

# **OPTIMISATION APPROFONDIE**

LICENCE MIASHS Semestre 6

2021–2022

# Table des matières

<b>1</b>	<b>Quelques rappels et notations</b>	<b>1</b>
1.1	Matrices symétriques réelles . . . . .	1
1.2	Fonctions différentiables de $\mathbb{R}^n$ en $\mathbb{R}^m$ . . . . .	1
1.3	Convexité et concavité . . . . .	2
1.4	Exercices . . . . .	3
<b>2</b>	<b>Optimisation sans contraintes</b>	<b>6</b>
2.1	Conditions nécessaires d'optimalité . . . . .	6
2.2	Conditions suffisantes d'optimalité . . . . .	6
2.3	Résolution d'un problème d'optimisation sans contraintes . . . . .	6
2.4	Exercices . . . . .	8
<b>3</b>	<b>Optimisation avec contraintes</b>	<b>10</b>
3.1	Qualification des contraintes . . . . .	10
3.2	Conditions nécessaires d'optimalité : théorème KKT . . . . .	11
3.3	Exercices . . . . .	13
<b>4</b>	<b>Quelques algorithmes de résolution</b>	<b>18</b>
4.1	Choix de la direction . . . . .	18
4.1.1	Méthodes de gradient . . . . .	19
4.1.2	Méthodes Newtonniennes . . . . .	24
4.2	Choix du pas . . . . .	26
4.2.1	Pas optimal . . . . .	26
4.2.2	Recherche linéaire . . . . .	27
4.3	Choix de la direction en présence de contraintes . . . . .	29
4.3.1	Méthode de pénalisation . . . . .	29
4.3.2	Construction d'une direction admissible . . . . .	29
4.4	Exercices . . . . .	30

# 1 Quelques rappels et notations

## 1.1 Matrices symétriques réelles

$\mathbb{R}^n$  est l'espace des vecteurs colonnes de longueur  $n$  et  $\mathcal{M}_n(\mathbb{R}) = \mathbb{R}^{n \times n}$  est l'espace des matrices carrées réelles de taille  $n$ . Soit  $A = [a_{ij}]_{1 \leq i \leq n, 1 \leq j \leq n} \in \mathcal{M}_n(\mathbb{R})$ . Alors  $A$  est

- une matrice symétrique si

$$\forall i, j : a_{ij} = a_{ji},$$

ou équivalent, si  $A = A^T$ .

- une matrice semi-définie négative (resp. semi-définie positive) si

$$\forall y \in \mathbb{R}^n : y^T A y \leq 0 \quad (\text{resp. } y^T A y \geq 0)$$

- une matrice définie négative (resp. définie positive) si

$$\forall y \in \mathbb{R}^n, y \neq 0 : y^T A y < 0 \quad (\text{resp. } y^T A y > 0)$$

$\mathcal{S}_n(\mathbb{R})$  est l'espace des matrices carrées réelles symétrique de taille  $n$ , donc

$$A \in \mathcal{S}_n(\mathbb{R}) \text{ ssi } (A \in \mathcal{M}_n(\mathbb{R}) \text{ et } A = A^T).$$

## 1.2 Fonctions différentiables de $\mathbb{R}^n$ en $\mathbb{R}^m$

La dérivée directionnelle d'une fonction  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  au point  $a \in \mathbb{R}^n$  dans la direction  $h \in \mathbb{R}^n$  est définie par :

$$D|_h f(a) := \lim_{t \rightarrow 0} \frac{f(a + th) - f(a)}{t}.$$

On dit qu'une fonction  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  est différentiable en  $a$  si sa dérivée directionnelle au point  $a$  existe pour tout  $h \in \mathbb{R}^n$  et si

$$D|_h f(a) = f'(a) \cdot h,$$

Ici,  $f'(a)$  est une application linéaire de  $\mathbb{R}^n \rightarrow \mathbb{R}^m$  qui peut être représentée par sa matrice dans les bases canoniques (*la matrice jacobienne de  $f$* )

$$[f'(a)]_{ij} = \frac{\partial f_i}{\partial x_j}(a) = [J_f(a)]_{ij} = [\nabla f_i(a)]_j, \quad 1 \leq i \leq m, 1 \leq j \leq n,$$

avec  $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$  et  $\nabla f_i(x)$  le gradient de  $f_i$  définie par :  $\nabla f_i(x) := \left( \frac{\partial f_i}{\partial x_1}(x), \dots, \frac{\partial f_i}{\partial x_n}(x) \right)$ .

Notez que

- le gradient est toujours un vecteur ligne;
- dans le cas particulier où  $m = 1$  on a  $f'(a) = \nabla f(a)$ .

On dit qu'une fonction  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  est deux fois différentiable en  $a$  si elle est différentiable au point  $a$  et s'il existe une matrice  $\nabla^2 f(a) \in \mathcal{S}_n(\mathbb{R})$ , appelée *la matrice hessienne de  $f$*  et définie par

$$\nabla^2 f(x) := \left[ \frac{\partial^2 f}{\partial x_i \partial x_j} \right]_{1 \leq i \leq n, 1 \leq j \leq n} = \nabla ( [\nabla f(x)]^T ),$$

telle que

$$f(a + h) = f(a) + \nabla f(a) \cdot h + h^T \nabla^2 f(a) h + \|h\|^2 \varepsilon(h),$$

où  $\varepsilon(h) \rightarrow 0$  pour  $h \rightarrow 0$ .

**Définition 1** Soit  $f : \Omega \rightarrow \mathbb{R}$  avec  $\Omega$  un ouvert de  $\mathbb{R}^n$ .

- $f$  est de classe  $\mathcal{C}^0$  sur  $\Omega$  si  $f$  est continue sur  $\Omega$ .
- $f$  est de classe  $\mathcal{C}^1$  sur  $\Omega$  si  $f$  est de classe  $\mathcal{C}^0$  et pour tout  $1 \leq i \leq n$  les fonctions dérivées partielles premières  $\frac{\partial f}{\partial x_i}$  sont définies et continues sur  $\Omega$ .
- $f$  est de classe  $\mathcal{C}^2$  sur  $\Omega$  si  $f$  est de classe  $\mathcal{C}^1$  et pour tout  $1 \leq i \leq n$  et pour tout  $1 \leq j \leq n$  les fonctions dérivées partielles secondes  $\frac{\partial^2 f}{\partial x_i \partial x_j}$  sont définies et continues sur  $\Omega$ .

### 1.3 Convexité et concavité

**Définition 2**

- Un ensemble  $D \subset \mathbb{R}^n$  est convexe ssi

$$\forall x, y \in D, \forall t \in [0, 1] : [(1 - t)x + ty] \in D.$$

- Une fonction  $f$  est concave ssi

$$\forall x, y \in D, \forall t \in [0, 1] : f((1 - t)x + ty) \geq (1 - t)f(x) + tf(y).$$

- Une fonction  $f$  est strictement concave ssi

$$\forall x, y \in D, x \neq y, \forall t \in ]0, 1[ : f((1 - t)x + ty) > (1 - t)f(x) + tf(y).$$

On obtiendra la définition pour une fonction (strictement) convexe en échangeant le sens de l'inégalité. Notez que  $f$  est (strictement) convexe ssi  $-f$  est (strictement) concave.

**Théorème 1** Soit  $f : D \rightarrow \mathbb{R}$  une fonction de classe  $\mathcal{C}^2$  où  $D$  est un ouvert convexe de  $\mathbb{R}^n$ . Alors :

- $f$  est concave ssi  $\forall x \in D, \nabla^2 f(x)$  est semi-définie négative.
- Si  $\forall x \in D, \nabla^2 f(x)$  est définie négative, alors  $f$  est strictement concave. La réciproque n'est pas forcément vraie (considérer p.ex. la fonction  $f(x) = -x^4$  sur  $\mathbb{R}$ ).

## 1.4 Exercices

### Exercice 1 Matrices symétriques

1. Soit  $A \in \mathcal{S}_n(\mathbb{R})$ . Montrer que  $A$  est semi-définie positive (resp. définie positive) ssi les valeurs propres de  $A$  sont positives ou nulles (resp. strictement positives).
2. Montrer que toute matrice symétrique réelle définie positive est inversible.
3. Donner un exemple de matrice symétrique réelle qui n'est ni semi-définie positive ni semi-définie négative.
4. Quelles sont les matrices symétriques réelles qui sont à la fois semi-définies négatives et semi-définies positives?

Exercice 2 Soit  $A = [a_{ij}] \in \mathcal{S}_2(\mathbb{R})$  une matrice symétrique réelle et  $F : \mathbb{R}^2 \rightarrow \mathbb{R}$  la fonction définie par  $F(x) = \frac{1}{2}x^T Ax$ . Calculer pour tout  $x \in \mathbb{R}^2$ ,  $\nabla F(x)$  (le gradient) et  $\nabla^2 F(x)$  (la matrice hessienne). Généraliser au cas  $A = [a_{ij}] \in \mathcal{S}_n(\mathbb{R})$ . Comment procéder dans le cas où  $A$  n'est pas symétrique ? *Indication: Remplacer  $x^T Ax$  par  $\frac{1}{2}x^T (A^T + A)x$ .*

### Exercice 3 Dérivées de fonctions dans $\mathbb{R}^n$

1. Soit  $f(x) = x^T Hx + hx + \gamma$ , avec  $H \in \mathcal{S}_n(\mathbb{R})$ ,  $h \in \mathbb{R}^{1 \times n}$  et  $\gamma \in \mathbb{R}$ . Montrer que

$$\nabla f(x) = 2x^T H + h, \quad \nabla^2 f(x) = 2H.$$

2. Soit  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $f(x) = Ax + b$ . Montrer que  $\nabla f(x) = A$  et  $\nabla^2 f(x) = 0$ .
3. Soit  $f(x) = h(Ax + b)$ . Montrer que  $\nabla f(x) = \nabla h(Ax + b)A$ .
4. Soit  $g : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $g(x) = \sum_{i=1}^m [r_i(x)]^2$ , où  $r_i : \mathbb{R}^n \rightarrow \mathbb{R}$  sont des fonctions deux fois dérivables. Calculer  $\nabla g(x)$  et  $\nabla^2 g(x)$ .

### Exercice 4 Ensembles convexes

1. Montrer que pour tout  $x \in \mathbb{R}^n$  et pour tout  $r > 0$ ,  $B(x, r)$  (la boule centrée en  $x$  de rayon  $r$ ) est convexe.
2. Montrer que  $D = \{x \in \mathbb{R}^2 : x_1 + x_2 \leq 1, x_1 \geq 0, x_2 \geq 0\}$  est un convexe de  $\mathbb{R}^2$ .
3. Montrer que  $D = \{x \in \mathbb{R}^n : \sum_{i=1}^n x_i = 1, x_i \geq 0, i = 1, \dots, n\}$  est un convexe de  $\mathbb{R}^n$ .
4. Quelles sont les parties convexes de  $\mathbb{R}$ ?

### Exercice 5 Fonctions convexes

1. Soient  $g_j : \mathbb{R}^n \rightarrow \mathbb{R}$  ( $1 \leq j \leq p$ ) des fonctions convexes. Montrer que

$$D = \{x \in \mathbb{R}^n : g_j(x) \leq 0, 1 \leq j \leq p\}$$

est un convexe de  $\mathbb{R}^n$ .

