLOAT ;

function f(x : in FLOAT ; y : in VARIABLE_2) return FLOAT is
begin
    return EVALUE(FF,x,y) ;
end f ;

package EQUADIFF_1Q_F is new EQUADIFF_1Q_GENERIQUE(F) ;
use EQUADIFF_1Q_F ;

begin
MODE_AFFICHAGE ;
PUT_LINE(" L'equation est ecrite sous la forme :") ;
PUT_LINE("      (q)                (q-1)") ;
PUT_LINE("  y      = f(x,y,...,y      )") ;
PUT_LINE(" ") ;
PUT_LINE("La derivee d'ordre i est notee yi, pour i = 0, ..., q - 1") ;
NEW_LINE ;
ENTRER_ENTIER_BORNE(1,MAX_ORDRE,q,"Ordre de l'equation : ") ;
NEW_LINE ;
declare
    ya : VARIABLE_2(0..q-1) ;
begin
    PUT("Equation f(x,y0,") ;
    for k in 1..q - 2
    loop
        PUT("y") ; PUT(k,2) ; PUT(",") ;
    end loop ;
    PUT("y") ; PUT(q-1,2) ; PUT_LINE(") = ") ;
    GET_LINE(FF," ") ;
    NEW_LINE ;
    ENTRER_REEL(a,"Valeur de a : ") ;
    ENTRER_REEL_BORNE(a,1.0E+6,b,"Valeur de b : ") ;
    for k in 0..q - 1
    loop
        PUT("Valeur initiale de la derivee d'ordre ") ; PUT(k,2) ;
        NEW_LINE ;
    end loop ;
end ;

```

```

        ENTRER_REEL(ya(k), "    ") ;
    end loop ;
    NEW_LINE ;
    AFFICHER := LIRE_REPONSE("Voulez-vous l'affichage des resultats? ") ;
    CONTROLE_DU_PAS := LIRE_REPONSE("Methode avec controle du pas? ") ;
    CLRSCR ;
    if (CONTROLE_DU_PAS) then
        RESOLUTION_RK4_1Q_CONTROLE(a,b,ya,AFFICHER,x,y) ;
    else
        RESOLUTION_RK4_1Q(a,b,ya,AFFICHER,x,y) ;
    end if ;
    PAUSE ; CLRSCR ;
if LIRE_REPONSE("Voulez vous le trace de la solution trouvee? ") then
    X_UNITE := (b - a)/20.0 ;
    X_FEN_MIN := a - X_UNITE ;
    X_FEN_MAX := b + X_UNITE ;
    Y_FEN_MIN := 1.0E+38 ;
    Y_FEN_MAX := -1.0E+38 ;
    for i in y'range
    loop
        Y_FEN_MIN := MIN(Y_FEN_MIN,y(i)) ;
        Y_FEN_MAX := MAX(Y_FEN_MAX,y(i)) ;
        if (x(i) > b) then
            n := i - 1 ;
            exit ;
        end if ;
    end loop ;
    Y_UNITE := (Y_FEN_MAX - Y_FEN_MIN)/20.0 ;
    Y_FEN_MIN := Y_FEN_MIN - Y_UNITE ;
    Y_FEN_MAX := Y_FEN_MAX + Y_UNITE ;
    INIT_GRAPHIQUE ;
    FENETRE_GRAPHIQUE(X_FEN_MIN,X_FEN_MAX,Y_FEN_MIN,Y_FEN_MAX) ;
    XY_AXES(a,Y_FEN_MIN + Y_UNITE,X_UNITE,Y_UNITE,TRUE) ;
    DEPLACE(x(x'FIRST),y(y'FIRST)) ;
    for i in x'FIRST + 1..n
    loop
        TRACE(x(i),y(i)) ;
    end loop ;
    TITRE("Graphe de la solution approchee y = y(x)") ;
    SORTIE_GRAPHIQUE ;
end if ;
end ;
CLOSE(IMP) ;
end DEM_RK1Q ;

```

## 8.8 Démonstration du paquetage EQUADIFF\_PQ\_GENERIQUE

```

with TEXT_IO, IIO, CRT, COMMON_MATH0, MATH0, MATH1, GRAPH,
    MATH2, MATH3, COMMON_MATRIX, EQUADIFF_PQ_GENERIQUE ;
use TEXT_IO, IIO, CRT, COMMON_MATH0, MATH0, MATH1, GRAPH,
    MATH2, MATH3, COMMON_MATRIX ;

```

```

procedure DEM_RKPQ is
N_MAX : constant INTEGER := 500 ;
MAX_FONCTIONS : constant INTEGER := 10 ;
subtype INDICE_FCT is INTEGER range 1..MAX_FONCTIONS ;
type SYSTEME_DIFF is array(INDICE_FCT range <>) of FONCTION ;
a, b : FLOAT ;
AFFICHER, CONTROLE_DU_PAS : BOOLEAN ;
p, q, n : INTEGER ;
X_UNITE, Y_UNITE, X_FEN_MIN, X_FEN_MAX, Y_FEN_MIN, Y_FEN_MAX : FLOAT ;
x : VECTEUR(0..N_MAX) ;
begin
  MODE_AFFICHAGE ;
  PUT_LINE(" Le systeme d'equations est ecrit sous la forme :") ;
  NEW_LINE ;
  PUT_LINE("      (q)                (q-1)                (q-1) ") ;
  PUT_LINE("  y      = f (x,y , ...,y ,.....,y ,...,y )") ;
  PUT_LINE("    j      j    1      1      p      p ") ;
  NEW_LINE ;
  PUT_LINE(" Pour j = 1, ..., p et i = 0, ..., q - 1, ") ;
  PUT_LINE("la derivee d'ordre i de y(j) est notee y(10*j + i)") ;
  NEW_LINE ;
  ENTRER_ENTIER_BORNE(1,MAX_FONCTIONS,p,"Nombre d'equations : ") ;
  NEW_LINE ;
  ENTRER_ENTIER_BORNE(1,MAX_FONCTIONS,q,"Ordre du systeme : ") ;
  NEW_LINE ;
  declare
    FF : SYSTEME_DIFF(1..p) ;
-- Pour p = q = 2, FF(j,x,y10,y11,y12,y20,y21) represente l'equation
-- y(j)'' = f(j)(x,y(1),y(1)',y(2),y(2)')
    ya : MATRICE(1..p,0..q-1) ;
    y : MATRICE(0..N_MAX,1..p) ;
    yy : VARIABLE_2(10..10*p + q - 1) ;
    function f(j : in POSITIVE ; x : in FLOAT ; y : in VARIABLE_2)
      return FLOAT is
    begin
      for j in 1..p
      loop
        for i in 0..q-1
        loop
          yy(10*j + i) := y(q*(j - 1) + i + 1) ;
        end loop ;
      end loop ;
      return EVALUE(FF(j),x,yy) ;
    end f ;

    package EQUADIFF_PQ_F is new EQUADIFF_PQ_GENERIQUE(F) ;
    use EQUADIFF_PQ_F ;

  begin
    for j in 1..p
    loop
      PUT("Equation Numero ") ; PUT(j,2) ; NEW_LINE ;
      PUT("f(x,") ;

```

```

    for k in 1..p
    loop
        for i in 0..q - 1
        loop
            PUT("y") ; PUT(10*k + i,2) ;
            if ((k /= p) or (i /= q - 1)) then
                PUT(",") ;
            else
                PUT_LINE(" = ") ;
            end if ;
        end loop ;
    end loop ;
    GET_LINE(ff(j),"  ") ;
end loop ;
NEW_LINE ;
PUT_LINE(" Bornes de l'intervalle") ;
NEW_LINE ;
ENTRER_REEL(a,"Valeur de a : ") ;
ENTRER_REEL_BORNE(a,1.0E+6,b,"Valeur de b : ") ;
for j in 1..p
loop
    for i in 0..q - 1
    loop
        PUT("Valeur initiale de la derivee d'ordre") ; PUT(i,2) ;
        PUT(" de y(") ; PUT(j,2) ; PUT_LINE(")") ;
        ENTRER_REEL(ya(j,i)," y = ") ;
    end loop ;
end loop ;
NEW_LINE ;
AFFICHER := LIRE_REPONSE("Voulez-vous l'affichage des resultats? ") ;
CONTROLE_DU_PAS := LIRE_REPONSE("Methode avec controle du pas? ") ;
CLRSCR ;
if (CONTROLE_DU_PAS) then
    RESOLUTION_RK4_PQ_CONTROLE(a,b,ya,AFFICHER,x,y) ;
else
    RESOLUTION_RK4_PQ(a,b,ya,AFFICHER,x,y) ;
end if ;
PAUSE ; CLRSCR ;
if LIRE_REPONSE("Voulez vous le trace de la solution trouvee? ") then
n := x'LAST ;
INIT_GRAPHIQUE ;
X_UNITE := (b - a)/20.0 ;
X_FEN_MIN := a - X_UNITE ;
X_FEN_MAX := b + X_UNITE ;
for j in 1..p
loop
    Y_FEN_MIN := 1.0E+38 ;
    Y_FEN_MAX := -1.0E+38 ;
    for i in y'range(1)
    loop
        Y_FEN_MIN := MIN(Y_FEN_MIN,y(i,j)) ;
        Y_FEN_MAX := MAX(Y_FEN_MAX,y(i,j)) ;
        if (x(i) > b) then

```

