

# Chapitre 8

## Optimisation sans contrainte

### 8.1 Introduction

On appelle problème d'optimisation un problème noté :

$$\mathcal{P} : \min_{x \in \mathcal{C}} f(x).$$

La fonction  $f$  est appelée fonction objectif et l'ensemble  $\mathcal{C}$  est l'ensemble des contraintes. Nous nous limitons dans ce cours au cas où  $\mathcal{C}$  est un sous-ensemble de  $\mathbb{R}^n$ .

**Exercice 8.1** *Différence entre dimension infinie et dimension finie sur un exemple. Soit*

$$\mathcal{P}_n : \min_{x \in \mathcal{C}_n \subset \mathbb{R}^n} f(x) = x^T x, \text{ où } \mathcal{C}_n = \{x \in \mathbb{R}^n, x_1 = \frac{1}{2}, \text{ et } \|x\|_2 \leq 1\}.$$

Soit

$$\mathcal{P}_\infty : \min_{x \in \mathcal{C}_\infty} f(x) = \int_0^1 x^2, \text{ où } \mathcal{C}_\infty = \{x, x \text{ est continue et } x(0) = \frac{1}{2}, \text{ et } \int_0^1 x^2 \leq 1\}.$$

Étudiez l'ensemble des solutions de  $\mathcal{P}_n$  et  $\mathcal{P}_\infty$ .

**Preuve 8.1** *Démonstration* : Le vecteur  $\bar{x} = (\frac{1}{2}, 0 \dots 0)^T$  est solution de  $\mathcal{P}_n$ . Pour tout  $x$ , on a  $f(x) > 0$ . Si

$$x_n(t) = \begin{cases} 0 & \text{si } t \in ]\frac{1}{n}, 1] \\ \frac{-n}{2}(t - \frac{1}{n}) & \text{si } t \in [0, \frac{1}{n}] \end{cases},$$

on a alors  $\lim_{n \rightarrow +\infty} f(x_n) = 0$ , mais il n'existe pas de fonction continue non nulle pour laquelle  $\int_0^1 x^2 = 0$ . Donc  $\mathcal{P}_\infty$  n'admet pas de solution.

□

Il arrive que l'on s'intéresse à l'existence, l'unicité et au calcul de points  $\bar{x}$  qui minimisent  $f$  sur  $\mathcal{C}$ , c'est à dire tels que

$$f(\bar{x}) \leq f(x) \text{ pour tout } x \text{ dans } \mathcal{C}.$$

c'est un problème d'optimisation globale. Dans les cas généraux, nous verrons qu'il est parfois possible de donner des conditions nécessaires, ou des conditions suffisantes d'optimalité, ou même quelquefois des conditions à la fois nécessaires et suffisantes. Les algorithmes rechercheront des points qui vérifient ces conditions. Le problème de moindres carrés linéaires vu précédemment est un exemple de problème d'optimisation.

**Exercice 8.2** Un fabricant de composants électroniques possède deux types de fabriques : A et B, notées  $A_i$ ,  $1 \leq i \leq m$  et  $B_j$ ,  $1 \leq j \leq n$ . Lors de la fabrication, chacun de ces composants doit tout d'abord passer par une des usines de type A puis par une de type B. Comme ces usines ne se trouvent pas dans le même lieu géographique, le fabricant doit étudier le meilleur moyen pour transporter ces composants à moindre coût des usines  $A_i$  vers les usines  $B_j$ . Connaissant la matrice des coûts  $C = [c_{ij}]$  où  $c_{ij}$  correspond au coût de transport d'une pièce de l'usine  $A_i$  vers l'usine  $B_j$ , ainsi que le nombre de pièces  $a_i$  produites par l'usine  $A_i$  et le nombre de pièces  $b_j$  que l'usine  $B_j$  doit recevoir, formuler le plan de transport optimal (en terme de coût de transport) sous la forme d'un problème d'optimisation. Données  $m = 2$ ,  $n = 3$ ,  $[a_1, a_2] = [10, 20]$ ,  $[b_1, b_2, b_3] = [5, 10, 15]$  et

