

# Mathématiques

Mazen SAAD

mazen.saad@ec-nantes.fr

Quelques notes sur les cours :

- 16h de cours, 44h de TD, 12h de TP,
- 7 chapitres,
- 2 DS (novembre & janvier): note EVI coef. 1/3 par DS,
- compte-rendus de TP (note EVC) coef 1/3.

Les chapitres de cette année :

- (1) Introduction à l'analyse numérique (méthodes directes de résolution de systèmes linéaires),
- (2) Méthodes itératives de résolution de systèmes linéaires d'approximation de valeurs propres,
- (3) Optimisation sans contrainte,
- (4) Optimisation avec contraintes,
- (5) Probabilités,
- (6) Statistiques,
- (7) Interpolation et approximation, intégration numérique.

# Table des matières

I. Introduction à l'analyse numérique, méthodes directes de résolution de systèmes linéaires .....	5
I.1. Arithmétiques flottantes .....	5
I.2. Calcul matriciel (document sur le serveur pédagogique) .....	6
I.3. Systèmes linéaires creux .....	6
I.3.a. Équation de la chaleur .....	6
I.3.b. Méthodes directes pour résoudre $Ax = b$ .....	8
I.3.b.i. Méthode de Cramer .....	8
I.3.b.ii. Élimination de Gauss .....	8
I.3.b.iii. Méthode de Choleski .....	8
I.4. Graphe associé à une matrice et inversement .....	9
I.5. Les matrices irréductibles .....	10
I.6. Localisation des valeurs propres .....	11
II. Méthodes itératives de résolution de systèmes linéaires .....	13
II.1. Convergence des méthodes itératives .....	13
II.2. Méthodes itératives classiques .....	14
II.2.a. Méthode de Jacobi .....	15
II.2.b. Méthode de Gauss-Seidel .....	15
II.2.c. Méthode de relaxation .....	16
II.2.d. Méthode de Richardson .....	16
II.2.e. Méthode du gradient à pas optimal .....	16
II.3. Calcul de valeurs et de vecteurs propres .....	18
II.3.a. Méthode de puissance itérées pour calculer la plus grande des valeurs propres en module .....	18
II.3.b. Méthode de la puissance inverse pour calculer la plus petite des valeurs propres en module .....	19
III. Optimisation sans contrainte .....	20
III.1. Introduction .....	20
III.2. Dérivabilité. ....	21
III.3. Existence des minimums sur $\mathbb{R}^N$ .....	22
III.4. Points stationnaires, et caractérisation d'un minimum .....	24
III.4.a. Étude de la nature d'un point stationnaire. ....	24
III.5. Convexité .....	25

IV. Optimisation non linéaire avec contraintes .....	28
IV.1. Existence d'un minimum .....	28
IV.2. Unicité du minimum .....	28
IV.3. Caractérisation du minimum sur un convexe .....	30
IV.4. Multiplicateurs de Lagrange .....	32
IV.4.a. Cas des contraintes d'égalité .....	32
IV.4.b. Multiplicateurs de Lagrange – cas général .....	33
V. Probabilités .....	35
V.1. Variables aléatoires réelles .....	37

# I. Introduction à l'analyse numérique, méthodes directes de résolution de systèmes linéaires

## I.1. Arithmétiques flottantes

Comment les nombres sont représentés dans une machine (un ordinateur) ? Un nombre est représenté par un nombre fini de caractères. Des nombres comme  $\sqrt{2}$  doivent être tronqués.

Un nombre réel est représenté par un un *nombre flottant*, où un nombre flottant est de la forme suivante :

$$\pm a \cdot 10^q$$

où  $q$  est un entier tel que  $-M \leq q \leq M$ , et  $a$  est la *mantisse* de la forme :

$$0.d_1 d_2 \dots d_t$$

avec  $d_1 \neq 0$  et  $t$  le nombre de chiffres de la mantisse.

Un nombre peut être représenté de plusieurs manières.

**Simple précision.** On utilise 4 octets (32 bits), on a  $M = 38$  et  $t = 7$ .

La précision est de 7 chiffres décimaux.

**Double précision.** On utilise 8 octets (64 bits), on a  $M = 308$  et  $t = 15$ .

La précision est de 15 chiffres décimaux. Ceci est le type le plus utilisé (type double en C).

Les conséquences de cette précision sont :

- (1) Il existe un plus petit nombre flottant machine (appelé *zéro machine*). Il est accessible avec `eps` en Matlab. Ce nombre est donc de la forme  $0,100 \dots 00 \cdot 10^{-M}$  (en valeur absolue).
- (2) Il existe un plus grand nombre :  $0, \underbrace{99 \dots 9}_{t \text{ fois}} \cdot 10^{+M}$ .
- (3) Tous les réels ne sont pas représentés :

$$\sqrt{2} \approx 0,1414234 \cdot 10^1 \quad \text{et} \quad \pi \approx 0,31415925 \cdot 10^1.$$

- (4) Toute opération élémentaire ( $+$ ,  $\times$ ,  $\div$ ) est en général entachée d'une erreur.

**Exemple.** Avec une machine avec  $t = 2$ , on pose  $a = 0,63 \cdot 10^1$  et  $b = 0,82 \cdot 10^{-4}$ . On fait la somme et l'ordinateur réduit au même exposant (le plus grand) :

$$\begin{aligned} a &: 0,6300000 \cdot 10^1 \\ b &: 0,0000082 \cdot 10^1 \\ a + b &: 0,6300082 \cdot 10^1 \end{aligned}$$

On a  $\text{fl}(a + b) = 0,63 \cdot 10^1$ , donc  $\text{fl}(a + b) = \text{fl}(a)$  et  $b \neq 0$ . Par contre,  $a + b > a$  car  $b > 0$ .

Si on note  $z = a + b$ , alors

$$\frac{1}{\text{fl}(z - a)} = \frac{1}{0} = +\infty \quad \text{au lieu de} \quad \frac{1}{b}.$$

Le réel  $\sqrt{2}$  est calculé  $10^{-7}$  près en simple précision. Si on l'utilise  $10^7$  fois, alors l'erreur est d'environ 1.

## I.2. Calcul matriciel (document sur le serveur pédagogique)

### I.3. Systèmes linéaires creux

#### I.3.a. Équation de la chaleur

On note  $u(x)$  la distribution de température d'une plaque chauffée aux deux extrémités (Figure 1).

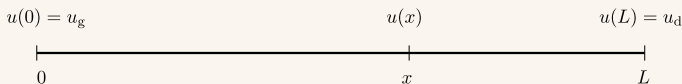


Figure 1: Plaque chauffée

On ajoute les conditions du système ci-dessous :

$$\begin{cases} -u''(x) = f(x) & 0 < x < L \\ u(0) = u_g \\ u(L) = u_d \end{cases}$$

où  $f$  est le *terme source*. Si  $f$  est intégrable, alors  $u$  est exacte. Sinon, on approxime la solution  $u$ .

**Analyse numérique.** On approxime la solution par une méthode d'approximation comme la *méthode des différences finies*. On décompose l'intervalle  $[0, L]$  en  $(N + 1)$  intervalles de pas  $h$ , où

$$h = \frac{L}{N + 1}$$

, comme montré dans la [Figure 2](#).

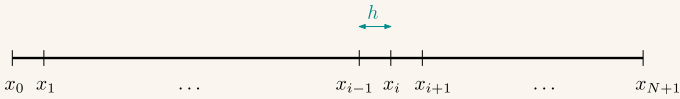


Figure 2: Approximation de la plaque chauffée

Ainsi, on a  $x_i = i \cdot h$ , pour  $i = 0, N + 1$ , et donc  $x_{i+1} = x_i + h$ .

On écrit le développement de Taylor :

$$-u''(x_i) = \frac{-u(x_{i-1}) + 2u(x_i) - u(x_{i+1}))}{h^2} + \mathcal{O}(h^2).$$

On note  $u_i \approx u(x_i)$  une approximation de la solution exacte au point  $x_i$ . La méthode des différences finies s'écrit comme le système :

$$(S_L) : \begin{cases} \frac{-u_{i+1} + 2u_i - u_{i-1}}{h^2} = f(x_i), & i = 1, N \\ u(0) = u_0 = u_g \\ u(L) = u_{N+1} = u_d \end{cases}$$

Le système  $(S_L)$  est un système linéaire. On cherche  $\mathbf{U}_h = (U_1, U_2, \dots, U_N)^\top$  solution de  $(S_L)$ . Le système linéaire  $(S_L)$  est équivalent à résoudre  $\mathbf{A}_h \cdot \mathbf{U}_h = \mathbf{b}_h$  avec

$$\mathbf{A}_h = \frac{1}{h^2} \begin{pmatrix} 2 & -1 & 0 & \cdots & 0 \\ -1 & 2 & -1 & \ddots & \vdots \\ 0 & -1 & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & -1 \\ 0 & \cdots & 0 & -1 & 2 \end{pmatrix} \in \mathcal{M}_{N,N}(\mathbb{R}) \quad \text{et} \quad \mathbf{b}_h = \begin{pmatrix} f(x_1) + \frac{u_g}{h^2} \\ \vdots \\ f(x_i) \\ \vdots \\ f(x_N) + \frac{u_d}{h^2} \end{pmatrix}.$$

### Remarques.

- (1)  $\mathbf{A}_h$  est une matrice creuse : le nombre de coefficients non nuls est  $3(N-2) + 4 \approx 3N$ , à comparer avec  $N^2$  coefficient pour une matrice pleine.
- (2)  $\mathbf{A}_h^{-1}$  est une matrice pleine.

## I.3.b. Méthodes directes pour résoudre $\mathbf{A}\mathbf{x} = \mathbf{b}$

### I.3.b.i. Méthode de Cramer

On pose  $x_i = \det \mathbf{A}_i / \det \mathbf{A}$  où  $\mathbf{A}_i$  est la matrice de  $\mathbf{A}$  en remplaçant la  $i$  ème colonne par  $\mathbf{b}$ . La *complexité de l'algorithme* est le nombre d'opérations nécessaires pour calculer la solution. On a besoin de  $(N+1) \times N!$  opérations pour la méthode de Cramer (le  $N!$  vient du calcul du déterminant). Pour  $N = 25$ , et la vitesse de calcul d'un ordinateur est 1 G/s, c'est à dire  $10^9$  opérations chaque seconde. Pour résoudre un système linéaire, il faut  $26 \times 25! \approx 10$  milliards d'années (deux fois l'âge de la Terre).