```
        n := i - 1 ;
        exit ;
    end if ;
end loop ;
Y_UNITE := (Y_FEN_MAX - Y_FEN_MIN)/20.0 ;
Y_FEN_MIN := Y_FEN_MIN - Y_UNITE ;
Y_FEN_MAX := Y_FEN_MAX + Y_UNITE ;
FENETRE_GRAPHIQUE(X_FEN_MIN,X_FEN_MAX,Y_FEN_MIN,Y_FEN_MAX) ;
XY_AXES(a,Y_FEN_MIN + Y_UNITE,X_UNITE,Y_UNITE,TRUE) ;
DEPLACE(x(x'FIRST),y(y'FIRST,j)) ;
for i in x'FIRST + 1..n
loop
    TRACE(x(i),y(i,j)) ;
end loop ;
if (j < p) then
    PAUSE_GRAPHIQUE ;
    CLEAR_SCREEN ;
else
    SORTIE_GRAPHIQUE ;
end if ;
end loop ;
end if ;
end ;
CLOSE(IMP) ;
end DEM_RKPQ ;
```

## CHAPITRE 8

# Méthode des différences finies

### 1. Problème de Dirichlet linéaire en dimension un

#### 1.1 Introduction

On s'intéresse ici aux équations différentielles linéaires d'ordre 2, avec conditions aux limites, du type (*problème de Dirichlet*) :

$$(1) \quad \begin{cases} y''(x) + \alpha(x) \cdot y'(x) + \beta(x) \cdot y(x) = f(x) & (a < x < b) \\ y(a) = y_a \text{ et } y(b) = y_b \end{cases}$$

où  $a, b, f$  sont des applications continues de  $[a, b]$  dans  $\mathbb{R}$  et la fonction inconnue  $y$  est de classe  $C^2$  sur  $[a, b]$ .

*Remarques* — (i) On sait que l'ensemble des solutions de l'équation différentielle  $y''(x) + \alpha(x) \cdot y'(x) + \beta(x) \cdot y(x) = f(x)$  est un espace affine de dimension deux. En particulier, une solution est uniquement déterminée par les conditions initiales  $y(a)$  et  $y'(a)$  donnés. Mais pour le problème qui nous intéresse, les conditions aux limites portent sur  $y(a)$  et  $y(b)$ , ce qui est différent et le résultat précédent n'assure pas l'existence et l'unicité de solutions.

*Exemple* — Considérons le problème :

$$\begin{cases} y''(x) + r^2 \cdot y(x) = 1 & (0 < x < 1) \\ y(0) = y_0 \text{ et } y(1) = y_1 \end{cases}$$

La solution générale est donnée par :

$$y(x) = c \cdot \cos(r \cdot x) + d \cdot \sin(r \cdot x) + \frac{1}{r^2}$$

où les constantes  $c$  et  $d$  sont définies par :



$$\begin{cases} c = y_0 - \frac{1}{r^2} \\ c \cdot \cos(r) + d \cdot \sin(r) = y_1 - \frac{1}{r^2} \end{cases}$$

Le déterminant de ce système étant  $D = \sin(r)$ , pour  $r = 2 \cdot k \cdot \pi$  avec  $k$  dans  $\mathbb{Z}$ , il y aura une infinité de solutions (si  $y_0 = y_1$ ) ou pas de solutions du tout (si  $y_0$  différent de  $y_1$ ).

(ii) Sauf dans des cas très particuliers, on ne sait pas résoudre de façon explicite les équations du type (1). On cherchera donc des procédés numériques qui permettent de donner des valeurs approchées d'une solution en certains points.

*Problèmes posés par (1)*

- Existence de solutions ;
- Elaboration de méthodes numériques ;
- Qualité de l'approximation d'une solution par une méthode donnée (conditions de convergence ; vitesse de convergence).

### **1.2 Théorème d'existence et d'unicité de solutions du problème de Dirichlet linéaire**

On peut montrer le :

*Théorème* : Si  $\beta(x) < 0$  pour tout  $x$  dans  $[a,b]$ , alors le problème de Dirichlet admet une unique solution  $y$  dans  $C^2([a,b])$ .

Contrairement au problème de Cauchy l'existence et l'unicité n'est pas assurée dans tous les cas.

On supposera donc, dans tout ce qui suit que  $\beta(x) < 0$  sur  $[a,b]$ .

### **1.3 Forme canonique de l'équation de Dirichlet. Le problème de Poisson**

On suppose que  $\alpha$  est de classe  $C^1$  et on note  $A(x)$  une primitive de  $\alpha/2$  sur  $[a,b]$ .

Soit  $y$  solution de (1). En posant, pour tout  $x$  dans  $[a,b]$  :

$$z(x) = y(x) \cdot e^{A(x)} \quad (\text{transformation de Liouville})$$

la fonction  $z$  est solution du problème de Poisson :

$$\begin{cases} z''(x) + p(x) \cdot z(x) = g(x) & (a < x < b) \\ z(a) = z_a \text{ et } z(b) = z_b \end{cases}$$

où :

$$\begin{cases} p(x) = \beta(x) - \frac{\alpha'(x)}{2} - \frac{\alpha^2(x)}{4} \\ g(x) = f(x) \cdot e^{A(x)} \\ z_a = y_a \cdot e^{A(a)} \text{ et } z_b = y_b \cdot e^{A(b)} \end{cases}$$

On peut donc se contenter d'étudier l'équation de Poisson, mais si on ne connaît pas de primitive de  $\alpha$  la transformation de Liouville n'a pas d'intérêt pratique.

### 1.4 Résolution approchée du problème de Dirichlet par discrétisation

En vue de discrétiser l'équation (1), on coupe l'intervalle  $[a,b]$  en  $n + 1$  intervalles égaux (on dit aussi qu'on fait un *maillage uniforme* de  $[a,b]$ ) en posant :

$$h = \frac{b-a}{n+1}, x_k = a + k \cdot h \text{ pour } k = 0, 1, \dots, n+1.$$

On dit que les  $x_k$  sont les *noeuds du maillage*.

On cherche alors, pour tout  $k = 1, 2, \dots, n$  une valeur approchée de  $y(x_k)$ . Ces valeurs seront consignées dans un vecteur

$$u = \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix}$$

Une fois ces valeurs approchées trouvées, on approxime la fonction  $y$  solution de (1) par la fonction affine par morceaux définie sur  $[a,b]$  par  $\varphi(x_k) = u_k$  et  $\varphi$  affine sur  $[x_k, x_{k+1}]$  pour tout  $k = 0, \dots, n$ , où on a posé  $u_0 = y_a$  et  $u_{n+1} = y_b$ .

Le graphe de  $\varphi$  sera alors une approximation de la courbe intégrale de (1).

On cherchera également à donner une majoration de l'erreur :

$$\text{Sup} \{ |y(x) - \varphi(x)|; 0 \leq x \leq 1 \}$$

## 2. Approximations des dérivées d'une fonction par différences finies

Soit  $f : [a,b] \rightarrow \mathbb{R}$  une fonction de classe  $C^4$  et  $x_0$  un point de  $[a,b]$ .

### 2.1 Approximation de $f'(x_0)$ par différence finie centrée

La première idée pour approximer  $f'(x_0)$  est de prendre :

$$f'(x_0) \cong \frac{f(x_0 + h) - f(x_0)}{h} \quad (\text{différence finie progressive})$$

ou

$$f'(x_0) \cong \frac{f(x_0) - f(x_0 - h)}{h} \quad (\text{différence finie régressive})$$

où  $h > 0$  est choisi assez petit.

Avec la formule de Taylor-Lagrange, on a :

$$\left| f'(x_0) - \frac{f(x_0 + h) - f(x_0)}{h} \right| = \left| \frac{f''(x_0 + \theta \cdot h)}{2} \right| \cdot h \quad (0 < \theta < 1)$$

c'est-à-dire que l'erreur commise est d'ordre 1 en  $h$ .

De même avec la différence finie régressive.

Mais en regardant de plus près, on peut obtenir une approximation qui sera d'ordre 2 en  $h$ .

Avec la formule de Taylor à l'ordre 3, on peut écrire :

$$\frac{f(x_0 + h) - f(x_0 - h)}{2 \cdot h} = f'(x_0) + \frac{f'''(x_0 + \theta_1 \cdot h) + f'''(x_0 - \theta_2 \cdot h)}{12} \cdot h^2 \quad (0 < \theta_1, \theta_2 < 1)$$

De plus, en appliquant le théorème des valeurs intermédiaires à la fonction :

$$\varphi: t \rightarrow \varphi(t) = f'''(x_0 + \theta_1 \cdot h) + f'''(x_0 - \theta_2 \cdot h) - 2 \cdot f'''(t)$$

on peut écrire :

$$\frac{f(x_0 + h) - f(x_0 - h)}{2 \cdot h} = f'(x_0) + \frac{f'''(x_0 + \theta \cdot h)}{6} \cdot h^2 \quad (-1 < \theta < 1)$$

Donc la différence finie centrée  $\frac{f(x_0 + h) - f(x_0 - h)}{2 \cdot h}$  donne une approximation d'ordre 2 en  $h$  de  $f'(x_0)$ .

## 2.2 Approximation de $f''(x_0)$ par différence finie centrée

Si, pour approximer la dérivée seconde, on applique deux fois ce qui précède on n'obtient qu'une approximation d'ordre 1 de cette dérivée.

Là encore, en regardant de plus près, on peut obtenir une approximation d'ordre 2.

Avec la formule de Taylor à l'ordre 4, on peut écrire :

$$\frac{f(x_0 + h) - 2 \cdot f(x_0) + f(x_0 - h)}{h^2} = f''(x_0) + \frac{f^{(4)}(x_0 + \theta \cdot h)}{12} \cdot h^2 \quad (-1 < \theta < 1)$$

Donc la différence finie centrée  $\frac{f(x_0 + h) - 2 \cdot f(x_0) + f(x_0 - h)}{h^2}$  donne une approximation d'ordre 2 en  $h$  de  $f''(x_0)$ .

## 3. Résolution approchée du problème de Dirichlet par la méthode des différences finies

### 3.1 Discrétisation du problème de Dirichlet

On reprend les notations du paragraphe (1.4).

En écrivant l'équation (1) en tous les points  $x_k$  ( $k = 1, \dots, n$ ) et en utilisant les résultats du (2.2), on voit que le vecteur  $y$  de coordonnées  $y_k = y(x_k)$  est solution du système :

$$A \cdot y = v + \varepsilon \quad (3)$$

où  $A$  est la matrice tridiagonale donnée par :

$$A = \begin{pmatrix} a_1 & c_1 & 0 & \dots & \dots & 0 \\ b_2 & a_2 & c_2 & 0 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \dots & 0 & b_{n-1} & a_{n-1} & c_{n-1} \\ 0 & \dots & \dots & 0 & b_n & a_n \end{pmatrix}$$

et les vecteurs  $v$  et  $\varepsilon$  sont donnés par :

$$v = \begin{pmatrix} v_1 \\ v_2 \\ \cdot \\ \cdot \\ v_{n-1} \\ v_n \end{pmatrix} = \begin{pmatrix} f(x_1) - b_1 \cdot y_a \\ f(x_2) \\ \cdot \\ \cdot \\ f(x_{n-1}) \\ f(x_n) - c_n \cdot y_b \end{pmatrix} \quad \text{et} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \cdot \\ \cdot \\ \varepsilon_{n-1} \\ \varepsilon_n \end{pmatrix}$$

avec, pour tout  $k = 1, 2, \dots, n$  :

$$\begin{aligned} a_k &= \beta(x_k) - \frac{2}{h^2} \\ b_k &= \frac{1}{h^2} - \frac{\alpha(x_k)}{2 \cdot h} \\ c_k &= \frac{1}{h^2} + \frac{\alpha(x_k)}{2 \cdot h} \\ \varepsilon_k &= \frac{f^{(4)}(\theta_k) + 2 \cdot \alpha(x_k) \cdot f'''(\eta_k)}{12} \cdot h^2 \end{aligned}$$

où  $\theta_k$  et  $\eta_k$  sont compris entre  $x_k$  et  $x_{k+1}$ .

On a donc  $\|\varepsilon\| < C \cdot h^2$  où  $C$  est une constante qui dépend de  $\alpha$  et des dérivées troisième et quatrième de  $f$ .

En négligeant le terme  $\varepsilon$ , le système discrétisé associé au problème de Dirichlet est alors le système tridiagonal :

$$A \cdot u = v \quad (4)$$

*Remarque* — Le système (4) sera à diagonale strictement dominante si :

$$|a_k| > |b_k| + |c_k| \quad (k = 1, \dots, n)$$

soit si :

$$\left| h^2 \cdot \beta(x_k) - 2 \right| > \left| 1 + \frac{\alpha(x_k) \cdot h}{2} \right| + \left| 1 - \frac{\alpha(x_k) \cdot h}{2} \right|$$

En prenant  $h > 0$  tel que  $h < \frac{2}{M}$ , où  $M = \text{Sup} \{ |\alpha(x)| ; x \in [a,b] \}$ , la condition ci-dessus s'écrit :  $\left| h^2 \cdot \beta(x_k) - 1 \right| > 2$ , ce qui est réalisé si  $\beta$  ne prend que des valeurs négatives sur  $[a,b]$ .

En résumé la condition suffisante du théorème (1.2) nous assure, que pour  $h$  assez petit, le système discrétisé associé au problème de Dirichlet est un système linéaire tridiagonal à diagonale strictement dominante, il admet donc une unique solution.

Cette dernière peut se trouver par la méthode du double balayage de Cholesky.

La programmation structurée ne pose pas de problème, il suffit de décrire une procédure qui donne les coefficients de la matrice tridiagonale  $A$ , puis de faire appel à une procédure de résolution d'un système tridiagonal.

### 3.2 Discrétisation de l'équation de Poisson

Dans le cas de l'équation de Poisson, on a  $\alpha = 0$  et la matrice du système est symétrique définie positive (car à diagonale strictement dominante et à termes diagonaux strictement positifs).

En multipliant par  $-h^2$ , on est en fait ramené à résoudre le système tridiagonal particulier  $A \cdot u = v$ , de matrice :

$$A = \begin{pmatrix} a_1 & -1 & 0 & \dots & \dots & 0 \\ -1 & a_2 & -1 & 0 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \dots & 0 & -1 & a_{n-1} & -1 \\ 0 & \dots & \dots & 0 & -1 & a_n \end{pmatrix}$$

avec :

$$a_k = 2 - h^2 \cdot \beta(x_k) > 2 \quad (k = 1, 2, \dots, n)$$

et de second membre :

$$v = \begin{pmatrix} v_1 \\ v_2 \\ \cdot \\ \cdot \\ v_{n-1} \\ v_n \end{pmatrix} = \begin{pmatrix} -f(x_1) \cdot h^2 + y_a \\ -f(x_2) \cdot h^2 \\ \cdot \\ \cdot \\ -f(x_{n-1}) \cdot h^2 \\ -f(x_n) \cdot h^2 + y_b \end{pmatrix}$$

On peut résoudre ce système en adaptant la méthode des pivots de Gauss, ce qui donne le système triangulaire :

$$\begin{cases} u_j - d_j \cdot u_{j+1} = e_j & (j = 1, \dots, n - 1) \\ u_n = e_n \end{cases}$$

de solution :

$$\begin{cases} u_n = e_n \\ u_k = e_k + d_k \cdot u_{k+1} \quad (k = n - 1, \dots, 1) \end{cases}$$

où on a posé :

$$\begin{cases} d_1 = \frac{1}{a_1}; e_1 = v_1 \cdot d_1 \\ d_k = \frac{1}{a_k - d_{k-1}} \quad (k = 2, \dots, n) \\ e_k = (v_k + e_{k-1}) \cdot d_k \end{cases}$$

*Remarque* — Comme  $a_k > 2$  pour tout  $k = 1, \dots, n$ , on en déduit que  $d_k$  est dans  $]0,1[$  pour tout  $k = 1, \dots, n$ .

En notant  $\|v\| = \text{Sup}\{|v_k|; k = 1, \dots, n\}$ , on en déduit alors que  $e_k < k \cdot \|v\|$ , pour tout  $k = 1, \dots, n$ . Puis :

$$u_k < \|v\| \cdot \frac{(n+1)^2}{2} \quad (k = 1, \dots, n).$$

### 3.3 Majoration de l'erreur

On reprend les notations du paragraphe (1.4) et le problème est de majorer l'erreur commise  $g(x) = y(x) - \varphi(x)$ , en remplaçant la solution exacte  $y(x)$  de l'équation de Dirichlet par la fonction affine par morceaux  $\varphi(x)$  qui passe par les points de coordonnées  $(x_k, u_k)$  pour tout  $k = 0, \dots, n + 1$ .

Grâce à la transformation de Liouville (Cf. paragraphe (1.3)), on peut se limiter au cas de l'équation de Poisson, et un changement de variable affine nous ramène facilement à l'intervalle  $[0,1]$ . On fait donc ces hypothèses dans ce paragraphe.

Notons  $\psi$  la fonction affine par morceaux qui passe par les points de coordonnées  $(x_k, y(x_k))$  pour  $k = 0, \dots, n + 1$ .

On écrit alors,  $g(x) = (y(x) - \psi(x)) + (\psi(x) - \varphi(x))$ .

#### 3.3.1 Majoration de $h(x) = y(x) - \psi(x)$

Pour  $k = 0, \dots, n + 1$ , on a  $h(x_k) = 0$  et  $h$  est deux fois continûment dérivable sur  $]x_k, x_{k+1}[$  avec  $h'' = y''$  car  $\psi$  est affine par morceaux (voir figure 8.1).

On peut alors montrer le :

*Lemme* — On notant  $M_2 = \text{Sup}\{|y''(x)|; x \in [a,b]\}$ , on a :

$$|h(x)| < \frac{(x - x_k) \cdot (x_{k+1} - x) \cdot M_2}{2} \quad \text{sur } [x_k, x_{k+1}].$$

Et avec  $(x - x_k) \cdot (x_{k+1} - x) \leq \frac{(x_{k+1} - x)^2}{4} \leq \frac{h^2}{4}$  sur  $[x_k, x_{k+1}]$ , on en déduit que :

$$|h(x)| \leq \frac{M^2 \cdot h^2}{8} \quad \text{sur } [0,1].$$

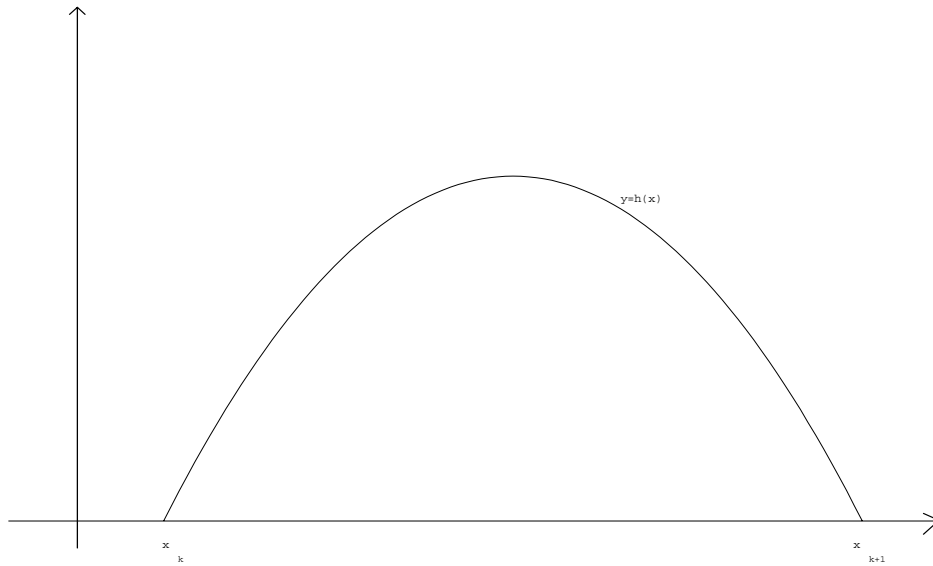


Figure 8.1

### 3.3.2 Majoration de $\theta(x) = \psi(x) - \phi(x)$

La fonction  $\theta$  est affine par morceaux, on a donc sur  $[0,1]$  :

$$|\theta(x)| \leq \text{Max}\{|\theta(x_k)|; k = 1, \dots, n\}$$

avec  $\theta(x_k) = y(x_k) - u_k$  pour  $k = 1, \dots, n$  ( $\theta(0) = \theta(1) = 0$ ).

Mais, avec :

$$\begin{cases} A \cdot y = v + \varepsilon \\ A \cdot u = v \end{cases}$$

on déduit que :

$$A \cdot (y - u) = \varepsilon$$

soit, en multipliant par  $-h^2$  :

$$-h^2 \cdot A \cdot (y - u) = -h^2 \cdot \varepsilon.$$

Avec la remarque du paragraphe (3.2), on déduit alors que :

$$\|y - u\| \leq h^2 \cdot \|\varepsilon\| \cdot \frac{(n+1)^2}{2}$$

puis avec  $\|\varepsilon\| \leq C \cdot h^2$ ,  $h = \frac{1}{n+1}$  et  $|\theta(x)| \leq \|y - u\|$ , on déduit que :

$$|\theta(x)| \leq M_4 \cdot \frac{h^2}{24} \text{ sur } [0,1], \text{ où } M_4 = \text{Sup}\{|y^{(4)}(x)|; x \in [0,1]\}$$

### 3.3.3 Majoration de $g(x) = y(x) - \phi(x)$

De ce qui précède, on déduit donc que :

$$|y(x) - \phi(x)| \leq K \cdot h^2 \text{ sur } [0,1], \text{ où } K = \frac{M_2 + M_4}{8}$$

Il en résulte que la méthode converge toujours et est d'ordre 2 en  $h$ .

*Remarque* — On a en fait montré que la suite de fonctions affines par morceaux  $\varphi_n$  converge uniformément vers la solution exacte  $y$  sur  $[0,1]$ .

## 4. Exemples d'application

### 4.1 Dissipation de la chaleur dans un disque

*Position du problème* — On se donne un disque plat de rayon  $R$  et d'épaisseur  $e$  négligeable devant  $R$ . On suppose qu'il est constitué d'un matériau de conductivité thermique  $\lambda$ , de masse volumique  $\mu$  et de capacité calorifique massique  $C$ .

La température en un point  $M(r)$ , situé à la distance  $r$  de l'axe du disque est supposée ne dépendre que de  $r$ , on la note  $T(r)$ .

Un corps cylindrique de rayon  $R_1 < R$ , maintenu à la température  $T_1$  est appliqué au centre du disque sur les deux faces. Donc les points à une distance  $r \in [0, R_1]$  sont à la température constante  $T_1$ .

Pour  $r \in ]R_1, R]$ , les deux faces du disque sont au contact avec un milieu ambiant à une température  $T_0$ .

La puissance thermique dissipée par unité de surface est  $A \cdot (T(r) - T_0)$  pour  $r \in [R_1, R]$ .

Enfin, le bord du disque est maintenu à la température constante  $T_0$ .

Les données numériques sont les suivantes :

$\lambda = 100 \cdot \text{W} \cdot \text{m}^{-1} \cdot \text{K}^{-1}$  ;  $A = 100 \cdot \text{W} \cdot \text{K}^{-1} \cdot \text{m}^{-2}$  ;  $R = 8 \cdot \text{cm}$  ;  $R_1 = 2 \cdot \text{cm}$  ;  $e = 2 \cdot \text{mm}$   
 $T_1 = 500 \cdot \text{K}$  ;  $T_0 = 300 \cdot \text{K}$ .

*Analyse du problème* — En régime permanent, le bilan thermique d'une couronne de rayons  $r$  et  $r + dr$  nous permet d'écrire :

$$2 \cdot r \cdot e \cdot \lambda \cdot T'(r) - 2 \cdot (r + dr) \cdot e \cdot \lambda \cdot T'(r + dr) + 2 \cdot (2 \cdot r \cdot dr) \cdot A \cdot (T(r) - T_0) = 0.$$

Soit en écrivant que  $T'(r + dr) \cong T'(r) + dr \cdot T''(r)$ , et en négligeant  $(dr)^2$  :

$$-2 \cdot dr \cdot e \cdot \lambda \cdot T'(r) - 2 \cdot r \cdot dr \cdot e \cdot \lambda \cdot T''(r) + 2 \cdot (2 \cdot r \cdot dr) \cdot A \cdot (T(r) - T_0) = 0.$$

Soit, en divisant par  $-2r \cdot dr \cdot e \cdot \lambda$  :

$$T''(r) + \frac{1}{r} \cdot T'(r) - \frac{2 \cdot A}{e \cdot \lambda} \cdot (T(r) - T_0) = 0.$$

En conclusion, l'équation différentielle vérifiée par  $y(x) = T(x) - T_0$ , où  $x = r$  et  $T_0 = 300$  degrés Kelvin est :

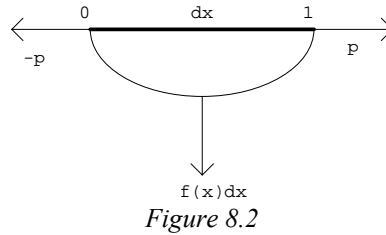
$$\begin{cases} y''(x) + \alpha(x) \cdot y'(x) + \beta(x) \cdot y(x) = f(x) & (0.02 < x < 0.08) \\ y(0.02) = 0; y(0.08) = 200 \end{cases}$$

avec :

$$\alpha(x) = \frac{1}{x} ; \beta(x) = -\frac{2 \cdot A}{\lambda \cdot e} = -1000 ; f(x) = 0.$$



## 4.2 Fléchissement d'une poutre



Considérons une poutre de longueur 1, appuyée à ses extrémités 0 et 1, soumise à une charge transversale  $f(x)dx$  par unité de longueur, et étirée selon son axe par une force  $P$ .

Alors, le moment fléchissant  $u(x)$ , au point d'abscisse  $x \in [0,1]$  est solution de l'équation de Poisson :

$$\begin{cases} -u''(x) + c(x) \cdot u(x) = f(x) & (0 < x < 1) \\ u(0) = u(1) = 0 \end{cases}$$

où  $c(x) = \frac{P}{E \cdot I(x)}$ , avec  $E$  désignant le module de Young du matériau constituant

la poutre et  $I(x)$  désignant le moment principal d'inertie de la section de la poutre au point d'abscisse  $x$ .

## 5. Résolution approchée d'équations aux dérivées partielles par la méthode des différences finies

### 5.1 Intervention des équations aux dérivées partielles en physique

Les équations aux dérivées partielles interviennent dans la description de nombreux phénomènes physiques. Ci-dessous on décrit certains exemples parmi les plus connus.

#### 5.1.1 L'équation des cordes vibrantes ou équation des ondes en dimension un

*Position du problème* — Considérons une corde de section constante, de longueur  $L$ , fixée en ses extrémités sur un axe  $Ox$ . La corde est tendue avec une tension  $T$ . On fait, de plus, les hypothèses suivantes :

- le poids par unité de longueur  $dx$  est  $w(x)$  ;
- la corde est soumise à une force verticale  $f(x,t)$  par unité de longueur  $dx$  ;
- la déformation de la corde influe de manière négligeable sur sa longueur, donc en tout point de la corde l'angle  $\alpha$  de cette dernière avec l'axe des  $x$  est supposé petit de sorte que  $\text{tg}(\alpha) \cong \alpha$  ;
- la position de la corde et la distribution des vitesses à  $t = 0$  sont supposées connues, c'est-à-dire qu'on connaît  $y(x,0) = y(x)$  et  $\frac{\partial y}{\partial t}(x,0) = y_1(x)$  (on aura  $y_1(x) = 0$  si la corde est lâchée à partir d'une position donnée).

On se propose alors de décrire les petits mouvements verticaux de la corde en fonction de l'abscisse  $x$  et du temps  $t$ . Pour ce faire, On note  $y(x,t)$  le déplacement vertical de la corde pour  $x \in [0,L]$  et  $t > 0$ .

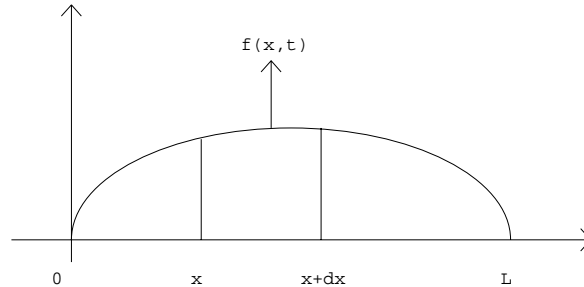


Figure 8.3

*Analyse du problème* — On travaille sur un élément de corde  $[x,x+dx]$  de masse

$$dm = \frac{w(x) \cdot dx}{g}$$

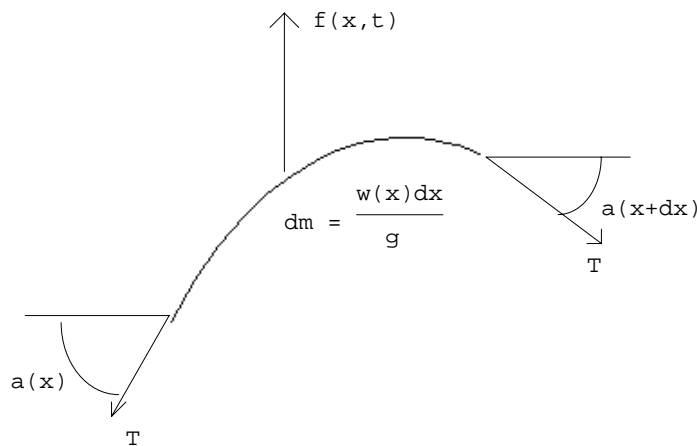


Figure 8.4

La tension en chacune de ses extrémités est dirigée suivant la tangente au segment et les composantes verticales de cette tension sont respectivement :

$$T \cdot \sin(\alpha(x+dx)) \cong T \cdot \text{tg}(\alpha(x+dx)) ; T \cdot \sin(\alpha(x)) \cong T \cdot \text{tg}(\alpha(x)).$$

L'accélération subie par  $dm$  étant  $\frac{\partial^2 y}{\partial t^2}$ , avec la loi de Newton on déduit que :

$$T \cdot \text{tg}(\alpha(x+dx)) - T \cdot \text{tg}(\alpha(x)) + f(x,t) \cdot dx = \frac{\partial^2 y}{\partial t^2} \cdot \frac{w(x) \cdot dx}{g}$$

En écrivant que  $\text{tg}(\alpha(x)) = \frac{\partial y}{\partial x}(x,t)$ , et en faisant tendre  $dx$  vers 0 on déduit alors que le déplacement vertical  $y(x,t)$  est solution de l'équation aux dérivées partielles :