$$C = \begin{pmatrix} 2 & 8 & 7 \\ 3 & 4 & 5 \end{pmatrix}$$

**Preuve 8.2** *Démonstration* : Soient les variables de décision suivantes :  $x_{ij}$  nombre de pièces allant de l'usine  $A_i$  vers l'usine  $B_j$  avec  $1 \leq i \leq 2$  et  $1 \leq j \leq 3$ . Le problème d'optimisation s'écrit : Minimiser  $z = 2x_{11} + 8x_{12} + 7x_{13} + 3x_{21} + 4x_{22} + 5x_{23}$  sous les contraintes

$$\begin{cases} x_{11} + x_{12} + x_{13} & = 10 \\ x_{21} + x_{22} + x_{23} & = 20 \\ x_{11} + x_{21} & = 5 \\ x_{12} + x_{22} & = 10 \\ x_{13} + x_{23} & = 15 \\ x_{11}; x_{12}; x_{13}; x_{21}; x_{22}; x_{23} & \geq 0 \end{cases}$$

□

**Exercice 8.3** *Principe de Fermat*. Soient  $a, b, c$  trois réels positifs. On suppose que l'on a deux milieux  $M_1 = \{(x, y), y > 0\}$  et  $M_2 = \{(x, y), y < 0\}$  et que la vitesse de propagation d'un rayon lumineux est  $c_i$  dans  $M_i$ . On considère que le rayon se propage en ligne droite dans chaque milieu et que le rayon suit un trajet de temps global de parcours minimum (principe de Fermat). Formuler le problème de la recherche du trajet entre  $A(0, a)$  et  $B(c, -b)$  sous la forme d'un problème d'optimisation. En utilisant une étude de fonction, montrez que le principe de Fermat se traduit par la loi de Snell

$$\frac{\sin \alpha_1}{c_1} = \frac{\sin \alpha_2}{c_2}.$$

**Preuve 8.3** *Démonstration* : Soit  $X(x_0, 0)$  le point où le rayon change de milieu. Le temps de trajet est

$$T(x) = \frac{1}{c_1} \sqrt{a^2 + x^2} + \frac{1}{c_2} \sqrt{(c-x)^2 + b^2},$$

et  $x_0$  minimise  $T(x)$ . Les minima d'une fonction réelle vérifient  $T'(x) = 0$ , ce qui donne

$$\frac{x}{c_1\sqrt{a^2 + x^2}} = \frac{c - x}{c_2\sqrt{(c - x)^2 + b^2}},$$

ce qui donne bien la loi de Snell puisque  $\sin(\alpha_1) = \frac{x}{\sqrt{a^2 + x^2}}$  et  $\sin(\alpha_2) = \frac{c - x}{\sqrt{(c - x)^2 + b^2}}$ . □

**Exercice 8.4** *Mission : désamorcer une bombe nucléaire sur un yacht. Yacht amarré à 50 mètres du rivage. James Bond se trouve à 100 mètres du point le plus proche du yacht sur la plage. Vitesses : course 18km/h, nage : 10km/h. Temps de désamorçage : 30 secondes. Explosion dans 65 secondes. Formaliser la faisabilité de cette mission sous la forme d'un problème d'optimisation.*

**Preuve 8.4** *Démonstration* : Le temps de parcours du héros est  $f(x) = \frac{x}{5} + 0.36\sqrt{50^2 + (100 - x)^2}$ . On veut donc que  $f(x)$  soit inférieur à  $65 - 30 = 35$  secondes, ce qui conduit au problème

$$\min f(x),$$

sous les contraintes

$$x \geq 0 \text{ et } x \leq 100.$$

Note :  $f(0) = 40.25$ ,  $f(100) = 38$  mais  $f(66) = 34.96$ . □

**Definition 8.1** Une fonction  $f$  est semicontinue inférieurement sur  $\mathbb{R}^n$  ssi

pour tout  $\alpha \in \mathbb{R}$ , l'ensemble  $\{x, f(x) \leq \alpha\}$  est fermé .

Une fonction continue est semicontinue inférieurement.

**Exercice 8.5** Nous supposons que  $\mathcal{C}$  est fermé et qu'il existe un point de  $\mathcal{C}$  en lequel  $f$  est finie. Supposons de plus que  $f$  est semicontinue inférieurement sur  $\mathcal{C}$ , et  $f$  est coercive ( $\lim_{\substack{\|x\| \rightarrow +\infty \\ x \in \mathcal{C}}} f(x) = +\infty$ ). La fonction  $f$  admet un minimum sur  $\mathcal{C}$ .