### I.3.b.ii. Élimination de Gauss

Elle consiste à trouver un système équivalent à  $\mathbf{A}\mathbf{x} = \mathbf{b}$  mais le système est triangulaire supérieur :

$\mathbf{A}$  matrice quelconque  $\rightsquigarrow \tilde{\mathbf{A}}$  triangulaire supérieure

$$\mathbf{A}\mathbf{x} = \mathbf{b} \iff \tilde{\mathbf{A}}\mathbf{x} = \tilde{\mathbf{b}}$$

### I.3.b.iii. Méthode de Choleski

La méthode de Choleski consiste à décomposer la matrice  $\mathbf{A}$  sous la forme  $\mathbf{A} = \mathbf{L} \cdot \mathbf{U}$  où  $\mathbf{L}$  est une matrice triangulaire supérieure avec  $\ell_{i,i} = 1$  et  $\mathbf{U}$  est une matrice triangulaire supérieure. Le calcul de ces



deux matrices est possible via un algorithme. Une fois ces matrices calculées, trouver la solution est facile à résoudre :

$$Ax = b \Leftrightarrow A \underbrace{Ux}_y = b \Leftrightarrow \begin{cases} Ly = b \\ Ux = y \end{cases}$$

#### I.4. Graphe associé à une matrice et inversement

Soit  $A = (a_{i,j})_{i,j}$  une matrice carrée d'ordre  $N$ . À chaque colonne, on fait correspondre un sommet  $S_i$  pour  $i = 1, N$ . Un arc relie deux sommets  $s_i$  et  $s_j$  si  $a_{i,j} \neq 0$ . Un graphe est formé de l'ensemble des sommets et des arcs.

**Exemple.** On considère la matrice

$$A = \begin{pmatrix} 4 & 3 & 0 & 0 \\ 0 & 2 & 1 & 2 \\ 0 & 1 & 0 & 3 \\ 6 & 5 & 0 & 0 \end{pmatrix}$$

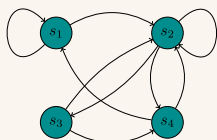


Figure 3: Graphe associé à la matrice  $A$

! Le graphe donne la « structure creuse » de la matrice  $A$ .

Un *chemin* allant de  $S_i$  à  $S_j$  est une suite d'arcs, si elle existe, telle que la suite

$$(s_i, s_{i_1}), (s_{i_1}, s_{i_2}), \dots, (s_{i_p}, s_j)$$

soient des arcs du graphe.

Un graphe est *fortement connexe* s'il existe au moins un chemin de tout sommet  $s_i$  à tout sommet  $s_j$ .

**Exemple.** Avec  $A = \begin{pmatrix} 3 & 2 & 5 \\ 4 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$ , le graphe n'est pas fortement connexe : il n'existe pas de chemin de  $s_3$  à  $s_1$ , comme le montre la Figure 4.

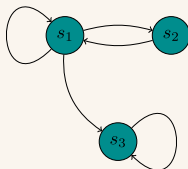


Figure 4: Graphe associé à la matrice  $A$

## I.5. Les matrices irréductibles

Soit  $A \in \mathcal{M}_{N,N}(\mathbb{R})$  de la forme

$$A = \begin{pmatrix} A_{1,1} & A_{1,2} \\ \mathbf{0} & A_{2,2} \end{pmatrix} \quad \text{et} \quad \mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}$$

Ainsi, il est équivalent :

$$A\mathbf{x} = \mathbf{b} \iff \begin{cases} A_{1,1}\mathbf{x}_1 = \mathbf{b}_1 - A_{1,2}\mathbf{x}_2 & \rightsquigarrow \text{on obtient } \mathbf{x}_1 \\ A_{1,2}\mathbf{x}_2 = \mathbf{b}_2 & \rightsquigarrow \text{on obtient } \mathbf{x}_2 \end{cases}$$

Le système est réduit à la résolution de deux systèmes de plus petite taille.

**Définition.** Une matrice  $A$  d'ordre  $N$  est réductible si, et seulement si, il existe  $P$  une matrice de permutation telle que :

$$B = P^\top A P = \begin{pmatrix} B_{1,1} & B_{1,2} \\ \mathbf{0} & B_{2,2} \end{pmatrix},$$

si, et seulement si, il existe une permutation  $\sigma : \{1, \dots, N\} \mapsto \{I, J\}$  où  $b_{i,j} = 0$  pour tout  $i \in I$ , et  $j \in J$ .

**Exemple.** Avec la matrice  $B$  ci-dessous, on applique la permutation  $\sigma$  définie ci-dessous pour trouver la matrice  $\tilde{B} = B_\sigma = (b_{\sigma(i),\sigma(j)})_{i,j}$ .

$$B = \begin{pmatrix} 30 & 0 & 20 & 0 \\ 11 & 10 & 6 & 5 \\ 13 & 0 & 11 & 0 \\ 4 & 2 & 3 & -1 \end{pmatrix} \quad \sigma : \begin{pmatrix} 1 \mapsto 4 \\ 2 \mapsto 2 \\ 3 \mapsto 3 \\ 4 \mapsto 1 \end{pmatrix} \quad \tilde{B} = \begin{pmatrix} \begin{pmatrix} -1 & 0 \\ 5 & 6 \end{pmatrix} & \begin{pmatrix} 3 & 4 \\ 10 & 11 \end{pmatrix} \\ \mathbf{0} & \begin{pmatrix} 11 & 13 \\ 20 & 30 \end{pmatrix} \end{pmatrix}$$

La matrice  $B$  est donc réductible.

**Proposition** (admis)  $A$  est irréductible si, et seulement si, le graphe de  $A$  est fortement connexe.

## I.6. Localisation des valeurs propres

**Théorème 1** (Gerschgorin-Hadamond) Soit  $\lambda$  une valeur propre de  $A$ , alors

$$\lambda \in \bigcup_{k=1}^N \bar{D}_k \quad \text{avec} \quad \bar{D}_k = \left\{ z \in \mathbb{C}^2 \mid |z - a_{k,k}| \leq \sum_{\substack{j=1 \\ j \neq k}}^N |a_{k,j}| \right\}$$

Les valeurs propres de  $A$  appartiennent à l'union de  $N$  disques  $\bar{D}_k$  de Gerschgorin.

**Preuve.** Soit  $(\lambda, \mathbf{u})$  un élément propre de  $A$  tel que  $\max_{i=1, \dots, N} |u_i| = |u_k| = 1$ .

On a  $A\mathbf{u} = \lambda\mathbf{u}$  donc  $(A\mathbf{u})_k = \lambda u_k$ .

Ainsi,  $\sum_{i=0}^N a_{k,j} u_k = \lambda u_k$ , ce qui est équivalent à,  $(\lambda - a_{k,k})u_k = \sum_{k \neq j} a_{k,j} u_k$ .

Donc  $|\lambda - a_{k,k}| \cdot |u_k| \leq \sum_{j \neq k} |a_{k,j}| \cdot |u_j|$ .

Et donc  $|\lambda - a_{k,k}| \leq \sum_{j \neq k} |a_{k,j}|$ , i.e.,  $\lambda \in \bar{D}_k$ .

**Exemple.**

$$A = \begin{pmatrix} 1+i & i & 2 \\ -3 & 2+i & 2 \\ 1 & i & 6 \end{pmatrix}$$

La matrice  $A$  a pour disques de Gerschgorin :

- $\bar{D}_1 = \{\lambda \mid |\lambda - (1 + i)| \leq |i| + 2 = 3\}$ ,
- $\bar{D}_2 = \{\lambda \mid |\lambda - (2 + i)| \leq |-3| + 1 = 4\}$ ,
- $\bar{D}_3 = \{\lambda \mid |\lambda - 6| \leq 1 + |i| = 2\}$ .

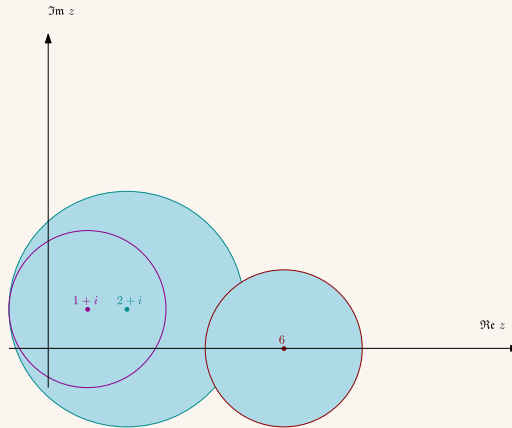


Figure 5: Disques de Gerschgorin pour la matrice  $A$

Dans la Figure 5, les valeurs propres se situent dans la zone bleue, correspondant à l'union des disques de Gerschgorin.

Si  $0$  n'est dans aucun disque, alors on sait que la matrice est inversible. On peut ainsi localiser les valeurs propres d'une matrice : on exclut tous les points en dehors des disques.

**Théorème 2** (Gerschgorin-Hadamond) Soit  $\lambda$  une valeur propre de  $A$  et  $A$  une matrice irréductible<sup>[1]</sup>. Si une valeur propre  $\lambda$  est située sur la frontière de la réunion des disques alors tous les cercles de Gerschgorin passent par  $\lambda$  :

$$\lambda \in \partial \left( \bigcup_{k=1}^N \bar{D}_k \right) \implies \lambda \in \bigcap_{k=1}^N \partial \bar{D}_k$$

On utilise souvent cela pour montrer que  $\lambda = 0$  n'est pas valeur propre.

<sup>[1]</sup> i.e. le graphe est fortement connexe

## II. Méthodes itératives de résolution de systèmes linéaires

### II.1. Convergence des méthodes itératives

Le principe est d'engendrer une suite de vecteurs  $(\mathbf{x}_k)_{k \in \mathbb{N}}$ , que l'on nomme les *itérés*, convergent vers la solution  $\mathbf{x}$  du système linéaire  $A\mathbf{x} = \mathbf{b}$ . La méthode itérative consiste à la procédure suivante.

Soit  $\mathbf{x}_0 \in \mathbb{R}^N$  (ou potentiellement  $\mathbb{C}^N$ ). On engendre une suite  $(\mathbf{x}_k)_{k \in \mathbb{N}}$  définies par :

$$(1) : \quad \mathbf{x}_{k+1} = B\mathbf{x}_k + \mathbf{c},$$

où  $B \in \mathcal{M}_{N \times N}$  est la *matrice d'itération* et  $\mathbf{c} \in \mathbb{R}^N$ .

**Définition.** Une méthode itérative de la forme (1) est convergente si, et seulement si,

$$\forall \mathbf{x}_0 \in \mathbb{R}^N, \quad \mathbf{x}_k \xrightarrow[k \rightarrow +\infty]{} \mathbf{x} \quad \text{et} \quad A\mathbf{x} = \mathbf{b}$$

La méthode est *consistante* si, et seulement si, si  $\mathbf{x}_k \rightarrow \mathbf{x}$ , alors  $\mathbf{x} = B\mathbf{x} + \mathbf{c} \Leftrightarrow A\mathbf{x} = \mathbf{b}$ .

**Définition.** L'erreur d'approximation à la  $k$ -ième itération est

$$\begin{aligned} e_k &= \mathbf{x}_k - \mathbf{x} = B\mathbf{x}_{k-1} + \mathbf{c} - B\mathbf{x} + \mathbf{c} \\ &= B(\mathbf{x}_{k-1} - \mathbf{x}) \\ &= B\mathbf{e}_{k-1} \\ &= B^2\mathbf{e}_{k-2} \end{aligned}$$

Donc  $e_k = B^k e_0$ , où  $e_0 = \mathbf{x}_0 - \mathbf{x}$ .