$$\frac{\partial^2}{\partial t^2}(x,t) = \frac{T \cdot g}{w(x)} \cdot \frac{\partial^2 y}{\partial x^2}(x,t) + \frac{g}{w(x)} \cdot f(x,t)$$

avec les conditions :

$$\begin{cases} y(0, t) = y(L, t) = 0 \text{ pour } t \geq 0 \text{ (conditions aux limites)} \\ y(x, 0) = y_0(x) \text{ et } \frac{\partial y}{\partial t}(x, 0) = y_1(x) \text{ sur } [0, L] \text{ (conditions initiales)} \end{cases}$$

L'équation obtenue est « l'équation des ondes à une dimension » encore appelée « équation des cordes vibrantes ».

Un tel problème est appelé *problème aux limites* (à cause des conditions aux limites de la corde) *d'évolution* (pour préciser que le temps intervient) *avec conditions initiales* (car on connaît le comportement initial de la corde).

*Remarque 1* — Si on note  $c(x)^2 = \frac{T \cdot g}{W(x)}$ , alors les dimensions de  $c(x)$  sont celles d'une vitesse.

Dans le cas particulier d'une corde homogène, la densité linéaire de masse  $w(x)$  est une constante  $\sigma$  et  $c = \sqrt{\frac{T \cdot g}{\sigma}}$  est la vitesse de propagation de l'onde (voir Barret et Wylie p. 465).

*Remarque 2* — Dans certains cas, les forces extérieures  $f(x, y)$  sont négligeables devant la tension  $T$ . De plus si la corde est homogène alors l'équation des ondes va s'écrire :

$$\frac{\partial^2 y}{\partial t^2} - c^2 \cdot \frac{\partial^2 y}{\partial x^2} = 0$$

Cette dernière est connue comme *l'équation de d'Alembert* et peut être résolue de façon explicite.

Il est facile de voir que si  $f$  et  $g$  sont des fonctions d'une variable réelle deux fois dérivables, alors  $y(x, t) = f(x - c \cdot t) + g(x + c \cdot t)$  est solution de l'équation des ondes. Inversement, en posant  $u = x - c \cdot t$  et  $v = x + c \cdot t$ , on transforme l'équation en

$\frac{\partial^2 y}{\partial u \partial v} = 0$  et les solutions de cette dernière s'écrivent  $f(u) + g(v)$  où  $f$  et  $g$  sont

deux fonctions d'une variable deux fois dérivables. Les conditions initiales permettent alors de déterminer  $f$  et  $g$  de façon explicite. En effet, on a :

$$\begin{cases} y(x, 0) = f(x) + g(x) = y_0(x) \\ \frac{\partial y}{\partial t}(x, 0) = -c \cdot f'(x) + c \cdot g'(x) = y_1(x) \end{cases}$$

donc :

$$\begin{cases} f(x) = \frac{1}{2} \cdot \left( y_0(x) - \int_0^x y_1(t) \frac{dt}{c} \right) + k \\ g(x) = \frac{1}{2} \cdot \left( y_0(x) + \int_0^x y_1(t) \frac{dt}{c} \right) - k \end{cases}$$

et :

$$y(x, t) = \frac{1}{2} \cdot (y_0(x - c \cdot t) + y_0(x + c \cdot t)) + \int_{x-c \cdot t}^{x+c \cdot t} y_1(t) \frac{dt}{2 \cdot c}$$

*Remarque 3* — Si on sait résoudre de façon explicite l'équation des ondes pour une corde homogène, dans le cas général il y a peu d'espoir d'obtenir une solution explicite. On cherchera alors une solution approchée à l'aide de la méthode des différences finies.

**5.1.2 L'équation des ondes en dimension deux**

De la même façon, on peut étudier les petits déplacements verticaux d'une membrane  $\Omega$  tendue le long d'une courbe fermée dans le plan des  $x, y$ .

Les hypothèses sont les suivantes :

- le poids par unité de surface  $dx dy$  est  $w(x, y)$  ;
- la membrane est tendue avec une même tension  $T$  en tout point et toute direction ;
- la membrane est soumise à une force verticale  $f(x, y, t)$  par unité de surface  $dx dy$ ;
- la déformation de la membrane influe de manière négligeable sur ses dimensions;
- la position de la membrane et la distribution des vitesses à  $t = 0$  sont supposées connues, c'est-à-dire qu'on connaît  $z(x, y, 0) = z_0(x, y)$  et  $\frac{\partial z}{\partial t}(x, y, 0) = z_1(x, y)$ .

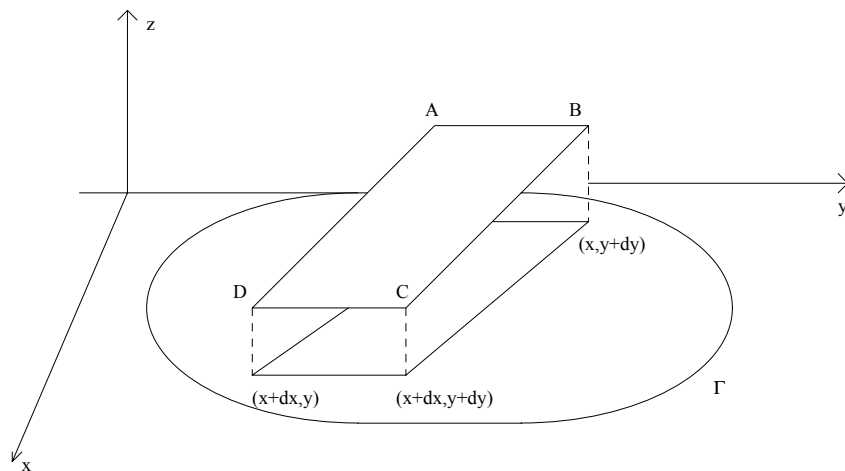


Figure 8.5

En travaillant sur un élément de surface, on déduit de la loi de Newton que le déplacement  $z(x, y, t)$  est solution de l'équation :

$$\frac{\partial^2 z}{\partial t^2} = \frac{T \cdot g}{w(x, y)} \cdot \left( \frac{\partial^2 z}{\partial x^2} + \frac{\partial^2 z}{\partial y^2} \right) + \frac{g}{w(x, y)} \cdot f(x, y, t)$$

avec les conditions initiales :

$$\begin{cases} z(x, y, 0) = z_0(x, y) \text{ sur } \Omega \\ \frac{\partial z}{\partial t}(x, y, 0) = z_1(x, y) \text{ sur } \Omega \\ z(x, y, t) = 0 \text{ sur } \Gamma \text{ pour tout } t \geq 0 \end{cases}$$

L'équation obtenue est l'équation des ondes en dimension deux.

On note  $\Delta = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2}$  (opérateur de Laplace ou *laplacien*).

*Remarque 1* — Si on s'intéresse aux solutions stationnaires, c'est-à-dire qui ne dépendent pas du temps, on obtient l'équation de Poisson :

$$\Delta z = h(x, y) \text{ sur } \Omega$$

qui pour  $h = 0$  est connue comme l'équation de Laplace.

*Remarque 2* — Pour  $w(x, y) = w$  constant et sans forces extérieures, l'équation devient :

$$\frac{\partial^2}{\partial t^2} = c^2 \cdot \Delta z$$

avec  $c = \sqrt{\frac{T \cdot g}{w}}$  qui a les dimensions d'une vitesse (c'est la vitesse de propagation de l'onde).

### 5.1.3 Equation de la chaleur en dimension un

On considère ici le problème de l'écoulement de la chaleur dans un milieu conducteur. On se donne donc un solide dans lequel la chaleur se transmet par conduction.

Comme toujours, on travaille sur un élément de volume  $dx dy dz$ . Si le poids volumique est  $\rho$ , alors la masse d'un tel élément est  $dm = \rho \cdot \frac{dx dy dz}{g}$ . De plus, si

$u(x, y, z, t)$  mesure la température en un point  $M(x, y, z)$  à l'instant  $t$  et si  $du$  est la variation de température durant  $dt$  alors la quantité de chaleur stockée dans cet élément de volume durant  $dt$  est :  $dH = c \cdot dm \cdot du$  (la quantité de chaleur gagnée ou perdue par un corps quand sa température change est proportionnelle à la masse du corps et au changement de température. La constante de proportionnalité est la *chaleur spécifique* du corps notée  $c$ ).

La vitesse de stockage de la chaleur est alors :

$$\frac{dH}{dt} = \frac{c \cdot \rho}{g} \cdot dx dy dz \cdot \frac{du}{dt}$$

Les changements de température sont dus à deux causes :

- (i) la chaleur peut être générée à l'intérieur du corps par effet électrique, chimique, ..., à une vitesse  $f(x, y, z, t)$  par unité de volume ;
- (ii) l'élément de volume reçoit de la chaleur par chacune de ses faces.

En particulier, la vitesse à laquelle la chaleur s'écoule à travers la face EFGH est approximativement (en prenant comme valeur moyenne du gradient de température sa valeur au centre de EFGH) :

$$-k \cdot dydz \cdot \frac{\partial u}{\partial x} \left( x, y + \frac{dy}{2}, z + \frac{dz}{2} \right)$$

(la vitesse à laquelle la chaleur s'écoule à travers une surface est proportionnelle à son aire et au gradient de température, en degrés par unité de longueur, dans une direction perpendiculaire à la surface. La constante de proportionnalité, notée  $k$ , est la *conductivité thermique*)

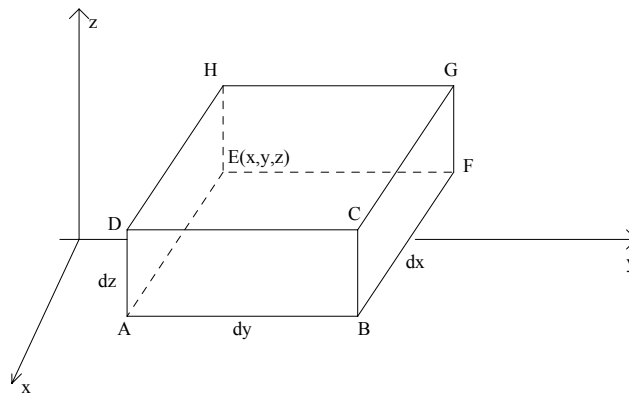


Figure 8.6

De même pour ABCD, on a approximativement :

$$k \cdot dydz \cdot \frac{\partial u}{\partial x} \left( x + dx, y + \frac{dy}{2}, z + \frac{dz}{2} \right)$$

La somme de ces deux quantités donne la vitesse de conduction de la chaleur suivant l'axe des  $x$ . On procède de même pour les deux autres axes.

Ensuite, on écrit que la vitesse à laquelle la chaleur est stockée par l'élément de volume est égale à la vitesse à laquelle elle est produite par cet élément (c'est le bilan thermique) plus la vitesse d'écoulement de la chaleur à travers ce dernier. Ce qui donne, au voisinage de l'élément de volume :

$$\begin{aligned} \frac{c \cdot \rho}{g} \cdot dx dy dz \cdot \frac{du}{dt} &= f(x, y, z, t) \cdot dx dy dz \\ &+ k \cdot dy dz \cdot \left( \frac{\partial u}{\partial x} \left( x + \frac{dx}{2}, y + \frac{dy}{2}, z + \frac{dz}{2} \right) - \frac{\partial u}{\partial x} \left( x, y + \frac{dy}{2}, z + \frac{dz}{2} \right) \right) \\ &+ k \cdot dx dz \cdot \left( \frac{\partial u}{\partial y} \left( x + \frac{dx}{2}, y + \frac{dy}{2}, z + \frac{dz}{2} \right) - \frac{\partial u}{\partial y} \left( x + \frac{dx}{2}, y, z + \frac{dz}{2} \right) \right) \\ &+ k \cdot dx dy \cdot \left( \frac{\partial u}{\partial z} \left( x + \frac{dx}{2}, y + \frac{dy}{2}, z + \frac{dz}{2} \right) - \frac{\partial u}{\partial z} \left( x + \frac{dx}{2}, y + \frac{dy}{2}, z \right) \right) \end{aligned}$$

En divisant par  $k \cdot dx dy dz$ , puis en faisant tendre  $dx$ ,  $dy$  et  $dz$  vers 0, on obtient l'équation de conduction de la chaleur (ou équation de Fourier) :

$$a^2 \cdot \frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} + \frac{\partial^2 u}{\partial z^2} + \frac{1}{k} \cdot f(x, y, z, t)$$

avec  $a^2 = \frac{C \cdot \rho}{k \cdot g}$

On a, de plus les conditions :

$$\begin{cases} u(x, y, z, 0) = u_0(x, y, z) \text{ (conditions initiales)} \\ u(x, y, z, t) = g(x, y, z, t) \text{ sur } \partial\Omega \text{ (conditions aux limites)} \end{cases}$$

Là encore, c'est un problème d'évolution avec conditions aux limites.

*Remarque* — Dans des cas particuliers importants, la chaleur n'est ni générée ni perdue dans le volume  $\Omega$ . On dit alors qu'on s'intéresse à la distribution de température en régime stationnaire. Dans ces conditions, on a  $f = 0$  et  $\frac{\partial u}{\partial t} = 0$ , de sorte que l'équation de la chaleur se réduit à l'équation de Laplace  $\Delta u = 0$ .

## 5.2 Classification des équations aux dérivées partielles d'ordre 2

### 5.2.1 Introduction

On peut se poser la question de savoir quelles sont les équations aux dérivées partielles qui se résolvent aussi simplement que celle des ondes à une dimension par une méthode analogue à celle de d'Alembert.

On s'intéresse donc ici aux équations du second ordre quasi-linéaires du type :

$$a \cdot \frac{\partial^2 u}{\partial x^2} + 2 \cdot b \cdot \frac{\partial^2 u}{\partial xy} + c \cdot \frac{\partial^2 u}{\partial y^2} = d \quad (1)$$

où  $a$ ,  $b$ ,  $c$  sont des fonctions continues de  $x$  et  $y$  et  $d$  est une fonction continue de  $x$ ,  $y$ ,  $u$ ,  $\frac{\partial u}{\partial x}$  et  $\frac{\partial u}{\partial y}$ .

*Remarque* — Une telle équation peut se ramener à un système de deux équations du premier ordre en posant  $p = \frac{\partial u}{\partial x}$  et  $q = \frac{\partial u}{\partial y}$ , ce qui donne :

$$\begin{cases} a \cdot \frac{\partial p}{\partial x} + b \cdot \frac{\partial p}{\partial y} + b \cdot \frac{\partial q}{\partial x} + c \cdot \frac{\partial q}{\partial y} = d \\ \frac{\partial q}{\partial x} - \frac{\partial p}{\partial y} = 0 \end{cases}$$

### 5.2.2 Cas particulier des équations à coefficients constants et sans second membre

On suppose ici que  $a, b, c$  sont des constantes et  $d = 0$ .

En utilisant la méthode de d'Alembert, on cherche des solutions de la forme  $u(x, y) = f(x + \lambda \cdot y)$  et on voit que nécessairement  $\lambda$  est solutions de l'équation caractéristique :

$$c \cdot \lambda^2 + 2 \cdot b \cdot \lambda + a = 0$$

Par analogie avec les coniques, on dira alors que l'équation est :

- *hyperbolique* si  $b^2 - a \cdot c > 0$  ;
- *parabolique* si  $b^2 - a \cdot c = 0$  ;
- *elliptique* si  $b^2 - a \cdot c < 0$  .

Si  $\lambda_1$  et  $\lambda_2$  sont les racines de l'équation caractéristique, les droites  $x + \lambda_j \cdot y = c_j$ , pour  $j = 1$  et  $j = 2$ , sont appelées les *caractéristiques* de l'équation.

### 5.2.3 Cas général

De manière générale on dira que l'équation (1) est :

- *hyperbolique* si  $b^2(x, y) - a(x, y) \cdot c(x, y) > 0$  sur  $\Omega$  ;
- *parabolique* si  $b^2(x, y) - a(x, y) \cdot c(x, y) = 0$  sur  $\Omega$  ;
- *elliptique* si  $b^2(x, y) - a(x, y) \cdot c(x, y) < 0$  sur  $\Omega$  ;

où  $\Omega$  est un ouvert sur lequel est définie l'équation.

Les courbes intégrales de l'équation différentielle :

$a(x, y) \cdot (y')^2 - 2 \cdot b(x, y) \cdot y' + c(x, y) = 0$  sont alors appelées les *caractéristiques* de (1). On écrira ces courbes intégrales sous la forme  $\varphi(x, y) = c_1$  et  $\psi(x, y) = c_2$ .

Les équations elliptiques, dont les exemples typiques sont l'équation de Laplace  $\Delta u = 0$  ou l'équation de Poisson  $\Delta u = g$ , décrivent en général des *phénomènes stationnaires* (c.a.d. indépendants du temps) comme le champ électrique dans un conducteur, l'écoulement d'un fluide non compressible, le transfert de chaleur par conduction, ...

De manière générale, le changement de variables :

$$r = \frac{\varphi(x, y) + \psi(x, y)}{2}, \quad s = \frac{\varphi(x, y) - \psi(x, y)}{2 \cdot i}$$

va transformer une équation elliptique sous la forme standard :

$$\Delta u(r, s) = g\left(u, \frac{\partial u}{\partial r}, \frac{\partial u}{\partial s}, r, s\right)$$

Les équations hyperboliques ont généralement pour origines des problèmes de vibrations ou de propagation d'ondes . L'exemple typique étant l'équation des



ondes à une dimension  $\frac{\partial^2 u}{\partial t^2} - v^2 \cdot \frac{\partial^2 u}{\partial x^2} = 0$  (ici  $t$  joue le rôle de la variable  $y$ ). Ce type d'équation apparaît souvent dans le cadre de *problèmes d'évolution* (i.e. le temps intervient) avec valeurs initiales.

Le changement de variable :

$$r = \varphi(x,y), s = \psi(x,y)$$

va alors transformer une équation hyperbolique sous la forme standard :

$$\frac{\partial^2 u}{\partial r \partial s} = g\left(u, \frac{\partial u}{\partial r}, \frac{\partial u}{\partial s}, r, s\right)$$

Les équations paraboliques interviennent aussi dans les *problèmes d'évolution*.

L'exemple typique étant l'équation de diffusion  $\frac{\partial u}{\partial t} + d \cdot \frac{\partial^2 u}{\partial x^2} = 0$ .

Le changement de variable :

$$r = x, s = \varphi(x,y)$$

va transformer une équation parabolique en :

$$\frac{\partial^2 u}{\partial r^2} = g\left(u, \frac{\partial u}{\partial r}, \frac{\partial u}{\partial s}, r, s\right)$$

### 5.2.4 Exemple

L'équation du potentiel des vitesses des écoulements irrotationnels d'un fluide compressible peut s'écrire, dans le cas stationnaire :

$$(u^2 - a^2) \cdot \frac{\partial^2 \varphi}{\partial x^2} + 2 \cdot u \cdot v \cdot \frac{\partial^2 \varphi}{\partial x \partial y} + (v^2 - a^2) \cdot \frac{\partial^2 \varphi}{\partial y^2} = 0$$

où  $x$  et  $y$  sont des variables d'espace,  $\varphi$  le potentiel des vitesses,  $u$  et  $v$  les composantes de la vitesse ( $u = \frac{\partial \varphi}{\partial x}$  et  $v = \frac{\partial \varphi}{\partial y}$ ) et  $a$  est la vitesse du son qui est reliée à  $u$  et  $v$  par :

$$\frac{1}{\gamma - 1} \cdot a^2 + \frac{1}{2} \cdot (u^2 + v^2) = \text{Const.}$$

en supposant le gaz parfait à chaleur spécifique constante de rapport  $\gamma$ .

Le type de cette équation dépend de  $\delta = a^4 \cdot (M^2 - 1)$  où  $M = \frac{\sqrt{u^2 + v^2}}{a}$  est le *nombre de Mach*.

En un point où  $M < 1$  (*écoulement subsonique*),  $\delta < 0$  et l'équation est elliptique, pour  $M = 1$  (*écoulement sonique*), l'équation est parabolique et pour  $M > 1$  (*écoulement supersonique*), l'équation est hyperbolique.

On a donc ainsi un exemple où le type de l'équation peut changer d'une région à l'autre de l'espace des  $(x,y)$ .

### 5.2.5 Remarque

Du point de vue numérique, il sera plutôt important de savoir si on a affaire à un problème d'évolution avec valeurs initiales (le temps intervient) ou si on a affaire à un problème avec conditions aux bords (de type « statique »).

### 5.3 Principe de la méthode des différences finies

Comme dans le cas des équations différentielles, l'idée est de trouver des valeurs approchées de la solution sur un ensemble fini de points d'un maillage donné du domaine  $\Omega$  sur lequel on étudie l'équation aux dérivées partielles.

En dimension deux, on pourra prendre un maillage rectangulaire qui recouvre le domaine considéré .

Par exemple, pour l'équation de Laplace sur un domaine borné  $\Omega \subset \mathbb{R}^2$  (problème avec conditions aux bords) on prendra un maillage du type suivant :

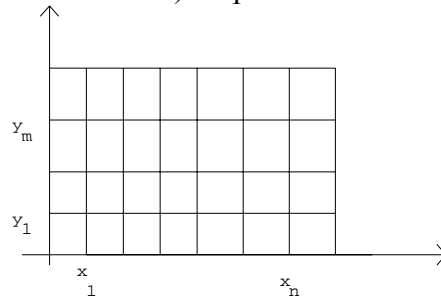


Figure 8.7

Pour l'équation des ondes à une dimension ou l'équation de la chaleur (problème d'évolution) on prendra un maillage infini du type suivant

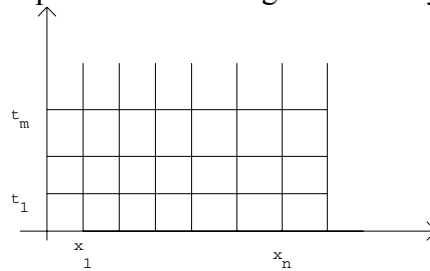


Figure 8.8

L'idée sera de remplacer chaque dérivée partielle par des différences finies et on aboutira en général à un système linéaire qui pourra être résolu soit directement soit par des méthodes itératives.

### 5.4 Cas d'un problème avec conditions au bord : l'équation de Poisson (équation de la chaleur en régime stationnaire)

On cherche donc une fonction de classe  $C^2$  sur un domaine  $\Omega$  de  $\mathbb{R}^2$  telle que :

$$\begin{cases} \Delta u(x, y) = f(x, y) \text{ à l'intérieur de } \Omega \\ u(x, y) = g(x, y) \text{ sur le bord de } \Omega \end{cases}$$

où les fonctions  $f$  et  $g$  sont données.

*Remarque* — Avec de bonnes hypothèses de régularité sur  $f$ ,  $g$  et  $\Gamma = \partial\Omega$ , on peut montrer que ce problème admet une unique solution.

*Notations* — On fait un maillage de  $\mathbb{R}^2$  de même pas  $h > 0$  dans les directions  $x$  et  $y$  dont les noeuds sont notés :

$$x_i = i \cdot h \text{ et } y_j = j \cdot h \text{ pour } i, j \text{ dans } Z$$

On note  $\Omega_h$  l'ensemble des noeuds qui sont dans  $\Omega$  et  $\Gamma_h$  l'ensemble des points d'intersection du bord  $\Gamma = \partial\Omega$  avec les lignes du maillage (a priori les points de  $\Gamma_h$  ne sont pas des noeuds du maillage).

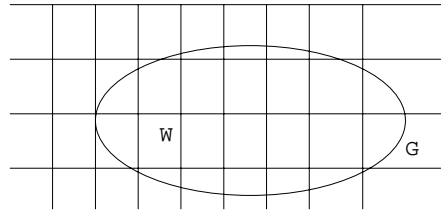


Figure 8.9

On posera aussi :

$$\begin{cases} f_{i,j} = f(x_i, y_j) \\ g_{i,j} = g(x_i, y_j) \end{cases}$$

et  $u_{i,j}$  est la valeur approchée cherchée de  $u(x_i, y_j)$  pour un point intérieur à  $\Omega$ .

Le calcul approché du laplacien d'une fonction  $\varphi$  de classe  $C^2$  dans un voisinage de  $(x_0, y_0)$  se fera à l'aide des différences finies, c'est-à-dire en posant :

$$\begin{cases} \frac{\partial^2 \varphi}{\partial x^2}(x_0, y_0) \cong \frac{\varphi(x_0 + h, y_0) - 2 \cdot \varphi(x_0, y_0) + \varphi(x_0 - h, y_0)}{h^2} \\ \frac{\partial^2 \varphi}{\partial y^2}(x_0, y_0) \cong \frac{\varphi(x_0, y_0 + h) - 2 \cdot \varphi(x_0, y_0) + \varphi(x_0, y_0 - h)}{h^2} \end{cases}$$

Le système discrétisé associé au problème de Poisson s'écrit alors :

$$\frac{u_{i+1,j} - 2 \cdot u_{i,j} + u_{i-1,j}}{h^2} + \frac{u_{i,j+1} - 2 \cdot u_{i,j} + u_{i,j-1}}{h^2} = f_{i,j}$$

pour tous les  $(i,j)$  tels que  $(x_i, y_j) \in \Omega_h$ .

En choisissant une façon de numéroter les points intérieurs, on aboutit ainsi à un système linéaire.

## 5.5 Equations elliptiques avec conditions aux bords sur un rectangle

### 5.5.1 Position du problème et notations

On considère une équation elliptique du type :

$$(1) \begin{cases} \alpha \cdot \frac{\partial^2 u}{\partial x^2} + \beta \cdot \frac{\partial u}{\partial x} + \gamma \cdot \frac{\partial^2 u}{\partial y^2} + \delta \cdot \frac{\partial u}{\partial y} + \varepsilon \cdot u = f \text{ sur } \Omega = ]0, a[ \times ]0, b[ \\ u(x, y) = g(x, y) \text{ sur } \partial\Omega \end{cases}$$

où  $\alpha, \beta, \gamma, \delta, \varepsilon, f$  et  $g$  sont des fonctions données assez régulières sur  $\Omega$  pour assurer l'existence et l'unicité d'une solution.

*Remarque 1* — Par un changement de variable toute équation elliptique d'ordre deux se ramènera à la forme ci-dessus.

*Remarque 2* — Les conditions aux bords peuvent aussi porter sur le gradient de  $u$ .

Les exemples les plus classiques de tels problèmes sont :

- *le problème de Dirichlet* — 
$$\begin{cases} \Delta u = f \text{ sur } \Omega \\ u = g \text{ sur } \partial\Omega \end{cases}$$
- *le problème de Neumann* — 
$$\begin{cases} \Delta u = f \text{ sur } \Omega \\ \frac{\partial u}{\partial n} = g \text{ sur } \partial\Omega \end{cases}$$
- *le problème mixte de Fourier* — 
$$\begin{cases} \Delta u = f \text{ sur } \Omega \\ \frac{\partial u}{\partial n} + k \cdot u = 0 \text{ sur } \partial\Omega \end{cases}$$

On suppose que les dimensions du rectangle  $\Omega$  vérifient :

$$\begin{cases} a = (n + 1) \cdot h \\ b = (m + 1) \cdot h \end{cases}$$

et on définit un maillage de  $\Omega$  en posant :

$$\begin{cases} x_i = i \cdot h \quad (i = 0, 1, \dots, n + 1) \\ y_j = j \cdot h \quad (j = 0, 1, \dots, m + 1) \end{cases}$$

Pour toute fonction  $\phi$  définie sur  $\Omega$ , on posera :

$$\phi_{ij} = \phi(x_i, y_j) \quad (0 \leq i \leq n + 1, 0 \leq j \leq m + 1)$$