2. Soient  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  une fonction convexe et  $\rho > 0$ . Montrer que

$$D = \left\{ x \in \mathbb{R}^n : \sum_{i=1}^n x_i = 1, \quad x_i \geq 0, 1 \leq i \leq n, \quad g(x) \leq \rho \right\}$$

est un convexe de  $\mathbb{R}^n$ .

3. Soient  $g_j : \mathbb{R}^n \rightarrow \mathbb{R}$  ( $1 \leq j \leq p$ ) des fonctions convexes et  $\mu_j$ ,  $1 \leq j \leq p$  des nombres réels positifs ou nuls. Montrer que  $g = \mu_1 g_1 + \mu_2 g_2 + \cdots + \mu_p g_p$  est une fonction convexe.

### Exercice 6 Soient $B = \mathbb{R}^p, B_1 = \mathbb{R}^q, B_2 = \mathbb{R}^n$ . Montrer les propriétés suivantes

1. Si  $C_1, C_2 \subset B$  sont des convexes alors l'ensemble  $C_1 + C_2 = \{x + y : x \in C_1, y \in C_2\}$  est convexe.
2. Soient  $C_1 \subset B_1, C_2 \subset B_2$  des ensembles convexes. Alors le produit cartésien

$$C_1 \times C_2 = \left\{ \begin{pmatrix} x \\ y \end{pmatrix} : x \in C_1, y \in C_2 \right\}$$

est convexe.

3. Soit  $C \subset B_1 \times B_2$  un convexe. Alors  $D = \left\{ x \in B_1 : \exists y \in B_2 \quad \begin{pmatrix} x \\ y \end{pmatrix} \in C \right\}$  est convexe.
4. Si  $C \subset B$  est un convexe alors  $\text{Int}(C)$  est un convexe et  $\text{Clos}(C)$  est un convexe.
5. Un ellipsoïde  $\mathcal{E}$  est défini par

$$\mathcal{E} = \{Ax + x_0 : \|x\|_2 \leq 1\} = \{x : (x - x_0)^T E^{-1} (x - x_0) \leq 1\},$$

avec  $A$  une matrice symétrique définie positive et  $E = A^2$ . Montrer que  $\mathcal{E}$  est un convexe.

6. Soit  $C \subset \mathbb{R}^n$  un convexe,  $x_1, x_2, \dots, x_k \in C$ ,  $\theta_1, \dots, \theta_k \geq 0$  tq  $\sum_{i=1}^k \theta_i = 1$ . Montrer que  $\theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_k x_k \in C$ .
7. Montrer que  $C$  est un convexe ssi l'intersection avec toute droite est un convexe.
8. Supposons que  $C$  vérifie la propriété de convexité du point milieu, i.e.,  $\forall a, b \in C, \quad \frac{1}{2}(a+b) \in C$ . Montrer que si  $C$  est de plus fermé alors  $C$  est convexe.

Exercice 7 Soit  $C \subset \mathbb{R}^n$  un ensemble convexe.

1. Soit  $h : C \rightarrow \mathbb{R}$  une fonction convexe et  $g : \mathbb{R} \rightarrow \mathbb{R}$  une fonction convexe et croissante. Montrer que  $f = g \circ h$  est convexe.
2. Soit  $f : C \rightarrow \mathbb{R}$ . On définit l'épigraphe de  $f$  par

$$\text{epi}(f) = \left\{ \begin{pmatrix} x \\ r \end{pmatrix} \in \mathbb{R}^n \times \mathbb{R} : x \in C, f(x) \leq r \right\}.$$

Montrer que  $f$  convexe  $\Leftrightarrow \text{epi}(f)$  est convexe.

3. Soit  $f$  une fonction convexe. Montrer que  $\forall r \in \mathbb{R}$ , l'ensemble de niveau

$$\{x \in C : f(x) \leq r\} \text{ est un convexe.}$$

Exercice 8 Soit  $C \subset \mathbb{R}^n$  l'ensemble solution de l'inégalité quadratique

$$C = \{x \in \mathbb{R}^n : x^T A x + b^T x + c \leq 0\}$$

avec  $A \in \mathcal{S}_n(\mathbb{R})$ ,  $b \in \mathbb{R}^n$  et  $c \in \mathbb{R}$ . Montrer que  $C$  est convexe si  $A$  est semi-définie positive.

### Exercice supplémentaire

Exercice 9 On considère la moyenne géométrique

$$f(x) = \left( \prod_{i=1}^n x_i \right)^{1/n} \quad \text{dans } \text{dom}(f) = \mathbb{R}_{++}^n = \{x \in \mathbb{R}^n : x_i > 0, i = 1, \dots, n\}.$$

1. Montrer que la matrice hessienne peut s'écrire sous la forme

$$\nabla^2 f(x) = \frac{f(x)}{n^2} \times \left[ \begin{pmatrix} 1/x_1 \\ \vdots \\ 1/x_n \end{pmatrix} \begin{pmatrix} 1/x_1 & \dots & 1/x_n \end{pmatrix} - n \times \begin{pmatrix} 1/x_1^2 & 0 & \dots & 0 \\ 0 & 1/x_2^2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & 1/x_n^2 \end{pmatrix} \right].$$

2. Montrer que  $f$  est concave. *Indication : On utilisera la propriété (sans démonstration) que  $\forall (a_1 \dots a_n)^T \in \mathbb{R}^n$ ,  $n \sum_{i=1}^n a_i^2 - (\sum_{i=1}^n a_i)^2 \geq 0$ .*
3. On considère maintenant la moyenne arithmétique  $g(x) = \sum_{i=1}^n x_i$  et on choisit  $\alpha \in ]0, 1[$ . Montrer que l'ensemble  $E = \{x \in \mathbb{R}_{++}^n : f(x) \geq \alpha g(x)\}$  est un convexe.

## 2 Optimisation sans contraintes

Le problème devient

$$\min_{x \in \Omega} f(x), \quad \Omega \text{ ouvert de } \mathbb{R}^n.$$

### 2.1 Conditions nécessaires d'optimalité

**Théorème 2** Soit le problème  $(P) \min_{x \in \Omega} f(x)$ ,  $\Omega$  ouvert de  $\mathbb{R}^n$  et  $f$  une fonction de classe  $\mathcal{C}^2$ . Si  $\bar{x}$  est minimum local de  $(P)$  alors

- i) Condition de premier ordre :  $\nabla f(\bar{x}) = 0$  (point stationnaire / point critique)
- ii) Condition de second ordre : la matrice hessienne  $\nabla^2 f(\bar{x})$  est semi-définie positive.

### 2.2 Conditions suffisantes d'optimalité

Les conditions dans Théorème 2 ne sont pas suffisantes.

**Contre-exemple :** Soit

$$\begin{aligned} f : \mathbb{R} &\rightarrow \mathbb{R} \\ x &\mapsto x^3 \end{aligned}$$

Alors les conditions dans Théorème 2 sont vérifiées pour  $\bar{x} = 0$ . Or  $\bar{x} = 0$  n'est pas un minimum (local) de  $f$ .

**Théorème 3** Soit le problème  $(P) \min_{x \in \Omega} f(x)$ ,  $\Omega$  ouvert de  $\mathbb{R}^n$  et  $f$  une fonction de classe  $\mathcal{C}^2$ . Soit  $\bar{x}$  vérifiant

- i)  $\nabla f(\bar{x}) = 0$
- ii)  $\nabla^2 f(\bar{x})$  est définie positive.

Alors  $\bar{x}$  est un minimum local de  $f$ .

**Théorème 4** Soit le problème  $(P) \min_{x \in \Omega} f(x)$ ,  $\Omega$  ouvert convexe de  $\mathbb{R}^n$  et  $f$  une fonction convexe de classe  $\mathcal{C}^1$ . Alors  $\bar{x}$  est un minimum local de  $f$  ssi  $\nabla f(\bar{x}) = 0$ .

### 2.3 Résolution d'un problème d'optimisation sans contraintes

**Exemple 1** Soit  $f(x_1, x_2) = -x_1^3 - x_2^3 + 3x_1 + 12x_2$ ,  $\Omega = \mathbb{R}^2$ .

- Détermination des candidats à minimum :

$$\nabla f(x_1, x_2) = (-3x_1^2 + 3, -3x_2^2 + 12),$$



donc

$$\nabla f(x_1, x_2) = 0 \Leftrightarrow \begin{cases} -3x_1^2 + 3 = 0 \\ -3x_2^2 + 12 = 0 \end{cases} \Leftrightarrow \begin{cases} x_1 = \pm 1 \\ x_2 = \pm 2 \end{cases}.$$

On a donc 4 points critiques :  $x^{(1)} = (1, 2)^T$ ,  $x^{(2)} = (1, -2)^T$ ,  $x^{(3)} = (-1, 2)^T$ ,  $x^{(4)} = (-1, -2)^T$ . Puis la matrice

$$\nabla^2 f(x_1, x_2) = \begin{pmatrix} -6x_1 & 0 \\ 0 & -6x_2 \end{pmatrix}$$

est semi-définie positive ssi les valeurs propres  $\lambda_i$ ,  $i = 1, 2$ , sont toutes positive. Il faut donc que

$$\begin{cases} \lambda_1 = -6x_1 \geq 0 \\ \lambda_2 = -6x_2 \geq 0 \end{cases} \Leftrightarrow \begin{cases} x_1 \leq 0 \\ x_2 \leq 0 \end{cases}.$$

Par conséquent le seul candidat est  $x^{(4)} = (-1, -2)^T$ , car  $\nabla^2 f(x^{(4)})$  est définie positive à valeurs propres  $\lambda_1 = 6 > 0$  et  $\lambda_2 = 12 > 0$ . On en déduit que  $x^{(4)}$  est un minimum local.