**Définition.** Une méthode itérative est convergente si, et seulement si,

$$\begin{aligned}
\forall \mathbf{x}_0, \quad \mathbf{e}_k \xrightarrow[k \rightarrow +\infty]{} \mathbf{0} &\iff \mathbf{B}^k \xrightarrow[k \rightarrow +\infty]{} \mathbf{0} \\
&\iff \|\mathbf{B}^k\| \xrightarrow[k \rightarrow +\infty]{} 0 \\
&\iff \mathbf{B}^k \mathbf{x} \xrightarrow[k \rightarrow +\infty]{} \mathbf{0}, \forall \mathbf{x} \in \mathbb{R}^N.
\end{aligned}$$

**Théorème** (Convergence des méthodes itératives).

$$\lim_{k \rightarrow +\infty} \mathbf{B}^k = \mathbf{0} \iff \rho(\mathbf{B}) < 1.$$

où  $\rho(\mathbf{B}) = \max |\lambda(\mathbf{B})|$ .

*Preuve.*

$\implies$ . Si  $\rho(\mathbf{B}) > 1$ , il existe  $\lambda$  une valeur propre de  $\mathbf{B}$  telle que  $|\lambda| \geq 1$ . Soit  $\mathbf{x} \neq \mathbf{0}$  un vecteur propre associé à la valeur propre  $\lambda$ . Alors,  $\mathbf{B}^k \mathbf{x} = \lambda^k \mathbf{x}$ . Et,  $|\lambda^k| \rightarrow 1$  ou vers  $+\infty$ . On en déduit donc que  $\mathbf{B}^k \mathbf{x} \rightarrow \mathbf{0}$ .

$\impliedby$ . On suppose  $\rho(\mathbf{B}) < 1$ , i.e.  $|\lambda| < 1$  pour  $\lambda \in \text{Sp } \mathbf{B}$ .

- Si  $\mathbf{B}$  est diagonalisable, alors on pose  $\mathbf{B} = \mathbf{P}\mathbf{A}\mathbf{P}^{-1}$  où  $\mathbf{A} = \text{diag}(\lambda_1, \dots, \lambda_n)$ . Ainsi,  $\mathbf{B}^k = \mathbf{P}\mathbf{A}^k\mathbf{P}^{-1}$ . Et, on sait que  $\mathbf{A}^k = \text{diag}(\lambda_1^k, \dots, \lambda_n^k)$ . Comme  $|\lambda_i| < 1$ , on sait que  $\lambda_i^k \rightarrow 0$  et donc  $\mathbf{B}^k \rightarrow \mathbf{0}$ .
- Cas général : toute matrice est semblable à une matrice de Jordan.

$$\tilde{\mathbf{B}} = \mathbf{S}\mathbf{B}\mathbf{S}^{-1} = \text{diag}(\mathbf{J}_1, \dots, \mathbf{J}_r)$$

où  $\mathbf{J}_k = \lambda_k \mathbf{I} + \mathbf{E}$ , où  $\mathbf{E} = \begin{pmatrix} 0 & \mathbf{I} \\ 0 & 0 \end{pmatrix}$  (matrice de surdiagonale de 1).  
Donc  $\mathbf{J}_k^n \rightarrow \mathbf{0}$  car  $\mathbf{E}^k = 0$  à partir d'un certain rang, car la matrice est nilpotente. Ainsi, comme  $\tilde{\mathbf{B}}^k = \mathbf{S}\mathbf{B}^k\mathbf{S}^{-1}$ , on a l'équivalence  $\tilde{\mathbf{B}}^k \rightarrow \mathbf{0} \iff \mathbf{B}^k \rightarrow \mathbf{0}$ .

**Corrolaire.** Si  $\|\mathbf{B}\| < 1$ , alors  $\rho(\mathbf{B}) < 1$  et la méthode converge donc.

(Ceci vient de  $\rho(\mathbf{B}) \leq \|\mathbf{B}\| < 1$ , montré en TD.)

## II.2. Méthodes itératives classiques

### II.2.a. Méthode de Jacobi

Soit  $A$  une matrice d'ordre  $n$  telle que  $a_{i,i} \neq 0$ , pour tout  $i$ . On décompose  $A$  sous la forme  $A = D - E - F$ , où

- $D$  est la diagonale,
- $-E$  est la partie inférieure stricte,<sup>[2]</sup>
- $-F$  est la partie supérieure stricte,

**Remarque.** On a l'équivalence suivante :

$$Ax = b \iff (D - E - F)x = b \iff Dx = (E + F)x + b.$$

Le principe de la *méthode de Jacobi* est :

$$\begin{cases} \mathbf{x}_0 \in \mathbb{R}^N \\ D\mathbf{x}_{k+1} = (E + F)\mathbf{x}_k + \mathbf{b} \end{cases} \iff \begin{cases} \mathbf{x}_0 \in \mathbb{R}^N \\ \mathbf{x}_{k+1} = \underbrace{D^{-1}(E + F)}_J \mathbf{x}_k = D^{-1}\mathbf{b}. \end{cases}$$

La matrice  $J = D^{-1}(E + F)$  est la matrice d'itération de Jacobi. On peut ré-écrire  $J$  comme  $D^{-1}(D - A) = I - D^{-1}A$ . Ainsi, la méthode de Jacobi converge si, et seulement si,  $\rho(J) < 1$ .

**Algorithme.** Pour calculer les itérés de la méthode de Jacobi :

$$D\mathbf{x}_{k+1} = (E + F)\mathbf{x}_k + \mathbf{b} \iff a_{i,i}(\mathbf{x}_{k+1})_i = - \sum_{j \neq i} a_{i,j}(\mathbf{x}_k)_j + b_i$$

Ainsi,  $(\mathbf{x}_{k+1})_i = \left( b_i - \sum_{j \neq i} (\mathbf{x}_k)_j \right)$  pour tout  $i$ .

### II.2.b. Méthode de Gauss-Seidel

**Remarque.** On décompose  $A = D - E - F$ . Il y a l'équivalence suivante :

$$(D - E - F)x = b \iff (D - E)x = Fx + b.$$

---

<sup>[2]</sup>On veut dire ici sans la diagonale.

Le principe de la méthode de Gauss-Seidel est : on pose  $\mathbf{x}_0 \in \mathbb{R}^N$ , et on itère avec  $(\mathbf{D} - \mathbf{E})\mathbf{x}_{k+1} = \mathbf{F}\mathbf{x}_k + \mathbf{b}$ . Sous forme matricielle, cela devient :  $\mathbf{x}_{k+1} = \underbrace{(\mathbf{D} - \mathbf{E})^{-1}\mathbf{F}}_{\mathbf{G}}\mathbf{x}_k + (\mathbf{D} - \mathbf{E})^{-1}\mathbf{b}$ .

La méthode de Gauss-Seidel converge si, et seulement si,  $\rho(\mathbf{G}) < 1$ . Pour calculer  $\mathbf{x}_{k+1}$ , il suffit d'appliquer l'algorithme de descente sur  $\mathbf{D} - \mathbf{E}$ . Pour  $i = 1, N$ , on a

$$a_{i,i}(\mathbf{x}_{k+1})_i = - \sum_{j=1}^{i-1} a_{i,j}(\mathbf{x}_{k+1})_j - \sum_{j=i+1}^N a_{i,j}(\mathbf{x}_k)_j + b_i.$$

### II.2.c. Méthode de relaxation

Soit  $\omega \neq 0$ . On décompose  $\mathbf{A}$  en  $\mathbf{A} = \left(\frac{\mathbf{D}}{\omega} - \mathbf{E}\right) + \left(\mathbf{D} - \frac{\mathbf{D}}{\omega} - \mathbf{F}\right)$ .

Le principe de la méthode de relaxation s'écrit :

$$\begin{aligned} \left(\frac{\mathbf{D}}{\omega} - \mathbf{E}\right)\mathbf{x}_{k+1} &= \left(\frac{1-\omega}{\omega}\mathbf{D} + \mathbf{F}\right)\mathbf{x}_k + \mathbf{b} \\ \Leftrightarrow (\mathbf{D} - \omega\mathbf{E})\mathbf{x}_{k+1} &= ((1-\omega)\mathbf{D} + \omega\mathbf{F})\mathbf{x}_k + \omega\mathbf{b}. \end{aligned}$$

Ainsi, la méthode de relaxation converge si, et seulement si,  $\rho(\mathbf{L}_\omega) < 1$ , où  $\mathbf{L}_\omega = (\mathbf{D} - \omega\mathbf{E})^{-1}((1-\omega)\mathbf{D} + \omega\mathbf{F})$ .

### II.2.d. Méthode de Richardson

On écrit  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha(\mathbf{b} - \mathbf{A}\mathbf{x}_k) = \underbrace{(\mathbf{I} - \alpha\mathbf{A})}_{\mathbf{R}}\mathbf{x}_k + \alpha\mathbf{b} = \mathbf{R}\mathbf{x}_k + \mathbf{b}$ .

C'est une méthode très simple, car il n'y a pas de système linéaire à résoudre. La recherche d'une valeur de  $\alpha$  pour que la méthode soit convergente sera faite en TD.

### II.2.e. Méthode du gradient à pas optimal

On écrit  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k(\mathbf{b} - \mathbf{A}\mathbf{x}_k) = \underbrace{(\mathbf{I} - \alpha_k\mathbf{A})}_{\mathbf{G}_k}\mathbf{x}_k + \alpha_k\mathbf{b}$ . On choisit  $\alpha_k$  de façon optimale.



La matrice d'itération  $\mathbf{G}_k = \mathbf{I} - \alpha_k \mathbf{A}$  dépend de  $k$ . On ne peut donc pas utiliser le théorème précédent pour montrer que la méthode converge (c.f. TD).

C'est la méthode la plus efficace dans le cas général.

**Définition.** Une matrice est dit

- à *diagonale strictement dominante* si, et seulement si

$$\forall i, \quad |a_{i,i}| > \sum_{j \neq i} |a_{i,j}|$$

- à *diagonale fortement dominante* si, et seulement si

$$\forall i, \quad |a_{i,i}| \geq \sum_{j \neq i} |a_{i,j}|$$

**et** il existe au moins un  $i_0$  tel que  $|a_{i_0,i_0}| > \sum_{j \neq i_0} |a_{i_0,j}|$ .

**Théorème.**

- (1) Soit  $\mathbf{A}$  une matrice réelle symétrique définie positive (ou hermitienne définie positive dans  $\mathbb{C}$ ). Alors, la méthode de relaxation converge si  $0 < \omega < 2$ .
- (2) Soit  $\mathbf{A}$  une matrice à diagonale strictement dominante (ou  $\mathbf{A}$  une matrice à diagonale fortement dominante **et** irréductible). Alors,
  - (a) la méthode de Jacobi converge.
  - (b) la méthode de Gauss-Seidel converge.
  - (c) la méthode de relaxation converge pour  $0 < \omega \leq 1$ .

**Preuve.** faite en TD.

**Remarque.** On peut montrer que  $\rho(\mathbf{G}) < \rho(\mathbf{J})$ . C'est à dire que, la méthode de Gauss-Seidel converge plus vite que celle de Jacobi. On pourra le vérifier en TP.