La fonction  $u$  solution de (1) est donc supposée connue aux points du bord  $(x_i, 0)$ ,  $(x_i, m+1)$ ,  $(0, y_j)$  et  $(n+1, y_j)$  pour  $i = 0, \dots, n + 1$  et  $j = 0, \dots, m + 1$ .

Le problème est alors de trouver des valeurs approchées, que nous noterons  $u_{ij}$  des  $u(x_i, y_j)$  pour tous les points intérieurs de  $\Omega$ , c'est-à-dire pour  $i = 1, 2, \dots, n$  et  $j = 1, 2, \dots, m$ .

Ces valeurs inconnues seront consignées dans le vecteur :

$$u = {}^t(u_{11}, u_{21}, \dots, u_{n1}, u_{12}, u_{22}, \dots, u_{n2}, \dots, u_{1m}, u_{2m}, \dots, u_{nm})$$

ce qui revient à numéroter les points de l'intérieur de  $\Omega$  de gauche à droite en démarrant par le bas.

Les dérivées partielles sont alors approximées par :

$$\begin{aligned}\frac{\partial u}{\partial x}(x_i, y_j) &\cong \frac{u_{i+1,j} - u_{i-1,j}}{2 \cdot h} \\ \frac{\partial u}{\partial y}(x_i, y_j) &\cong \frac{u_{i,j+1} - u_{i,j-1}}{2 \cdot h} \\ \frac{\partial^2 u}{\partial x^2}(x_i, y_j) &\cong \frac{u_{i+1,j} - 2 \cdot u_{i,j} + u_{i-1,j}}{h^2} \\ \frac{\partial^2 u}{\partial y^2}(x_i, y_j) &\cong \frac{u_{i,j+1} - 2 \cdot u_{i,j} + u_{i,j-1}}{h^2}\end{aligned}$$

On dit qu'on a des approximations à 5 points des dérivées partielles.

### 5.5.2 Discrétisation du problème (1)

Le problème se discrétise donc de la façon suivante :

$$a_{i,j} \cdot u_{i+1,j} + b_{i,j} \cdot u_{i-1,j} + c_{i,j} \cdot u_{i,j+1} + d_{i,j} \cdot u_{i,j-1} + e_{i,j} \cdot u_{i,j} = f_{i,j}$$

pour  $i = 1, \dots, n, j = 1, \dots, m$ , avec :

$$\begin{aligned}a_{i,j} &= \frac{\alpha_{i,j}}{h^2} + \frac{\beta_{i,j}}{2 \cdot h} \\ b_{i,j} &= \frac{\alpha_{i,j}}{h^2} - \frac{\beta_{i,j}}{2 \cdot h} \\ c_{i,j} &= \frac{\gamma_{i,j}}{h^2} + \frac{\delta_{i,j}}{2 \cdot h} \\ d_{i,j} &= \frac{\gamma_{i,j}}{h^2} - \frac{\delta_{i,j}}{2 \cdot h} \\ e_{i,j} &= -\frac{2 \cdot \alpha_{i,j}}{h^2} - \frac{2 \cdot \gamma_{i,j}}{h^2} + \varepsilon_{i,j}\end{aligned}$$

Ce qui donne un système de  $n \cdot m$  équations à  $n \cdot m$  inconnues :

$$(2) \quad A \cdot u = v$$

où la matrice  $A$  est tridiagonale par blocs de la forme :

$$A = \begin{pmatrix} A_1 & C_1 & 0 & \dots & \dots & 0 \\ B_2 & A_2 & C_2 & 0 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \dots & 0 & B_{m-1} & A_{m-1} & C_{m-1} \\ 0 & \dots & \dots & 0 & B_m & A_m \end{pmatrix}$$

où :

$$\begin{aligned}
 A_j &= \begin{pmatrix} e_{1,j} & a_{1,j} & 0 & \dots & \dots & 0 \\ b_{2,j} & e_{2,j} & a_{2,j} & 0 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \dots & 0 & b_{n-1,j} & e_{n-1,j} & a_{n-1,j} \\ 0 & \dots & \dots & 0 & b_{n,j} & e_{n,j} \end{pmatrix} \\
 B_j &= \begin{pmatrix} d_{1,j} & 0 & 0 & \dots & \dots & 0 \\ 0 & d_{2,j} & 0 & 0 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \dots & 0 & 0 & d_{n-1,j} & 0 \\ 0 & \dots & \dots & 0 & 0 & d_{n,j} \end{pmatrix} \\
 c_j &= \begin{pmatrix} c_{1,j} & 0 & 0 & \dots & \dots & 0 \\ 0 & c_{2,j} & 0 & 0 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \dots & 0 & 0 & c_{n-1,j} & 0 \\ 0 & \dots & \dots & 0 & 0 & c_{n,j} \end{pmatrix}
 \end{aligned}$$

Le vecteur  $v$  contient les  $f_{i,j}$  et les  $u_{i,j}$  qui sont connus sur le bord.

De manière plus précise, la méthode des différences finies sur la ligne  $j$  va se traduire par :

$$(B_j, A_j, C_j) \cdot \begin{pmatrix} u_{1,j-1} \\ \cdot \\ \cdot \\ u_{n,j-1} \\ u_{1,j} \\ \cdot \\ \cdot \\ u_{n,j} \\ u_{1,j+1} \\ \cdot \\ \cdot \\ u_{n,j+1} \end{pmatrix} = \begin{pmatrix} f_{1,j} - b_{1,j} \cdot u_{0,j} \\ f_{2,j} \\ \cdot \\ \cdot \\ f_{n-1,j} \\ f_{n,j} - a_{n,j} \cdot u_{n+1,j} \end{pmatrix}$$

pour  $j = 2, \dots, m - 1$ .

Et pour  $j = 1, j = m$ , on a :

$$\begin{aligned}
 (A_1, C_1) \cdot \begin{pmatrix} u_{1,1} \\ \cdot \\ \cdot \\ u_{n,1} \\ u_{1,2} \\ \cdot \\ \cdot \\ u_{n,2} \end{pmatrix} &= \begin{pmatrix} f_{1,1} - b_{1,1} \cdot u_{0,1} - d_{1,1} \cdot u_{1,0} \\ f_{2,1} - d_{2,1} \cdot u_{2,0} \\ \cdot \\ \cdot \\ f_{n-1,1} - d_{n-1,1} \cdot u_{n-1,0} \\ f_{n,1} - a_{n,1} \cdot u_{n+1,1} - d_{n,1} \cdot u_{n,0} \end{pmatrix} \\
 (B_m, A_m) \cdot \begin{pmatrix} u_{1,m-1} \\ \cdot \\ \cdot \\ u_{n,m-1} \\ u_{1,m} \\ \cdot \\ \cdot \\ u_{n,m} \end{pmatrix} &= \begin{pmatrix} f_{1,m} - b_{1,m} \cdot u_{0,m} - d_{1,m} \cdot u_{1,m-1} \\ f_{2,m} - d_{2,m} \cdot u_{2,m+1} \\ \cdot \\ \cdot \\ f_{n-1,m} - d_{n-1,m} \cdot u_{n-1,m+1} \\ f_{n,m} - a_{n,m} \cdot u_{n+1,m} - d_{n,m} \cdot u_{n,m+1} \end{pmatrix}
 \end{aligned}$$

*Remarque* — Le système obtenu est tridiagonal par blocs avec une matrice très creuse (sur chaque ligne, il y a au plus cinq termes non nuls). Si pour résoudre ce système, on emploie une méthode directe, par exemple la méthode L-R de Crout (ou de Cholesky quand la matrice est symétrique définie positive), le calcul des coefficients des matrices L et R va demander environ  $(n \cdot m)^2$  opérations élémentaires, ce qui est prohibitif pour  $n$  et  $m$  grands. Pour cette raison on préférera utiliser une méthode itérative (Gauss-Seidel ou relaxation) qui va conserver la structure creuse de la matrice de départ et demander un nombre de calculs beaucoup moins important.

### 5.5.3 Méthode SOR (Simultaneous Over-Relaxation) pour résoudre un système tridiagonal par blocs

Dans les méthodes itératives de Gauss-Seidel et de relaxation, on avait utilisé les notations suivantes :

$$A = D + E + F$$

où  $D$  est la diagonale de  $A$ ,  $E$  est le triangle inférieur strict et  $F$  est le triangle supérieur strict.

L'algorithme de Gauss-Seidel pour résoudre  $A \cdot u = v$  est alors le suivant :

$$\begin{cases} u^{(0)} \text{ donné} \\ (D + E) \cdot u^{(k)} = -F \cdot u^{(k-1)} + v \quad (k \geq 1) \end{cases}$$

La méthode de relaxation consiste à "corriger" l'algorithme de Gauss-Seidel à l'aide d'un facteur  $\omega$  de la façon suivante :

$$\begin{cases} u^{(0)} \text{ donné} \\ (D + \omega \cdot E) \cdot u^{(k)} = (1 - \omega) \cdot D \cdot u^{(k-1)} + \omega \cdot (-F \cdot u^{(k-1)} + v) \quad (k \geq 1) \end{cases}$$

Ce qui donne :

$$u_i^{(k)} = \frac{1}{a_{ii}} \cdot \left\{ (1-\omega) \cdot a_{ii} \cdot u_i^{(k-1)} - \omega \cdot \sum_{j=1}^{i-1} a_{i,j} \cdot u_j^{(k)} - \omega \cdot \sum_{j=i+1}^n a_{i,j} \cdot u_j^{(k-1)} + \omega \cdot v_i \right\}$$

Ce qui peut aussi s'écrire, pour  $k \geq 1$  :

$$u_i^{(k)} = u_i^{(k-1)} - \omega \cdot \xi_i^{(k)} \quad (1 \leq i \leq n)$$

où  $\xi_i^{(k)}$  est le terme résiduel défini par :

$$\xi_i^{(k)} = \sum_{j=1}^{i-1} a_{i,j} \cdot u_j^{(k)} + \sum_{j=i+1}^n a_{i,j} \cdot u_j^{(k-1)} - v_i$$

Pour le type de matrice qui nous intéresse, on a alors les résultats suivants :

*Théorème* : Si A est une matrice à diagonale strictement dominante telle que toutes les valeurs propres de la matrice de Jacobi correspondante soient réelles, alors la méthode de relaxation converge si, et seulement si  $\omega \in ]0,2[$  et le paramètre de relaxation optimal est donné par :

$$\omega_0 = \frac{2}{1 + \sqrt{1 - \rho(J)^2}}$$

Dans ces conditions, le rayon spectral de la matrice  $L_{\omega_0}$  qui intervient dans la méthode de relaxation est donné par :

$$\rho(L_{\omega_0}) = \omega_0 - 1 = \left( \frac{\rho(J)}{1 + \sqrt{1 - \rho(J)^2}} \right)^2$$

*Démonstration* — Voir Ciarlet p. 106 .

Si on n'a pas d'idée sur  $\rho(J)$  (le rayon spectral de la matrice d'itérations dans la méthode de Jacobi) on pourra essayer plusieurs valeurs de  $\omega$  dans  $]1,2[$  (on dit alors qu'on utilise une méthode de sur-relaxation).

Pour le problème qui nous concerne, on utilise plutôt la matrice A d'ordre  $((n+2) \cdot (m+2))^2$  qui fait intervenir tous les coefficients  $u_{i,j}$  pour  $i = 0, \dots, n+1$  et  $j = 0, \dots, m+1$ . Ce qui donne, pour le terme résiduel à l'ordre  $k-1$  :

$$\xi_{i,j}^{(k-1)} = a_{i,j} \cdot u_{i+1,j}^{(k-1)} + b_{i,j} \cdot u_{i-1,j}^{(k)} + c_{i,j} \cdot u_{i,j+1}^{(k-1)} + d_{i,j} \cdot u_{i,j-1}^{(k)} + e_{i,j} \cdot u_{i,j}^{(k-1)} - f_{i,j}$$

Et alors :

$$u_{i,j}^{(k)} = u_{i,j}^{(k-1)} - \omega \cdot \frac{\xi_{i,j}^{(k-1)}}{e_{i,j}}$$

Comme test d'arrêt, on pourra prendre :

$$\|\xi^{(k-1)}\| < \varepsilon \cdot \|f\|$$

où f désigne le vecteur de coordonnées  $f_{i,j}$  et  $\varepsilon > 0$  est une précision donnée.

*Remarque 1* — Si on sépare les coordonnées du vecteur  $u^{(k)}$  en deux groupes « pair » et « impair » (la « parité » de  $u_{i,j}^{(k)}$  étant celle de  $i+j$ ), alors les



coordonnées paires de  $u^{(k)}$  ne dépendent que des coordonnées impaires de  $u^{(k-1)}$  et vice versa. Il nous suffit donc, à chaque étape, de calculer seulement la moitié des coordonnées.

*Remarque 2 : Procédé d'accélération de la convergence de Tchébycheff* — Quand on a la possibilité de le connaître, le choix de  $\omega$  optimal dès le début des itérations n'est pas le plus judicieux, un minimum de  $n$  itération sera nécessaire à la convergence. L'idée de Tchébycheff est de changer le paramètre  $\omega$  à chaque itération de la manière suivante : on part de  $\omega^{(0)} = 1$ , puis on prend  $\omega^{(1)} =$

$\frac{1}{1 - \frac{\rho(J)^2}{2}}$  et à l'étape  $n \geq 1$ , on prendra :

$$\omega^{(n+1)} = \frac{1}{1 - \frac{\rho(J)^2 \cdot \omega^{(n)}}{4}}$$

La suite  $(\omega^{(n)})$  converge en fait vers  $\omega$  optimal et le procédé de Tchébycheff diminuera, en général, le nombre d'itérations.

#### 5.5.4 Programmation structurée de la méthode S.O.R. avec accélération de Tchébycheff

Au préalable on a initialisé les  $u_{i,0}$ ,  $u_{i,m+1}$ ,  $u_{0,j}$  et  $u_{n+1,j}$ .

PROCEDURE S.O.R. (Entrée :  $n, m$  : Entier ;  $\rho$  : Réel ;  $A, B, C, D, E, F$  : Matrice ;  
Entrée\_Sortie  $U$  : Matrice) ;

Début

$\omega = 1$  ;

$k := 1$  ;

$\|f\| = 0$  ;

Pour  $i$  allant de 1 à  $n$  faire

Début

Si  $|f_{i,j}| > \|f\|$

Alors faire  $\|f\| = |f_{i,j}|$  ;

Fin ;

Répéter

NormeRésidu = 0 ;

Pour  $i$  allant de 1 à  $n$  faire

Début

Pour  $j$  allant de 1 à  $m$  faire

Début

Si  $(i + j$  et  $k$  de même parité)

Alors faire

Début

Résidu =  $a_{i,j}u_{i+1,j} + b_{i,j}u_{i-1,j} + c_{i,j}u_{i,j+1} + d_{i,j}u_{i,j-1} + e_{i,j}u_{i,j} - f_{i,j}$  ;

Si  $|Résidu| > NormeRésidu$

Alors faire NormeRésidu =  $|Résidu|$  ;

$$u_{i,j} = u_{i,j} - \omega \frac{\text{Résidu}}{e_{i,j}} ;$$

Fin ;

Fin ;

Fin ;

Si  $k = 1$

$$\text{Alors } \omega = \frac{1}{1 - \frac{\rho(J)^2}{2}} ;$$

$$\text{Sinon faire } \omega = \frac{1}{1 - \frac{\rho(J)^2 \cdot \omega}{4}} ;$$

$k = k + 1 ;$

Jusqu'à ( $k = \text{MaxItération}$ ) ou ( $\text{NormeRésidu} < \varepsilon \cdot \|f\|$ ) ;

Si  $k = \text{MaxItération}$

Alors Afficher('Méthode inadaptée : convergence trop lente ou non convergence. ') ;

Fin ;

## 5.6 Cas particulier de l'équation de Poisson

Dans ce cas, on a :

$$a_{i,j} = b_{i,j} = c_{i,j} = d_{i,j} = 1, e_{i,j} = -4 \text{ et } f_{i,j} = h_2 \cdot f(x_i, y_j)$$

De plus, on connaît explicitement  $\rho(J)$ , il est donné par :

$$\rho(J) = \frac{\text{Cos}\left(\frac{\pi}{n+1}\right) + \text{Cos}\left(\frac{\pi}{m+1}\right)}{2}$$

De manière plus générale, pour des pas  $\delta x$  et  $\delta y$  différents, on a :

$$\rho(J) = \frac{\text{Cos}\left(\frac{\pi}{n+1}\right) + \left(\frac{\delta x}{\delta y}\right)^2 \cdot \text{Cos}\left(\frac{\pi}{m+1}\right)}{1 + \left(\frac{\delta x}{\delta y}\right)^2}$$

(Cf. Stoer et Burlisch p. 558) .

*Remarque* — Pour  $f = 0$ , on a l'équation de Laplace  $\Delta u = 0$  et la méthode des différences finies s'exprime en disant que  $u_{i,j}$  est la moyenne des quatre points qui l'entourent  $u_{i-1,j}$ ,  $u_{i+1,j}$ ,  $u_{i,j-1}$  et  $u_{i,j+1}$ . C'est simplement la version discrète de la propriété de la moyenne pour les fonctions harmoniques (une fonction est harmonique si, et seulement si elle vérifie  $\Delta u = 0$ , ce qui équivaut à dire qu'elle vérifie la propriété de la moyenne, et de cette propriété on déduit que le maximum de  $u$  ne peut être atteint qu'en un point du bord. — cf. Cartan ch. 7 S 4 —).

### 5.7 Exemple d'application : distribution de la température en régime stationnaire dans une cheminée

Le problème de la distribution de la température, en régime stationnaire, dans une cheminée nous conduit à résoudre le problème de Dirichlet :

$$\begin{cases} \Delta u = 0 \text{ sur } \Omega = R - R' \\ u = g \text{ sur } \partial R \\ u = h \text{ sur } \partial R' \end{cases}$$

où  $R' \subset R$  sont deux rectangles.

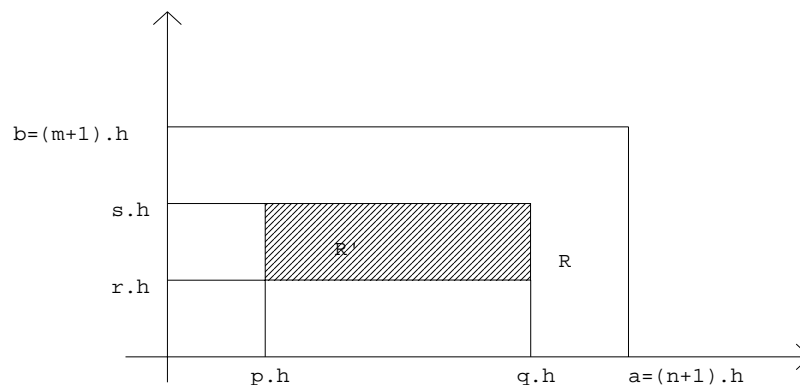


Figure 8.10

En modifiant la procédure ci-dessus, on écrit facilement un programme qui permet de calculer la température en tout point intérieur de  $R-R'$ .

### 5.8 Premier exemple d'équation parabolique : l'équation de la chaleur à une dimension

#### 5.8.1 Introduction

Dans ce paragraphe, on s'intéresse au problème de la propagation de la chaleur par conduction dans une barre isolée de sorte que les échanges de chaleur avec l'extérieur soient limités aux extrémités maintenues à température constante, la distribution de température à un instant initial  $t = 0$  étant donnée.

Avec les notations du paragraphe (5.1.3), on est donc ramené à trouver une fonction  $u$  des deux variables  $t \geq 0$  et  $x \in [0, L]$  solution du problème d'évolution avec conditions initiales :

$$(1) \quad \begin{cases} \frac{\partial u}{\partial t}(x, t) = k \cdot \frac{\partial^2 u}{\partial x^2}(x, t), \text{ pour } t > 0 \text{ et } 0 < x < L \\ u(x, 0) = f(x), \text{ pour } 0 \leq x \leq L \\ u(0, t) = u(L, t) = 0, \text{ pour } t \geq 0 \end{cases}$$

où la constante  $k$  est le coefficient de diffusion et  $L > 0$  est la longueur de la barre.  
*Remarque* — Par continuité, on a nécessairement  $f(0) = f(L) = 0$ .

**5.8.2 Existence et unicité de solutions de (1)**

Pour simplifier, on fait  $k = 1$ .

*Théorème* : Une solution de (1) est donnée par :

$$u(x,t) = \frac{1}{2 \cdot \sqrt{\pi \cdot t}} \cdot \int_{\mathbb{R}} e^{-\frac{(x-y)^2}{4t}} \cdot f(y) dy$$

*Démonstration* — On suppose que les fonctions  $f, x \mapsto u(x,t), x \mapsto \frac{\partial u}{\partial t}(x,t)$  et  $x \mapsto \frac{\partial^2 u}{\partial x^2}(x,t)$  sont dans  $L^1(\mathbb{R})$  (i.e. absolument intégrables) avec  $\left| \frac{\partial u}{\partial t}(x,t) \right| \leq g(x)$  où  $g \in L^1(\mathbb{R})$ .

En notant :

$$\hat{u}(v,t) = \int_{\mathbb{R}} e^{-2ivx} \cdot u(x,t) dx$$

la transformée de Fourier en  $x$  de  $u$ , l'équation (1) devient :

$$(1') \quad \begin{cases} \frac{\partial \hat{u}}{\partial t} + 4 \cdot \pi^2 \cdot v^2 \cdot \hat{u} = 0 \\ \hat{u}(v,0) = f(v) \\ \hat{u}(0,t) = \hat{u}(L,t) = 0 \end{cases}$$

En intégrant par rapport à  $t$ , on déduit que :

$$\hat{u}(v,t) = f(v) \cdot e^{-4\pi^2 \cdot v^2 \cdot t}$$

Mais  $e^{-4\pi^2 \cdot v^2 \cdot t}$  est la transformée de Fourier de  $\frac{1}{\sqrt{4 \cdot \pi \cdot t}} \cdot e^{-\frac{x^2}{4t}}$  et la transformation de Fourier transformant le produit de convolution en produit ordinaire, avec l'injectivité de cette transformation, on déduit le résultat.

*Remarque* — Du théorème précédent, on déduit que :

$$\lim_{t \rightarrow +\infty} u(x,t) = 0, \text{ pour tout } x \in [0,L],$$

si  $u$  est solution de (1) (ce qui est évident d'un point de vue physique si  $u$  représente une température).

Pour montrer l'unicité de  $u$ , on utilise la quantité :

$$E(t) = \int_0^L u^2(x,t) \frac{dx}{2}$$

qui correspond à une *énergie cinétique* de la barre à l'instant  $t$ .

*Théorème* : L'énergie cinétique à l'instant  $t$ , vérifie les propriétés suivantes :

- (i)  $\forall t \geq 0, E(t) \geq 0$  ;
- (ii)  $E$  est décroissante ;

(iii)  $E$  est continue à droite en  $0$  ;  
 (iv)  $E = 0 \Leftrightarrow u = 0$  ;  
 cette dernière propriété assurant l'unicité de la solution de (1).

*Démonstration* — On a,  $E'(t) = \int_0^L u(x,t) \cdot \frac{\partial u}{\partial t}(x,t) dx = \int_0^L u(x,t) \cdot \frac{\partial^2 u}{\partial x^2}(x,t) dx$  et une intégration par parties donne  $E'(t) = -\int_0^L \left( \frac{\partial u}{\partial x}(x,t) \right)^2 dx \leq 0$ . Donc  $E$  est décroissante.

Avec un théorème de convergence dominée de Lebesgue, on déduit que  $\lim_{t \rightarrow 0} E(t) = E(0)$ , c'est-à-dire que  $E$  est continue à droite en  $0$ .

Si  $E = 0$ , avec la continuité de  $u$  on déduit que  $u = 0$ . D'où l'unicité.

*Conclusion* — Le problème (1) admet une unique solution nulle à l'infini.

*Remarque* — La solution ainsi obtenue n'est pas satisfaisante d'un point de vue numérique. Et la situation sera encore moins satisfaisante si on ajoute un second membre non nul à l'équation (1), ou si le coefficient de diffusion  $\sqrt{k}$  est non constant ou encore si on s'intéresse à l'équation de la chaleur en dimension deux.

On utilisera donc une méthode de différences finies pour obtenir des valeurs approchées de  $u$ .

### 5.8.3 Résolution de l'équation de la chaleur à une dimension par des méthodes de différences finies

On se donne un maillage de  $[0, +\infty[ \times ]0, L]$  en faisant : une discrétisation en temps de pas  $\Delta t > 0$  : c'est-à-dire en posant,  $t_j = j \cdot \Delta t$ , pour tout  $j \in \mathbb{N}$  ; et une discrétisation en espace de pas  $\Delta x = \frac{L}{n+1}$ , où  $n \geq 1$  : c'est-à-dire en posant,  $x_i = i \cdot \Delta x$ , pour tout  $i = 0, 1, \dots, n+1$ .

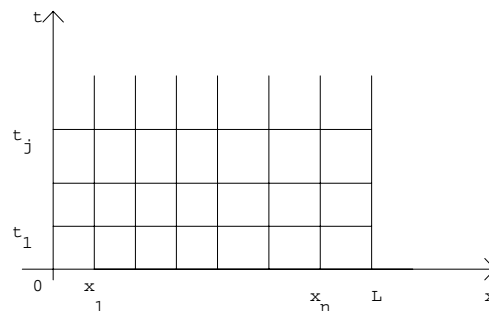


Figure 8.11

On notera  $u_{i,j}$  une valeur approchée (inconnue) de  $u(x_i, t_j)$  pour  $1 \leq i \leq n$  et  $j \geq 1$ . Les valeurs connues étant les :

$$\begin{cases} u_{i,0} = u(x_i, 0) = f(x_i) = f_i, \text{ pour } 0 \leq i \leq n + 1 \\ u_{0,j} = u(0, t_j) = 0, \text{ pour } j \geq 0 \\ u_{n+1,j} = u(L, t_j) = 0, \text{ pour } j \geq 0 \end{cases}$$

On consignera les  $u_{i,j}$ , à  $j$  fixé, dans un vecteur  $u_j = \begin{pmatrix} u_{1,j} \\ \vdots \\ u_{n,j} \end{pmatrix}$ .

Les dérivées partielles d'ordre 2 par rapport à  $x$ , seront approximées par différences finies centrées, soit :

$$\frac{\partial^2 u}{\partial x^2}(x_i, t_j) \cong \frac{u_{i+1,j} - 2 \cdot u_{i,j} + u_{i-1,j}}{\Delta x^2}$$

Pour calculer les  $u_j$  par une récurrence d'ordre 1 ( $u_0$  étant connu), on a le choix pour approximer  $\frac{\partial u}{\partial t}(x_i, t_j)$  entre les différences finies progressives ou régressives.

En utilisant les *différences finies progressives*, on a le schéma :

$$\frac{u_{i,j+1} - u_{i,j}}{\Delta t} = k \cdot \frac{u_{i+1,j} - 2 \cdot u_{i,j} + u_{i-1,j}}{\Delta x^2} \quad (1 \leq i \leq n \text{ et } j \geq 0)$$

soit, pour tout  $j \geq 0$  :

$$u_{i,j+1} = u_{i,j} + k \cdot \frac{\Delta t}{\Delta x^2} \cdot (u_{i+1,j} - 2 \cdot u_{i,j} + u_{i-1,j}) \quad (1 \leq i \leq n)$$

En posant  $r = \frac{\Delta t}{\Delta x^2}$  et en tenant compte du fait que  $u_{0,j} = u_{n+1,j} = 0$ , on peut écrire les choses sous la forme matricielle suivante :

$$u_{j+1} = u_j + k \cdot r \cdot A \cdot u_j, \text{ pour } j \geq 0$$

où on a posé :

$$A = \begin{pmatrix} -2 & 1 & 0 & \dots & \dots & 0 \\ 1 & -2 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & \dots & \dots & 1 & -2 \end{pmatrix}$$

Donc les  $u_j$  sont déterminés par le *schéma explicite* suivant :

$$(2) \quad \begin{cases} u_0 \text{ donné} \\ u_j = B \cdot u_{j-1} \quad (j \geq 1) \end{cases}$$

où  $B = I_d + k \cdot r \cdot A$ .

En utilisant les *différences finies régressives*, on a le schéma suivant :

$$\frac{u_{i,j} - u_{i,j-1}}{\Delta t} = k \cdot \frac{u_{i+1,j} - 2 \cdot u_{i,j} + u_{i-1,j}}{\Delta x^2} \quad (1 \leq i \leq n \text{ et } j \geq 0)$$

soit, pour  $j \geq 1$  :

$$u_{i,j} - k \cdot r \cdot (u_{i+1,j} - 2 \cdot u_{i,j} + u_{i-1,j}) = u_{i,j-1} \quad (1 \leq i \leq n)$$

Toujours avec  $u_{0,j} = u_{n+1,j} = 0$  et en notant  $C = I_d - k \cdot r \cdot A$ , les  $u_j$  sont définis par le schéma implicite :

$$(3) \quad \begin{cases} u_0 \text{ donné} \\ C \cdot u_j = u_{j-1} \quad (j \geq 1) \end{cases}$$

Dans le schéma implicite,  $u_j$  est solution d'un système linéaire tridiagonal de matrice constante  $C$ . Il semble donc, a priori, que ce schéma soit moins intéressant car plus coûteux en temps de calcul. En fait, on verra par la suite que ce coût est compensé par une meilleure « stabilité » de la méthode.

*Remarque* — Pour chacune des deux méthodes, l'erreur par rapport à la solution exacte est un  $O(h)$ . Une idée pour améliorer la précision est alors de faire la moyenne entre les méthodes implicites et explicites, ce qui donne le schéma de Crank-Nicholson :