**Preuve 8.5** *Démonstration* : Faisons la démonstration dans le cas où  $f$  est continue. Soit  $x_0 \in \mathcal{C}$  en lequel  $f$  est finie. Une conséquence de la coercivité est que il existe  $\alpha$  tel que  $\|x\| > \alpha$  entraîne  $f(x) \geq f(x_0)$ . Alors le problème d'optimisation revient à la minimisation de la fonction continue  $f$  sur le compact  $K = \{x \in \mathcal{C}, \|x\| \leq \alpha\}$ . Comme l'image continue d'un compact est un compact,  $f(K)$  est un compact en dimension finie, donc c'est un fermé borné. Donc le réel  $\inf_{x \in K} f(x)$  qui appartient à l'adhérence de  $f(K)$  appartient à  $f(K)$ , ainsi il existe  $x^* \in K$  tel que  $f(x^*) = \inf_{x \in K} f(x) \leq f(x)$  pour tout  $x \in \mathbb{R}^n$ . □

Le résultat ci-dessus peut-être utilisé pour montrer que le problème d'optimisation de l'exemple 8.3 admet au moins une solution.

**Definition 8.2** Une partie  $\mathcal{C}$  est dite convexe ssi pour tout  $(x, y) \in \mathcal{C}^2$ , et pour tout  $\alpha \in [0, 1]$ ,  $\alpha x + (1 - \alpha)y \in \mathcal{C}$ . Une fonction  $f$  définie sur une partie  $\mathcal{C}$  convexe est une fonction convexe ssi  $(x, y) \in \mathcal{C}^2$ , et pour tout  $\alpha \in [0, 1]$  on a  $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$ . Une fonction  $f$  définie sur une partie  $\mathcal{C}$  convexe est une fonction strictement convexe ssi  $(x, y) \in \mathcal{C}^2$ ,  $x \neq y$ , et pour tout  $\alpha \in ]0, 1[$  on a  $f(\alpha x + (1 - \alpha)y) < \alpha f(x) + (1 - \alpha)f(y)$ .

**Exercice 8.6** Si  $\mathcal{C}$  est convexe et si  $f$  est strictement convexe sur  $\mathcal{C}$ , alors  $f$  admet au plus un minimum sur  $\mathcal{C}$ .

**Preuve 8.6** *Démonstration* : Supposons qu'il existe deux minima  $x_0$  et  $x_1$  de  $f$  dans  $\mathcal{C}$  (i.e.  $f(x_0) = f(x_1) \leq f(x)$  pour tout  $x \in \mathcal{C}$ ). D'après la stricte convexité de  $f$  sur  $\mathcal{C}$ , on a

$$f\left(\frac{x_0}{2} + \frac{x_1}{2}\right) < \frac{1}{2}f(x_0) + \frac{1}{2}f(x_1) = f(x_0),$$

ce qui est impossible d'après la définition même du minimum.

□

## 8.2 Rudiments en calcul différentiel

**Definition 8.3** Une fonction  $f$  définie sur un ouvert  $\mathcal{O} \subset \mathbb{R}^n$  est dite dérivable (au sens de Fréchet) en  $x$  ssi il existe un vecteur ligne  $f'(x)$  tel que

$$f(x + h) = f(x) + f'(x)h + o(h),$$

où l'on a posé  $o(h) = \|h\|\epsilon(h)$ , avec  $\lim_{\|h\| \rightarrow 0} \epsilon(h) = 0$ . Le vecteur colonne  $f'(x)^T$  s'appelle gradient de  $f$  en  $x$  et est noté  $\nabla f(x)$ . Notez que cette notion généralise la notion de dérivabilité d'une fonction de  $\mathbb{R}$  dans  $\mathbb{R}$  et que  $f'(x)$  ne dépend pas de la norme considérée.

**Exercice 8.7** Montrez que si  $f$  est dérivable en  $x$ , alors

1.  $f$  est continue en  $x$ ,
2.  $f$  admet des dérivées partielles en  $x$  et  $f'(x) = \left[\frac{\partial f(x)}{\partial x_1}, \dots, \frac{\partial f(x)}{\partial x_n}\right] \in \mathbb{R}^{1 \times n}$ .