## II.3. Calcul de valeurs et de vecteurs propres

### II.3.a. Méthode de puissance itérées pour calculer la plus grande des valeurs propres en module

On pose une matrice exemple :  $\mathbf{A} = \begin{pmatrix} 10 & 0 \\ -0 & 1 \end{pmatrix}$ . Ses valeurs propres sont 10 et 1. Deux vecteurs propres associés sont  $(1 \ -1)^\top$  et  $(0 \ 1)^\top$ . Soit  $\mathbf{x}_0 = (1 \ 1)^\top$ . On calcule

- (1)  $\mathbf{x}_1 = \mathbf{A}\mathbf{x}_0 = (10 \ -8)^\top$ ,
- (2)  $\mathbf{x}_2 = \mathbf{A}\mathbf{x}_1 = (100 \ -98)^\top$ ,
- (3)  $\mathbf{x}_3 = \mathbf{A}\mathbf{x}_2 = (1000 \ -998)^\top$ ,
- (4) ...

On voit que  $\mathbf{x}_{k+1} \approx 10\mathbf{x}_k$  et  $\mathbf{x}_{k+1} \parallel (1 \ -1)^\top$ .

**Algorithme.** Pour le calcul des puissances itérées : soit  $\mathbf{q}_0 \in \mathbb{C}^N$  tel que  $\|\mathbf{q}_0\| = 1$  ; alors, pour  $k = 1, 2, \dots$ , on pose  $\mathbf{x}_k = \mathbf{A}\mathbf{x}_{k-1}$  et  $\mathbf{q}_k = \mathbf{x}_k / \|\mathbf{x}_k\|$ .

**Théorème.** Soit  $\mathbf{A}$  une matrice d'ordre  $N$ , diagonalisable dont la valeur propre du plus grand module  $\lambda_1$  est unique et vérifie :

$$|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_N|.$$

Soit  $\mathbf{q}_0 \in \mathbb{C}^N$  choisi *convenablement*. Alors,

- (1)  $\lim_{k \rightarrow +\infty} (\lambda_i / |\lambda_1|)^k \mathbf{q}_k = \mathbf{q}$  est un vecteur propre de norme 1 associé à  $\lambda_1$ .
- (2)  $\lim_{k \rightarrow +\infty} \|\mathbf{A}\mathbf{q}_k\| = \lim_{k \rightarrow +\infty} \|\mathbf{x}_k\| = |\lambda_1|$
- (3)  $\lim_{k \rightarrow +\infty} (\mathbf{x}_{k+1})_j / (\mathbf{q}_k)_i = \lambda_1$  pour  $1 \leq j \leq 1$  (?) tel que  $(\mathbf{q}_k)_j \neq 0$ .

**Preuve.** On suppose  $\mathbf{A}$  diagonalisable, et  $(\mathbf{u}_i)_i$  une base de vecteurs propres. On pose  $\mathbf{q}_0 = \alpha_1 \mathbf{u}_1 + \sum_{i=2}^N \alpha_i \mathbf{u}_i$  en supposant  $\alpha_k \neq 0$ .

- (1) On pose  $\mathbf{q}_1 = \mathbf{x}_1 / \|\mathbf{x}_1\| = \mathbf{A}\mathbf{q}_0 / \|\mathbf{A}\mathbf{q}_0\|$ . Par récurrence on montre que  $\mathbf{q}_k = (\mathbf{A}^k \mathbf{q}_0) / \|\mathbf{A}^k \mathbf{q}_0\|$ . Ainsi,

$$\begin{aligned} \mathbf{A}^k \mathbf{q}_0 &= \mathbf{A}^k \left( \alpha_1 \mathbf{u}_1 + \sum_{i \geq 2} \alpha_i \mathbf{u}_i \right) \\ &= \alpha_1 \mathbf{A}^k (\mathbf{u}_1 + \mathbf{e}_k) \end{aligned}$$

$$\text{où } e_k = \sum_{i \geq 2} (\alpha_i / \alpha_1) (\lambda_i / \lambda_1)^k \mathbf{u}_i.$$

(2) Exercice

(3) Exercice

### II.3.b. Méthode de la puissance inverse pour calculer la plus petite des valeurs propres en module

Les valeurs propres de  $\mathbf{A}^{-1}$  sont  $\mu_i = 1/\lambda_i$ . On suppose que l'on peut ordonner les valeurs propres comme  $|\lambda_1| \geq \dots \geq |\lambda_{N-1}| > |\lambda_N|$ . Ainsi, par passage à l'inverse,  $1/|\lambda_1| \leq \dots \leq 1/|\lambda_{N-1}| < 1/|\lambda_N|$ . Autrement dit,  $|\mu_1| \leq \dots < |\mu_N|$ .

Le réel  $|\mu_N| = 1/|\lambda_N|$  est la plus grande des valeurs propres de  $\mathbf{A}^{-1}$ . On applique l'algorithme des puissances itérées sur  $\mathbf{A}^{-1}$ .

# III. Optimisation sans contrainte

## III.1. Introduction

L'optimisation consiste à traiter le ou les éléments d'un ensemble « admissible » de données, pour trouver la meilleure solution, nommée solution optimale, qui minimise un certain critère.

Un *problème d'optimisation* s'écrit sous la forme suivante.

Soit  $E$  un espace vectoriel normé (evn) de dimension finie  $n$  (ou non<sup>[3]</sup>).  
Soit  $K$  un sous-ensemble (en général convexe et fermé) de  $E$ , et soit  
 $J : E \rightarrow \mathbb{R}$  une fonction.

On nomme  $E$  l'espace des états admissibles,  $K$  l'ensemble des constantes, et  $J$  la fonction de coût, ou fonction objectif.

Dans ce cours, on supposera  $K = E = \mathbb{R}^N$ . On dit que c'est un problème sans contraintes.

### **Exemple. (Problème du sac à dos)**

On souhaite maximiser l'utilité du sac à dos sous la contrainte du poids. On liste  $n$  objets : une bouteille d'eau, une loupe, un briquet, ...  
On suppose que le sac à dos ait une contrainte maximale de poids  $P$ , par exemple 15 kg. On note

- $p_1, p_2, \dots, p_n$  les poids respectifs des objets ;
- $u_1, u_2, \dots, u_n$  les utilités respectives des objets ;
- $x_i$  valant 1 si on a l'objet  $i$  dans le sac à dos, et 0 sinon.

Le problème d'optimisation est

- $J(\mathbf{x}) = x_1 u_1 + x_2 u_2 + \dots + x_n u_n$ ,
- $K = \{\mathbf{x} \in \mathbb{R}^n \mid x_1 p_1 + x_2 p_2 + \dots + x_n p_n \leq P\}$ .

On cherche  $\mathbf{x}^*$  tel que  $J(\mathbf{x}^*) = \max_{\mathbf{x} \in K} J(\mathbf{x})$ .

---

<sup>[3]</sup>Espace fonctionnel  $\mathcal{C}^0(\mathbb{R})$

### Exemple. (Détection d'un objet)

Un objet est situé sur une droite et sa position est inconnue. On dispose de deux radars donnant les distances *entachées d'erreurs*.

Les radars mesurent approximativement les distances  $a - s_1 \approx d_1$  et  $s_2 - a \approx d_2$ . Pour résoudre ce problème, on cherche à minimiser  $J(a) = \min_{x \in \mathbb{R}} J(x)$ , où

$$J(x) = |x - s_1 d_1|^2 + |s_2 - x - d_2|^2.$$

On dit que c'est un *problème quadratique*.

## III.2. Dérivabilité.

Soient  $E$  et  $F$  deux espaces vectoriels normés, et  $f : E \rightarrow F$ .

### Définition.

On dit que  $f$  est dérivable en  $\mathbf{x}$  s'il existe  $\ell_{\mathbf{x}} \in \mathcal{L}(E, F)$  telle que :

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \ell_{\mathbf{x}}(\mathbf{h}) + \|\mathbf{h}\|\varepsilon(\mathbf{h})$$

où  $\varepsilon(\mathbf{h}) \rightarrow \mathbf{0}$  quand  $\|\mathbf{h}\| \rightarrow 0$ . On note alors  $f'(\mathbf{x}) = \ell_{\mathbf{x}}$ .

### Exemple.

Dans le cas où  $f : \mathbb{R}^N \rightarrow \mathbb{R}$  est définie par  $f(\mathbf{x}) = \langle \mathbf{x} | \mathbf{x} \rangle = \sum_{i=1}^n x_i^2$ . Alors  $f'(\mathbf{x})\mathbf{h} = 2\langle \mathbf{x} | \mathbf{h} \rangle$ .

### Exemple.

Dans le cas où  $E = \mathbb{R}^N$ , on note  $\frac{\partial f}{\partial x_i}$

### Définition.

Pour  $f : \mathbb{R}^N \rightarrow \mathbb{R}$  deux fois dérivable, on a

$$f''(\mathbf{x}) = \nabla^2 f(\mathbf{x})$$

est une matrice symétrique de  $\mathcal{M}_{N,N}(\mathbb{R})$ . On la note  $\mathbf{H}_f$ . À compléter

### Formule de Taylor Lagrange.

Il existe  $\theta \in ]0, 1[$  tel que

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}) \mid \mathbf{h} \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x} + \theta \mathbf{h}) \mathbf{h} \mid \mathbf{h} \rangle.$$

### Formule de Taylor Young.

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}) \mid \mathbf{h} \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}) \mathbf{h} \mid \mathbf{h} \rangle + \|\mathbf{h}\|^2 \varepsilon(\mathbf{h}),$$

où  $\varepsilon(\mathbf{h}) \rightarrow \mathbf{0}$  quand  $\|\mathbf{h}\| \rightarrow 0$ .

### Formule de Taylor avec reste intégral.

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}) \mid \mathbf{h} \rangle + \int_0^1 (1-t) \langle \nabla^2 f(\mathbf{x} + t\mathbf{h}) \mid \mathbf{h} \rangle dt.$$

### Formule de Taylor à l'ordre 1.

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \int_0^1 \langle \nabla f(\mathbf{x} + t\mathbf{h}) \mid \mathbf{h} \rangle dt.$$

## III.3. Existence des minimums sur $\mathbb{R}^N$

On note  $V$  un espace vectoriel normé de dimension finie, dans ce cas  $V = \mathbb{R}^N$ . On suppose  $J : V \rightarrow \mathbb{R}$  continue. On cherche  $\min_{\mathbf{v} \in V} J(\mathbf{v})$ .