- $x^{(4)}$  n'est pas un minimum global car

$$\lim_{(x_1, x_2) \rightarrow (+\infty, +\infty)} f(x_1, x_2) = -\infty, \quad \text{or} \quad f(x^{(4)}) = -18 > -\infty.$$

**Exemple 2** (problème convexe). Soit  $f(x_1, x_2) = x_1^2 + 2x_2^2 - 2x_1 - 4x_2$ ,  $\Omega = \mathbb{R}^2$ . Alors

$$\nabla f(x_1, x_2) = (2x_1 - 2, 4x_2 - 4) \quad \text{et} \quad \nabla^2 f(x_1, x_2) = \begin{pmatrix} 2 & 0 \\ 0 & 4 \end{pmatrix}.$$

Comme  $f$  est strictement convexe sur un domaine convexe, il s'agit d'un problème convexe. D'après Théorème des problèmes convexes (voir poly d'Optimisation) tout minimum local est minimum global et de plus il y a au plus un minimum. D'après Théorème 4,  $\bar{x}$  est minimum ssi  $\nabla f(\bar{x}) = 0$

$$\Leftrightarrow \begin{cases} 2x_1 - 2 = 0 \\ 4x_2 - 4 = 0 \end{cases} \Leftrightarrow \begin{cases} x_1 = 1 \\ x_2 = 1 \end{cases} \Leftrightarrow \bar{x} = (1, 1)^T.$$

**Exemple 3** Soit  $f(x_1, x_2) = -x_1^3 - x_1x_2^2 + x_1^2x_2 + x_2^3$ ,  $\Omega = \mathbb{R}^2$ . Points critiques :

$$\begin{aligned} \nabla f(x_1, x_2) = 0 &\Leftrightarrow \begin{cases} -3x_1^2 - x_2^2 + 2x_1x_2 = 0 \\ +3x_2^2 + x_1^2 - 2x_1x_2 = 0 \end{cases} \Leftrightarrow \begin{cases} 3x_1^2 + x_2^2 - 2x_1x_2 = 0 \\ +3x_2^2 + x_1^2 - 2x_1x_2 = 0 \end{cases} \\ &\Leftrightarrow \begin{cases} 2x_1^2 + (x_1 - x_2)^2 = 0 \\ 2x_2^2 + (x_1 - x_2)^2 = 0 \end{cases} \Leftrightarrow x_1 = x_2 = 0. \end{aligned}$$

Puis on a :

$$\nabla^2 f(x_1, x_2) = \begin{pmatrix} -6x_1 + 2x_2 & -2x_2 + 2x_1 \\ -2x_2 + 2x_1 & 6x_2 - 2x_1 \end{pmatrix} \Rightarrow \nabla^2 f(0, 0) = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}.$$

Donc on ne peut rien dire quant à la nature de  $\bar{x} = (0, 0)^T$  par l'utilisation des théorèmes précédents. Pour savoir si  $\bar{x}$  est optimum local on remarque que

$$f(x_1, x_2) = (x_2 - x_1)(x_1^2 + x_2^2) \begin{cases} > 0 \text{ si } x_2 > x_1 \\ < 0 \text{ si } x_2 < x_1 \\ = 0 \text{ si } x_2 = x_1 \end{cases}$$

Par conséquent  $\forall r > 0, \exists x', x'' \in B(\bar{x}, r)$  tq  $f(x') > 0 = f(\bar{x})$  et  $f(x'') < 0 = f(\bar{x})$ . Donc  $f$  ne possède pas d'optimum.

## 2.4 Exercices

Exercice 1 Soit  $A \in \mathcal{S}_n(\mathbb{R})$  et  $C = (c_1, \dots, c_n)$ . On considère la fonction  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  définie par

$$f(x) = \frac{1}{2}x^T Ax + Cx$$

1. Montrer que  $f$  est convexe si et seulement si  $A$  est semi-définie positive.
2. Supposons que  $A$  est définie positive et considérons le problème  $\min_{x \in \mathbb{R}^n} f(x)$ . Montrer que se problème possède un et un seul optimum global que l'on exprimera en fonction de  $A$  et  $C$ .

Exercice 2 Soit  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  la fonction définie par  $f(x) = 3x_1^4 - 4x_1^2x_2 + x_2^2$  et le problème  $(P) \min_{x \in \mathbb{R}^2} f(x)$

1. Montrer que le seul candidat à optimum de  $(P)$  est  $\tilde{x} = (0, 0)^T$ .
2. Montrer que l'on ne peut pas actionner les conditions suffisantes d'optimalité. Montrer que  $f(x) = (x_1^2 - x_2)(3x_1^2 - x_2)$  et en déduire que  $\tilde{x}$  n'est pas optimum local de  $(P)$ .
3. Pour tout  $t \in \mathbb{R}$  on note  $D_t = \{x \in \mathbb{R}^2 : x_2 = tx_1\}$  et  $(P_t) \min_{x \in D_t} f(x)$ . Montrer que pour tout  $t \in \mathbb{R}$   $\tilde{x}$  est optimum local de  $(P_t)$ .

Exercice 3 On considère le problème

$$(P) \min_{x \in \Omega} (-x_1^2 - x_2^2) \text{ avec } \Omega = \{x \in \mathbb{R}^2 : x_1^2 + x_2^2 < 1\}.$$

Montrer que la fonction possède un seul point critique  $\tilde{x} = (0, 0)^T$ . La matrice hessienne de la fonction est-elle semi-définie positive? Que peut-on en conclure?

## Exercices supplémentaires

Exercice 4 On considère sur  $\mathbb{R}^2$  la fonction

$$f(x_1, x_2) = 2x_1^2 + x_2^2 - 2x_1x_2 + 2x_1^3 + x_1^4.$$

1. Déterminer les candidats à minimum local.
2. Y-a-t-il des minimums locaux? Justifier votre réponse.
3. La fonction a-t-elle un minimum global? Pourquoi?

Exercice 5 Soient  $n \geq 2$  et  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  la fonction définie par

$$f(x) = (1 + x_n)^3 \sum_{i=1}^{n-1} x_i^2 + x_n^2.$$

1. Déterminer les candidats à minimum local.
2. Y-a-t-il des minimums locaux? Justifier votre réponse.
3. La fonction a-t-elle un minimum global? Pourquoi?

Exercice 6 On considère le problème  $(P)$   $\min_{x \in \mathbb{R}^2} f(x)$  où

$$f(x) = \frac{1}{2} x^T A x + b x, \quad \text{avec } A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \text{ et } b = \begin{pmatrix} 1 & 0 \end{pmatrix}.$$

1. Montrer que  $(P)$  est un problème strictement convexe.
2. Déterminer les candidats à minimum local.
3. Y-a-t-il une solution optimale? Justifier votre réponse.
4. Quelle est la valeur optimale?

Exercice 7 On considère sur  $\mathbb{R}^2$  la fonction

$$f(x_1, x_2) = x_1^3 - 2x_1^2 x_2 + x_2^2$$

1. Déterminer les candidats à minimum local.
2. Y-a-t-il des minimums locaux? Justifier votre réponse.
3. La fonction a-t-elle un minimum global? Pourquoi?

### 3 Optimisation avec contraintes

On rappelle la formulation du problème

$$(P) \quad \begin{array}{ll} \min f(x) \\ \text{s.c.} & g(x) = 0 \quad \text{contraintes d'égalité} \\ & h(x) \leq 0 \quad \text{contraintes d'inégalité} \\ & x \in \Omega \end{array}$$

avec  $g = (g_i)_{i=1,\dots,p}$  et  $h = (h_j)_{j=1,\dots,q}$ .

**Définition 3** On dit qu'une contrainte  $h_j(x) \leq 0$  est saturée (ou active) en  $x^*$  ssi  $h_j(x^*) = 0$ .

On note

$$\mathcal{A} = \{x \in \Omega : g_i(x) = 0, i = 1, \dots, p, h_j(x) \leq 0, j = 1, \dots, q\}$$

l'ensemble admissible (le domaine) et on remarque que (*exercice 1*)

- si les fonctions  $g_i$  et  $h_j$  sont continues, alors  $\mathcal{A}$  est fermé;
- si les fonctions  $g_i$  sont linéaires et  $h_j$  sont convexes, alors  $\mathcal{A}$  est convexe.

Pour pouvoir caractériser la solution du problème on a besoin d'une propriété sur le domaine – la *qualification* des contraintes en rapport avec la géométrie du domaine.

#### 3.1 Qualification des contraintes

La qualification des contraintes au point  $x \in \mathcal{A}$  garantit que la géométrie de l'ensemble admissible au voisinage du point n'est pas trop atypique. Donnons juste quelques conditions suffisantes qui garantissent que les contraintes sont qualifiées en un point admissible  $x$  :

1. *Mangasarian-Fromovitz* :  $g_i, i = 1, \dots, p$  et  $h_j, j = 1, \dots, q$  sont de classe  $\mathcal{C}^1$

$$\left\{ \begin{array}{l} * \quad \nabla g_i(x), i = 1, \dots, p \text{ sont linéairement indépendants} \\ * \quad \exists v \in \mathbb{R}^n : \begin{cases} \langle \nabla g_i(x)^T, v \rangle = 0 & i = 1, \dots, p \\ \langle \nabla h_j(x)^T, v \rangle < 0 & \text{si } h_j(x) = 0 \end{cases} \end{array} \right.$$

2. *Fiacco-McCormick* : les gradients de toutes les contraintes saturées en  $x$  (c'est-à-dire,  $\nabla g_i(x), i = 1, \dots, p$  et  $\nabla h_j(x)$  si  $h_j(x) = 0$ ) sont linéairement indépendants.

Ensuite quelques conditions suffisantes qui garantissent que les contraintes sont qualifiées en tout point admissible :

1. les fonctions  $g_i$  et  $h_j$  sont affines

2. *Slater* :

$$\left\{ \begin{array}{l} * \quad f \text{ est convexe sur l'ouvert convexe } \Omega \\ * \quad g : x \mapsto (Ax - b) \text{ où } A \text{ est surjective;} \\ \quad \quad \quad \text{c-à-d., les vecteurs lignes de } A \text{ sont linéairement indépendants} \\ * \quad h_j : \Omega \rightarrow \mathbb{R}, j = 1, \dots, q, \text{ convexes de classe } \mathcal{C}^1 \\ * \quad \exists x_0 \in \Omega \text{ admissible vérifiant } h_j(x_0) < 0 \text{ pour toute contrainte d'inégalité non linéaire} \end{array} \right.$$

### 3.2 Conditions nécessaires d'optimalité : théorème KKT

**Théorème 5** (Karush-Kuhn-Tucker). *Si  $x^*$  est solution du problème*

$$(P) \quad \begin{aligned} \min & f(x) \\ \text{s.c.} & \quad g(x) = 0 \\ & \quad h(x) \leq 0 \\ & \quad x \in \Omega \end{aligned}$$

*avec  $f$ ,  $g$  et  $h$  dérivables en  $x^*$ , et les contraintes sont qualifiées en  $x^*$ , alors il existe un vecteur  $\Lambda = (\lambda_1, \dots, \lambda_p)^T \in \mathbb{R}^p$  et  $M = (\mu_1, \dots, \mu_q)^T \in \mathbb{R}^q$  tels que*

$$(KKT) \quad \begin{cases} (I) & \nabla f(x^*) - \Lambda^T \nabla g(x^*) - M^T \nabla h(x^*) = 0 \\ (II) & M \leq 0 \\ (III) & \langle M, h(x^*) \rangle = M^T \cdot h(x^*) = \sum_{j=1}^q \mu_j h_j(x^*) = 0 \quad (*) \end{cases}$$

*où  $\nabla g$  et  $\nabla h$  sont les matrices jacobiniennes définies par*

$$\nabla g = \begin{pmatrix} \frac{\partial g_1}{\partial x_1} & \cdots & \frac{\partial g_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial g_p}{\partial x_1} & \cdots & \frac{\partial g_p}{\partial x_n} \end{pmatrix} \quad \text{et} \quad \nabla h = \begin{pmatrix} \frac{\partial h_1}{\partial x_1} & \cdots & \frac{\partial h_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial h_q}{\partial x_1} & \cdots & \frac{\partial h_q}{\partial x_n} \end{pmatrix}.$$

Les  $\lambda_i$  et  $\mu_j$  mis en évidence dans ce théorème de KKT sont appelés *multiplicateurs de Lagrange* (ou *de Lagrange-KKT*). Les (in-)égalités (I)–(III) du théorème de KKT sont appelées *conditions de KKT*.

Les relations de complémentarité (\*) peuvent aussi s'écrire  $\mu_j h_j(x^*) = 0$ ,  $j = 1, \dots, q$ , car on a une somme de termes tous positifs qui doit être nul.

Comment procéder pour résoudre un problème?

- (i) Vérifier la qualification des contraintes dans les points admissibles.