**Preuve 8.7** *Démonstration* : Par définition de la différentiabilité en  $x$ ,  $f(x + h) = f(x) + f'(x)h + o(h)$  donc  $\lim_{h \rightarrow 0} f(x + h) = \lim_{h \rightarrow 0} f(x) + f'(x)h + o(h) = f(x)$ , ce qui est bien la définition de la continuité. En ce qui concerne les dérivées partielles, posons  $h = \delta e_i$ , où  $\delta \neq 0$  et  $e_i$  est le  $i$ -ème vecteur de la base canonique de  $\mathbb{R}^n$ . On a alors en considérant la norme Euclidienne,  $f(x + \delta e_i) - f(x) = \delta \cdot f'(x)e_i + |\delta| \cdot \epsilon(\delta e_i)$ , c'est à dire, pour  $\delta \neq 0$ ,

$$\lim_{\delta \rightarrow 0} \frac{f(x + \delta \cdot e_i) - f(x)}{\delta} = \lim_{\delta \rightarrow 0} f'(x)e_i + \frac{|\delta|}{\delta} \epsilon(\delta e_i) = f'(x)e_i.$$

On obtient donc  $\frac{\partial f(x)}{\partial x_i} = f'(x)e_i$ .

□

**Exercice 8.8** On considère la fonction quadratique définie sur  $\mathbb{R}^n$  par  $f(x) = \frac{1}{2}x^T Ax - x^T b$ , où  $A$  est carrée et symétrique montrez que  $\nabla f(x) = Ax - b$ .

**Preuve 8.8** *Démonstration* : On a

$$\begin{aligned} f(x+h) &= \frac{1}{2}(x+h)^T A(x+h) - (x+h)^T b \\ &= \frac{1}{2}x^T Ax + \frac{1}{2}h^T Ah + \frac{1}{2}x^T Ah + \frac{1}{2}h^T Ax - (x+h)^T b \\ &= f(x) + (Ax - b)^T h + \frac{1}{2}h^T Ah. \end{aligned}$$

De plus,  $0 \leq \frac{|h^T Ah|}{\|h\|_2} \leq \frac{\|h\|_2^2 \|A\|_2}{\|h\|_2} = \|A\|_2 \|h\|_2$ , donc  $\lim_{\|h\|_2 \rightarrow 0} \frac{|h^T Ah|}{\|h\|_2} = 0$ , ce qui montre que  $h^T Ah = o(h)$ .

□

**Définition 8.4** Une fonction  $f$  est dite deux fois dérivable si chaque dérivée partielle  $\frac{\partial f(x)}{\partial x_i}$  est dérivable. Une fonction est  $k$  fois dérivable si elle est  $k-1$  fois dérivable et si les dérivées partielles à l'ordre  $k-1$  sont dérivables.

**Exercice 8.9** (Dérivation d'une composée) Soit  $f$ , définie sur un ouvert  $\mathcal{O} \subset \mathbb{R}^n$ , dérivable en tout  $x \in \mathcal{O}$ . Soit  $d \in \mathbb{R}^n$ . On définit localement en  $x$  la fonction de la variable réelle  $t$  par  $\phi : t \mapsto \phi(t) = f(x+td)$ . Montrez que  $\phi$  est dérivable en 0 et que

$$\phi'(0) = \nabla f(x)^T d = \sum_{i=1}^n \frac{\partial f(x)}{\partial x_i} d_i.$$

On suppose chaque composante  $f'_i(x) = \frac{\partial f(x)}{\partial x_i}$  de  $f'(x)$  est dérivable en  $x$ . Montrez que  $\phi$  est deux fois dérivable en 0 et montrez que

$$\phi''(0) = \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 f(x)}{\partial x_i \partial x_j} d_i d_j.$$

**Preuve 8.9** *Démonstration* : En utilisant la définition de la différentiabilité de  $f$  en  $x+td$ , on trouve que  $f(x+(t+\delta t)d) = f(x+td) + f'(x)(d \cdot \delta t) + \|\delta t \cdot d\| \epsilon(\delta t \cdot d) = f(x) + (f'(x+td)d)\delta t + |\delta t| \epsilon_1(\delta t)$ , avec  $\epsilon_1(\delta t) = \|d\| \epsilon(\delta t \cdot d)$  et  $\lim_{\delta t \rightarrow 0} \epsilon_1(\delta t) = 0$ , d'où le résultat obtenu en posant  $t=0$ . On applique ce résultat à la fonction  $\psi_i(t) = \frac{\partial f(x+td)}{\partial x_i}$ , pour obtenir que  $\psi'_i(0) = \sum_{j=1}^n \frac{\partial^2 f(x)}{\partial x_i \partial x_j} d_j$ . On finit en remarquant que  $\phi'(t) = \sum_{i=1}^n \psi_i(t) d_i$ , et donc que  $\phi''(t) = \sum_{i=1}^n \psi'_i(t) d_i$ .