### Remarque.

On cherche toujours le minimum. En effet, pour trouver le maximum, on peut se ramener à la recherche d'un minimum par :

$$\max_{\mathbf{v} \in V} J(\mathbf{v}) = - \min_{\mathbf{v} \in V} (-J(\mathbf{v})).$$

### Définition (minimum local et global).

- On dit que  $\mathbf{x}$  est un point de minimum *global* de  $J$  si  $J(\mathbf{y}) \geq J(\mathbf{x})$ , quel que soit  $\mathbf{y} \in V$ .
- On dit que  $\mathbf{x}$  est un point de minimum *local* de  $J$  si  $J(\mathbf{y}) \geq J(\mathbf{x})$ , quel que soit  $\mathbf{y}$  voisin de  $\mathbf{x}$ . C'est à dire :

$$\forall \mathbf{y} \in V, \quad [\exists \varepsilon > 0, \|\mathbf{y} - \mathbf{x}\| \leq \varepsilon] \implies J(\mathbf{y}) \geq J(\mathbf{x}).$$

- On dit que  $\mathbf{x}$  est un point de *selle* de  $J$  si
  - $\mathbf{x}$  est un minimum local dans une direction,
  - $\mathbf{x}$  est un maximum local dans une autre direction.

### Définition (coercivité).

Soit  $J : \mathbb{R}^n \rightarrow \mathbb{R}$  continue. On dit que  $J$  est *coercive* sur  $\mathbb{R}^N$  si,  $J(\mathbf{x}) \rightarrow +\infty$  quand  $\|\mathbf{x}\| \rightarrow +\infty$ . On dira que  $J$  est *infinie* à l'infini.

La définition précédente en précise pas la « direction » de la limite. Par exemple, dans  $\mathbb{R}$ , une fonction est coercive si, et seulement si, sa limite en  $+\infty$  et en  $-\infty$  est  $+\infty$ .

### Exemples.

- La fonction  $J$  définie par  $J(\mathbf{x}) = x_1^3 + x_2^2$  est non coercive. En effet, on a  $J(x_1, 0) \rightarrow -\infty$  quand  $x_1 \rightarrow -\infty$ .
- La fonction  $J$  définie par  $J(\mathbf{x}) = x_1^2 + x_2^2$  est coercive.

### Théorème. (Existence d'un minimum)

Soit  $J$  une fonction continue et coercive sur  $\mathbb{R}^n$ . Alors  $J$  admet au moins un point de minimum global.

La preuve est sur le serveur pédagogique. C'est le seul résultat d'existence.

### Remarque.

Si  $J$  n'est pas coercive, on ne peut pas conclure. En effet, voici les deux cas :

- $J$  n'est pas coercive et admet un minimum,

- $J$  n'est pas coercive et n'admet pas de minimum.

### III.4. Points stationnaires, et caractérisation d'un minimum

#### Définition (point stationnaire).

Soit  $J$  dérivable en  $\mathbf{x}^*$ . On dit que  $\mathbf{x}^*$  est un *point stationnaire* de  $J$  si  $\nabla J(\mathbf{x}^*) = \mathbf{0}$ .

#### Théorème. (Caractérisation d'un minimum)<sup>[4]</sup>

Soit  $J$  dérivable en  $\mathbf{x}^*$ , où  $\mathbf{x}^*$  est un point d'extremum local de  $J$ , alors  $\nabla J(\mathbf{x}^*) = \mathbf{0}$ . Un *extremum local* est un minimum ou un maximum local.

*Preuve.*

Soit  $\varphi : t \mapsto J(\mathbf{x}^* + t\mathbf{y})$  est définie au voisinage de  $\mathbf{x}^*$ . On a  $\varphi'(t) = \langle \nabla J(\mathbf{x}^* + t\mathbf{y}) \mid \mathbf{y} \rangle$ . Ainsi, en  $t = 0$ , on a donc

$$\varphi'(0) = \langle \nabla J(\mathbf{x}^*) \mid \mathbf{y} \rangle.$$

- Si  $\mathbf{x}^*$  est un minimum local, on a donc  $J(\mathbf{x}^* + t\mathbf{y}) \geq J(\mathbf{x}^*)$  et donc  $J(\mathbf{x}^* + t\mathbf{y}) - J(\mathbf{x}^*) \geq 0$ . À finir . . .
- De même si  $\mathbf{x}^*$  est un maximum local.

#### III.4.a. Étude de la nature d'un point stationnaire.

Soit  $\mathbf{x}^*$  un point stationnaire. On suppose que  $J$  est de classe  $\mathcal{C}^2$ . On écrit le théorème de Taylor Young à l'ordre 2. Soit  $\mathbf{h} \in \mathbb{R}^n$ . On a :

$$J(\mathbf{x}^* + \mathbf{h}) = J(\mathbf{x}^*) + \underbrace{\langle \nabla J(\mathbf{x}^*) \mid \mathbf{h} \rangle}_{=0} + \frac{1}{2} \langle \nabla^2 J(\mathbf{x}^*) \mathbf{h} \mid \mathbf{h} \rangle + \|\mathbf{h}\|^2 \varepsilon(\mathbf{h}).$$

On a  $\nabla^2 J(\mathbf{x}^*)$  qui est symétrique et réelle. D'après le quotient de Rayleigh, on a

$$\lambda_{\min}^* \|\mathbf{h}\|^2 \leq \langle \nabla^2 J(\mathbf{x}^*) \mathbf{h} \mid \mathbf{h} \rangle \leq \lambda_{\max}^* \|\mathbf{h}\|^2,$$

<sup>[4]</sup>On nomme aussi ce théorème « l'équation d'Euler ».



où  $\lambda_{\min}^*$  et  $\lambda_{\max}^*$  sont la plus petite et la plus grande valeur propre des valeurs propres de  $\nabla^2 J(\mathbf{x}^*)$ . On en déduit que

$$(E) : \|\mathbf{h}\|^2 \left( \frac{1}{2} \lambda_{\min}^* + \varepsilon(\mathbf{h}) \right) \leq J(\mathbf{x}^* + \mathbf{h}) - J(\mathbf{x}^*) \leq \|\mathbf{h}\|^2 \left( \frac{1}{2} \lambda_{\max}^* + \varepsilon(\mathbf{h}) \right).$$

### Conclusion.

- (1) Si  $\lambda_{\min}^* > 0$ , alors toutes les valeurs propres sont positives, et pour  $\mathbf{h}$  petit, on a  $\frac{1}{2} \lambda_{\min}^* + \varepsilon(\mathbf{h}) > 0$  et donc  $J(\mathbf{x}^* + \mathbf{h}) - J(\mathbf{x}^*) \geq 0$ . C'est donc un minimum local.
- (2) Si  $\lambda_{\max}^* < 0$ , alors toutes les valeurs propres sont négatives, et pour  $\mathbf{h}$  petit, on a  $\frac{1}{2} \lambda_{\max}^* + \varepsilon(\mathbf{h}) < 0$  et donc  $J(\mathbf{x}^* + \mathbf{h}) - J(\mathbf{x}^*) \leq 0$ . C'est donc un maximum local.
- (3) Si les valeurs propres sont de signes différents, et non nulles, que l'on notera  $\lambda_+$  et  $\lambda_-$ . Alors, pour  $\lambda_+ > 0$ , soit  $\mathbf{u}$  un vecteur propre associé à  $\lambda_+$ . Alors, on pose  $\mathbf{h} = t\mathbf{u}$  pour  $t$  petit. Ainsi,  $J(\mathbf{x}^* + t\mathbf{u}) - J(\mathbf{x}^*) = \frac{t^2}{2} \lambda_+ \|\mathbf{u}\|^2 + t^2 \|\mathbf{u}\|^2 \varepsilon(t\mathbf{u}) > 0$  pour  $t$  petit. Ainsi  $\mathbf{x}^*$  est un minimum local.
- (4) Si  $\lambda = 0$ , alors  $J(\mathbf{x}^* + t\mathbf{u}) - J(\mathbf{x}^*) = t^2 \|\mathbf{u}\|^2 \varepsilon(t\mathbf{u})$ , qui n'a pas de signe. On ne peut donc pas conclure.

On a donc montré le théorème suivant.

### Théorème.

Soit  $J$  de classe  $\mathcal{C}^2$  sur  $\mathbb{R}^n$ , et  $\mathbf{x}^*$  un point stationnaire de  $J$ . On a alors les situations suivantes.

Signe des valeurs propres de $\nabla^2 J(\mathbf{x}^*)$	Nature du point stationnaire
$\forall i \in \llbracket 1, n \rrbracket, \lambda_i > 0$	$\mathbf{x}^*$ est un minimum local
$\forall i \in \llbracket 1, n \rrbracket, \lambda_i < 0$	$\mathbf{x}^*$ est un maximum local
$\forall i \in \llbracket 1, n \rrbracket, \lambda_i \neq 0$ de signes différents	$\mathbf{x}^*$ est un point de selle
$\exists i \in \llbracket 1, n \rrbracket, \lambda_i = 0$	$\mathbf{x}^*$ est un point dégénéré

## III.5. Convexité

Les définitions de la convexité ci-dessous sont équivalentes.

**Définition.**

Soit  $J$  continue sur  $\mathbb{R}^N$

(1) On dit que  $J$  est *convexe* si

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \forall t \in ]0, 1[, J(t\mathbf{x} + (1-t)\mathbf{y}) \leq tJ(\mathbf{x}) + (1-t)J(\mathbf{y}).$$

(1) On dit que  $J$  est *convexe* si l'inégalité ci-dessus est stricte :

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \forall t \in ]0, 1[, J(t\mathbf{x} + (1-t)\mathbf{y}) < tJ(\mathbf{x}) + (1-t)J(\mathbf{y}).$$

Autrement dit,

(1)  $J$  est convexe si elle est au dessous de toutes ses cordes.

(2)  $J$  est strictement convexe si elle est strictement au dessous de toutes ses cordes.

**Définition.**

Soit  $J$  de classe  $\mathcal{C}^1$  sur  $\mathbb{R}^N$ . On dit que  $J$  est *convexe* si

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad \langle \nabla J(\mathbf{y}) - \nabla J(\mathbf{x}) \mid \mathbf{y} - \mathbf{x} \rangle \geq 0$$

i.e.

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad J(\mathbf{y}) \geq J(\mathbf{x}) + \langle \nabla J(\mathbf{x}) \mid \mathbf{y} - \mathbf{x} \rangle$$

**Définition.**

Soit  $J$  de classe  $\mathcal{C}^1$  sur  $\mathbb{R}^N$ . On dit que  $J$  est  $\alpha$ -*convexe* s'il existe  $\alpha > 0$  tel que :

$$J(\mathbf{y}) \geq J(\mathbf{x}) + \langle \nabla J(\mathbf{x}) \mid \mathbf{y} - \mathbf{x} \rangle + \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|^2$$

ou

$$\langle \nabla J(\mathbf{y}) - \nabla J(\mathbf{x}) \mid \mathbf{y} - \mathbf{x} \rangle \geq \alpha \|\mathbf{x} - \mathbf{y}\|^2.$$

**Définition.**

Soit  $J$  de classe  $\mathcal{C}^2$  sur  $\mathbb{R}^N$ . On dit que :

- (1)  $J$  est convexe si  $\langle \nabla^2 J(\mathbf{x})\mathbf{y} \mid \mathbf{y} \rangle \geq 0$ , quels que soient  $\mathbf{x}$  et  $\mathbf{y}$ , i.e. que toutes les valeurs propres de  $\nabla^2 J(\mathbf{x})$  sont positives ou nulles ;
- (2)  $J$  est strictement convexe si  $\langle \nabla^2 J(\mathbf{x})\mathbf{y} \mid \mathbf{y} \rangle > 0$ , quels que soient  $\mathbf{x}$  et  $\mathbf{y}$ , i.e. que toutes les valeurs propres de  $\nabla^2 J(\mathbf{x})$  sont strictement positives ;
- (3)  $J$  est  $\alpha$ -convexe si  $\langle \nabla^2 J(\mathbf{x})\mathbf{y} \mid \mathbf{y} \rangle \geq \alpha\|\mathbf{y}\|^2$ , quels que soient  $\mathbf{x}$  et  $\mathbf{y}$ , i.e. que toutes les valeurs propres de  $\nabla^2 J(\mathbf{x})$  sont supérieures ou égales à  $\alpha$ .