*Tous les points pour lesquels on n'a pas réussi à montrer une condition suffisante de qualification des contraintes doivent être considérés comme potentiels candidats à optimum.*

- (ii) Résoudre le système (KKT) qui donnera, si qualification des contraintes, les candidats à optimum.
- (iii) Utiliser les arguments de compacité du domaine, coercivité de la fonction objectif, convexité (analyse de la matrice hessienne) pour conclure.

**Théorème 6** *Si (P) est un problème convexe les conditions KKT sont suffisantes d'optimalité; c-à-d.,  $x^*$  vérifie les conditions KKT  $\Rightarrow x^*$  solution optimale.*

Deux grandes classes de problèmes :

A) Problèmes linéaires

$$(P) \quad \begin{array}{ll} \max & f^T x \\ \text{s.c.} & Ax = b \\ & Cx \leq d \end{array}$$

C'est toujours un problème convexe. La structure très particulière conduit à des méthodes spécifiques de résolution (voir poly d'Optimisation : Programmation linéaire).

B) Problèmes quadratiques

$$(P) \quad \begin{array}{ll} \min & [\frac{1}{2}\langle x, Qx \rangle - \langle r, x \rangle] \\ \text{s.c.} & Ax = b \\ & Cx \leq d \end{array}$$

avec  $Q \in \mathcal{S}_n(\mathbb{R})$  et  $r \in \mathbb{R}^n$ . C'est un problème où la fonction objectif est quadratique et les contraintes sont linéaires. Comme la matrice hessienne  $\nabla^2 f(x) = Q$ , c'est un problème convexe si  $Q$  est semi-définie positive et strictement convexe si  $Q$  est définie positive.

**Corollaire 1** *Tout problème quadratique avec  $Q$  définie positive admet au plus une solution optimale.*

**Exemple 4** (Solution de portefeuille). *Les rendements annuels de 3 actifs financiers sont modélisés par des variables aléatoires indépendantes dont on a estimé les espérances : 5, 10 et 15% et les écarts types : 2, 8 et 10% respectivement à partir de historique des cours. La variance du rendement d'un actif financier est une mesure du risque associé à cet actif. Comment constituer avec les 3 actifs un portefeuille de variance minimale (et donc le moins risqué)?*

*Le problème s'écrit*

$$(P) \quad \begin{array}{ll} \min & (4x_1^2 + 64x_2^2 + 100x_3^2) \\ \text{s.c.} & x_1 + x_2 + x_3 = 1 \\ & x_i \geq 0, \quad 1 \leq i \leq 3 \end{array}$$

avec  $x_i =$  part de l'actif  $i$ ,  $i = 1, \dots, 3$ . C'est un problème quadratique elliptique.

*(Une stratégie naïve consiste à n'investir que dans le premier actif pour un rendement de 5% et un risque de 2%. La résolution du problème permet de constituer un portefeuille dont le rendement est strictement supérieur à 5% avec un risque inférieur à 2%.)*

*D'abord on note qu'il s'agit d'un problème d'optimisation strictement convexe, donc*

- 1. d'après Théorème 6, si  $x^*$  vérifie les conditions KKT, alors  $x^*$  est solution optimal;*
- 2. d'après Corollaire 1, cette solution  $x^*$  est unique.*

*Le système (KKT) devient*

$$\begin{cases} 8x_1 - \lambda + \mu_1 = 0 \\ 128x_2 - \lambda + \mu_2 = 0 \\ 200x_3 - \lambda + \mu_3 = 0 \\ \mu_j \leq 0, \quad 1 \leq j \leq 3 \\ \mu_j x_j = 0, \quad 1 \leq j \leq 3 \\ x_1 + x_2 + x_3 - 1 = 0 \\ -x_j \leq 0, \quad 1 \leq j \leq 3 \end{cases}$$

En résolvant ce système on obtient  $M = 0$ ,  $\Lambda = \frac{3200}{441}$ ,  $x^* = \frac{1}{441}(400, 25, 16)^T$  et  $f(x^*) = \frac{1600}{441}$ . Le rendement est donc  $(5, 10, 15) \cdot x^* = \frac{830}{147} \approx 5,65\%$  avec un risque de  $\sqrt{f(x^*)} = \frac{40}{21} \approx 1,90\%$ .

### 3.3 Exercices

Exercice 1 Considérons le domaine suivant

$$\mathcal{A} = \{x \in \Omega : g_i(x) = 0, i = 1, \dots, p; \quad h_j(x) \leq 0, j = 1, \dots, q\}.$$

Montrer que:

- Si  $g_i$  et  $h_j$  sont des fonctions continues alors le domaine est fermé.
- Si les fonctions  $g_i$  sont linéaires et  $h_j$  sont convexes alors le domaine est convexe.

Exercice 2 Etudier la classification des contraintes pour le domaine suivant

$$\mathcal{A}_1 = \{(x, y)^T : x \geq 0, y \geq 0, (1 - x)^3 - y \geq 0\}.$$

Résoudre  $\min_{x \in \mathcal{A}_1} f(x, y)$  avec  $f(x, y) = -x^2 - (y + 1)^2$ .

Exercice 3 Résoudre les problèmes suivants

(a)

$$\min_{x_1^2 + x_2^2 - 2 = 0} (-x_1 x_2)$$

(b)

$$\min_{(x_1 - 1)^3 - x_2^2 = 0} (x_1^2 + x_2^2)$$

(c)

$$\min_{x \in \mathcal{A}} \frac{1}{2} (x_1^2 + 2x_2^2 + 3x_3^2) \text{ avec } \mathcal{A} = \{x \in \mathbb{R}^3 : x_1 + x_2 + x_3 = 1, \quad x_1 - x_2 = 18\}$$

(d)

$$\min_{x_1^2 + x_2^2 \leq 1} (x_2 - x_1^2)$$

(e)

$$\min_{x \in \{x_1 + x_2 \leq 2; x_1 \geq 0; x_2 \geq 0\}} (x_1^2 + x_2^2 - 8x_1 + 4x_2)$$

(f)

$$\min_{x \in \mathcal{A}} \frac{1}{2} (x_1^2 + (x_2 - 2)^2 + (x_3 + 2)^2) \text{ avec } \mathcal{A} = \{x \in \mathbb{R}^3 : x_1 + x_2 + x_3 = 0; x_1^2 + x_2^2 + x_3^2 \leq 2; x_2 \geq 0\}$$

Exercice 4 Soit  $f \in \mathcal{C}^1(\mathbb{R}^n)$ . Montrer que si  $\bar{x}$  est solution optimale pour le problème

$$\min\{f(x) : Ax = a, x \geq 0\} \text{ alors}$$

$$\exists \lambda \in \mathbb{R}^m : \quad \nabla f(\bar{x}) + \lambda^T A \geq 0, \quad A\bar{x} = a, \quad \bar{x} \geq 0, \quad \nabla f(\bar{x})\bar{x} + \lambda^T a = 0.$$

Exercice 5 Soit  $\mathcal{P}$  le polyèdre défini par

$$\mathcal{P} = \{x \in \mathbb{R}^n : \quad Bx \geq b, \quad Cx = c\}$$

Montrer que une condition nécessaire pour que  $\bar{x}$  soit solution optimale pour le problème

$$\min\{f(x), \quad x \in \mathcal{P}\} \text{ avec } f \in \mathcal{C}^1 \text{ est}$$

$$\exists \mu \leq 0 \quad \exists \lambda \text{ tels que } \begin{cases} \nabla f(\bar{x}) + \mu^T B + \lambda^T C = 0 \\ \mu^T (B\bar{x} - b) = 0, \quad B\bar{x} - b \geq 0, \quad C\bar{x} = c. \end{cases}$$

Exercice 6 On se donne  $n$  fonctions dérivables  $f_j : \mathbb{R} \rightarrow \mathbb{R}$  et on considère le problème

$$\begin{cases} \min \sum_{j=1}^n f_j(x_j) \\ \sum_{j=1}^n x_j = 1, \quad x_j \geq 0, \quad j = 1, \dots, n \end{cases}$$

Montrer que si  $\bar{x}$  est solution de ce problème alors  $\exists \lambda$  tel que:

$$\begin{aligned} \bar{x}_j > 0 &\Rightarrow f'_j(\bar{x}_j) = \lambda \\ \bar{x}_j = 0 &\Rightarrow f'_j(\bar{x}_j) \geq \lambda \quad j = 1, \dots, n \end{aligned}$$

Résoudre  $\min\{x_1^3 + 3x_1 + x_2^2 : x_1 + x_2 = 1, x_1, x_2 \geq 0\}$ .

Exercice 7

1. Maximiser les fonctions

$$\begin{aligned} f(x) &= x_1^4 + 4x_1^2x_2^2 - 2x_1^2 + 2x_2^2 - 1, \\ g(x) &= 2x_1^2 + x_1x_2 + x_2^2 + x_2x_3 + x_3^2 - 6x_1 - 7x_2 - 8x_3 + 9. \end{aligned}$$

dans le disque d'unité  $\{x : \|x\| \leq 1\}$ .

2. Résoudre le problème suivant

$$\max\{f(x), \quad x \in \mathcal{A}\} \text{ avec}$$

$$f(x) = -x_1^2 - x_2^2 - x_3^2, \quad \mathcal{A} = \{x \in \mathbb{R}^3 : \quad -x_1 + x_2 - x_3 \geq -10, \quad x_1 + x_2 + 4x_3 \geq 20\}$$



Exercice 8 On considère le problème

$$(P) \quad \min \left[ \frac{1}{2}((x_1 - 1)^2 + (x_2 + 1)^2 + (x_3 - 1)^2) \right]$$
$$\text{s.c. } x_1 \leq x_2 \leq x_3$$

1. Vérifier que (P) est un problème d'optimisation convexe.
2. Montrer que (P) possède un optimum et un seul.
3. Résoudre (P).

Exercice 9 Soit le problème

$$(P) \quad \begin{cases} \min(x_1^2 + x_2^2 - 4x_1x_2 - 2x_1 - 4x_2) \\ ax_1 + bx_2 \leq 4 \\ x_1 \geq 0; x_2 \geq 0 \end{cases}$$

où  $a$  et  $b$  sont 2 paramètres réels strictement positifs. Déterminer  $a$  et  $b$  pour que  $\bar{x} = (1/2, 1/2)^T$  soit un optimum local de (P).

Exercice 10 On considère le problème

$$(P) \quad \begin{cases} \min f(x) \\ a_i \leq x_i \leq b_i \quad (1 \leq i \leq n) \end{cases}$$

avec  $f$  de classe  $\mathcal{C}^2$  et  $a_i < b_i$ ,  $i = 1, \dots, n$ . Soit  $\bar{x} = (\bar{x}_i)$  un optimum local de (P). Montrer que l'on a

$$\begin{cases} \frac{\partial f}{\partial x_i}(\bar{x}) \geq 0 & \text{si } \bar{x}_i = a_i \\ \frac{\partial f}{\partial x_i}(\bar{x}) \leq 0 & \text{si } \bar{x}_i = b_i \\ \frac{\partial f}{\partial x_i}(\bar{x}) = 0 & \text{si } \bar{x}_i \in ]a_i, b_i[ \end{cases}$$

## Exercices supplémentaires

Exercice 11 On veut résoudre le problème suivant:

$$\min(x_1^2 + x_2^2 - 3x_1 - 5x_2)$$
$$(P) \quad \begin{cases} \text{s.c. } & x_1 + 2x_2 \leq 4 \\ & 3x_1 + x_2 \leq 7 \\ & x_1, x_2 \geq 0 \end{cases}$$

1. Ecrire les conditions de KKT pour ce problème.

2. Sont-elles vérifiées pour le point  $x = (2, 1)^T$ ?
3. Déterminer les candidats à optimum.
4. Résoudre le problème  $(P)$ .