□

**Exercice 8.10** Supposons que  $f$  est une fonction définie sur un ouvert convexe  $\mathcal{O}$  et 3 fois continûment dérivable en tout  $x \in \mathcal{O}$ . Montrez qu'alors la matrice carrée symétrique  $\nabla^2 f(x) = [\frac{\partial^2 f(x)}{\partial x_i \partial x_j}]$  appelée Hessian de  $f$  en  $x$ , est telle que

$$f(x+h) = f(x) + \nabla f(x)^T h + \frac{1}{2} h^T \nabla^2 f(x) h + o(h^2), \quad (8.1)$$

où l'on a posé  $o(h^2) = \|h\|^2 \epsilon(h)$ , avec  $\lim_{\|h\| \rightarrow 0} \epsilon(h) = 0$ .

**Preuve 8.10** *Démonstration* : Soit  $h$  tel que  $x+h \in \mathcal{O}$ . On pose  $\phi(t) = f(x+th)$ . On a alors  $\phi'(t) = \sum_{j=1}^n \frac{\partial f(x+th)}{\partial x_j} h_j$ ,  $\phi''(t) = \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 f(x+th)}{\partial x_i \partial x_j} h_i h_j = h^T \nabla^2 f(x+th) h$ , et

$$\phi'''(t) = \sum_{k=1}^n \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^3 f(x+th)}{\partial x_i \partial x_j \partial x_k} h_i h_j h_k.$$

D'après la formule de Taylor avec reste intégral on a  $\phi(1) = \phi(0) + \phi'(0) + \frac{1}{2} \phi''(0) + \frac{1}{2} \int_0^1 (1-s)^2 \phi'''(s) ds$ , et il reste à montrer que  $\lim_{\|h\| \rightarrow 0} \frac{\int_0^1 (1-s)^2 \phi'''(s) ds}{\|h\|^2} = 0$ . En utilisant l'équivalence des normes en dimension finie, on peut travailler avec n'importe quelle norme. Choisissons la norme infinie. Notons tout d'abord que comme  $h \mapsto \frac{\partial^3 f(x+h)}{\partial x_i \partial x_j \partial x_k}$  est continue, il existe  $M$  et  $\delta$  tel que  $|\frac{\partial^3 f(x+h)}{\partial x_i \partial x_j \partial x_k}| < M$  pour tout  $h$  tel que  $\|h\|_\infty \leq \delta$  (pour tout  $i$  et  $j$ ). Alors, comme  $|h_i| \leq \|h\|_\infty$ , on a  $|\int_0^1 (1-s)^2 \phi'''(s) ds| \leq \int_0^1 |\phi'''(s)| ds \leq \|h\|_\infty^3 \sum_{k=1}^n \sum_{i=1}^n \sum_{j=1}^n M = Mn^3 \|h\|_\infty^3$ . On a donc  $0 \leq \frac{|\int_0^1 (1-s)^2 \phi'''(s) ds|}{\|h\|^2} \leq \frac{Mn^3 \|h\|_\infty^3}{\|h\|_\infty^2}$ , d'où le résultat. □

**Exercice 8.11** On considère la fonction quadratique définie sur  $\mathbb{R}^n$  par  $f(x) = \frac{1}{2} x^T A x - x^T b$ , où  $A$  est carrée et symétrique montrez que  $\nabla^2 f(x) = A$ .