Attention,  $\alpha$  de dépend pas de  $\mathbf{x}$ .

### Propositions.

Soit  $J$  une fonction  $\alpha$ -convexe sur  $\mathbb{R}^n$ . Alors,  $J$  est strictement convexe et coercive.

### Théorème.

- (1) Soit  $J$  une fonction convexe. Alors, si  $\mathbf{x}^*$  est un minimum de  $J$ , alors  $\mathbf{x}^*$  est un minimum global.

*« Si  $J$  est convexe, alors tout minimum local est global. »*

- (2) Soit  $J$  une fonction strictement convexe. Alors,  $J$  a au plus un minimum global.

*« Si  $J$  est strictement convexe, alors tout minimum local est unique et global. »*

- (3) Soit  $J$  une fonction strictement convexe et coercive. Alors,  $J$  a un, et un seul, minimum global.

- (4) Soit  $J$  une fonction  $\alpha$ -convexe. Alors,  $J$  a un, et un seul, minimum global.

Attention, on n'a pas l'implication  $J$  convexe  $\not\Rightarrow \mathbf{x}^*$  existe. Par exemple,  $\exp : x \mapsto e^x$  est convexe mais n'a pas de minimum global.

## IV. Optimisation non linéaire avec contraintes

La forme générale du problème est :

$$\inf_{\mathbf{v} \in K} J(\mathbf{v})$$

où  $K \subset V$  est l'ensemble des contraintes,  $J : K \rightarrow \mathbb{R}$  et  $V = \mathbb{R}^n$ .<sup>[5]</sup>

### IV.1. Existence d'un minimum

**Théorème** (Existence dans  $K \subset \mathbb{R}^n$ ).

Soient  $K \subset \mathbb{R}^n$ , où  $K$  est fermé et non vide, et  $J : K \rightarrow \mathbb{R}$  continue.

On suppose, de plus,

- (1) ou bien  $K$  est borné,
- (2) ou bien  $J$  est coercive sur  $K$  (i.e.  $J(\mathbf{x}) \rightarrow +\infty$  quand  $\|\mathbf{x}\| \rightarrow +\infty$ ).

Alors, il existe au moins un minimum sur  $K$ .

*Preuve.*

- (1) Soit  $K$  fermé et bornée, donc  $K$  est compact. Comme  $J$  est continue, alors  $J$  atteint ses bornes sur  $K$ . Il y a donc au moins un minimum et un maximum de  $J$  sur  $K$ .
- (2) On suppose  $K$  fermé et  $J$  coercive. Soit  $(\mathbf{u}_n)_{n \in \mathbb{N}}$  une suite minimisante de  $J$  sur  $K$ . Alors,  $\mathbf{u}_n \in K$  et  $J(\mathbf{u}_n) \rightarrow \inf_K J$ . La suite  $(\mathbf{u}_n)$  est bornée car sinon, il existe une sous suite  $(\mathbf{u}_{\varphi(n)})$  telle que  $J(\mathbf{u}_{\varphi(n)}) \rightarrow +\infty$ , ce qui contredit que  $J(\mathbf{u}_n) \rightarrow \inf_K J$ . D'où, la suite  $(\mathbf{u}_n)$  est bornée et dans un fermé. On peut en extraire une sous-suite qui converge vers  $\mathbf{u}$  dans  $K$  (car  $K$  fermé). Par continuité de  $J$ , on a  $J(\mathbf{u}_n) \rightarrow J(\mathbf{u}) = \inf_K J$  quand  $n$  tend vers  $+\infty$ .

□

### IV.2. Unicité du minimum

On s'intéresse à la convexité de l'ensemble  $K$ .

---

<sup>[5]</sup>On peut choisir un espace vectoriel muni d'un produit scalaire et complet.

**Définition** ( $K$  convexe).

On dit que  $K$  est convexe dans  $V = \mathbb{R}^n$  si, pour tout couple de vecteurs  $(\mathbf{x}, \mathbf{y}) \in K^2$ , alors  $t\mathbf{x} + (1-t)\mathbf{y} = \mathbf{z} \in K$  pour  $t \in [0, 1]$ .

**Théorème 2** (minimum local  $\rightarrow$  minimum global).

Soient  $K$  convexe sur  $V$ ,  $J : K \rightarrow \mathbb{R}$  continue et convexe sur  $K$ . Alors, un minimum local de  $J$  sur  $K$  est un minimum global de  $J$  sur  $K$ .

*Preuve.*

Soit  $\mathbf{x}$  un minimum local de  $J$  sur  $K$ , i.e.  $J(\mathbf{z}) \geq J(\mathbf{x})$  pour  $\mathbf{z} \in K$  et  $\mathbf{z}$  proche de  $\mathbf{x}$  (i.e.  $\exists \eta > 0, \|\mathbf{z} - \mathbf{x}\| \leq \eta$ ). Soit  $\mathbf{y} \in K$ , alors

$$\mathbf{z} = \mathbf{x} + t(\mathbf{y} - \mathbf{x}) = t\mathbf{y} + (1-t)\mathbf{x} \in K,$$

pour  $t \in [0, 1]$ , car  $K$  est convexe. Pour  $t$  petit, on considère que  $\mathbf{z}$  est un voisin de  $\mathbf{x}$ , et donc

$$J(\mathbf{x}) \leq J(\mathbf{z}) = J(t\mathbf{y} + (1-t)\mathbf{x}) \leq tJ(\mathbf{y}) + (1-t)J(\mathbf{x}),$$

car  $J$  est convexe. Ainsi,

$$tJ(\mathbf{y}) + (1-t)J(\mathbf{x}) - J(\mathbf{x}) = t(J(\mathbf{y}) - J(\mathbf{x})) \geq 0$$

et on a  $t > 0$ , car  $\mathbf{y} \neq \mathbf{x}$ . Ceci est vrai quel que soit  $\mathbf{y}$  :

$$\forall \mathbf{y} \in K, J(\mathbf{y}) \geq J(\mathbf{x}).$$

On en déduit que  $\mathbf{x}$  est un minimum global. □

**Attention.** Ce théorème ne montre ni l'unicité, ni l'existence.

**Théorème 3** (unicité).

Soient  $K$  convexe sur  $V$ ,  $J : K \rightarrow \mathbb{R}$  continue et strictement convexe sur  $K$ . Alors,  $J$  admet au plus un minimum sur  $K$ . C'est à dire, un minimum local est unique.

*Preuve.*

On suppose  $J$  strictement convexe. Supposons que  $J$  admet deux mi-

nima globaux  $x_1$  et  $x_2$ . Alors,  $J(x_1) = J(x_2) = m$ . Si  $x_1 \neq x_2$ , alors  $z = tx_1 + (1-t)x_2 \in K$ , pour  $t \in [0, 1]$ . Ainsi,

$$\begin{aligned} J(x_1) = m &\leq J(z) = J(tx_1 + (1-t)x_2) \\ &< tJ(x_1) + (1-t)J(x_2) \\ &= tm + (1-t)m \\ &= m \end{aligned}$$

Ceci démontre  $m < m$ , ce qui est *absurde*.

On en déduit  $x_1 = x_2$ . □

**Proposition** (Existence et unicité).

Soit  $K$  un convexe, fermé, non vide de  $V = \mathbb{R}^n$ . Soit  $J : K \rightarrow \mathbb{R}$  continue. Si de plus,

- (1) ou bien  $J$  est strictement convexe et coercive sur  $K$ ,
- (2) ou bien  $J$  est  $\alpha$ -convexe sur  $K$ .

Alors  $J$  admet un unique minimum global sur  $K$ .

*Preuve.*

On regroupe les théorèmes précédents. □

### IV.3. Caractérisation du minimum sur un convexe

**Exemple.** On considère  $f : x \mapsto x^2$  sur  $\mathbb{R}$ , avec  $K = [1, 2]$ . On a  $f'(x) = 2x$  et donc  $f''(x) = 2 > 0$ , d'où  $f$  est  $\alpha$ -convexe. Ainsi  $x^* = 0$  est le minimum de  $f$  sur  $\mathbb{R}$ . Sur  $K = [1, 2]$ , on a  $\min_K f$  est atteint en  $x^* = 1$ . On peut remarquer que  $f'(1) = 2 \neq 0$ . L'égalité d'Euler  $\nabla f(x^*) = 0$  n'est pas valide sur  $K$ .

**Théorème** (condition d'optimalité sur un convexe).

Soient  $K$  convexe non vide de  $V$ , et  $J : V \rightarrow \mathbb{R}$  continue (ou sur un ouvert contenant  $K$ ), et dérivable en  $x^* \in K$ .

- (1) *Condition nécessaire.* Si  $x^*$  est un minimum local de  $J$  sur  $K$  alors  $x^*$  vérifie l'inégalité d'Euler :

$$\forall \mathbf{x} \in K, \langle \nabla J(\mathbf{x}^*) \mid \mathbf{x} - \mathbf{x}^* \rangle \geq 0.$$

- (2) *Condition nécessaire et suffisante.* Si, de plus,  $K$  est convexe sur  $K$ , alors  $\mathbf{x}^*$  est un minimum si, et seulement si, l'inégalité d'Euler est vérifiée. Dans ce cas,  $\mathbf{x}^*$  est un minimum global.