$$\frac{u_{i,j+1} - u_{i,j}}{\Delta t} = \frac{k}{2} \cdot \frac{u_{i+1,j+1} - 2 \cdot u_{i,j+1} + u_{i-1,j+1} + u_{i+1,j} - 2 \cdot u_{i,j} + u_{i-1,j}}{\Delta x^2} \quad (1 \leq i \leq n \text{ et } j \geq 0)$$

Ce qui donne un schéma implicite, écrit sous forme matricielle :

$$(4) \quad \begin{cases} u_0 \text{ donné} \\ \left( I_d - \frac{k}{2} \cdot r \cdot A \right) \cdot u_j = \left( I_d + \frac{k}{2} \cdot r \cdot A \right) \cdot u_{j-1} \quad (j \geq 1) \end{cases}$$

Dans ce cas, non seulement la méthode est stable, mais elle est aussi d'ordre 2, c'est-à-dire que l'erreur est un  $O(h^2)$ .

De manière plus générale, en prenant les moyennes des schémas implicites et explicites pondérées par  $\omega$  et  $1 - \omega$ , avec  $0 \leq \omega \leq 1$ , on a le schéma :

$$\frac{u_{i,j+1} - u_{i,j}}{\Delta t} = \omega \cdot \frac{k}{\Delta x^2} \cdot \{ u_{i+1,j+1} - 2 \cdot u_{i,j+1} + u_{i-1,j+1} \} + (1 - \omega) \cdot \frac{k}{\Delta x^2} \cdot \{ u_{i+1,j} - 2 \cdot u_{i,j} + u_{i-1,j} \}$$

soit, sous forme matricielle :

$$(5) \quad \begin{cases} u_0 \text{ donné} \\ (I_d - \omega \cdot k \cdot r \cdot A) \cdot u_j = (I_d + (1 - \omega) \cdot k \cdot r \cdot A) \cdot u_{j-1} \quad (j \geq 1) \end{cases}$$

Pour  $\omega = 0$ , on a le schéma explicite donné par les différences finies progressives, pour  $\omega = 1$ , on a le schéma implicite donné par les différences finies régressives et pour  $\omega = \frac{1}{2}$ , on a le schéma implicite de Crank-Nicholson.

La matrice  $(I_d - \omega \cdot k \cdot r \cdot A)$  est toujours symétrique définie positive car symétrique, à diagonale strictement dominante et de coefficients diagonaux tous strictement positifs. On aura donc intérêt, pour résoudre le système linéaire donnant  $u_j$ , dans les méthodes implicites, à utiliser la méthode de décomposition de Cholesky, qui va se simplifier à cause du caractère tridiagonal des matrices.

#### 5.8.4 Problèmes de stabilité

Dans ce paragraphe, on va se contenter d'une présentation succincte de la notion de stabilité d'une méthode numérique pour résoudre l'équation de la chaleur. Pour plus de détails sur ce sujet, on pourra se reporter à Euvrard (p. 71 à 84) ou Numerical Recipes (p. 625 à 640)

On a vu que si  $u$  est solution de (1), alors  $\lim_{t \rightarrow +\infty} u(x, t) = 0$ , pour tout  $x \in [0, L]$ .

Une condition nécessaire de stabilité, pour une méthode itérative donnant les valeurs approchées  $u_{i,j}$  des  $u(x_i, t_j)$ , sera alors :  $\lim_{j \rightarrow +\infty} u_j = 0$ , pour tout  $n \geq 1$ .

Si on considère le schéma (5), avec  $0 \leq \omega \leq 1$ , en notant :

$$B_\omega = I_d - \omega \cdot k \cdot r \cdot A \text{ et } C_\omega = I_d + (1 - \omega) \cdot k \cdot r \cdot A$$

on a :

$$u_j = B_\omega^{-1} \cdot C_\omega \cdot u_{j-1} = (B_\omega^{-1} \cdot C_\omega)^j \cdot u_0, \text{ pour } j \geq 0$$

et la condition de stabilité va se traduire par :

$$\lim_{j \rightarrow +\infty} (B_\omega^{-1} \cdot C_\omega)^j = 0, \text{ pour tout } n \geq 1$$

ce qui équivaut à :

$$\sigma(B_\omega^{-1} \cdot C_\omega) < 1, \text{ pour tout } n \geq 1$$

où  $\sigma$  désigne le rayon spectral.

*Lemme* : Les valeurs propres de la matrice  $A = \begin{pmatrix} -2 & 1 & 0 & \dots & 0 \\ 1 & -2 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & 1 & -2 \end{pmatrix}$  sont

les :

$$\lambda_k = -4 \cdot \text{Sin}^2(\theta_k), \text{ où } \theta_k = \frac{k \cdot \pi}{2 \cdot (n+1)} \text{ (} k = 1, \dots, n \text{)}$$

*Démonstration* — Dire que  $\lambda$  est valeur propre de  $A$ , équivaut à dire qu'il existe un vecteur non nul  $x$  tel que  $A \cdot x = \lambda \cdot x$ , ce qui équivaut au système :

$$x_{i-1} - 2 \cdot x_i + x_{i+1} = \lambda \cdot x_i \quad (i = 1, \dots, n)$$

en posant  $x_0 = x_{n+1} = 0$ .

Si on pose  $x_i = \text{Sin}(i \cdot \theta)$ , pour  $i = 0, 1, \dots, n + 1$ , on a :

$$x_{i-1} + x_{i+1} = 2 \cdot \cos(\theta) \cdot x_i$$

et il vient :  $2 \cdot (\cos(\theta) - 1) \cdot x_i = \lambda \cdot x_i$ . Ce qui donne :

$$\lambda = 2 \cdot (\cos(\theta) - 1) = -4 \cdot \text{sin}^2(\theta/2)$$

Enfin, en tenant compte de  $x_{n+1} = \text{sin}((n+1) \cdot \theta) = 0$ , on déduit que les seules valeurs possibles de  $\theta$  sont les  $\frac{k \cdot \pi}{n+1}$ , pour  $k = 1, \dots, n$ . D'où le résultat.

On déduit donc que les valeurs propres de  $B_\omega^{-1} \cdot C_\omega$  sont les :

$$\mu_k = \frac{1 + (1 - \omega) \cdot k \cdot r \cdot \lambda_k}{1 - \omega \cdot k \cdot r \cdot \lambda_k} = \frac{1 - 4 \cdot (1 - \omega) \cdot k \cdot r \cdot \text{Sin}^2(\theta_k)}{1 + 4 \cdot \omega \cdot k \cdot r \cdot \text{Sin}^2(\theta_k)}$$

Et la condition de stabilité va se traduire par :

$$2 \cdot k \cdot r \cdot \text{sin}^2(\theta_k) \cdot (1 - 2 \cdot \omega) < 1 \quad (1 \leq k \leq n \text{ et pour tout } n \geq 1)$$

Pour  $\omega \in [\frac{1}{2}, 1]$ , cette condition est toujours vérifiée, on dit alors que la méthode est *inconditionnellement stable*.



Si  $\omega \in [0, \frac{1}{2}[$ , on a une condition de stabilité sur  $r = \frac{\Delta t}{\Delta x^2}$  donnée par :

$$\sin^2(\theta_k) \cdot r < \frac{1}{2 \cdot k \cdot (1 - 2 \cdot \omega)}, \text{ pour } k = 1, \dots, n \text{ et } n \geq 1$$

ce qui équivaut à :

$$\sin^2(\theta_n) \cdot r < \frac{1}{2 \cdot k \cdot (1 - 2 \cdot \omega)}, \text{ pour tout } n \geq 1$$

et en faisant  $n \rightarrow +\infty$ , on a une condition nécessaire de stabilité :

$$r < \frac{1}{2 \cdot k \cdot (1 - 2 \cdot \omega)}$$

*En résumé :*

- Pour  $\frac{1}{2} < \omega < 1$ , la méthode est *inconditionnellement stable* ;
- pour  $0 < \omega < \frac{1}{2}$ , la méthode est *conditionnellement stable*.

En particulier, les méthodes implicites pour  $\omega = 1$  et  $\omega = \frac{1}{2}$  sont inconditionnellement stables et la méthode explicite pour  $\omega = 0$  est conditionnellement stable, la condition de stabilité étant  $r < \frac{1}{2 \cdot k}$ .

### 5.9 Equations paraboliques linéaires à une variable d'espace avec conditions initiales

On considère ici le cas d'une équation parabolique, avec conditions initiales du type :

$$\begin{cases} \frac{\partial u}{\partial t}(x, t) = a(x) \cdot \frac{\partial^2 u}{\partial x^2}(x, t) + b(x) \cdot \frac{\partial u}{\partial x}(x, t) + c(x) \cdot u(x, t) + d(x) & (0 \leq x \leq L, t \geq 0) \\ u(x, 0) = h(x), \text{ pour } 0 \leq x \leq L \\ u(0, t) = f(t); u(L, t) = g(t) & (t \geq 0) \end{cases}$$

On considère le même maillage qu'en (5.8) et on garde les mêmes notations.

Le schéma (5) du paragraphe (5.8.3), avec paramètre  $\omega \in [0, 1]$  s'écrit alors :

$$\begin{aligned} \frac{u_{i,j+1} - u_{i,j}}{\Delta t} = \omega \cdot \left\{ a_i \cdot \frac{u_{i+1,j+1} - 2 \cdot u_{i,j+1} + u_{i-1,j+1}}{\Delta x^2} + b_i \cdot \frac{u_{i+1,j+1} - u_{i-1,j+1}}{2 \cdot \Delta x} + c_i \cdot u_{i,j+1} + d_i \right\} \\ + (1 - \omega) \cdot \left\{ a_i \cdot \frac{u_{i+1,j} - 2 \cdot u_{i,j} + u_{i-1,j}}{\Delta x^2} + b_i \cdot \frac{u_{i+1,j} - u_{i-1,j}}{2 \cdot \Delta x} + c_i \cdot u_{i,j} + d_i \right\} \end{aligned}$$

pour  $j \geq 0$  et  $i = 1, \dots, n$ .

Ce qui peut s'écrire :

$$\beta_i \cdot u_{i-1,j+1} + \alpha_i \cdot u_{i,j+1} + \gamma_i \cdot u_{i+1,j+1} = \beta'_i \cdot u_{i-1,j} + \alpha'_i \cdot u_{i,j} + \gamma'_i \cdot u_{i+1,j} + d_i$$

pour  $j \geq 0$  et  $i = 1, \dots, n$ ,

avec :

$$\begin{cases} \beta_i = -\omega \cdot \Delta t \cdot \left\{ \frac{a_i}{\Delta x^2} - \frac{b_i}{2 \cdot \Delta x} \right\} \\ \alpha_i = 1 - \omega \cdot \Delta t \cdot \left\{ c_i - \frac{2 \cdot a_i}{\Delta x^2} \right\} \\ \gamma_i = -\omega \cdot \Delta t \cdot \left\{ \frac{a_i}{\Delta x^2} + \frac{b_i}{2 \cdot \Delta x} \right\} \end{cases}$$

et :

$$\begin{cases} \beta'_i = (1 - \omega) \cdot \Delta t \cdot \left\{ \frac{a_i}{\Delta x^2} - \frac{b_i}{2 \cdot \Delta x} \right\} \\ \alpha'_i = 1 + (1 - \omega) \cdot \Delta t \cdot \left\{ c_i - 2 \cdot \frac{a_i}{\Delta x^2} \right\} \\ \gamma'_i = (1 - \omega) \cdot \Delta t \cdot \left\{ \frac{a_i}{\Delta x^2} + \frac{b_i}{2 \cdot \Delta x} \right\} \end{cases}$$

ce qui donne encore  $u_{j+1}$  comme solution d'un système tridiagonal.

*Remarque* — Pour  $\omega = 0$ , on a une méthode explicite, pour  $\omega = \frac{1}{2}$ , c'est la méthode de Crank-Nicholson et pour  $\omega \neq 0$  c'est une méthode implicite.

Pour  $0 \leq \omega < \frac{1}{2}$ , la méthode est conditionnellement stable et pour  $\frac{1}{2} \leq \omega \leq 1$ , la méthode est inconditionnellement stable.

## 6. Exemple d'application : Phénomène de sustentation par utilisation d'une source de pression. Le patin hydrostatique

### 6.1 Présentation du problème

Les dispositifs de lubrification (ex. glissières de machines) utilisent l'effet de sustentation hydrostatique : laminage par utilisation d'un fluide sous pression.

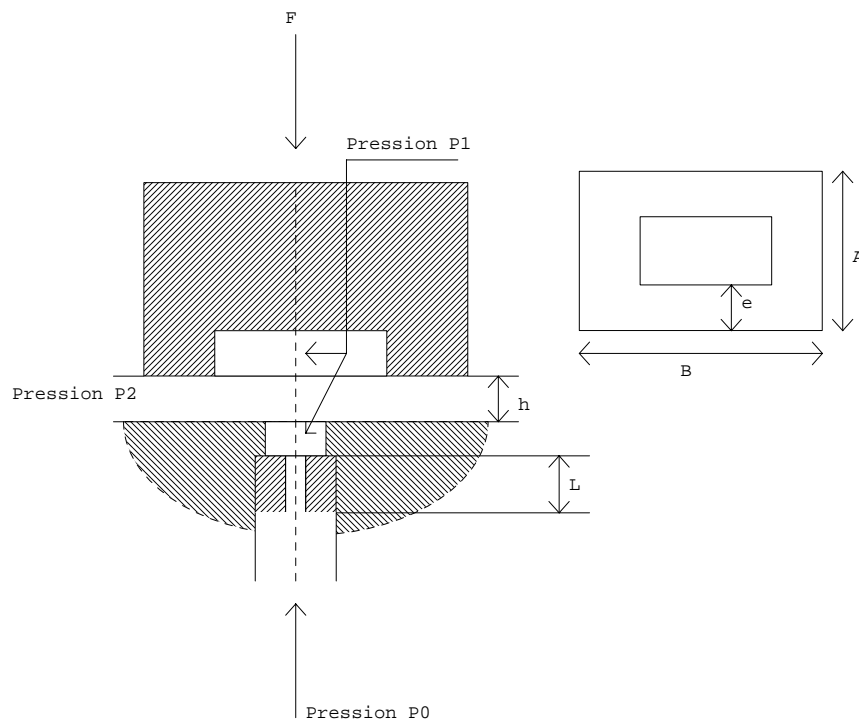


Figure 8.12

Les données mécaniques sont les suivantes : un patin rectangulaire est alimenté par un capillaire de diamètre  $d_0$  et on désigne par :

- $P_0$  la pression d'alimentation ;
- $P_1$  la pression dans la chambre ;
- $P_2$  la pression à la sortie :  $P_2 = P_{at}$  (pression atmosphérique) ;
- la pression utile est la pression effective notée  $P_{eff} = P - P_{at}$  ;

On utilise le rapport de pression :  $r = \frac{P_{eff}}{P_0}$ , pour  $0 < r < 1$

On note :

$$\left\{ \begin{array}{l} Kq \text{ le coefficient de débit} \\ Ks \text{ le coefficient de surface} \\ Kf \text{ le coefficient d'effort} \end{array} \right.$$

Pour une application numérique, on pourra prendre les valeurs suivantes :  
 $A = 50 \cdot \text{mm}$ ,  $B = 70 \cdot \text{mm}$  et  $C = 10 \cdot \text{mm}$ .

Les lois gouvernant l'équilibre du patin sont :

$$(1) \quad P_0 - P_1 = \frac{128 \cdot \mu}{\pi} \cdot \frac{L}{d_0^4} \cdot Q_{v_0} \quad (\text{Loi des écoulements laminaires})$$

$$(2) \quad P_1 - P_2 = \frac{1}{Kq} \cdot \frac{\mu \cdot Q_{v_0}}{h^3}$$

## 6.2 Détermination du coefficient $Ks$

Pour chaque hauteur de stabilisation  $h$ , on mesure :  $P_{0\text{eff}}$ ,  $P_{1\text{eff}}$  et  $F = f - P$ , où :

- $f$  désigne la charge due à la pression ;
- $F$  désigne la charge mesurée (compensation du poids du coulisseau).

On écrit alors :  $f = P_{1\text{eff}} \cdot S_{\text{eff}}$ , où  $S_{\text{eff}}$  est une surface effective inférieure à la surface réelle, pour tenir compte de la chute de pression.

On pose  $S_{\text{eff}} = Ks \cdot A \cdot B$ , avec  $0 < Ks < 1$  et on écrit  $F = Ks \cdot A \cdot B \cdot P_{1\text{eff}} + b_1$ , soit en posant  $a_1 = Ks \cdot A \cdot B$  :

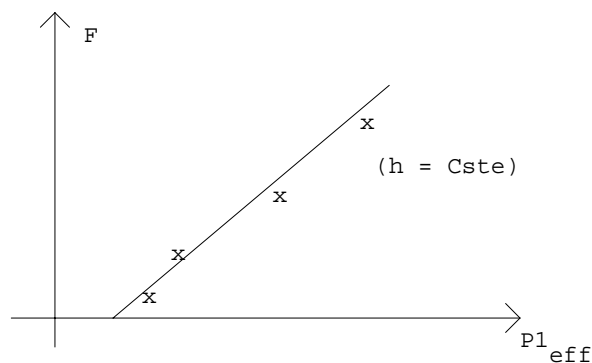


Figure 8.13

$$F = a_1 \cdot P_{1\text{eff}} + b_1$$

Le coefficient  $Ks$  peut alors s'obtenir à l'aide d'une régression affine.

*Remarques* — On aura intérêt à calculer  $Ks$  pour chaque groupe de valeurs mesurées à  $h$  constante. Dans le cas de 4 groupes de valeurs, on obtiendra des valeurs différentes de  $Ks$ . La comparaison de ces valeurs permet de juger, si la loi de comportement à adopter dans le laminage, est indépendante ou non de la hauteur.

Si les différentes valeurs des  $Ks$  varient peu, on pourra prendre la valeur moyenne  $Ks_{\text{moy}}$ , pour modéliser le comportement du patin dans son ensemble. L'extrapolation pour des valeurs de  $h$  inférieures à la valeur minimale des essais, devra cependant être envisagée avec prudence.

### 6.3 Détermination du coefficient $Kq$

$Kq$  peut être déduit des tableaux de mesures, en effet,  $\frac{(1)}{(2)}$  donne :

$$\frac{P0 - P1}{P2 - P1} = \frac{\frac{128 \cdot \mu \cdot L}{\pi \cdot d_0^4} \cdot Q_{v_0}}{\frac{1}{Kq} \cdot \frac{\mu}{h^3} \cdot Q_{v_0}}$$

$$\frac{1-r}{r} = \frac{128 \cdot Kq \cdot L \cdot h^3}{\pi \cdot d_0^4}$$

$$h = \left( \frac{\pi \cdot d_0^4}{128 \cdot Kq \cdot L} \right)^{\frac{1}{3}} \cdot \left( \frac{1-r}{r} \right)^{\frac{1}{3}}$$

Soit en posant  $X = \left( \frac{1-r}{r} \right)^{\frac{1}{3}}$  :  $h = a \cdot X$ .

Si  $Kq$  est constant, alors  $a$  est constant et une régression affine donnera  $Kq$ .

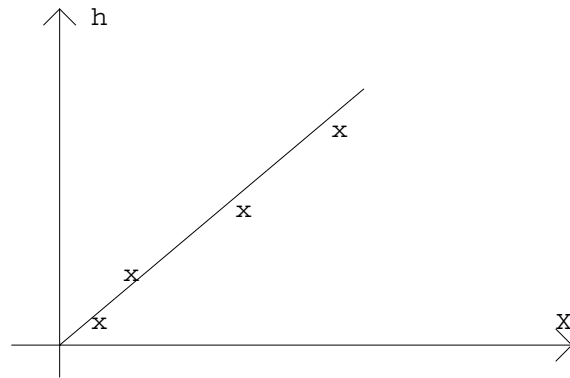


Figure 8.14

### 6.4 Raideur du patin

Il s'agit d'un paramètre mécanique d'une grande importance, notamment dans l'application aux glissières ; elle traduit l'aptitude à la variation de  $h$  sous l'action d'une variation de  $f$ .

La force portante est  $f = Ks \cdot A \cdot B \cdot P1_{\text{eff}}$  et  $r = \frac{P1_{\text{eff}}}{P0_{\text{eff}}}$ , donc :

$$f = (Ks \cdot A \cdot B \cdot P0_{\text{eff}}) \cdot r = Kf \cdot r$$

où  $Kf$  est le *coefficient d'effort*.

La raideur est définie par :

$$C = -\frac{df}{dh} = -\frac{df}{dr} \cdot \frac{dr}{dh} = -Kf \cdot \frac{dr}{dh}$$

or  $h = a \cdot \left( \frac{1-r}{r} \right)^{\frac{1}{3}}$ , donc  $\frac{dh}{dr} = -\frac{3}{a} \cdot r^{\frac{4}{3}} \cdot (1-r)^{\frac{2}{3}}$ , d'où :

$$C = \frac{3 \cdot Kf}{a} \cdot r^{\frac{4}{3}} \cdot (1-r)^{\frac{2}{3}}$$

La raideur maximale du patin est obtenue pour  $\frac{dC}{dr} = 0$ , soit  $r = \frac{2}{3}$ .

## 6.5 Modélisation du champ de pression

### 6.5.1 Introduction

Le patin rectangulaire étant chargé uniformément (charge P), une alimentation d'huile à la partie inférieure assurant le maintien d'un film d'huile de hauteur h.

Pour les calculs, on admettra que la pression  $P_1$  est uniforme dans tout le volume intérieur. La pression est  $P_2$  à la sortie du film.

### 6.5.2 Les équations de laminage

- Bilan de masse :  $\frac{\partial \rho}{\partial t} + \text{Div}(\rho \cdot \vec{V}) = 0$ , où  $\rho$  est la masse volumique, supposée constante et  $\vec{V}$  est le champ des vitesses.

Dans la phase de mouvement stationnaire, on a  $\frac{\partial \rho}{\partial t} = 0$ . L'équation se réduit donc à :

$$\text{Div}(\vec{V}) = 0 \quad (1)$$

- Bilan de quantité de mouvement (forme locale) :  $\rho \cdot \gamma_i = \rho \cdot g_i + \frac{\partial \sigma_{ij}}{\partial x_j}$  avec :

$$\sigma_{ij} = -p \cdot \delta_{ij} + 2 \cdot \mu \cdot D_{ij} + \lambda \cdot \theta \cdot \delta_{ij} ;$$

$\sigma_{ij}$  : Tenseur des contraintes ;

$$D_{ij} : \text{Tenseur des taux de déformation} : D_{ij} = \frac{1}{2} \cdot \left( \frac{\partial V_i}{\partial x_j} + \frac{\partial V_j}{\partial x_i} \right)$$

$$q = \text{Div}(\vec{V}) = 0 \quad (\text{d'après (1)})$$

$p$  est la pression en un point

$\mu$  est le coefficient de viscosité dynamique supposé constant est un autre coefficient de viscosité.

Finalement, cette équation se réduit à :

$$\rho \cdot \gamma_i = \rho \cdot g_i + \mu \cdot \sum_{j=1}^3 \frac{\partial^2 V_i}{\partial x_j^2} - \frac{\partial p}{\partial x_i}$$

soit encore :

$$\rho \cdot \vec{\gamma} = \rho \cdot \vec{g} - \text{Grad}(p) + \mu \cdot \Delta \vec{V} \quad (2)$$

où  $\Delta \vec{V}$  est le laplacien de la vitesse.

### 6.5.3 Résolution dans le cas où les termes d'accélération sont très petits

Dans ce cas le système se réduit à :

$$\begin{cases} \text{Div}(\vec{V}) = \vec{0} \\ -\rho \cdot \vec{g} + \text{Grad}(p) + \mu \cdot \Delta \vec{V} \end{cases}$$

Notations et propriétés

En posant :  $\hat{P} = P + \rho \cdot g \cdot z$ , le premier membre s'écrit  $\text{Grad}(\hat{P})$  et le système :

$$\begin{cases} \text{Div}(\vec{V}) = 0 \\ \text{Grad}(\hat{P}) = \mu \cdot \Delta \vec{V} \end{cases}$$

Mais,  $\Delta \vec{V} = -\text{Rot}(\text{Rot}(\vec{V}))$  et  $\text{Div}(\text{Rot}(\vec{V})) = 0$ , d'où avec (2) :  $\Delta \hat{P} = 0$

*Remarque* — Dans le problème posé, le terme  $\rho \cdot g \cdot z$  est négligeable devant  $P$  et l'équation se réduit à :  $\Delta P = 0$ .

## 7. Exercices

### 7.1 Un problème de thermique

Dans un problème d'échange thermique, on est amené à rechercher deux fonctions  $f$  et  $g$  définies sur  $[0, +\infty[ \times [0, +\infty[$  et telles que :

$$(1) \quad \begin{cases} \frac{\partial f}{\partial y}(x, y) = g(x, y) - f(x, y), \text{ pour tous } x \geq 0 \text{ et } y \geq 0 \\ \frac{\partial g}{\partial x}(x, y) = f(x, y) - g(x, y), \text{ pour tous } x \geq 0 \text{ et } y \geq 0 \\ f(x, 0) = \varphi(x), \text{ pour tout } x \geq 0 \\ g(0, y) = \psi(y), \text{ pour tout } y \geq 0 \end{cases}$$

où  $\varphi$  et  $\psi$  sont des fonctions données.

On choisit une valeur limite, pour  $x$  et  $y$ , qui sera notée  $a > 0$ .

En vue de résoudre (1), par discrétisation, on définit un maillage du carré  $\Omega = [0, a] \times [0, a]$  en posant  $h = \frac{a}{n}$ ,  $x_i = i \cdot h$  pour  $0 \leq i \leq n$  et  $y_j = j \cdot h$  pour  $0 \leq j \leq n$

On notera  $f_{ij}$  et  $g_{ij}$  des valeurs approchées de  $f(x_i, y_j)$  et  $g(x_i, y_j)$ .

1°) En approchant les dérivées partielles par des différences finies progressives, traduire le problème (1).

2°) Décrire un algorithme de résolution du problème discrétisé obtenu en 1°). Ecrire ensuite la programmation structurée correspondante.

3°) Application numérique : On prend  $\varphi(x) = 1$  et  $\psi(y) = 0$ . Pour  $n = 3$  et  $a = 1$ , donner les valeurs approchées  $f_{ij}$  et  $g_{ij}$ .

4°) Quels problèmes pose l'utilisation de différences finies régressives ?

## 7.2 Méthode des éléments finis

On considère le problème aux limites sur  $[0,1]$  :

$$(1) \quad \begin{cases} -(p(x) \cdot u'(x))' + q(x) \cdot u(x) = f(x) & (0 < x < 1) \\ u(0) = u(1) = 0 \end{cases}$$

où les fonctions  $p$ ,  $q$ , et sont données telles que :

$$\begin{cases} p \in C^1([0,1]); \forall x \in [0,1], p(x) \geq p_0 > 0 \\ q \in C^0([0,1]); \forall x \in [0,1], q(x) \geq 0 \\ f \in C^0([0,1]) \end{cases}$$

et la fonction inconnue  $u$  est cherchée dans  $C^2([0,1])$ .

On admettra, qu'avec nos hypothèses, une telle fonction  $u$  existe.

La méthode étudiée pour résoudre ce problème est la *méthode variationnelle de Rayleigh, Ritz et Galerkin*, ou *méthode des éléments finis*.

On note  $E = \{v : [0,1] \rightarrow \mathbb{R} ; \text{continue et } C^1 \text{ par morceaux ; } v(0) = v(1) = 0\}$ .

Pour tout  $n \geq 1$  dans  $\mathbb{N}$ , on pose :  $h = \frac{1}{n+1}$ ,  $x_i = i \cdot h$ , pour  $i = 0, 1, \dots, n+1$ .

On notera  $u$  la solution exacte de (1).

1°) En utilisant une intégration par parties, montrer que  $u$  est aussi solution du problème :

$$(2) \quad \forall v \in E, \int_0^1 (p \cdot u' \cdot v' + q \cdot u \cdot v)(x) dx = \int_0^1 (f \cdot v)(x) dx$$

Dans la suite du problème, on pose, pour tous  $v, w$  dans  $E$  :

$$a(v, w) = \int_0^1 (p \cdot w' \cdot v' + q \cdot w \cdot v)(x) dx \text{ et } b(v) = \int_0^1 (f \cdot v)(x) dx$$

2°) Montrer que  $a$  est une forme bilinéaire symétrique définie positive.

3°) En déduire que (1) admet une unique solution.

Pour résoudre (2) de façon approchée, l'idée de la méthode des éléments finis est de remplacer l'espace vectoriel de dimension infinie  $E$  par un s. e. v.  $E_h$  de dimension finie, le paramètre  $h > 0$  étant lié à la dimension de  $E_h$  et destiné à tendre vers 0.

Le problème (2) sera alors remplacé par le problème discret :

$$(2)_h \quad \begin{cases} \text{Trouver } u_h \text{ dans } E_h \text{ tel que} \\ \forall v \in E_h, a(u_h, v) = b(v) \end{cases}$$

On espère que  $\lim_{h \rightarrow 0} (u_h) = 0$ . Et ce résultat sera atteint avec une précision qui dépendra du choix du s. e. v.  $E_h$ . On admet donc que la méthode converge.

Si  $\{v_1, \dots, v_n\}$  est une base de  $E_h$ , on écrira :  $u_h = \sum_{i=1}^n \alpha_i \cdot v_i$ .

3°) Montrer que le problème  $(2)_h$  admet une unique solution.

4°) On désigne par  $E_h$  le s. e. v. de  $E$  engendré par les fonctions  $v_1, \dots, v_n$ , où  $v_i$  est la fonction affine par morceaux dont le graphe est représenté sur la figure 8.15 :



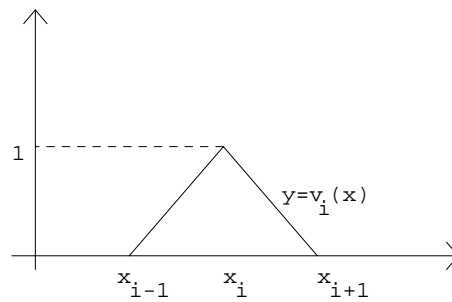


Figure 8.15

- (a) Expliciter le système linéaire dont est solution  $\alpha = (\alpha_1, \dots, \alpha_n)$ .
- (b) Décrire une procédure de résolution du système obtenu en (a), en supposant que l'on dispose d'une procédure de calcul des intégrales définies.
- (c) En prenant  $u_h$  comme valeur approchée de  $u$ , expliquer pourquoi les coefficients  $\alpha_i$  donnent des valeurs approchées des  $u(x_i)$ , pour  $i = 1, \dots, n$ .
- (d) Comparer cette méthode à celle utilisant les différences finies centrées.

### 7.3

Ecrire les procédures de résolution d'une équation aux dérivées partielles de type parabolique.

## 8. Programmation Ada

### 8.1 Spécification du paquetage DIRICHLET

Dans ce paquetage on trouve des procédures de résolution d'équations différentielles d'ordre 2 avec conditions aux limites. On pourra également en écrire une version générique.