Exercice 12 On veut résoudre le problème suivant:

$$\begin{aligned} & \min(6x_1 - 6x_2 - 3x_1^2 - 3x_2^2 - x_1x_2) \\ (P) \quad & \begin{cases} \text{s.c.} & 3x_1 + 4x_2 \leq 12 \\ & x_1, x_2 \geq 0 \end{cases} \end{aligned}$$

1. Ce problème a-t-il une solution? Justifier.
2. Pourquoi peut-on appliquer KKT pour le résoudre?
3. Résoudre le problème  $(P)$ .

Exercice 13 Vérifier si les contraintes

$$\begin{cases} (x_1 - 1)^2 + (x_2 - 1)^2 \leq 2 \\ (x_1 - 1)^2 + (x_2 + 1)^2 \leq 2 \\ x_1 \geq 0 \end{cases}$$

sont qualifiées au point  $x = (0, 0)^T$ .

Exercice 14 On veut déterminer le point de la parabole  $y = \frac{1}{5}(x - 1)^2$  qui est le plus proche de  $(1, -2)^T$ .

1. Montrer que ce problème peut s'écrire sous la forme suivante

$$(P) \quad \begin{cases} \min(x - 1)^2 + (y + 2)^2 \\ \text{s.c.} & (x - 1)^2 = 5y \end{cases}$$

2. Le problème admet-il une solution? Pourquoi?
3. Montrer que la solution est unique.
4. Calculer la solution.

Exercice 15 Résoudre

$$(P) \quad \begin{cases} \max x_1x_2 \\ \text{s.c.} & x_1^2 + x_2^2 \leq 1 \end{cases}$$

Exercice 16 On veut résoudre le problème suivant :

$$(P) \quad \max_{x_1^3 \leq 8x_2 \leq 2x_1^2} (x_1 - 1)^2 + (x_2 - 1)^2.$$

1. Les contraintes sont-elles qualifiées en tout point admissible ? Justifier.
2. Ecrire les conditions KKT pour le problème  $(P)$ .
3. Montrer que les conditions KKT sont vérifiées pour  $x = (2, 1)^T$ .
4. Déterminer la solution optimale et la valeur optimale du problème  $(P)$  sachant que  $x = (2, 1)^T$  est le seul point admissible qui vérifie les conditions KKT. Vérifier vos résultats à l'aide d'une représentation graphique du problème.

Exercice 17 On veut déterminer les points de

$$\mathcal{S} = \{(x, y, z)^T \in \mathbb{R}^3 : z^2 = xy + 4\}$$

les plus proches de l'origine  $(0, 0, 0)^T$ .

1. Formuler ce problème comme un problème d'optimisation.
2. Ce problème a-t-il une solution ? Justifier.
3. Montrer que l'on peut appliquer KKT pour le calcul de(s) solution(s).
4. Calculer la (les) solution(s) en utilisant la méthode KKT.

## 4 Quelques algorithmes de résolution

### Algorithme général

Le système d'équations KKT que l'on doit résoudre pour trouver les candidats à optimum de

$$(P) \quad \begin{array}{ll} \min & f(x) \\ \text{s.c.} & g(x) = 0 \\ & h(x) \leq 0 \\ & x \in \Omega \end{array}$$

est presque toujours très difficile à résoudre (sauf dans le cas de petite dimension). Il faut donc donner des *méthodes itératives* (qui génèrent une suite de vecteurs) qui permettent de calculer une approximation de la solution du problème. La structure générale d'un algorithme sera

$x \leftarrow x_0$  (point initial)  
Tant que *critère d'arrêt* = *faux* faire  
\* déterminer une direction  $d$  de déplacement  
\* déterminer un pas  $\rho$  dans cette direction  
\* actualiser le point courant :  $x \leftarrow x + \rho d$

Différentes méthodes correspondent à différents choix de la direction et différents choix du pas. Ces choix doivent être faits de façon à garantir la convergence de la méthode; c'est-à-dire, générer une suite qui converge vers la solution. Le critère d'arrêt correspond à une précision exigée pour l'approximation de la solution.

### 4.1 Choix de la direction

**Définition 4** Une direction  $d$  est une direction de descente pour  $f$  au point  $x$  si

$$\exists T > 0 \text{ tq } \forall \rho \in ]0, T[ : f(x + \rho d) < f(x).$$

C'est-à-dire, on peut diminuer la valeur de la fonction objectif avec un pas  $\rho$  arbitraire petit dans la direction  $d$ .

**Théorème 7** Si  $f$  est dérivable au point  $x$ , alors  $d$  est une direction de descente au point  $x$  si

$$\nabla f(x) \cdot d < 0.$$

*Démonstration.* Si on pose  $\varphi(t) = f(x + td)$ , alors  $\varphi'(0) = \nabla f(x) \cdot d < 0$ . Donc  $\exists L > 0$  tq  $\forall t \in ]0, L[ : f(x + td) < f(x)$ . ■

**Théorème 8** Si  $f$  est convexe sur  $[x, y]$  et  $f(y) < f(x)$ , alors  $d = y - x$  est une direction de descente pour  $f$  en  $x$ .

*Démonstration.*  $\forall t \in ]0, 1[$ ,

$$f(x + td) = f(x + t(y - x)) = f(ty + (1 - t)x) \leq tf(y) + (1 - t)f(x) < f(x). \quad \blacksquare$$

#### 4.1.1 Méthodes de gradient

Pour cette classe de méthodes le choix de la direction est  $d = -\nabla f(x)^T$ . Ceci correspond au choix de la direction qui, en norme constante, minimise l'approximation linéaire de  $f(x + d) : f(x) + \nabla f(x) \cdot d$ .

##### 1. Méthode de gradient à pas fixe

Problème	$\min_{x \in \Omega} f(x)$	$\Omega =$ ouvert de $\mathbb{R}^n$
		$f$ dérivable

##### Algorithme 1

```

GradPasFixe( $f, x_0, \rho, tol$ )
 $x \leftarrow x_0, d \leftarrow -\nabla f(x_0)^T$ 
Tant que  $\|d\| > tol$  faire
     $x \leftarrow x + \rho d$ 
     $d \leftarrow -\nabla f(x)^T$ 

```

**Théorème 9** (Convergence). *Sous les hypothèses*

(H0) l'ensemble de niveau  $S_0 = \{x \in \Omega : f(x) \leq f(x_0)\}$  est fermé dans  $\mathbb{R}^n$  et  $f$  est minorée sur  $\Omega$ ;

(H1)  $f$  est 2 fois dérivable en  $x$  et  $\exists \kappa, c > 0$  constantes pour tout  $x \in S_0$  tq

$$\forall u \in \mathbb{R}^n : c \|u\|^2 \leq u^T \nabla^2 f(x) u \leq \kappa \|u\|^2 ;$$

(H2)  $x_0$  n'est pas un point critique de  $f$ ;

la suite  $\{x_k\}_{k \geq 0}$  définie par  $x_{k+1} = x_k - \rho \nabla f(x_k)^T$ , avec  $\rho$  fixé vérifiant  $0 < \rho < 2/\kappa$ , soit  $\lim_{k \rightarrow \infty} \|x_k\| = +\infty$ , soit converge vers un minimiseur local  $x^*$  de  $f$ . La convergence est linéaire.

Pour démontrer le théorème précédent, on a besoin du lemme suivant :

**Lemma 1** Soit  $S_\nu = \{x \in \Omega : f(x) \leq \nu\}$ ,  $x \in S_\nu$  et  $d$  une direction de descente avec  $\alpha := -\nabla f(x) \cdot d > 0$ . Supposons que

(H0( $\nu$ ))  $S_\nu$  est fermé dans  $\mathbb{R}^n$ ;

(H1( $\nu$ ))  $f$  est 2 fois dérivable sur  $S_\nu$  et  $\exists \kappa > 0$  tq  $\forall y \in S_\nu$  et  $\forall u \in \mathbb{R}^n : u^T \nabla^2 f(y) u \leq \kappa \|u\|^2$ .

On pose  $\varphi(\rho) = f(x + \rho d)$ . Alors  $\forall \rho$  tq  $0 < \rho < \frac{2\alpha}{\kappa \|d\|^2}$  on a

$$\varphi(\rho) \leq f(x) - \alpha \rho + \frac{\kappa}{2} \rho^2 \|d\|^2.$$

(C'est-à-dire, le graphe de  $\varphi$  est dominé dans un voisinage de zero par le graphe d'une parabole tangente au graphe de  $\varphi$  en  $(0, \varphi(0))$ , décroissante pour  $\rho > 0$  dont la courbure est contrôlée par un majorant  $\kappa$  du rayon spectral de la hessienne.)

*Démonstration.*

$$\varphi(\rho) = f(x + \rho d) \Rightarrow \varphi'(\rho) = \nabla f(x + \rho d)d \Rightarrow \varphi''(\rho) = d^T \nabla^2 f(x + \rho d)d \leq \kappa \|d\|^2$$

si  $(x + \rho d) \in S_\nu$ . Alors  $\varphi'(0) = -\alpha < 0$ , donc  $\varphi(\rho) < \varphi(0)$  pour  $\rho$  suffisamment petit. Soit  $\sigma > 0$  tq

$$\sigma = \sup\{s > 0 : 0 \leq \rho \leq s \Rightarrow \varphi(\rho) \leq \varphi(0)\}.$$

Pour  $0 < \rho < \sigma$ ,  $\varphi(\rho) \leq \varphi(0) + \rho\varphi'(0) + \frac{\kappa}{2}\rho^2 \|d\|^2$ . Si  $\sigma < +\infty$  on a  $(x + \sigma d) \in S_\nu$  et  $\varphi(\sigma) = \varphi(0)$ . Donc

$$\sigma\varphi'(0) + \frac{\kappa}{2}\sigma^2 \|d\|^2 \geq 0 \Rightarrow \varphi'(0) + \frac{\kappa}{2}\sigma \|d\|^2 \geq 0 \Rightarrow \sigma \geq \frac{2\alpha}{\kappa \|d\|^2}. \quad \blacksquare$$

*Démonstration du Théorème 9.*

- La suite  $f_k := f(x_k)$  est strictement décroissante. En effet on applique le lemme avec  $\nu = f_k$  et  $d = -\nabla f(x_k)^T$  :

$$f_{k+1} \leq f_k - \rho \|\nabla f(x_k)\|^2 + \frac{\kappa\rho^2}{2} \|\nabla f(x_k)\|^2 = f_k - \frac{\kappa\rho}{2} \left(\frac{2}{\kappa} - \rho\right) \|\nabla f(x_k)\|^2 < f_k.$$