*Preuve.*

- (1) Soit  $\mathbf{x}^* \in K$ . On a  $J(\mathbf{z}) \geq J(\mathbf{x}^*)$ , quel que soit  $\mathbf{z}$  voisin de  $\mathbf{x}^*$ . Soit  $\mathbf{x} \in K$ , alors  $\mathbf{z} = \mathbf{x}^* + t(\mathbf{x} - \mathbf{x}^*)$ , pour  $0 < t < 1$ . La fonction  $J$  est dérivable en  $\mathbf{x}^*$ . On applique la formule de Taylor-Young :

$$\begin{aligned} J(\mathbf{z}) &= J(\mathbf{x}^*) + \langle \nabla J(\mathbf{x}^*) \mid \mathbf{z} - \mathbf{x}^* \rangle + \|\mathbf{z} - \mathbf{x}^*\| \cdot \varepsilon(\mathbf{z} - \mathbf{x}^*) \\ &= J(\mathbf{x}^*) + t \cdot \langle \nabla J(\mathbf{x}^*) \mid \mathbf{x} - \mathbf{x}^* \rangle + \|\mathbf{z} - \mathbf{x}^*\| \cdot \varepsilon(t(\mathbf{x} - \mathbf{x}^*)) \end{aligned}$$

Alors,

$$\begin{aligned} \langle \nabla J(\mathbf{x}^*) \mid \mathbf{x} - \mathbf{x}^* \rangle &= \frac{1}{t} \overbrace{(J(\mathbf{z}) - J(\mathbf{x}^*))}^{\geq 0} - \|\mathbf{x} - \mathbf{x}^*\| \cdot \varepsilon(t(\mathbf{x} - \mathbf{x}^*)) \\ &\geq -\|\mathbf{x} - \mathbf{x}^*\| \cdot \varepsilon(t(\mathbf{x} - \mathbf{x}^*)) \end{aligned}$$

On en déduit que

$$\forall \mathbf{x} \in K, \langle \nabla J(\mathbf{x}^*) \mid \mathbf{x} - \mathbf{x}^* \rangle \geq 0.$$

- (2) Supposons  $J$  convexe. Ainsi,  $J(\mathbf{x}) \geq J(\mathbf{x}^*) + \overbrace{\langle \nabla J(\mathbf{x}^*) \mid \mathbf{x} - \mathbf{x}^* \rangle}^{\geq 0}$ . Alors, si  $\langle \nabla J(\mathbf{x}^*) \mid \mathbf{x} - \mathbf{x}^* \rangle \geq 0$ , alors  $J(\mathbf{x}) \geq J(\mathbf{x}^*)$  quel que soit  $\mathbf{x} \in K$ . D'où,  $\mathbf{x}^*$  est un minimum global de  $J$  sur  $K$ .

□

### Remarques.

- (1) *Minimum de  $f(x) = x^2$  sur  $[1, 2]$ .* L'inégalité d'Euler s'écrit  $f'(x^*) \cdot (x - x^*) \geq 0$ , quel que soit  $x \in [0, 1]$  et pour  $x^* \in [1, 2]$ . Cette inégalité est vérifiée si, et seulement si  $2x^*(x - x^*) \geq 0$ , d'où  $x^* = 1$ .
- (2) Si  $K = V = \mathbb{R}^n$ , alors  $\nabla J(\mathbf{x}^*) = \mathbf{0}$ .
- (3) Si  $\mathbf{x}^* \in \overset{\circ}{K}$ , l'intérieur de  $K$ , alors  $\langle \nabla J(\mathbf{x}^*) \mid \mathbf{x} - \mathbf{x}^* \rangle = 0$ .

(4) Si  $\mathbf{x}^* \in \partial K$ , le bord de  $K$ , alors

$$\langle \nabla J(\mathbf{x}^*) \mid \mathbf{x} - \mathbf{x}^* \rangle = \|\nabla \dots\|.$$

## IV.4. Multiplicateurs de Lagrange

### IV.4.a. Cas des contraintes d'égalité

On pose  $V = \mathbb{R}^n$  et  $K = \{\mathbf{x} \in \mathbb{R}^n \mid \forall i \in \{1, \dots, p\}, h_i(\mathbf{x}) = 0\}$ , il y a donc  $p$  contraintes. Pour chaque contrainte, on associe un multiplicateur de Lagrange  $\lambda_i$ , on définit le lagrangien par

$$L(\mathbf{x}, \boldsymbol{\lambda}) = J(\mathbf{x}) + \sum_{i=1}^n \lambda_i h_i(\mathbf{x}).$$

**Théorème KKT** (Contrainte d'égalité).

On suppose que  $J$  et les  $h_i$  sont  $\mathcal{C}^1$  sur  $\mathbb{R}^n$ . On supposera également que, pour  $\mathbf{x}^* \in K$ , les  $(\nabla h_i(\mathbf{x}^*))$  sont linéairement indépendants.

(1) Si  $\mathbf{x}^*$  est un minimum global de  $J$  sur  $K$  alors il existe un vecteur de multiplicateurs  $\boldsymbol{\lambda}^* = (\lambda_i^*) \in \mathbb{R}^n$  tel que

$$(\mathbf{KKT}) : \begin{cases} \nabla J(\mathbf{x}^*) + \sum_{i=1}^n \lambda_i^* \cdot \nabla h_i(\mathbf{x}^*) = 0 \\ \forall i \in \{1, \dots, p\}, \quad h_i(\mathbf{x}^*) = 0 \end{cases}$$

Il faut résoudre un système de  $n + p$  équations.

(2) Si  $J$  est convexe sur  $K$  et  $K$  convexe. Les conditions **(KKT)** sont suffisantes.

*Preuve.*

On fait la démonstration dans le cas linéaire :

$$h_i(\mathbf{x}) = \sum_{j=1}^n a_{i,j} x_j - b_i = (\mathbf{A}\mathbf{x} - \mathbf{b})_i \iff h(\mathbf{x}) = \mathbf{A}\mathbf{x} - \mathbf{b}.$$

où  $\mathbf{A} = (a_{i,j}) \in \mathcal{M}_{p,n}$  et  $i \in \{1, \dots, p\}$ , avec  $p \leq n$ . Ainsi,



$$\nabla h_i(\mathbf{x}) = \frac{\partial h_i}{\partial \mathbf{x}_j} = \begin{pmatrix} a_{i,1} \\ \vdots \\ a_{i,n} \end{pmatrix}$$

On a  $K = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{b}\}$ .

...

On a  $\mathbf{x}^*$  qui vérifie  $\langle \nabla J(\mathbf{x}^*) \mid \mathbf{x} - \mathbf{x}^* \rangle \geq 0$ , que l que soit  $\mathbf{x} \in K$ . On a  $\mathbf{A}\mathbf{x}^* = \mathbf{b}$  et  $\mathbf{A}\mathbf{x} = \mathbf{b}$  d'où  $\mathbf{A}(\mathbf{x} - \mathbf{x}^*) = \mathbf{0}$ . On a  $\mathbf{x} - \mathbf{x}^* \in K_0$  quel que soit  $\mathbf{x} \in K$ . Alors,  $\langle \nabla J(\mathbf{x}^*) \mid \mathbf{y} \rangle \geq 0$ , quel que soit  $\mathbf{y} \in K_0$ . Or,  $-\mathbf{y} \in K_0$  et donc  $\langle \nabla J(\mathbf{x}^*) \mid \mathbf{y} \rangle \leq 0$ . On en déduit  $\langle \nabla J(\mathbf{x}^*) \mid \mathbf{y} \rangle = 0$  quel que soit  $\mathbf{y} \in K_0 = \ker \mathbf{A}$ . Ainsi,  $\nabla J(\mathbf{x}^*) \perp \ker \mathbf{A}$ . D'où,  $\nabla J(\mathbf{x}^*) \in (\ker \mathbf{A})^\perp$ . Or,  $(\ker \mathbf{A})^\perp = \text{im } \mathbf{A}^\top$ . Et donc  $\nabla J(\mathbf{x}^*) \in \text{im } \mathbf{A}^\top$  d'où,  $\mathbf{A}^\top \in \mathcal{M}_{p,n}$ . Il existe  $\tilde{\boldsymbol{\lambda}} \in \mathbb{R}^p$  tel que  $\nabla J(\mathbf{x}^*) = \mathbf{A}^\top \tilde{\boldsymbol{\lambda}}$ . Soit encore  $\nabla J(\mathbf{x}^*) + \mathbf{A}^\top \boldsymbol{\lambda}^* = \mathbf{0}$  et  $\boldsymbol{\lambda}^* = \tilde{\boldsymbol{\lambda}}$ . On a

$$(\mathbf{A}^\top \boldsymbol{\lambda}^*)_j = \sum_{i=1}^p a_{j,i}^\top \lambda_i^* = \sum_{i=1}^p a_{i,j}^\top \lambda_i^* = \sum_{i=1}^n \frac{\partial h_i}{\partial x_j} \lambda_i^*$$

Enfin  $\mathbf{A}^\top \boldsymbol{\lambda}^* = \sum_{i=1}^p \lambda_i^* \nabla h_i$ . On a un système non linéaire à résoudre : trouver  $\boldsymbol{\lambda}^* \in \mathbb{R}^p$  et  $\mathbf{x}^* \in K$  solution du système (**KKT**).  $\square$

**Exemple.** On calcule  $\min_{x+y=1} x^2 + y^2 - 4$ . On pose  $h(x, y) = x + y + 1 = 0$ ; On a un seul multiplicateur de Lagrange : on cherche  $\boldsymbol{\lambda}^*$  et  $\mathbf{x}^* = (x^*, y^*)$  tel que

$$\begin{cases} \nabla J(\mathbf{x}^*) + \boldsymbol{\lambda}^* \nabla h(\mathbf{x}^*) = 0 \\ h(\mathbf{x}^*) = 0 \end{cases} \implies x = y = -\frac{1}{2}$$

et  $\lambda = 1$ .

#### IV.4.b. Multiplicateurs de Lagrange – cas général

On considère des contraintes d'égalité et d'inégalité :

$$K = \{\mathbf{x} \in \mathbb{R}^n \mid \forall i \in \{1, \dots, p\}, h_i(\mathbf{x}) = 0 \text{ et } \forall j \in \{1, \dots, q\}, g_j(\mathbf{x}) \leq 0\}.$$

On suppose

- $J$ , les  $h_i$  et les  $g_j$  de classe  $\mathcal{C}^1$  ;

- pour  $\mathbf{x}^* \in K$ ,  $(\nabla h_i(\mathbf{x}^*))$  sont linéairement indépendants ;
- *constantes qualifiées*, il existe  $\mathbf{d} \neq \mathbf{0}$  un vecteur de  $\mathbb{R}^n$  tel que l'on ait  $\langle \nabla h_i(\mathbf{x}^*) \mid \mathbf{d} \rangle = 0$  et  $\langle \nabla g_j(\mathbf{x}^*) \mid \mathbf{d} \rangle < 0$ .

Le multiplicateur de Lagrange est donc

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = J(\mathbf{x}) + \sum_{i=1}^p \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^q \mu_j g_j(\mathbf{x}).$$

Le nom **KKT** vient des noms Karuch, Kuhn et Tucker.

### **Théorème KKT.**

- (1) Si  $\mathbf{x}^*$  est un minimum local de  $J$  sur  $K$ , alors il existe  $\boldsymbol{\lambda}^*$  et  $\boldsymbol{\mu}^*$  qui vérifie :
  - (a)  $L(\mathbf{x}^*, \boldsymbol{\lambda}, \boldsymbol{\mu}) = 0$ ,
  - (b)  $h_i(\mathbf{x}^*) = 0$ ,
  - (c)  $\mu_j^* g_j(\mathbf{x}^*) = 0$ ,
  - (d)  $\mu_j^* \geq 0$ ,
  - (e)  $g_j(\mathbf{x}^*) \geq 0$ .

On a  $n + p + q$  équations à résoudre.

- (2) Si, de plus,  $J$  est convexe, et les  $g_j$  sont convexes. Alors les conditions sur des conditions nécessaires et suffisantes, et  $\mathbf{x}^*$  est un minimum global.