```
with CURVE, COMMON_MATH2 ;
use CURVE, COMMON_MATH2 ;
package DIRICHLET is
  procedure RESOLUTION_DIRICHLET(p, q, f : in FONCTION ; -- § 3.1
    a, b, ya, yb : in FLOAT ; C : in out COURBE) ; -- § 3.2
  procedure RESOLUTION_POISSON(p, f : in FONCTION ; y0, y1 : in FLOAT ;
    C : in out COURBE) ;
end DIRICHLET ;
```

### 8.2 Démonstration du paquetage DIRICHLET

```
with TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, CURVE,
  MATH2, MATH3, DIRICHLET ;
use TEXT_IO, IIO, FIO, CRT, COMMON_MATH0, MATH0, CURVE,
  MATH2, MATH3, DIRICHLET ;
procedure DEM_DIRI is
  a, b, ya, yb, y0, y1 : FLOAT ;
  p, q, f : FONCTION ;
```

```

C : COURBE(300) ;
CHOIX : INTEGER range 1..3 ;

procedure DEMO_POISSON is
begin
  MODE_AFFICHAGE ;
  PUT_LINE("Resolution de y'' + p(x).y(x) = f(x) sur [0,1] ") ;
  PUT_LINE("avec p negative et y(0) et y(1) donnees.") ;
  GET_LINE(p,"Entrer p(x) : ") ;
  GET_LINE(f,"Entrer f(x) : ") ;
  ENTRER_REEL(y0,"y(0) = ") ;
  ENTRER_REEL(y1,"y(1) = ") ;
  RESOLUTION_POISSON(p,f,y0,y1,C) ;
  CLOSE(IMP) ;
end DEMO_POISSON ;

procedure DEMO_DIRICHLET is
begin
  MODE_AFFICHAGE ;
  PUT_LINE("Resolution d'une equation differentielle lineaire ") ;
  PUT_LINE("du second ordre avec conditions initiales.") ;
  PUT_LINE(" y''(x) + p(x).y'(x) + q(x).y(x) = f(x) (a < x < b)") ;
  PUT_LINE(" avec y(a) et y(b) donnees.") ;
  GET_LINE(p,"Entrer p(x) : ") ;
  GET_LINE(q,"Entrer q(x) : ") ;
  GET_LINE(f,"Entrer f(x) : ") ;
  ENTRER_REEL(a,"a = ") ;
  ENTRER_REEL(ya,"ya = ") ;
  ENTRER_REEL_BORNE(a,1.0E12,b,"b = ") ;
  ENTRER_REEL(yb,"yb = ") ;
  RESOLUTION_DIRICHLET(p,q,f,a,b,ya,yb,C) ;
  CLOSE(IMP) ;
end DEMO_DIRICHLET ;

begin
  loop
    PUT_LINE(" -1- Equation du type de Poisson") ;
    PUT_LINE(" -2- Equation du type de Dirichlet") ;
    PUT_LINE(" -3- FIN") ;
    ENTRER_ENTIER_BORNE(1,3,CHOIX," Votre choix : ") ;
    case CHOIX is
      when 1 => DEMO_POISSON ;
      when 2 => DEMO_DIRICHLET ;
      when 3 => exit ;
    end case ;
    if LIRE_REPONSE("Voulez vous l'affichage des resultats? ") then
      for k in 0..C.n
        loop
          PUT("x = ") ; PUT(C.u(k),4,6,0) ;
          PUT(" y = ") ; PUT(C.v(k),4,6,0) ;
          NEW_LINE ;
          if ((k + 1) mod 20 = 0) then
            PAUSE ;
          end if ;
        end loop ;
      end for ;
    end if ;
  end loop ;
end ;

```

```

        end if ;
    end loop ;
end if ;
NEW_LINE ;
if LIRE_REPONSE("Voulez vous le trace graphique? ") then
    TRACE_UNE_COURBE(C) ;
    SORTIE_GRAPHIQUE ;
else
    CLRSCR ;
end if ;
end loop ;
end DEM_DIRI ;

```

### 8.3 Spécification du paquetage NIVEAUX

Ce paquetage est utilisé pour tracer les courbes de niveaux d'une fonction  $u(x,y)$  solution d'une équation elliptique. Sa conception est inspirée du logiciel Modulog. Pour plus de détails, on se reportera à l'exercice 10 du chapitre 2.