- Comme  $\inf f > -\infty$  la suite  $f_k$  étant décroissante minorée converge. Alors la suite  $\nabla f(x_k)$  doit converger vers 0.
- On extrait de  $\{x_k\}_{k \geq 0}$ , si elle est bornée, un sous-suite convergente vers un point  $x^*$  qui par continuité vérifie  $\nabla f(x^*) = 0$ . Si  $x_0$  n'est pas un point critique,  $x^* \in \text{Int}(S_0)$  car  $f(x^*) < f(x_0)$ . Comme  $\nabla^2 f(x)$  est symétrique définie positive dans  $S_0$ ,  $x^*$  est un minimiseur local.
- On pose  $\phi(x) = x - \rho \nabla f(x)^T$  et donc  $x_{k+1} = \phi(x_k)$ . La jacobienne de  $\phi$  est  $J_\phi(x) = I - \rho \nabla^2 f(x)$ , donc le rayon spectral est majoré par

$$\theta(\rho) = \max_{c \leq \lambda \leq \kappa} |1 - \rho\lambda| = \max(|1 - c\rho|, |1 - \kappa\rho|) < 1$$

et donc

$$\|x_{k+1} - x_k\| = \|\phi(x_k) - \phi(x_{k-1})\| \leq \theta(\rho) \|x_k - x_{k-1}\| \leq \theta(\rho)^k \|x_1 - x_0\|.$$

On en déduit que

$$\|x_k - x_0\| \leq \frac{1}{1 - \theta(\rho)} \|x_1 - x_0\|,$$

donc convergence de la suite (qui est de Cauchy) et

$$\|x_{k+1} - x^*\| = \|\phi(x_k) - \phi(x^*)\| \leq \theta(\rho) \|x_k - x^*\|. \quad \blacksquare$$

Cas d'un critère quadratique : Supposons que

$$f(x) = \frac{1}{2} \langle Ax, x \rangle - \langle b, x \rangle$$

avec  $A \in \mathcal{S}_n(\mathbb{R})$ ,  $b \in \mathbb{R}^n$ ,  $\nabla f(x)^T = Ax - b$  et  $\nabla^2 f(x) = A$  constante. Les points calculés par l'algo vérifient

$$x_{k+1} = x_k - \rho(Ax_k - b).$$

Tout minimiseur  $x^*$  de  $f$  vérifie  $\nabla f(x^*) = 0 \Leftrightarrow Ax^* - b = 0$ . Donc

$$x_{k+1} - x^* = x_k - \rho(Ax_k - b) - x^* = (I - \rho A)(x_k - x^*).$$

On écrit  $x_0 - x^* = \sum_{i=1}^n \beta_i v_i$  avec  $v_i$  vecteur propre normé de  $I - \rho A$  associé à la valeur propre  $\mu_i = (1 - \rho \lambda_i)$ ,  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ . Alors

$$x_k - x^* = \sum_{i=1}^n \mu_i^k \beta_i v_i \Rightarrow \|x_k - x^*\|^2 = \sum_{i=1}^n \mu_i^{2k} \beta_i^2.$$

Ceci ne converge vers 0 que si tous les  $\mu_i$  correspondant aux composantes  $\beta_i v_i$  non nulles de la décomposition sont strictement inférieurs à 1 en valeur absolue. C'est-à-dire,  $x_0 - x^*$  appartient au sous-espace engendré par les vecteurs propres de  $I - \rho A$  associés à une valeur propre appartenant à  $\text{sp}(I - \rho A) \cap ]-1, 1[$  et la convergence est linéaire. En particulier l'algo converge pour toute initialisation  $x_0 \in \mathbb{R}^n$  si  $\text{sp}(I - \rho A) \subset ]-1, 1[$ . Le taux de convergence sera  $\max(|1 - \rho \lambda_1|, |1 - \rho \lambda_n|)$ . Le taux sera minimal pour  $\rho = \frac{2}{\lambda_1 + \lambda_n}$  et il vaut  $\frac{\gamma-1}{\gamma+1}$  avec  $\gamma = \text{cond}_2(A)$ .

## 2. Méthode de gradient à pas optimal

Problème	$\min_{x \in \Omega} f(x)$	$\Omega =$ ouvert de $\mathbb{R}^n$
		$f$ dérivable

### Algorithme 2

```

GradPasOpt( $f, x_0, tol$ )
 $x \leftarrow x_0, d \leftarrow -\nabla f(x_0)^T$ 
Tant que  $\|d\| > tol$  faire
     $\rho = \text{ArgMin}\{f(x + td), t > 0\}$ 
     $x \leftarrow x + \rho d$ 
     $d \leftarrow -\nabla f(x)^T$ 

```

**Théorème 10** (Convergence). *Sous les hypothèses*

(H0)  $S_0 = \{x \in \Omega : f(x) \leq f(x_0)\}$  est un convexe compact de  $\mathbb{R}^n$ ;

(H1)  $f$  est 2 fois dérivable en  $S_0$ ,  $\nabla^2 f(x)$  est définie positive pour tout  $x \in S_0$  et  $\exists \kappa > 0$  constante tq

$$\forall x \in S_0 \text{ et } \forall u \in \mathbb{R}^n : u^T \nabla^2 f(x) u \leq \kappa \|u\|^2 ;$$

la suite  $\{x_k\}_{k \geq 0}$  définie par  $x_{k+1} = x_k - \rho_k \nabla f(x_k)^T$ , avec

$$\rho_k = \text{ArgMin}\{f(x_k - t \nabla f(x_k)^T), t > 0\},$$

converge vers l'unique minimiseur de  $f$  dans  $\Omega$ .

**Remarque :** Cet algorithme appliqué à une fonction elliptique a un comportement en zig-zag. On a  $\varphi(t) = f(x_k + t d_k)$  avec  $d_k = -\nabla f(x_k)^T$ , et  $\varphi'(t) = \nabla f(x_k + t d_k) d_k$ . Pour  $t = \rho_k$  paramètre optimal,  $x_k + \rho_k d_k = x_{k+1}$  et  $\varphi'(\rho_k) = 0$ , donc

$$\nabla f(x_k + \rho_k d_k) d_k = 0 \Leftrightarrow d_{k+1}^T d_k = 0.$$

Deux directions de descente successives sont donc toujours orthogonales.

### 3. *Le gradient conjugué*

Le problème à résoudre est

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{avec} \quad f(x) = \frac{1}{2} \langle Ax, x \rangle - \langle b, x \rangle$$

$A$  définie positive ( $f$  elliptique)

Donc  $f$  est quadratique. Le problème est strictement convexe, donc il admet une solution unique  $x^*$  qui doit vérifier

$$\nabla f(x^*) = 0 \Leftrightarrow Ax^* - b = 0 \Leftrightarrow Ax^* = b.$$

Donc ça équivaut à résoudre ce système linéaire, ce qui fait du GC une des méthodes les plus efficace de résolution de systèmes linéaires symétriques définis positifs.

**Définition 5** Deux directions  $d$  et  $d'$  sont dites  $A$ -conjuguées ssi  $\langle Ad, d' \rangle = \langle d, Ad' \rangle = 0$ ; donc orthogonales pour le produit scalaire associée à la matrice  $A$  symétrique définie positive  $\langle x, y \rangle_A = x^T A y$ .

Détaillons l'algorithme. On va construire une suite  $\{x_k\}_{k \geq 0}$  définie par

$$x_{k+1} = x_k + \rho_k d_k, \quad x_0 \text{ donné,}$$

avec

(i)  $\rho_k$  le pas optimal dans la direction  $d_k$  :

$$\rho_k = \text{ArgMin}\{f(x_k + t d_k), t > 0\}.$$

On pose

$$\varphi(t) = f(x_k + t d_k) \Rightarrow \varphi'(t) = \nabla f(x_k + t d_k) d_k.$$



Alors le minimum est atteint pour

$$\begin{aligned}\varphi'(t) &= [A(x_k + td_k) - b]^T d_k = 0 \Leftrightarrow \langle Ax_k - b, d_k \rangle + t \langle Ad_k, d_k \rangle = 0 \Leftrightarrow \\ \rho_k &= t = \frac{\langle b - Ax_k, d_k \rangle}{\langle Ad_k, d_k \rangle}.\end{aligned}$$

En posant  $r_k = b - Ax_k$ , on peut démontrer (*exercice 1*) que  $\langle r_{k+1}, d_k \rangle = 0$ . Donc avec ce choix optimal de paramètre on a :

$$\langle Ax_{k+1} - b, d_k \rangle = 0 \Leftrightarrow \langle A(x_{k+1} - x^*), d_k \rangle = 0.$$

La direction optimale de recherche  $(x_{k+1} - x^*)$  est donc  $A$ -conjuguée de  $d_k$ .

(ii) Choix de la direction :

$d_0 = -\nabla f(x_0)^T$ , et à l'étape  $k+1$ , si  $x_{k+1}$  n'est pas solution, on pose  $r_{k+1} = b - Ax_{k+1}$  (résidu) et on calcule

$$d_{k+1} = r_{k+1} + \alpha_k d_k$$

avec  $\alpha_k$  un paramètre de sorte que  $d_{k+1}$  et  $d_k$  soient  $A$ -conjuguées. Donc

$$\langle d_{k+1}, Ad_k \rangle = \langle r_{k+1}, Ad_k \rangle + \alpha_k \langle d_k, Ad_k \rangle = 0 \Leftrightarrow \alpha_k = -\frac{\langle r_{k+1}, Ad_k \rangle}{\langle d_k, Ad_k \rangle}.$$

Ces formules peuvent être simplifiées (*exercice 1*) :

$$\begin{aligned}& * \quad \rho_k = \| r_k \|^2 / \langle d_k, Ad_k \rangle ; \\& * \quad \forall k \geq 0 : \langle r_{k+1}, r_k \rangle = 0 ; \\& * \quad \alpha_k = \| r_{k+1} \|^2 / \| r_k \|^2 .\end{aligned}$$

On obtient donc l'algorithme suivant :

### Algorithme 3

```

GradConj( $A, b, x_0, tol$ )
 $x \leftarrow x_0, r \leftarrow b - Ax, d \leftarrow r, s \leftarrow \| r \|^2$ 
Tant que  $s > tol$  faire
     $\rho \leftarrow s / \langle d, Ad \rangle$ 
     $x \leftarrow x + \rho \cdot d$ 
     $r \leftarrow b - Ax$ 
     $\mu \leftarrow \| r \|^2$ 
     $d \leftarrow r + \frac{\mu}{s} d$ 
     $s \leftarrow \mu$ 

```

**Théorème 11** (Convergence). *Pour tout  $k$ ,  $x_{k+1}$  calculé par l'algorithme du gradient conjugué est l'unique minimiseur de  $f$  sur le sous-espace affine  $E_k = x_0 + \text{vect}\{d_0, \dots, d_k\}$  de  $\mathbb{R}^n$ .*

*Démonstration.* (Par récurrence.) C'est évident pour  $k = 0$  car  $x_1$  est l'unique minimiseur de  $f$  sur  $x_0 + td_0$ ,  $t \in \mathbb{R}$ .