**Remarque.**  $\mu_j^* g_j(\mathbf{x}^*) = 0$ , alors  $g_j(\mathbf{x}^*) = 0$  alors la contrainte est active ; et si  $g_j(\mathbf{x}^*) < 0$  alors elle est inactive car  $\mu_j^* = 0$ .

## V. Probabilités

Ce chapitre sera vu en deux CM :

- (1) Mesure de probabilités, lois usuelles, variables aléatoires, densité de probabilité
- (2) Vecteur aléatoire, convergence d'une suite aléatoire (convergence presque sûre, convergence en loi)

**Exemple** ( $\Omega$  fini : dé).

Soit  $\Omega = \{1, 2, 3, 4, 5, 6\}$ . On considère les événements  $A_1 = \{1, 3, 5\}$  « lancé de dé impair »,  $A_2 = \Omega$  « lancé de dé »,  $A_3 = \{1, 3\}$  « lancé de dé impair, inférieur à 5 ».

On dit que  $M = \wp(\Omega)$  est l'ensemble des événements de  $\Omega$ . On peut définir

$$P : M \rightarrow [0, 1]$$

$$A \mapsto P(A) = \sum_{a \in A} P(a) = \frac{\text{card } A}{6}$$

On note  $\text{card } A = \#A$  le nombre d'éléments de  $A$ .

**Définition** (Mesure de probabilité uniforme discrète sur  $\Omega$  fini).

On pose  $\Omega = \{\omega_0, \dots, \omega_n\}$  et  $M = \wp(\Omega)$ . Pour tout  $i$ , on pose  $P(\omega_i) = 1/n$ , équiprobabilité. Ainsi,

$$\forall A \in M, \quad P(A) = \frac{\text{card } A}{n}.$$

**Exemple** (Tirage d'un nombre dans  $\Omega = [0, 1]$ , infini).

On pose  $M = \{[0, 1], \emptyset, [0, 1/4], (1/4, 1]\}$ . On définit  $P : M \rightarrow [0, 1]$  par  $P(\Omega) = 1$ ,  $P([0, 1/4]) = \alpha$  et  $P((1/4, 1]) = 1 - \alpha$ , pour  $\alpha \in [0, 1]$ .

**Définition.**

On dit que :

- $\Omega$  est l'événement certain,
- $\emptyset$  est l'événement impossible,
- $A \neq \Omega$  et  $P(A) = 1$  est un événement presque sûr,
- $A \neq \emptyset$  et  $P(A) = 0$  est un événement négligeable.

**Définition** (Espace mesurable).

Soient  $\Omega$  un ensemble, et  $M$  un ensemble de parties de  $\Omega$  (i.e.  $M \subseteq \wp(\Omega)$ ). On dit que  $M$  est une *tribu* si

- $\Omega \in M$ ,
- si  $A \in M$ , alors  $A^c \in M$  (le complémentaire de  $A$  dans  $\Omega$ ),
- si  $\{A_i\}_{i \in \mathbb{N}}$  est une famille d'éléments de  $M$ , alors  $\bigcup_{i \in \mathbb{N}} A_i \in M$ .

On dira que  $(\Omega, M)$  est un *espace mesurable* ou *espace prbabilisable*.

De plus, pour  $A \in M$ , on dit que  $A$  est un *événement mesurable*.

- (1) L'ensemble des parties  $M = \wp(\Omega)$  est une tribu.
- (2) L'ensemble  $\{\Omega, \emptyset, [0, 1/4], (1/4, 1]\}$  est une tribu engendrée pour  $[0, 1/4]$ .
- (3) Pour  $\Omega = \{0, 1, 2\}$ , les ensembles  $C_1 = \{\emptyset, \Omega, \{0\}, \{1, 2\}\}$ , et  $C_2 = \{\emptyset, \Omega, \{1\}, \{0, 2\}\}$  sont tous deux des tribus, mais  $C_1 \cup C_2$  n'en n'est pas une.  
En effet,  $A = \{0, 1\} \in C_1 \cup C_2$ , mais  $A^c = \{2\} \notin C_1 \cup C_2$ .

**Définition** (Tribu borélienne, Borel 1871–1956).

- (1) On appelle la *tribu borélienne* de  $\mathbb{R}$ , notée  $\mathcal{B}(\mathbb{R})$ , la tribu engendré par les ouverts de  $\mathbb{R}$  (ou les fermés). La tribu  $\mathcal{B}(\mathbb{R})$  est engendrée par tous les ouverts de la forme  $(-\infty, x]$ , pour  $x \in \mathbb{R}$ .
- (2) La tribu  $\mathcal{B}(\mathbb{R}^n)$  est engendrée par tous les ouverts de la forme  $\prod_{i=1}^n (-\infty, x_i)$ , pour  $(x_1, \dots, x_n) \in \mathbb{R}^n$ .

**Définition** (Espace probabilisé, Kolmogorov 1930).

Soit  $(\Omega, M)$  un espace mesurable. On appelle *probabilité* toute application  $P : M \rightarrow [0, 1]$  telle que :

- $P(\Omega) = 1$ ,
- si  $\{A_n\}_{n \in \mathbb{N}}$  est une famille finie ou dénombrable d'événements de  $M$  disjoints deux à deux, alors  $P\left(\bigcup_{n \in \mathbb{N}} A_n\right) = \sum_{n \in \mathbb{N}} P(A_n)$ .

Le triplet  $(\Omega, M, P)$  s'appelle *espace probabilisé*.

- (1) L'équiprobabilité est  $P(\{\omega\}) = 1/(\text{card } \Omega)$ , pour  $\omega \in \Omega$ .
- (2) La loi de poisson  $\mathcal{P}(\lambda)$  est définie sur  $\mathbb{N}$  par  $P(n) = e^{-\lambda} \lambda^n / n!$ .
- (3) Pour  $\Omega = [0, 1]$ , et  $M = \mathcal{B}(\Omega)$ , on définit  $P$  sur  $M$  par  $P(A) = |A|$ , pour tout intervalle  $A$  de  $\Omega$ .

### Proposition/Définition.

- (1)  $P(\emptyset) = 0, P(\Omega) = 1, P(A^c) = 1 - P(A)$
- (2)  $P(A \cup B) = P(A) + P(B) - P(A \cap B), P(A | B) = P(A \cap B) / P(B)$
- (3) Deux événements  $A$  et  $B$  sont indépendants si  $P(A \cap B) = P(A) \cdot P(B)$ .
- (4) Des événements  $(A_i)_{i \in I}$  sont mutuellement indépendants si

$$\forall p \leq |I|, \forall (i_1, \dots, i_p) \in I^p, \quad P(A_{i_1} \cap \dots \cap A_{i_p}) = P(A_{i_1}) \cdots P(A_{i_p})$$

## V.1. Variables aléatoires réelles

On définit une fonction  $X : \Omega \rightarrow \mathbb{R}, \omega \mapsto X(\omega)$ , où  $\omega$  est une éventualité sur  $\Omega$ . C'est une fonction du hazard.

### Définition.

Soit  $(\Omega, M, P)$  un espace probabilisé.

On dit que  $X$  est une variable aléatoire réelle si  $\forall x \in \mathbb{R}, P(X \in (-\infty, x])$  est défini.

Ceci est équivalent à  $\forall B \in \mathcal{B}(\mathbb{R}), P(X^{-1}(B))$  est bien défini, où

$$X^{-1}(B) = \{\omega \in \Omega \mid X(\omega) \in B\}.$$

On note  $P(X \in (-\infty, x]) = P(X \leq x) = P(\omega \in \Omega \mid X(\omega) \leq x)$ .

**Définition** (Fonction de répartition).

La fonction de répartition de  $X$  est définie par  $F_X : \mathbb{R} \rightarrow \mathbb{R}$  telle que  $F_X(x) = P(X \leq x)$ , quel que soit  $x \in \mathbb{R}$ .

**Exemple 1** (Cas discret).

On pose  $\Omega = \{1, 2, 3, 4, 5, 6\}$ , et  $M = \wp(\Omega)$ . La probabilité  $P$  est discrète uniforme :  $P(i) = 1/6$ . On pose  $X : \Omega \rightarrow \mathbb{R}$ ,  $X(i) = i$ , la valeur du dé. Ainsi, la fonction de répartition de  $X$  est donnée par :

$$F_X(x) = P(X \leq x) = \begin{cases} 0 & \text{si } x < 1 \\ \frac{1}{6} & \text{si } 1 \leq x < 2 \\ \frac{2}{6} & \text{si } 2 \leq x < 3 \\ \frac{3}{6} & \text{si } 3 \leq x < 4 \\ \frac{4}{6} & \text{si } 4 \leq x < 5 \\ \frac{5}{6} & \text{si } 5 \leq x < 6 \\ 1 & \text{si } x \geq 6 \end{cases}$$

**Exemple 2** (Loi uniforme sur  $[0, 1]$ ).

On pose  $\Omega = [0, 1]$ , et  $M = \mathcal{B}(\Omega)$ . On pose  $P$  la probabilité continue uniforme, et  $P(A) = |A|$ . Soit  $X : \Omega \rightarrow \mathbb{R}$  une variable aléatoire. Alors,

$$\forall x \in \mathbb{R}, F_X(x) = P(X \leq x) = \begin{cases} 0 & \text{si } x < 0 \\ x & \text{si } 0 \leq x \leq 1 \\ 1 & \text{si } x > 1 \end{cases}$$

Dans les deux exemples, on remarque que la fonction de répartition est une fonction *croissante*.

**Proposition.**

$$P(X \in (a, b]) = F_X(b) - F_X(a).$$

**Définition.** On dit que deux variables  $X$  et  $Y$  suivent la même loi si leurs fonction de répartition sont égales  $F_X = F_Y$ , i.e.

$$\forall x \in \mathbb{R}, P(X \leq x) = P(Y \leq x).$$

**Définition** (Densité de probabilité).

Soit  $f : \mathbb{R} \rightarrow \mathbb{R}^+$  une fonction positive et continue (par morceaux) telle que  $\int_{-\infty}^{+\infty} f(x) dx = 1$ . Alors  $f$  est une densité de probabilité.

On dit que  $X$  est une variable aléatoire à densité s'il existe  $f_X$  une fonction de probabilité et  $F_X(x) = \int_{-\infty}^x f_X(t) dt = P(X \leq x)$ .

**Exemples.**

(1) Loi uniforme  $X \sim \mathcal{U}(a, b)$  de densité  $f_X(x) = \mathbb{1}_{[a,b]}(x)/(b-a)$ .

(2) Loi normale  $X \sim \mathcal{N}(m, \sigma^2)$  de densité

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{(x-m)^2}{2}\sigma^2}.$$

(3) Loi exponentielle  $X \sim \mathcal{E}(\lambda)$  de densité  $f_X(x) = \lambda e^{-\lambda x} \mathbb{1}_{\{X \geq 0\}}(x)$

**Définition.** Soit  $X$  une variable aléatoire à valeurs dans l'ensemble  $\{x_1, x_2, \dots, x_n, \dots\}$ . On a  $X = \sum_{i=1}^{\infty} x_i \mathbb{1}_{\{X=x_i\}}$ .