**Preuve 8.11** *Démonstration* : D'après l'exercice 8.8,  $\frac{\partial f(x)}{\partial x_i} = (Ax - b)_i = \sum_{j=1}^n a_{ij} x_j - b_i$ . Donc  $\frac{\partial^2 f(x)}{\partial x_i \partial x_j} = a_{ij}$ . □

**Exercice 8.12** (Taylor avec reste intégral) Supposons que  $f$  est une fonction définie sur un ouvert convexe  $\mathcal{O}$  et 1 fois continûment dérivable sur  $\mathcal{O}$ . Montrez qu'alors pour tout  $x$  et  $y$  de  $\mathcal{O}$ ,

$$f(y) = f(x) + \int_0^1 \nabla f(x + s(y-x))^T (y-x) ds.$$

Si  $f$  est 2 fois continûment dérivable sur  $\mathcal{O}$ ,

$$\nabla f(y) = \nabla f(x) + \int_0^1 \nabla^2 f(x + s(y-x))(y-x) ds.$$

**Preuve 8.12 Démonstration :** 1) Soit  $\phi$  la fonction continûment différentiable sur  $[0, 1]$ ,  $\phi(t) = f(x + t(y - x))$ . Alors le premier résultat n'est autre que

$$\phi(1) = \phi(0) + \int_{s=0}^1 \phi'(s) ds.$$

2) Soit  $\phi_i$  la fonction continûment différentiable sur  $[0, 1]$ ,  $\phi(t) = \frac{\partial f(x + t(y - x))}{\partial x_i}$ .

Alors 1) s'écrit  $\frac{\partial f(y)}{\partial x_i} = \frac{\partial f(x)}{\partial x_i} + \int_0^1 \left( \sum_{j=1}^n \frac{\partial^2 f(x + s(y - x))}{\partial x_j \partial x_i} (y - x)_j \right) ds$ , ce qui mis sous forme matricielle donne 2), puisque par définition du Hessien,  $\nabla^2 f(x + s(y - x)) = \left[ \frac{\partial^2 f(x + s(y - x))}{\partial x_i \partial x_j} \right]$ .

□

**Definition 8.5** Soit  $f$  définie sur un ouvert  $\mathcal{O} \subset \mathbb{R}^n$  à valeurs dans  $\mathbb{R}^m$ . On dit que  $f$  est dérivable (au sens de Fréchet) en  $x$ , si chacune des composantes  $f_i$  est dérivable (au sens de Fréchet) en  $x$ . On a alors

$$f(x + h) = f(x) + f'(x)h + o(h),$$

où l'on a posé  $f'(x) = [f_1(x)'; \dots; f_m(x)'] \in \mathbb{R}^{m \times n}$  ainsi que  $o(h) = \|h\|\epsilon(h) \in \mathbb{R}^m$ , avec  $\lim_{\|h\| \rightarrow 0} \epsilon(h) = 0 \in \mathbb{R}^m$ . La matrice

$$f'(x) = D_f(x) = \begin{pmatrix} \frac{\partial f_1(x)}{\partial x_1} & \cdots & \frac{\partial f_1(x)}{\partial x_n} \\ \vdots & \vdots & \vdots \\ \frac{\partial f_m(x)}{\partial x_1} & \cdots & \frac{\partial f_m(x)}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \nabla f_1(x)^T \\ \vdots \\ \nabla f_m(x)^T \end{pmatrix} \in \mathbb{R}^{m \times n}$$

est appelée matrice Jacobienne de  $f$  en  $x$ .

**Exercice 8.13** Dérivation d'une composée. Soit  $f$  définie sur un ouvert  $\mathcal{O} \subset \mathbb{R}^n$ , différentiable en  $x \in \mathcal{O}$ , à valeurs dans  $\mathbb{R}^m$ . Soit  $g$  définie sur un ouvert  $\mathcal{V} \subset \mathbb{R}^m$ , différentiable en  $f(x) \in \mathcal{V}$ , à valeurs dans  $\mathbb{R}^p$ . Alors la fonction  $x \mapsto g \circ f(x) = g(f(x))$  définie sur l'ouvert  $\mathcal{O}$  est différentiable en  $x$  et vérifie  $(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$ , où  $f'(x) \in \mathbb{R}^{m \times n}$  et  $g'(f(x)) \in \mathbb{R}^{p \times m}$ .