Supposons que  $x_k$  est l'unique minimiseur de  $f$  sur  $E_{k-1}$ . Soit  $x \in E_k$ , alors  $x$  s'écrit  $x = y + \rho d_k$  avec  $y \in E_{k-1}$ . Comme  $f$  est quadratique on a :

$$\begin{aligned} f(x) &= f(y) + \rho \langle \nabla f(y)^T, d_k \rangle + \frac{\rho^2}{2} \langle d_k, \nabla^2 f(y) d_k \rangle \\ &= f(y) + \rho \langle Ay - b, d_k \rangle + \frac{\rho^2}{2} \langle d_k, Ad_k \rangle. \end{aligned}$$

Comme  $y - x_k \in \text{vect}\{d_0, \dots, d_{k-1}\}$  et  $d_k$  est  $A$ -conjuguée à tout direction de  $\text{vect}\{d_0, \dots, d_{k-1}\}$  (*exercice 1*), on a

$$\langle \nabla f(y)^T - \nabla f(x_k)^T, d_k \rangle = \langle A(y - x_k), d_k \rangle = 0.$$

On en déduit que

$$\begin{aligned} \forall x \in E_k : f(x) &= f(y) + \rho \langle \nabla f(x_k)^T, d_k \rangle + \frac{\rho^2}{2} \langle d_k, Ad_k \rangle \\ &\geq f(x_k) + \rho \langle Ax_k - b, d_k \rangle + \frac{\rho^2}{2} \langle d_k, Ad_k \rangle = f(x_k + \rho d_k). \end{aligned}$$

Donc l'unique minimiseur de  $f$  sur  $E_k$  est le point  $x_{k+1} = x_k + \rho_k d_k$  avec

$$\rho_k = \text{ArgMin}\{f(x_k + \rho d_k)\}. \quad \blacksquare$$

**Corollaire 2** *L'algorithme du gradient conjugué converge en au plus  $n$  itérations.*

#### 4.1.2 Méthodes Newtonniennes

La direction de Newton dans un point  $x$  s'obtient en résolvant le problème de minimisation qui consiste à remplacer  $f$  par son approximation de Taylor d'ordre 2 :

$$f(x + d) \approx f(x) + \nabla f(x) \cdot d + \frac{1}{2} \langle d, \nabla^2 f(x) \cdot d \rangle = g(d).$$

On résoud donc le problème de minimisation  $\min_{d \in \mathbb{R}^n} g(d)$ . Si  $\nabla^2 f(x)$  est définie positive, la solution est donnée par

$$\nabla g(d) = 0 \Leftrightarrow \nabla f(x)^T + \nabla^2 f(x) \cdot d = 0 \Leftrightarrow d = -[\nabla^2 f(x)]^{-1} \nabla f(x)^T.$$

#### Algorithme 4

```

Newton( $f, x_0, tol$ )
 $x \leftarrow x_0$ ,  $H \leftarrow \nabla^2 f(x)$ ,  $d \leftarrow -H^{-1} \nabla f(x)^T$ 
Tant que  $|\langle d, \nabla f(x)^T \rangle| > tol$  faire
     $\rho \leftarrow \text{ArgMin}\{f(x + td), t > 0\}$ 
     $x \leftarrow x + \rho \cdot d$ 
     $H \leftarrow \nabla^2 f(x)$ 
     $d \leftarrow -H^{-1} \nabla f(x)^T$ 

```

**Théorème 12** (Convergence). *Si  $f$  est coercive sur  $\mathbb{R}^n$  et  $\nabla^2 f(x)$  est définie positive  $\forall x \in \mathbb{R}^n$ , alors l'algo converge pour toute initialisation  $x_0 \in \mathbb{R}^n$ . La convergence est au moins quadratique.*

### 1. Méthode de Newton à pas fixe

C'est l'algorithme précédent avec  $\rho = 1$ .

**Théorème 13** (Convergence). *Sous les hypothèses suivantes*

(H0) *l'ensemble de niveau  $S_0 = \{x \in \Omega : f(x) \leq f(x_0)\}$  est fermé dans  $\mathbb{R}^n$ ;*

(H1)  *$f$  est 2 fois dérivable en  $x$  et  $\exists c > 0$  constante tq*

$$\forall x \in S_0, \forall u \in \mathbb{R}^n : u^T \nabla^2 f(x) u \geq c \|u\|^2 ;$$

(H2)  *$\nabla^2 f(x)$  est lipschitzienne de constante  $L$  sur  $S_0$ ;*

(H3)  *$\|\nabla^2 f(x)^{-1} \nabla f(x)^T\| < 2c/L$ ;*

*la suite des itérés  $\{x_k\}_{k \geq 0}$  calculés par l'algorithme de Newton à pas fixe converge vers un minimiseur local  $x^*$  de  $f$ . La convergence est au moins quadratique et*

$$\|x_{k+1} - x^*\| \leq \frac{L}{2c} \|x_k - x^*\|^2 .$$

Ce théorème garantit la convergence dès qu'un des points calculés par l'algo est suffisamment proche d'un minimiseur local.

### 2. Méthode de type Quasi-Newton

Le calcul de la hessienne et de son inverse à chaque étape rendent la méthode de Newton trop coûteuse. Alors l'idée est de remplacer l'inverse par une approximation actualisée à chaque étape. On cherche une matrice  $B_k$  symétrique définie positive pour approcher l'inverse de l'hessienne vérifiant :

$$x_k - x_{k-1} = B_k (\nabla f(x_k)^T - \nabla f(x_{k-1})^T) .$$

Soient

$$\begin{cases} w_k = \nabla f(x_k)^T - \nabla f(x_{k-1})^T \\ v_k = x_k - x_{k-1} \end{cases}$$

alors

a) l'algorithme de Davidon-Fletcher-Powell (DFP) consiste à choisir

$$B_k = B_{k-1} - \frac{B_{k-1} w_k w_k^T B_{k-1}}{w_k^T B_{k-1} w_k} + \frac{v_k v_k^T}{v_k^T w_k}$$

b) l'algorithme de Broyden-Fletcher-Goldfarb-Shanno (BFGS) consiste à choisir

$$B_k = B_{k-1} + \frac{1}{w_k^T v_k} \left( 1 + \frac{w_k^T B_{k-1} w_k}{w_k^T v_k} \right) v_k v_k^T - \frac{B_{k-1} w_k v_k^T + v_k w_k^T B_{k-1}}{w_k^T v_k}$$

## 4.2 Choix du pas

Une fois fixée la direction de déplacement on doit fixer le pas. Pour certains algorithmes ce pas est fixé (Newton, gradient à pas fixé), pour d'autres il faut par une procédure de *recherche linéaire*. Considérons différentes procédures.

### 4.2.1 Pas optimal

$$\rho = \text{ArgMin}_{t \geq 0} f(x + td)$$

nécessite une procédure de minimisation unidirectionnelle.

**Définition 6**  $\varphi : [0, T] \rightarrow \mathbb{R}$  est unimodale sur  $[0, T]$  si elle admet un minimum strict  $\rho \in ]0, T[$ , décroît strictement de 0 à  $\rho$  et croît strictement de  $\rho$  à  $T$ .

Dans le cas d'une fonction unimodale l'algorithme de la section dorée calcule une valeur approchée du minimum  $\rho$  de  $\varphi$  sur  $[0, T]$  à une précision  $tol$  près.

### Algorithme 5

```

SectionDorée( $\varphi, T, tol$ )
 $\alpha \leftarrow \frac{\sqrt{5}-1}{2}$   (inverse du nombre d'or)
 $a \leftarrow 0, b \leftarrow T, c \leftarrow \alpha a + (1 - \alpha)b$ 
 $d \leftarrow a + b - c$   (symétrique de c par rapport au milieu de [a, b])
Tant que  $|b - a| > 2 \cdot tol$  faire
    si  $\varphi(c) < \varphi(d)$  alors
         $b \leftarrow d, d \leftarrow c, c \leftarrow a + b - d$ 
    sinon
         $a \leftarrow c, c \leftarrow d, d \leftarrow a + b - c$ 
retourner  $(a + b)/2$ 

```

**Explication :** Supposons que

$$c_k - a_k = (1 - \alpha)(b_k - a_k) = b_k - d_k$$

et que  $a_{k+1} = a_k, b_{k+1} = d_k, d_{k+1} = c_k$  et  $c_{k+1} = a_{k+1} + b_{k+1} - d_{k+1}$ . Alors

$$b_{k+1} - a_{k+1} = d_k - a_k = b_k - (1 - \alpha)(b_k - a_k) - a_k = \alpha(b_k - a_k)$$

et

$$b_{k+1} - d_{k+1} = d_k - c_k = (d_k - b_k) + b_k + (a_k - c_k) - a_k = (2\alpha - 1)(b_k - a_k);$$

donc

$$\frac{b_{k+1} - d_{k+1}}{b_{k+1} - a_{k+1}} = \frac{2\alpha - 1}{\alpha}.$$

Or  $\alpha$  est l'inverse du nombre d'or et donc vérifie l'équation  $x^2 - x - 1 = 0$  :

$$\frac{1}{\alpha^2} - \frac{1}{\alpha} - 1 = 0 \Leftrightarrow 1 - \alpha - \alpha^2 = 0 \Leftrightarrow 2\alpha - 1 = \alpha(1 - \alpha).$$

On en déduit que

$$\frac{b_{k+1} - d_{k+1}}{b_{k+1} - a_{k+1}} = 1 - \alpha,$$

donc

$$c_{k+1} - a_{k+1} = (1 - \alpha)(b_{k+1} - a_{k+1}) = b_{k+1} - d_{k+1}.$$

De plus puisque  $b_{k+1} - a_{k+1} = \alpha(b_k - a_k)$  avec  $0 < \alpha < 1$ , l'intervalle est réduit à chaque étape.

**Remarque :** On espère que la fonction  $\varphi(t) = f(x + td)$  est unimodale sur un intervalle  $[0, T]$ , où on peut déterminer  $T$  par l'algorithme

**Algorithme 6**

```

T ← 1
Tant que  $\varphi(T) < \varphi(0)$  faire
  T ← 2T

```

#### 4.2.2 Recherche linéaire

Si l'évaluation de la fonction objectif est trop coûteuse on procède autrement : on essaie de trouver un pas  $\rho$  qui fait suffisamment décroître  $f$  pour garantir la convergence. Le principe commun à ces méthodes est

- (a)  $\rho$  ne doit pas être choisi trop grand sinon l'algo risque d'avoir un comportement oscillatoire;
- (b)  $\rho$  ne doit pas être choisi trop petit sinon l'algo risque de converger prématurément.

##### 1. Règle d'Armijo

On se donne  $m \in ]0, 1[$  et  $M > 1$ , et on choisit  $\rho$  tq

- (i)  $f(x + \rho d) < f(x) + m\rho \nabla f(x)d$  (condition d'Armijo)
- (ii)  $f(x + M\rho d) > f(x) + mM\rho \nabla f(x)d$

La condition (i) est vérifiée pour  $\rho$  suffisamment petit :

$$\begin{aligned} f(x + \rho d) &= f(x) + \rho \nabla f(x)d + \mathcal{O}(\rho^2) \\ &= f(x) + \rho [\nabla f(x)d + \mathcal{O}(\rho)] \\ &< f(x) + \rho [m \nabla f(x)d] \quad \forall \rho \in ]0, T[. \end{aligned}$$

Pour que la condition (ii) est vérifiée aussi, il faut que  $M\rho > T$ , donc  $\rho \in ]\frac{T}{M}, T[$ .