```

with MATH2, COMMON_MATRIX ;
use MATH2, COMMON_MATRIX ;
package NIVEAUX is
procedure DISCRETISE_FONCTION(f : in FONCTION ;
    xMin, xMax, yMin, yMax : in FLOAT ; M : in out MATRICE) ;
procedure TRACE_NIVEAU(XMin,XMax,YMin,YMax : in FLOAT ;
    M : in MATRICE ; f : in FLOAT) ;
procedure TRACE_EQUIPOTENTIELLES(M : in MATRICE ;
    xMin, xMax, yMin, yMax : in FLOAT ;
    u0 : in VECTEUR ; AXES : in BOOLEAN := FALSE) ;
end NIVEAUX ;

```

### 8.4 Spécification du paquetage EDP\_ELLIPTIQUES

Dans ce paquetage on décrit des procédures de résolution de l'équation de Poisson  $\Delta u = f$  et d'une équation plus générale de type elliptique. On pourra également en écrire une version générique.

```

with MATH2, COMMON_MATRIX ;
use MATH2, COMMON_MATRIX ;
package EDP_ELLIPTIQUES is
procedure RESOLUTION_POISSON(f, g1, g2, g3, g4 : in FONCTION ;
    RHO_JACOBI, h : in FLOAT ; U : in out MATRICE) ; -- § 5.6
PROCEDURE RESOLUTION_ELLIPTIQUE(f, g1, g2, g3, g4,
    Alpha, Beta, Gama, Delt, Epsi : in FONCTION ;
    Omega, h : in FLOAT ; U : in out MATRICE) ;
procedure AFFICHER_SOLUTION(U : in MATRICE) ; -- § 5.5.4
procedure RESOLUTION_POISSON_2(f, g1, g2, g3, g4,
    h1, h2, h3, h4 : in FONCTION ;
    p, q, r, s : in INTEGER ; RHO_JACOBI, h : in FLOAT ;
    U : in out MATRICE) ;
procedure AFFICHER_SOLUTION_2(p, q, r, s : in INTEGER ;
    U : in MATRICE) ;

```

```
ERREUR_CONVERGENCE exception ;
end EDP_ELLIPTIQUES ;
```

### 8.5 Démonstration du paquetage EDP\_ELLIPTIQUES

```
with TEXT_IO, IIO, CRT, COMMON_MATH0, MATH0, COMMON_MATH1, MATH1,
     MATH2, COMMON_MATRIX, NIVEAUX, EDP_ELLIPTIQUES ;
use TEXT_IO, IIO, CRT, COMMON_MATH0, MATH0, COMMON_MATH1, MATH1,
     MATH2, COMMON_MATRIX, NIVEAUX, EDP_ELLIPTIQUES ;
procedure DEM_ELLI is
  CHOIX : INTEGER range 1..4 ;

  procedure DEMO_EDP_POISSON is
    f, g1, g2, g3, g4 : FONCTIONION ;
    n, m, NB_COURBES : INTEGER ;
    h, RHO_JACOBI : FLOAT ;
  begin
    MODE_AFFICHAGE ;
    PUT_LINE("   Resolution de l'equation de Poisson : ") ;
    PUT_LINE("   Delta(u) = f sur un rectangle [0,a]x[0,b] ;") ;
    PUT_LINE("   avec les conditions au bord : u(x,y) = g(x,y) ;") ;
    PUT_LINE("   ou g(x,0) = g1(x), g(x,b) = g2(x),
              g(0,y) = g3(y), g(a,y) = g4(y).") ;
    PUT_LINE("On utilise les notations suivantes : ") ;
    PUT_LINE(" a = (n + 1).h ; b = (m + 1).h ,
              ( h > 0 pas de discretisation donne ) ;") ;
    PUT_LINE("un maillage du rectangle etant defini par : ") ;
    PUT_LINE(" (x(i),y(j)) = (i.h,j.h) pour i = 0, ..., n + 1
              et j = 0, ..., m + 1 . ") ;
    PAUSE ; CLRSCR ;
    PUT_LINE("ENTREE DES DONNEES POUR L'EQUATION DE POISSON") ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,n,
                        "Nombre de points sur l'axe des x : ") ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,m,
                        "Nombre de points sur l'axe des y : ") ;
    ENTRER_REEL(h,"h = ") ;
    GET_LINE(f,"f(x,y) = ") ;
    GET_LINE(g1,"Condition a y = 0 , g1(x) = ") ;
    GET_LINE(g2,"Condition a y = b , g2(x) = ") ;
    GET_LINE(g3,"Condition a x = 0 , g3(y) = ") ;
    GET_LINE(g4,"Condition a x = a , g4(y) = ") ;
    RHO_JACOBI :=0.5*(COS(PI/FLOAT(n + 1)) + COS(PI/FLOAT(m + 1))) ;
    declare
      U : MATRICE(0..n + 1,0..m + 1) ;
    begin
      RESOLUTION_POISSON(f,g1,g2,g3,g4,RHO_JACOBI,h,U) ;
      if LIRE_REPONSE("Voulez vous l'affichage des resultats? ") then
        CLRSCR ;
        AFFICHER_SOLUTION(U) ;
        PAUSE ; CLRSCR ;
      end if ;
      if LIRE_REPONSE("Voulez vous tracer les equipotentiellles? ") then
        ENTRER_ENTIER(NB_COURBES,"Nombre de courbes a tracer : ") ;
```

```

declare
    xMin, xMax, yMin, yMax : FLOAT ;
    U0 : VECTEUR(1..NB_COURBES) ;
begin
    for i in 1..NB_COURBES
    loop
        PUT(" Courbe numero ") ; PUT(i) ; NEW_LINE ;
        ENTRER_REEL(U0(i),"Entrer u0 : ") ;
    end loop ;
    xMin := 0.0 ; xMax := FLOAT(n + 1)*h ;
    yMin := 0.0 ; yMax := FLOAT(m + 1)*h ;
    TRACE_EQUIPOTENTIELLES(U,xMin,xMax,yMin,yMax,U0) ;
end ;
end if ;
end ;
CLOSE(IMP) ;
end DEMO_EDP_POISSON ;

procedure DEMO_EDP_POISSON_2 is
f, g1, g2, g3, g4, h1, h2, h3, h4 : FONCTION ;
n, m, p, q, r, s, NB_COURBES : INTEGER ;
h, RHO_JACOBI : FLOAT ;
begin
    MODE_AFFICHAGE ;
    PUT_LINE(" Resolution de l'equation de Poisson : ") ;
    PUT_LINE(" Delta(u) = f sur un rectangle troue
        [0,a]x[o,b] - [c,d]x[e,f]") ;
    PUT_LINE("avec les conditions au bord : u(x,y) = g(x,y) sur le bord"
        & " exterieur et u(x,y) = h(x,y) sur le bord interieur.") ;
    NEW_LINE ;
    PUT_LINE("Avec g(x,0) = g1(x), h(x,e) = h1(x), g(x,b) = g2(x), "
        & "h(x,f) = h2(x), g(0,y) = g3(y), h(c,y) = h3(y),
        g(a,y) = g4(y), h(d,y) = h4(y).") ;
    PUT_LINE("On utilise les notations suivantes : ") ;
    PUT_LINE(" a = (n + 1).h , b = (m + 1).h, c = p.h, d = q.h, "
        & "e = r.h, f = s.h (h > 0 pas de discretisation donne) ;") ;
    PUT_LINE("un maillage du rectangle etant defini par :") ;
    PUT_LINE(" (x(i),y(j)) = (i.h,j.h) pour i = 0, ..., n + 1
        et j = 0, ..., m + 1 . ") ;
    PAUSE ; CLRSCR ;
    PUT_LINE("ENTREE DES DONNEES") ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,n,
        "Nombre de points sur l'axe des x, n = ") ;
    ENTRER_ENTIER_BORNE(2,DIM_MAX,m,
        "Nombre de points sur l'axe des y, m = ") ;
    ENTRER_ENTIER_BORNE(2,n - 2,p,"Entrer p (entre 2 et n - 2) : ") ;
    ENTRER_ENTIER_BORNE(p + 1,n - 1,q,
        "Entrer q (entre p + 1 et n - 1) : ") ;
    ENTRER_ENTIER_BORNE(2,m - 2,r,"Entrer r (entre 2 et m - 2) : ") ;
    ENTRER_ENTIER_BORNE(r + 1,m - 1,s,
        "Entrer s (entre r + 1 et m - 1) : ") ;
    ENTRER_REEL(h,"h = ") ;
    GET_LINE(f,"f(x,y) = ") ;

```

```

GET_LINE(g1,"Condition a y = 0 , g1(x) = " ) ;
GET_LINE(h1,"Condition a y = r.h , h1(x) = " ) ;
GET_LINE(h2,"Condition a y = s.h , h2(x) = " ) ;
GET_LINE(g2,"Condition a y = b , g2(x) = " ) ;
GET_LINE(g3,"Condition a x = 0 , g3(y) = " ) ;
GET_LINE(h3,"Condition a x = p.h , h3(y) = " ) ;
GET_LINE(h4,"Condition a x = q.h , h4(y) = " ) ;
GET_LINE(g4,"Condition a x = a , g4(y) = " ) ;
RHO_JACOBI := 0.5*(COS(PI/FLOAT(n + 1)) + COS(Pi/FLOAT(m + 1))) ;
CLRSCR ;
declare
  U : MATRICE(0..n + 1,0..m + 1) ;
begin
  RESOLUTION_POISSON_2(f,g1,g2,g3,g4,h1,h2,h3,h4,p,q,r,s,
    RHO_JACOBI,h,U) ;
  if LIRE_REPONSE("Voulez vous l'affichage des resultats? ") then
    CLRSCR ;
    AFFICHER_SOLUTION_2(p,q,r,s,U) ;
    PAUSE ; CLRSCR ;
  end if ;
  if LIRE_REPONSE("Voulez vous tracer les equipotentiellles? " then
    ENTREER_ENTIER(NB_COURBES,"Nombre de courbes a tracer : ") ;
    declare
      xmin, xmax, ymin, ymax : FLOAT ;
      U0 : VECTEUR(1..NB_COURBES) ;
    begin
      for i in 1..NB_COURBES
        loop
          PUT(" Courbe numero ") ; PUT(i) ; NEW_LINE ;
          ENTREER_REEL(U0(i),"Entrer U0 : ") ;
        end loop ;
        xmin := 0.0 ; xmax := FLOAT(n + 1)*h ;
        ymin := 0.0 ; ymax := FLOAT(m + 1)*h ;
        TRACE_EQUIPOTENTIELLES(U,xmin,xmax,ymin,ymax,U0) ;
      end ;
    end if ;
  end ;
  CLOSE(IMP) ;
end DEMO_EDP_POISSON_2 ;

procedure DEMO_EDP_ELLIPTIQUE is
Alpha, Beta, Gama, Delt, Epsi, f, g1, g2, g3, g4 : FONCTION ;
n, m, NB_COURBES : INTEGER ;
h, Omega : FLOAT ;
begin
  MODE_AFFICHAGE ;
  PUT_LINE(" Resolution d'une equation aux derivees partielles
    de type elliptique : ") ;
  PUT_LINE("sur un rectangle [0,a]x[o,b] ;") ;
  PUT_LINE(" avec les conditions au bord : u(x,y) = g(x,y) ;") ;
  PUT_LINE(" ou g(x,0) = g1(x), g(x,b) = g2(x),
    g(0,y) = g3(y), g(a,y) = g4(y) .") ;
  PUT_LINE("On utilise les notations suivantes : ") ;

```

```

PUT_LINE(" a = (n + 1).h ; b = (m + 1).h ,
          ( h > 0 pas de discretisation donne ) ;") ;
PUT_LINE("un maillage du rectangle etant defini par : ") ;
PUT_LINE(" (x(i),y(j)) = (i.h,j.h) pour i = 0, ..., n + 1
          et j = 0, ..., m + 1 . ") ;
PAUSE ; CLRSCR ;
PUT_LINE("ENTREE DES DONNEES POUR UNE EQUATION ELLIPTIQUE") ;
ENTRER_ENTIER_BORNE(2,DIM_MAX,n,
                    "Nombre de points sur l'axe des x : ") ;
ENTRER_ENTIER_BORNE(2,DIM_MAX,m,
                    "Nombre de points sur l'axe des y : ") ;
ENTRER_REEL(h,"h = ") ;
GET_LINE(Alpha,"Alpha(x,y) = ") ;
GET_LINE(BETA,"Beta(x,y) = ") ;
GET_LINE(GAMA,"Gamma(x,y) = ") ;
GET_LINE(DELT,"Delta(x,y) = ") ;
GET_LINE(EPSI,"Epsilon(x,y) = ") ;
GET_LINE(f,"f(x,y) = ") ;
GET_LINE(g1,"Condition a y = 0 , g1(x) = ") ;
GET_LINE(g2,"Condition a y = b , g2(x) = ") ;
GET_LINE(g3,"Condition a x = 0 , g3(y) = ") ;
GET_LINE(g4,"Condition a x = a , g4(y) = ") ;
ENTRER_REEL(OMEGA,"Coeeficient de relaxation omega = ") ;
declare
  U : MATRICE(0..n + 1,0..m + 1) ;
begin
  RESOLUTION_ELLIPTIQUE(f,g1,g2,g3,g4,
                        Alpha,Beta,Gama,Delt,Epsi,Omega,h,U) ;
  if LIRE_REPONSE("Voulez vous l'affichage des resultats? ") then
    CLRSCR ;
    AFFICHER_SOLUTION(U) ;
    PAUSE ; CLRSCR ;
  end if ;
  if LIRE_REPONSE("Voulez vous tracer les equipotentiellles? ") then
    ENTRER_ENTIER(NB_COURBES,"Nombre de courbes a tracer : ") ;
    declare
      xMin, xMax, yMin, yMax : FLOAT ;
      U0 : VECTEUR(1..NB_COURBES) ;
    begin
      for i in 1..NB_COURBES
      loop
        PUT(" Courbe numero ") ; PUT(i) ; NEW_LINE ;
        ENTRER_REEL(U0(i),"Entrer u0 : ") ;
      end loop ;
      xMin := 0.0 ; xMax := FLOAT(n + 1)*h ;
      yMin := 0.0 ; yMax := FLOAT(m + 1)*h ;
      TRACE_EQUIPOTENTIELLES(U,xMin,xMax,yMin,yMax,U0) ;
    end ;
  end if ;
end ;
CLOSE(IMP) ;
end DEMO_EDP_ELLIPTIQUE ;

```

```

begin
  loop
    CLRSCR ;
    PUT_LINE("-1- Resolution d'une equation de Poisson : Delta(u) = f") ;
    PUT_LINE("-2- Resolution de Delta(u) = f sur un rectangle troue") ;
    PUT_LINE("-3- Resolution d'une e.d.p. de type elliptique") ;
    PUT_LINE("-4- FIN") ;
    NEW_LINE ;
    ENTRER_ENTIER_BORNE(1,4,CHOIX,"Votre choix : ") ;
    case CHOIX is
      when 1 => DEMO_EDP_POISSON ;
      when 2 => DEMO_EDP_POISSON_2 ;
      when 3 => DEMO_EDP_ELLIPTIQUE ;
      when 4 => exit ;
    end case ;
  end loop ;
end DEM_ELLI ;

```

## 8.6 Spécification du paquetage *EDP\_PARABOLIQUES*

Dans ce paquetage on décrit des procédures de résolution de l'équation de diffusion de la chaleur dans une barre. On pourra le compléter avec des procédures de résolution d'une équation de type parabolique plus générale.

```

with MATH2, COMMON_MATRIX ;
use MATH2, COMMON_MATRIX ;
package EDP_PARABOLIQUES is
  NX_MAX : constant := 200 ;
  NT_MAX : constant := 30 ;
  procedure EDP_CHALEUR_EXPLICITE(k, dx, dt : in FLOAT ;
    f : in FONCTION ; U : in out MATRICE) ; § 5.8.3
  procedure EDP_CHALEUR_IMPLICITE(w, k, dx, dt : in FLOAT ;
    f : in FONCTION ; U : in out MATRICE) ; § 5.8.3
  procedure AFFICHER_TEMPERATURES(U : in MATRICE ; dt : in FLOAT ;
    p : in INTEGER) ;
  procedure COURBES_CHALEUR(U : in MATRICE ; dx : in FLOAT ;
    p : in INTEGER) ;
end EDP_PARABOLIQUES ;

```

## 8.7 Démonstration du paquetage *EDP\_PARABOLIQUES*

```

with TEXT_IO, CRT, COMMON_MATH0, MATH0, MATH2, COMMON_MATRIX,
  EDP_PARABOLIQUES ;
use TEXT_IO, CRT, COMMON_MATH0, MATH0, MATH2, COMMON_MATRIX,
  EDP_PARABOLIQUES ;

procedure DEM_CHAL is
  Eps : constant FLOAT := 1.0E-12 ;
  nx, nt, p : INTEGER ;
  k, w, L, dx, dt, dtb : FLOAT ;
  f : FONCTION ;
  CORRECT_F : BOOLEAN ;

```



```

procedure PRESENTATION_CHALEUR is
begin
  PUT_LINE("la barre est a temperature nulle aux bords") ;
  PUT_LINE("le reste n'a pas d'echange de chaleur avec l'exterieur") ;
  PUT_LINE("la longueur de la barre est notee L = (nx + 1)*dx") ;
  PUT_LINE("tf = nt*dt est la valeur finale du temps.") ;
  PUT_LINE("le coefficient de conduction est note k") ;
  PUT_LINE("w est le coefficient de la methode de relaxation") ;
  PUT_LINE("la temperature a t = 0 en x est f(x)") ;
  -- Pour des raisons de compatibilite, f doit etre nulle en 0 et L
end PRESENTATION_CHALEUR ;

begin
  MODE_AFFICHAGE ;
  PRESENTATION_CHALEUR ;
  PAUSE ; CLRSCR ;
  ENTRER_REEL(L," Entrez L : ") ;
  ENTRER_ENTIER_BORNE(2,NX_MAX,nx,
    " Nombre de points sur l'axe des x : ") ;
  dx := L/FLOAT(nx + 1) ;
  ENTRER_REEL(k," Entrez la valeur de k = ") ;
  loop
    GET_LINE(f," Entrez f(x) : ") ;
    CORRECT_F := (abs(EVALUE(f,0.0)) <= Eps)
      and (ABS(EVALUE(f,L)) <= Eps) ;
    if not Correct_f then
      BIP ;
      PUT_LINE("* Erreur * f doit etre nulle en 0 et L") ;
    end if ;
    exit when CORRECT_F ;
  end loop ;
  ENTRER_REEL_BORNE(0.0,1.0,w,"Entrez w : ") ;
  -- Choix du pas de discretisation en t :
  -- Si  $0 \leq w < 1/2$ , une condition de stabilite est donnee par :
  --  $r = dt/(dx^2) < 1/(2*k*(1 - 2*w))$ 
  -- La valeur maximale de dt est alors  $dtb = (dx**2)/(2*k*(1 - w))$ 
  if (w < 0.5) then
    dtb := (dx*dx)/(2.0*k*(1.0 - 2.0*w)) ;
  else
    dtb := 10.0 ;
  end if ;
  ENTRER_REEL_BORNE(0.0,dtb,dt,"Entrez la valeur de dt : ") ;
  ENTRER_ENTIER_BORNE(2,NT_MAX,nt,
    "Entrez le nombre de points sur l'axe des t : ") ;
  ENTRER_ENTIER_BORNE(1,nt,p,
    "Entrez la periode des traces de courbes : ") ;
  declare
    U : MATRICE(0..nx + 1,0..nt) ;
  begin
    if (w = 0.0) then
      EDP_CHALEUR_EXPLICITE(k,dx,dt,f,U) ;
    else
      exit ;
    end if ;
  end ;
end ;

```