**Preuve 8.13 Démonstration :** Par définition de la différentiabilité en  $x$  de  $f$ , on a  $f(x + h) = f(x) + f'(x)h + \|h\|\epsilon_1(h)$  et

$$\begin{aligned} g(f(x + h)) &= g(f(x) + f'(x)h + \|h\|\epsilon_1(h)) \\ g(f(x + h)) &= g(f(x)) + g'(f(x)) (f'(x)h + \|h\|\epsilon_1(h)) \\ &\quad + \|f'(x)h + \|h\|\epsilon_1(h)\| \epsilon_2(f'(x)h + \|h\|\epsilon_1(h)) \end{aligned}$$

Posons  $\epsilon_3(h) = \frac{1}{\|h\|} (f'(x)h + \|h\|\epsilon_1(h)) \epsilon_2(f'(x)h + \|h\|\epsilon_1(h)) + g'(f(x))\|h\|\epsilon_1(h)$ . Alors  $0 \leq \epsilon_3(h) \leq (\|f'(x)\| + \|\epsilon_1(h)\|) \|\epsilon_2(f'(x)h + \|h\|\epsilon_1(h))\| + \|g'(f(x))\|\epsilon_1(h)$ . Le membre droit de cette inégalité tend bien vers 0 lorsque  $h$  tend vers 0 par définition de  $\epsilon_1$  et  $\epsilon_2$ . Ainsi  $\lim_{\|h\|_2 \rightarrow 0} \epsilon_3(h) = 0$  et

$$g(f(x+h)) = g(f(x)) + g'(f(x))f'(x)h + o(h).$$

□

**Exercice 8.14** *Dérivation numérique.* Pour une fonction différentiable, nous avons vu que le calcul de la dérivée se ramène au calcul de dérivées partielles, donc de dérivées de fonctions  $\phi$  de  $\mathbb{R}$  dans  $\mathbb{R}$ . On suppose que  $\phi$  est deux fois dérivable et que  $|\phi''(x)| \leq M$ . Sur un ordinateur, l'évaluation de  $\phi$  se fait à  $\epsilon$  près : à la place de  $\phi(x)$ , on calcule  $\tilde{\phi}(x) = \phi(x) + \delta(x)$ , avec  $|\delta(x)| \leq \epsilon$ . Posons  $\Delta_h^{\tilde{\phi}}(x) = \frac{\tilde{\phi}(x+h) - \tilde{\phi}(x)}{h}$ . Montrez que  $|\phi'(x) - \Delta_h^{\tilde{\phi}}(x)| \leq \frac{Mh}{2} + 2\frac{\epsilon}{h}$ . En déduire que un choix "raisonnable" pour  $h$  est  $h_0 = 2\sqrt{\frac{\epsilon}{M}}$ , pour lequel  $|\phi'(x) - \Delta_{h_0}^{\tilde{\phi}}(x)| \leq 2\sqrt{M\epsilon}$ .

**Preuve 8.14** *Démonstration :* D'après le théorème de Taylor Lagrange, il existe  $\theta$ ,  $0 < \theta < 1$ , tel que  $\phi(x+h) = \phi(x) + \phi'(x)h + \frac{h^2}{2}\phi''(x+\theta h)$ , ce qui montre que  $|\phi'(x) - \Delta_h^{\phi}(x)| \leq Mh/2$ . Cette erreur est une erreur d'approximation de la dérivée par une formule de différence finie. De plus,

$$|\Delta_h^{\phi}(x) - \Delta_h^{\tilde{\phi}}(x)| = \frac{|\delta(x+h) - \delta(x)|}{h} \leq 2\frac{\epsilon}{h}.$$

Cette erreur est une erreur numérique due au calcul de la différence finie sur ordinateur. On a alors

$$|\phi'(x) - \Delta_h^{\tilde{\phi}}(x)| \leq |\phi'(x) - \Delta_h^{\phi}(x)| + |\Delta_h^{\phi}(x) - \Delta_h^{\tilde{\phi}}(x)| \leq \frac{Mh}{2} + 2\frac{\epsilon}{h}.$$

Une idée pour choisir  $h$  est de minimiser pour  $h > 0$  la borne de l'erreur  $\frac{Mh}{2} + 2\frac{\epsilon}{h}$ . La dérivée vaut  $\frac{M}{2} - 2\frac{\epsilon}{h^2}$  et s'annule en  $h = 2\sqrt{\frac{\epsilon}{M}}$ , qui est bien le minimum (pour le voir étudier la fonction  $h \mapsto \frac{Mh}{2} + 2\frac{\epsilon}{h}$ ).

□

### 8.3 Minimisation locale

**Definition 8.6** Soit  $f$  une fonction définie sur un ouvert  $\mathcal{O}$  de  $\mathbb{R}^n$ . Un point  $\bar{x}$  pour lequel il existe  $\epsilon > 0$  tel que

$$f(\bar{x}) \leq f(x) \text{ pour tout } x \text{ tel que } \|\bar{x} - x\| < \epsilon$$

est un minimum local de  $f$ .