Pour trouver  $\rho$  vérifiant la condition d'Armijo on peut appliquer la procédure suivante en supposant que  $m < \frac{1}{2}$  :

## Algorithme 7

```

Armijo( $f, x, d, m$ )
 $t \leftarrow 1$ 
Tant que  $f(x + td) \geq f(x) + mt\nabla f(x)d$  faire
     $t \leftarrow \frac{-t^2\nabla f(x)d}{2[f(x+td)-f(x)-t\nabla f(x)d]}$ 
retourner  $x + td$ 

```

**Explication :** Si  $f(x + td) < f(x) + mt\nabla f(x)d$  alors on s'arrête. Sinon on interpole la fonction  $\varphi(\rho) = f(x + \rho d)$  par la parabole  $p(\rho)$  passant aux points  $(0, \varphi(0))$  et  $(t, \varphi(t))$  et dont la tangente au point  $(0, \varphi(0))$  à une pente  $\varphi'(0)$  :

$$p(\rho) = \varphi(0) + \rho\varphi'(0) + \rho^2 \frac{\varphi(t) - \varphi(0) - t\varphi'(0)}{t^2}.$$

Puis on détermine l'abscisse correspondant au minimum de cette parabole :

$$2\rho \frac{\varphi(t) - \varphi(0) - t\varphi'(0)}{t^2} + \varphi'(0) = 0 \Leftrightarrow \rho = \frac{-t^2\varphi'(0)}{2[\varphi(t) - \varphi(0) - t\varphi'(0)]}.$$

Or

$$\varphi(t) - \varphi(0) \geq tm\varphi'(0) \Leftrightarrow \varphi(t) - \varphi(0) - t\varphi'(0) \geq t(m-1)\varphi'(0).$$

Par conséquent

$$\rho \leq \frac{-t^2\varphi'(0)}{t(m-1)\varphi'(0)} = \frac{t}{2(1-m)}.$$

De plus  $\frac{t}{2(1-m)} < t$  si  $m < 1/2$  et on réduit donc la taille de l'intervalle.

## 2. Règle de Goldstein

Fixer  $m_1 \in ]0, 1[$  et  $m_2 \in ]m_1, 1[$ . Choisir  $\rho$  vérifiant

- (i)  $f(x + \rho d) \leq f(x) + m_1\rho\nabla f(x)d$
- (ii)  $f(x + \rho d) \geq f(x) + m_2\rho\nabla f(x)d$

## Algorithme 8

```

Goldstein( $f, x, d, m_1, m_2, \alpha_m, \alpha_M$ )
 $test \leftarrow faux$ 
Tant que ( $test$  est faux) faire
     $\rho \leftarrow (\alpha_m + \alpha_M)/2$ 
    calculer  $g = f(x + \rho d)$ 
    si  $g \leq f(x) + m_1\rho\nabla f(x)d$  alors
        si  $g \geq f(x) + m_2\rho\nabla f(x)d$  alors
             $test \leftarrow vrai$ 
        sinon
             $\alpha_m \leftarrow \rho$ 
    sinon
         $\alpha_M \leftarrow \rho$ 

```

### 4.3 Choix de la direction en présence de contraintes

Pour pouvoir se déplacer dans la direction  $d$  il faut qu'elle soit admissible (c-à-d.  $\exists T > 0$  tq  $\forall t \in [0, T] : x + td \in \mathcal{A}$ ) ce qui peut ne pas être le cas avec les choix précédents. On peut alors procéder de différentes façons :

#### 4.3.1 Méthode de pénalisation

On se débarrasse des contraintes en les passant dans la fonction objectif et pénalisant fortement les points qui ne sont pas admissibles :

$$(P_\mu) \quad \min f(x) + \frac{1}{\mu} \sum_{i=1}^p g_i^2(x) + \frac{1}{\mu} \sum_{j=1}^q \{\max(0, h_j(x))\}^2$$

avec  $\mu$  un paramètre “petit”. Sous des bonnes hypothèses on peut montrer que lorsque  $\mu \rightarrow 0+$  les solutions de ces problèmes convergent vers la solution de  $(P)$ .

#### 4.3.2 Construction d’une direction admissible

Si  $f$  est convexe et les contraintes sont linéaires on peut trouver une direction de descente admissible en procédant de la façon suivante :

1. On remplace la fonction objectif par une approximation linéaire au point  $x$

$$L(y) = f(x) + \nabla f(x) \cdot (y - x).$$

2. On résout le problème de programmation linéaire suivant

$$(P_L(x)) \quad \begin{array}{ll} \min & \nabla f(x)y \\ \text{s.c.} & g(y) = 0 \\ & h(y) \leq 0 \end{array} \quad (\text{méthode du simplexe}).$$

3. Si  $y$  est une solution de ce problème auxiliaire alors

- $d = y - x$  est une direction admissible car

$$\begin{aligned} g(x + td) &= (1 - t)g(x) + tg(y) = 0 \quad \forall t > 0 \\ h(x + td) &= (1 - t)h(x) + th(y) \leq 0 \quad \forall t \in [0, 1] \end{aligned}$$

- $d = y - x$  est une direction de descente car

$$\nabla f(x)d = \nabla f(x)(y - x) = \nabla f(x)y - \nabla f(x)x < 0.$$

## 4.4 Exercices

Exercice 1 Soit  $\{x_k\}_{k \geq 0}$  la suite définie par

$$\begin{cases} x_0 \text{ donné} \\ x_{k+1} = x_k + \rho_k d_k, \quad k \geq 0 \end{cases}$$

avec

$$\rho_k = \frac{\langle r_k, d_k \rangle}{\langle Ad_k, d_k \rangle}, \quad r_k = b - Ax_k$$

et  $A \in \mathcal{M}_n(\mathbb{R})$  symétrique définie positive.

1. Montrer que  $\forall k \geq 0 : \langle r_{k+1}, d_k \rangle = 0$ .

Supposons dans la suite que la suite des directions  $\{d_k\}_{k \geq 0}$  est définie par

$$\begin{cases} d_0 \text{ donné} \\ d_{k+1} = r_{k+1} + \alpha_k d_k, \quad k \geq 0 \end{cases}$$

avec

$$\alpha_k = -\frac{\langle r_{k+1}, Ad_k \rangle}{\langle Ad_k, d_k \rangle}.$$

2. Montrer que  $\forall k \geq 0 :$

(a)  $\rho_k = \|r_k\|^2 / \langle d_k, Ad_k \rangle$

(b)  $\langle r_{k+1}, r_k \rangle = 0$

(c)  $\alpha_k = \|r_{k+1}\|^2 / \|r_k\|^2$

3. Montrer que si  $d_0 = r_0$ , alors  $\forall k > 0$ ,  $d_k$  est  $A$ -conjuguée à tout direction de vect  $\{d_0, \dots, d_{k-1}\}$ .

## Exercices supplémentaires

Exercice 2 On veut minimiser la fonction

$$f(x_1, x_2) = 2x_1^2 + x_2^2 - 2x_1x_2 + 2x_1^3$$

sur  $\mathbb{R}^2$  par une méthode itérative. On part du point  $x = (2, -1)^T$  dans la direction  $d = (-1, 4)^T$ .

1. Cette direction est-elle de pente? Justifier.
2. Quel est le nouveau point que l'on obtient si l'on fait une minimisation exacte dans cette direction?
3. Supposons que l'on veut faire une minimisation approchée (recherche linéaire) en utilisant le critère d'Armijo avec paramètres  $m = 0.5$  et  $M = 2$ . Quelles conditions doit vérifier la valeur  $\rho$  définissant le nouvel itéré  $x + \rho d$ ? Vérifier si  $\rho = 1$  satisfait ces conditions.



4. Dire comment devrait-on procéder si on voulait faire une itération de la méthode du gradient à pas optimal.

Exercice 3 On veut minimiser la fonction

$$f(x_1, x_2) = 5x_1^2 + x_2^2 - 2x_1 - x_2 - 12$$

sur  $\mathbb{R}^2$  par une méthode itérative de gradient.

1. Expliciter l'algorithme du gradient à pas constant et égal à 0.25 appliqué à cette fonction avec point initial  $x_0 = (1, 2)^T$ .
2. Quel serait la valeur optimale du pas pour cette fonction?
3. Supposons maintenant que l'on choisit le pas de façon à minimiser la fonction dans la direction de déplacement.
  - (a) Calculer explicitement le pas à chaque itération.
  - (b) Quelle est la propriété vérifiée par deux directions consécutives? Justifier.
  - (c) Quelle méthode obtient-on?
4. En quoi consistent les méthodes de recherche linéaire? Pourquoi on les utilise?

Exercice 4 On veut déterminer le minimum de la fonction  $f(x_1, x_2) = (x_1 + x_2^2)^2$  sur  $\mathbb{R}^2$  en utilisant une méthode itérative. On part du point  $\tilde{x} = (1, 0)^T$  dans la direction  $\tilde{d} = (-1, 1)^T$ .

1. Montrer que la direction  $\tilde{d}$  est une direction de descente.
2. Calculer l'itéré suivant par  $\min_{\alpha \geq 0} f(\tilde{x} + \alpha \tilde{d})$
3. Par quoi peut-on remplacer le choix précédent de  $\alpha$ ?
4. Si on applique la méthode de Newton, quel sera le nouvel itéré?

Exercice 5 On veut déterminer le minimum de la fonction  $f(x_1, x_2) = x_1^4 + 3x_1^2x_2 - x_1 + 3x_2^2$  sur  $\mathbb{R}^2$  en utilisant une méthode itérative. On part du point  $x^{(0)} = (1/2, 1)^T$  dans la direction  $d_0 = (0, -1)^T$ .

1. Montrer que la direction  $d_0$  est une direction de descente.
2. Quel est le nouveau point  $x^{(1)}$  que l'on obtient si l'on fait une minimisation exacte dans cette direction?
3. Supposons que l'on veut faire une minimisation approchée (recherche linéaire) en utilisant la règle de Goldstein avec paramètres  $m_1 = 1/3$  et  $m_2 = 1/2$ . Quelles conditions doit vérifier la valeur  $\rho_0$  définissant le nouvel itéré  $x^{(1)}$ ? Montrer que  $\rho_0 = 5/4$  satisfait ces conditions.

4. Si on applique la méthode de Newton à pas fixe, quel sera le nouvel itéré  $x^{(1)}$  ?

Exercice 6 On applique la méthode de gradient à pas optimal à la minimisation en  $\mathbb{R}^2$  de la fonction

$$f(x_1, x_2) = x_1^2 + 2x_2^2$$

au départ du point  $y_0 = (2, 1)^T$ .

1. Montrer que le premier itéré que l'on obtient est  $y_1 = (2/3, -1/3)^T$ .
2. Montrer que la suite obtenue vérifie  $y_k = (2/3^k, (-1)^k/3^k)^T$ .
3. Montrer que  $f(y_{k+1}) = f(y_k)/9$ . En déduire la valeur optimale.

Exercice 7 On veut appliquer la méthode de gradient conjugué à la minimisation en  $\mathbb{R}^3$  de la fonction

$$f(x) = x_1^2 + x_2^2 + x_3^2 + x_1x_3 - x_2.$$

1. Vérifier que  $f(x) = \frac{1}{2}x^T Ax - b^T x$  avec

$$A = \begin{pmatrix} 2 & 0 & 1 \\ 0 & 2 & 0 \\ 1 & 0 & 2 \end{pmatrix}, \quad \text{et} \quad b = (0, 1, 0)^T.$$

2. Vérifier que  $\lambda$  est une valeur propre de  $A$  ssi  $\lambda \in \{1, 2, 3\}$ . En déduire que  $A$  est définie positive.
3. Soit  $x^{(0)} = (0, 0, 0)^T$ . Calculer  $\{x^{(k)}\}_{k \geq 0}$  en utilisant la méthode de gradient conjugué et arrêter les itérations dès que la suite a convergé. Quelle est la solution optimale et la valeur optimale ?