```
      EDP_CHALEUR_IMPLICITE(w,k,dx,dt,f,U) ;
    end if ;
    if LIRE_REPONSE("Voulez vous l'affichage des resultats? ") then
      AFFICHER_TEMPERATURES(U,dt,p) ;
    end if ;
    if LIRE_REPONSE("Voulez vous le trace des courbes? ") then
      COURBES_CHALEUR(U,dx,p) ;
    end if ;
  end ;
  CLOSE(IMP) ;
end DEM_CHAL ;
```

# INDEX

## A

access, 4  
Ada, 1; 8  
ADA.AUX, 8  
ADA.LIB, 8  
agrégat, 5  
Aitken, 164  
ALGEBLIN, 94  
approximations successives, 119  
ArcCos, 20  
ArcSin, 20  
ArcTan, 20  
array, 4

## B

B-Splines, 149; 157; 160  
B\_SPLINE, 183  
Bairstow, 115  
BAIRSTOW, 125  
BAMP, 8  
Bernouilli, 119; 195  
Bernstein, 149; 150  
Bézier, 149; 152; 155  
BEZIER, 181  
BIN, 7  
BIP, 16  
Brun, 206

## C

C.A.O., 149  
carreau de Bézier, 156  
CARTE\_GRAPHIQUE, 23  
CARTESIENNE, 30  
case, 5  
Cauchy, 238  
Cauchy-Lipschitz, 241  
Cayley-Hamilton, 86  
CENTIEME, 14  
cercle des moindres carrés, 146  
CERCLE, 27  
CHAINE\_ENTIERE, 16  
CHAINE\_REELLE, 16  
CHOIX\_AFFICHAGE, 14  
Cholesky, 62; 65; 97; 98  
CHRONOMETRE, 16  
CLEAR\_SCREEN, 24  
CLOSE\_GRAPH, 24  
CLREOL, 13  
CLRSCR, 13  
coefficient de corrélation, 131  
coefficients de Fourier, 224  
COLONNE\_TEXTE, 12  
COMMON\_CRT, 12  
COMMON\_FOURIER, 232

COMMON\_GRAPH, 23  
COMMON\_MATH0, 14  
COMMON\_MATH1, 20  
COMMON\_MATH3, 25  
COMMON\_MATRIX, 93  
COMMON\_POLY, 124  
COMPLEXE, 124  
composite, 194; 197  
conditionnement, 43; 45  
conditions aux limites, 256  
contrôle du pas, 253  
CONVERSION, 26  
convolution, 226  
Cooley, 219  
COORDONNEE, 28  
cordes vibrantes, 288  
Cotes, 192  
COULEUR, 23  
COURBE, 28  
courbes à pôles, 153  
covariance, 131  
Cramer, 43; 49  
Crank-Nicholson, 309  
CREER\_FICHER, 17  
creuse, 43  
CROIX, 27  
Crout, 58  
CRT, 11; 12  
CURSOR, 13  
CURVE, 28

## D

d'Alembert, 290  
D.O.D., 1  
Darboux-Christoffel, 204  
De Boor, 170  
De Casteljeau, 154; 157  
décomposition L-R, 58; 61  
définie positive, 47  
déflation, 78; 103  
dégénéré, 40  
delta, 4  
DEPLACE, 26  
déterminant, 57  
DETRUIT\_FONCTION, 23  
DEUX\_PI, 20  
diagonalisable, 46  
dichotomie, 108  
DICHOTOMIE, 120  
différence finie centrée, 281  
différence finie progressive, 281  
différence finie régressive, 281  
différences divisées, 166  
différences finies, 110; 114; 133; 281; 296

- digits, 4
- DIM\_MAX, 93
- Dirichlet, 240; 279; 282; 298
- DIRICHLET, 319
- Dirichlet-Neumann, 240
- DONNEES\_ALGEBLIN, 95
- DONNEES\_POINTS, 179
  
- E*
- E, 20
- écart quadratique, 129
- écart type, 131
- échantillonnage simple, 213
- échantillonnage stratifié, 213
- EDP\_ELLIPTIQUES, 321
- EDP\_PARABOLIQUES, 326
- éléments finis, 227; 318
- elliptique, 295
- EQUADIFF\_11\_GENERIQUE, 269
- EQUADIFF\_1Q\_GENERIQUE, 269
- EQUADIFF\_P1\_GENERIQUE, 269
- EQUADIFF\_PQ\_GENERIQUE, 270
- équation de la chaleur, 292; 306
- équation des ondes, 288
- équation parabolique, 306
- équations différentielles linéaires, 242
- équations différentielles, 239
- équations elliptiques, 298
- équations normales, 134
- équations paraboliques, 312
- EQUATIONS\_GENERIQUE, 120
- erreur de consistance, 243
- erreur de discrétisation, 243
- ERREUR\_ARGUMENT, 19; 20
- ERREUR\_CONVERGENCE, 120; 122; 125
- Euler, 246
- Euler-Cauchy, 246
- EVALUE, 22
- exception, 2
  
- F*
- F.F.T, 137; 220
- fausse position, 120
- FENETRE\_GRAPHIQUE, 26
- FICHER\_EXISTE, 17
- FLOAT, 2
- FONCTION, 21
- fonctions mélanges, 150; 158
- fonctions pondérantes, 150; 158
- Fourier discrète, 218
- Fourier rapide, 216
- Fourier, 136; 293; 298
- FOURIER\_GENERIQUE, 233
- FRAC, 17
- fractale, 34
  
- fractions continues, 32; 177
- Fubini, 187
  
- G*
- Gamma, 21
- Gauss, 51; 95; 200
- Gauss-Jordan, 55; 97
- Gauss-Seidel, 67
- generic, 6
- générique, 5
- GET\_LINE, 16; 23
- GET\_TIME, 15
- Givens, 91
- GOTOXY, 13
- GRAPH, 12; 23
- Gronwall, 245
  
- H*
- Hermite, 209
- Heun, 247
- HEURE, 14
- Hilbert, 41
- Householder, 91; 92
- hyperbolique, 294
- hyperboliques inverses, 21
  
- I*
- Ichbiach, 1
- if, 5
- IMP, 14
- in out, 5
- in, 5
- INIT\_GRAPH, 24
- INIT\_GRAPHIQUE, 25
- intégrales elliptiques, 227
- INTEGRATION\_GENERIQUE, 229
- interpolation, 162
- INTERPOLATION\_SPLINE, 185
- INTERRUPT, 12
- inversion binaire, 222
- inversion, 57
- is, 4
  
- J*
- Jacobi, 66; 81; 99; 105
- Jordan, 76
- Julia, 34
  
- K*
- Krylov, 87
  
- L*
- L-R, 96
- L.R.M., 3
- Lagrange, 163
- LAGRANGE, 184

Laguerre, 208  
 laplacien, 292  
 LECTURE\_REPERTOIRE, 16  
 Legendre, 205  
 Leverrier, 87; 89  
 LIGNE\_TEXTE, 12  
 LINE, 24  
 LINE\_TO, 24  
 LIO, 12  
 Liouville, 188; 280  
 Lipschitzienne, 240  
 LIRE\_FICHER, 17  
 LIRE\_REPONSE, 16  
 LnGamma, 21  
 loi Gamma, 177  
 LONGUEUR\_CHAINE, 18  
 loop, 5

*M*

maillage uniforme, 281  
 maillage, 298  
 MAJUSCULE, 15  
 Mandelbrot, 34  
 Markov, 205  
 MATH\_LIB, 11; 19  
 MATH0, 11; 14  
 MATH1, 11  
 MATH2, 12; 21  
 MATH3, 12; 24  
 matrice jacobienne, 113  
 MATRICE, 93  
 MATRICE\_JACOBIENNE, 122  
 MATRICE\_TRIDIAGONALE, 93  
 matrices tridiagonales, 61  
 MATRIX, 93  
 MAX\_VARIABLES, 22  
 MESSAGE\_ERREUR\_GRAPHIQUE, 26  
 méthode à un pas, 243  
 méthode variationnelle, 318  
 méthodes directes, 43  
 méthodes itératives, 43; 65  
 MINUTE, 14  
 MODE\_AFFICHAGE, 17  
 MODE\_ECRAN, 25  
 MODE\_TEXTE, 25  
 moindres carrés, 127  
 Monte-Carlo, 210  
 MOVE\_TO, 24  
 Muller, 120  
 multiprécision, 33

*N*

NB\_POINTS, 28  
 Neumann, 298  
 Neville, 164

new, 6  
 NEWLIB.BAT, 8  
 Newton, 86; 90; 165; 192  
 Newton-Raphson, 109; 114; 116; 120; 122;  
 257  
 NEWTON\_MAEHLY, 125  
 Newton\_Maehly, 111  
 NIVEAUX, 321  
 noeuds du maillage, 281  
 nombre de Mach, 296  
 nombres rationnels, 32  
 NOMBRES\_COMPLEXES, 124  
 norme matricielle, 44  
 normes, 44  
 Nyquist, 216

*O*

of, 4  
 Oslo, 161  
 Ostrowski, 71  
 out, 5

*P*

PACLIB, 7  
 paquetages, 5  
 parabolique, 294  
 paramétriques, 36  
 PAUSE, 16  
 PAUSE\_GRAPHIQUE, 26  
 PERIPHERIQUE\_DE\_SORTIE, 14  
 PI., 20  
 PI\_SUR\_DEUX., 20  
 PI\_SUR\_QUATRE, 20  
 pivots partiels, 52  
 pivots, 51  
 Plancherel, 218  
 POINT, 27  
 pointeur, 4  
 points de contrôles, 152  
 Poisson, 280; 284; 292; 297; 304  
 polaires, 36  
 pôles, 152  
 POLY, 124  
 polyèdre descripteur, 155  
 polygone descripteur, 152  
 polynôme caractéristique, 46; 85  
 POLYNOMES, 34  
 primitives, 187  
 private, 4  
 problème de Neumann, 240  
 problèmes aux limites, 240  
 problèmes d'évolution, 295  
 procedure, 4  
 produit scalaire, 47  
 puissance itérée, 75

puissances, 20  
 PUT\_CHAR, 13  
 PUT\_LINE\_STRING, 13  
 PUT\_PIXEL, 24  
 PUT\_STRING, 13

*Q*

Q-R, 92; 100  
 quasi-linéaires, 294

*R*

RANDOM, 15  
 RANDOMIZE, 15  
 range, 4  
 rayon de courbure, 148  
 rayon spectral, 46; 75  
 READKEY, 13  
 RECTANGLE, 24  
 régression affine, 130  
 régression polynomiale, 133  
 régression trigonométrique, 136  
 REGRESSIONS, 179  
 Reich, 71  
 relaxation, 70; 100  
 résidu, 129  
 Richardson, 198  
 RK2, 249  
 RK4, 248  
 Romberg, 198  
 Runge, 163; 166  
 Runge-Kutta, 246  
 Rutishauser, 79; 104

*S*

schéma explicite, 309  
 schéma implicite, 309  
 schémas à un pas, 242  
 Schoenberg, 168  
 SECONDE, 14  
 separate, 5  
 séries de Fourier, 225  
 signature, 64  
 Simpson, 196  
 SOR, 302  
 SORTIE\_GRAPHIQUE, 27  
 Souriau, 85; 106  
 spécification, 5  
 SPECTRE, 102  
 spectre, 46  
 spline cubique, 168  
 splines naturels, 170  
 splines périodiques, 170  
 stabilité, 310  
 stationnaires, 295  
 Stone-Weirstrass, 189

Stormer-Cowell, 267  
 surcharger, 5  
 symétrique, 47  
 système différentiel, 239  
 SYSTEME\_EQUATIONS\_GENERIQUE, 122  
 systèmes triangulaires, 49

*T*

tâche, 4  
 TAILLE\_CHARACTERE, 23  
 Taylor-Lagrange, 282  
 Tchébychev, 72; 167; 207; 303  
 TEST\_MODE\_GRAPHIQUE, 26  
 TEXT\_MOVE\_TO, 26  
 TEXTE\_DE\_FONCTION, 23  
 TEXTE\_HORIZONTAL, 23  
 TEXTE\_VERTICAL, 23  
 TIME, 14  
 tir, 256; 268  
 TITRE, 27  
 TRACE, 27  
 TRACE\_SEGMENT, 27  
 transformations antithétiques, 214  
 transposée, 47  
 trapèzes dichotomique, 195  
 trapèzes, 193  
 trois corps, 266  
 TRUNC, 17  
 Tukey, 219  
 type dérivé, 5  
 type, 3  
 types énumératifs, 3  
 types privés, 4

*U*

UNISURF, 150

*V*

VAL, 17  
 valeurs propre, 46; 75  
 Van Der Monde, 89; 102; 164; 202  
 VAN KOCH, 34  
 Varga, 72  
 VARIABLE\_2, 22  
 variance, 131  
 vecteur nodal, 157  
 VECTEUR, 93  
 vecteurs propres, 46; 75

*W*

Weirstrass, 149  
 WHERE\_X, 13  
 WHERE\_Y, 13

*X*

X\_AXE, 27  
XY\_AXES, 27

Y  
Y\_AXE, 27  
Young, 72

# Bibliographie

## Chapitre 1

A un premier niveau on pourra consulter :

- [1] *Manuel de référence du langage de programmation Ada*. Alsys.
- [2] J. BARNES. *Programmer en Ada*. InterEditions — Un ouvrage de base.
- [3] B. LEGUY. *Ada : Guide d'utilisation*. Eyrolles — Un excellent ouvrage d'initiation. Très pédagogique.
- [4] R. OGOR, R. RANNOU. *Langage Ada et algorithmique*. Hermes — On trouvera de nombreux exemples. C'est aussi une bonne initiation au langage.

En deuxième lecture :

- [5] G. BOOCH. *Ingénierie du logiciel avec Ada*. Addison Whesley.
- [6] CH. CARREZ. *Des structures aux bases de données*. Dunod.
- [7] N. H. COHEN. *Ada as a second language*. Mc Graw Hill.
- [8] F. CULWIN. *Ada a developmental approach*. Prentice Hall.
- [9] D. J. DAVID. *Le langage Ada*. Editions du PSI.
- [10] B. FORD, J. KOK, M. W. ROGERS. *Scientific Ada*. Cambridge University Press.
- [11] M. GAUTHIER. *Ada, un apprentissage*. Dunod.
- [12] N. GEHANI. *Ada. Introduction avancée*. Eyrolles.
- [13] P. LIGNELET. *Structures de données avec Ada. Conception orientée objet*. 2 Vol. Masson.
- [14] R. MITCHELL. *Abstract data types and Ada*. Prentice Hall.

## Chapitre 2

- [1] Logiciel Modulog : Edil Belin : Atelier/Courrier, 20 , Rue Gossin 92 120 Montrouge.
- [2] R. DONY. *Graphisme dans le plan et dans l'espace avec Turbo-Pascal 4.0*. Masson.
- [3] M. DUCAMP ET A. REVERCHON. *Mathématiques en Turbo-Pascal. Vol. 1 : Analyse*. Eyrolles.
- [4] A. LEVINE. *Calcul formel*. Ellipses.

## Chapitre 3

- [1] J. BARANGER. *Analyse numérique*. Hermann.
- [2] R. BURLISCH ET J. STOER. *Introduction to Numerical Analysis*. Springer Verlag.
- [3] CEA-EDF-INRIA. *Résolution numérique des grands systèmes linéaires*. Eyrolles.
- [4] P.G. CIARLET. *Introduction à l'analyse numérique matricielle et à l'optimisation*. Masson.



- [5] P.G CIARLET B. MIARA J.M. THOMAS. *Exercices d'analyse numérique matricielle et d'optimisation*. Masson.
- [6] A.S. HOUSEHOLDER. *The theory of matrices in numerical analysis*. Dover.
- [7] P. LASCAUX R. THEODOR. *Analyse numérique matricielle appliquée à l'art de l'ingénieur. Vol. 1 et 2*. Masson.
- [8] J.P. NOUGIER. *Méthodes de calcul numérique*. Masson.
- [9] C. NOWAKOWSKI. *Méthodes de calcul numérique. Programmes en Basic et en Pascal. Vol. 1 et 2*. Editions du PSI..
- [10] PRESS, FLANNERY, TEUKOLSKY, VETTERLING. *Numerical Recipes*. Cambridge Univ. Press.
- [11] C. REINSCH J.H. WILKINSON. *Linear Algebra*. Springer Verlag.
- [12] J. VIGNES. *Algorithmes numériques. Analyse et mise en oeuvre*. Vol. 1 et 2. Technip.

#### Chapitre 4

- [1] J. BARANGER. *Analyse numérique*. Hermann.
- [2] R. BURLISCH ET J. STOER. *Introduction to Numerical Analysis*. Springer Verlag.
- [3] D. MONASSE. *La méthode de Newton et les racines des polynômes*. RMS, Janvier 1987.
- [4] J.P. NOUGIER. *Méthodes de calcul numérique*. Masson.
- [5] J. M. ORTEGA, W. C. RHEINBOLDT. *Iterative solution of nonlinear equations in several variables*. Academic Press.
- [6] PRESS, FLANNERY, TEUKOLSKY, VETTERLING. *Numerical Recipes*. Cambridge Univ. Press.
- [7] R. THEODOR. *Initiation à l'analyse numérique*. Masson.

#### Chapitre 5

- [1] AHLBERG, NILSON, WALSH. *The theory of Splines and their applications. Mathematics in science and engineering*. Vol. 38. Academic Press.
- [2] R. H. BARTELS, C. BEATY, A. BARSKY. *Mathématiques et CAO*. Vol 6. Hermes.
- [3] R. BURLISCH, J. STOER. *Introduction to Numerical Analysis*. Springer Verlag.
- [4] CHUEN YEN CHOW. *An introduction to computational fluid mechanics*. Seminole Publishing Company.
- [5] PH. DAVIS. *Interpolation and approximation*. Dover Publication.
- [6] R. DONY. *Calcul des parties cachées, approximation des courbes par la méthode de Bézier et des B\_Splines*. Masson.
- [7] ENCYCLOPEADIA UNIVERSALIS. Article : *Représentation et approximation des fonctions*.
- [8] S. A. MORSE, A. J. ALEXANDER. *An investigation of particule trajectories in two phase flow systems*. Coughborough university of technology.
- [9] PICHAT, DI CRESCENZO, WOLF. *Mathématiques pour l'informatique. Exercices et problèmes*. Armand Colin.

- [10] PRESS, FLANNERY, TEUKOLSKY, VETTERLING. *Numerical Recipes*. Cambridge.
- [11] R. THEODOR. *Initiation à l'analyse numérique*. Masson.

### Chapitre 6

- [1] AYANT ET BORG. *Fonctions spéciales*. Dunod.
- [2] M. CROUZEIX ET A. MIGNOT. *Analyse numérique des équations différentielles*. Masson.
- [3] PH. DAVIS. *Interpolation and approximation*. Dover Publication.
- [4] P. DEHEUVELS. *L'intégrale*. PUF.
- [5] N. N. LEBEDEV. *Special functions and their applications*. Prentice-Hall.
- [6] P. PELLETIER. *Techniques numériques appliquées au calcul scientifique*. Masson.
- [7] PRESS, FLANNERY, TEUKOLSKY, VETTERLING. *Numerical recipes*. Cambridge.
- [8] B. RANDE. *Primitives élémentaires*. RMS, Mai 89. Vuibert.
- [9] J. STOER ET R. BURLISCH. *Introduction to numerical analysis*. Springer Verlag.
- [10] A. WARUSFEL. *A propos de la méthode d'intégration de Romberg*. RMS, Janvier 86. Vuibert.

### Chapitre 7

- [1] R. BURLISCH ET J. STOER. *Introduction to numerical Analysis*. Springer-Verlag.
- [2] CROUZEIX ET MIGNOT. *Analyse numérique des équations différentielles*. Masson.
- [3] J. P. DEMAILLY. *Analyse numérique des équations différentielles*. Presses Universitaires de Grenoble.
- [4] M. DUCAMP ET A. REVERCHON. *Mathématiques et Turbo-Pascal*. Eyrolles.
- [5] G. HACQUES. *Mathématiques pour l'informatique. Vol. 3 : Algorithmique numérique*. Armand Colin.
- [6] P. HENRICE. *Discrete variable methods in ordinary differential equations*. Wiley.
- [7] P. HENRICE. *Error propagation for difference methods*. Wiley.
- [8] A. HILLION. *Les théories mathématiques des populations*. PUF..
- [9] J. P. NOUGIER. *Méthodes de calcul numérique*. Masson.
- [10] PRESS, FLANNERY, TEUKOLSKY, VETTERLING. *Numerical recipes*. Cambridge.

### Chapitre 8

- [1] L. C. BARRET, C. R. WYLIE. *Advanced engineering mathematics*. Mac Graw Hill.
- [2] R. BURLISCH, J. STOER. *Introduction to numerical Analysis*. Springer Verlag.
- [3] H. CARTAN. *Théorie élémentaire des fonctions analytiques d'une ou plusieurs variables complexes*. Masson.
- [4] P. G. CIARLET. *Introduction à l'analyse numérique matricielle et à l'optimisation*. Masson.

- [5] D. EUVRARD. *Résolution numérique des équations aux dérivées partielles*. Masson.
- [6] FLANNERY, PRESS, TEUKOLSKY, VETTERLING. *Numerical Recipes*. Cambridge.
- [7] J. P. NOUGIER. *Méthodes de calcul numérique*. Masson.
- [8] H. REINHARD. *Equations aux dérivées partielles*. Dunod.