**Exercice 8.15** Si  $f$  est différentiable en  $\bar{x}$  et si  $\bar{x}$  est un minimum local  $f$  alors

$$\nabla f(\bar{x}) = 0. \quad (8.2)$$

Notez qu'en présence de contraintes, ce résultat ne tient plus (considérer  $\min_{x \in [0,1]} x$ ).

**Preuve 8.15** *Démonstration* : Supposons qu'il existe  $d$  tel que  $\nabla f(\bar{x})^T d < 0$ . Soit  $\phi : t \mapsto f(\bar{x} + td)$ , on a alors  $\phi'(0) = \nabla f(\bar{x})^T d < 0$ . On a alors  $\phi(t) = \phi(0) + \phi'(0)t + |t|\epsilon(t)$ . Comme  $\epsilon$  tend vers 0 en 0, il existe  $t_0$  tel que si  $t \leq t_0$ ,  $\epsilon(t) \leq \frac{|\phi'(0)|}{2} = \frac{-\phi'(0)}{2}$ . Mais alors, pour  $t > 0$ ,  $\phi(t) - \phi(0) \leq t \frac{\phi'(0)}{2} < 0$ , ce qui contredit que  $\bar{x}$  est un minimum local de  $f$ . □

**Exercice 8.16** Loi de Snell. En reprenant l'exercice 8.3, montrez que le principe de Fermat se traduit par la loi de Snell

$$\frac{\sin \alpha_1}{c_1} = \frac{\sin \alpha_2}{c_2}.$$

**Preuve 8.16** *Démonstration* : En appliquant l'exercice 8.5 à  $T(x)$ , il apparaît que la fonction  $T(x)$  admet au moins un minimum global. Ces minima vérifient  $T'(x) = 0$ , ce qui donne

$$\frac{x}{c_1 \sqrt{a^2 + x^2}} = \frac{c - x}{c_2 \sqrt{(c - x)^2 + b^2}},$$

ce qui donne bien la loi de Snell puisque  $\sin(\alpha_1) = \frac{x}{\sqrt{a^2 + x^2}}$  et  $\sin(\alpha_2) = \frac{c - x}{\sqrt{(c - x)^2 + b^2}}$ . □

Supposons  $f : \mathcal{O} \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  deux fois dérivable au point  $\bar{x} \in \mathcal{O}$ .

**Exercice 8.17** Si  $\bar{x}$  est un minimum local de  $f$  et si  $f$  est deux fois différentiable en  $\bar{x}$ , alors

$$\nabla f(\bar{x}) = 0 \text{ et } \nabla^2 f(\bar{x}) \text{ est semi-définie positive.} \quad (8.3)$$

Considérer  $\min_{x \in \mathbb{R}} x^3$  pour montrer que (8.3) n'est pas une condition suffisante de minimum local.

**Preuve 8.17** *Démonstration* : Pour  $d \in \mathbb{R}^n$  fixé, et  $\phi : t \mapsto f(\bar{x} + td)$ , on a  $\phi'(0) = d^T \nabla f(\bar{x})$ , et  $\phi''(0) = d^T \nabla^2 f(\bar{x}) d$ . Comme 0 est un minimum local de la fonction  $\phi$  (de la variable réelle  $t$ ) on a  $\phi'(0) = 0$ . Donc  $\phi(t) = \phi(0) + \phi''(0)t^2/2 + t^2\epsilon(t)$ . Comme  $\epsilon(t)$  tend vers 0 en 0, il existe  $t_0$  tel que si  $0 < t < t_0$ ,  $\epsilon(t) \leq |\phi''(0)|/4$ . Mais alors si  $\phi''(0) < 0$ ,  $\phi(t) - \phi(0) = \phi''(0)t^2/2 + t^2\epsilon(t) \leq \phi''(0)t^2(1/2 - 1/4) = \phi''(0)t^2/2 < 0$  pour tout  $t < t_0$ . Cela contredit le fait que 0 est un minimum local de  $\phi$ . Donc pour tout  $d$ ,  $\phi''(0) = d^T \nabla^2 f(\bar{x}) d \geq 0$ . □

**Exercice 8.18** Si  $f$  est deux fois différentiable en  $\bar{x}$  et si

$$\nabla f(\bar{x}) = 0 \text{ et } \nabla^2 f(\bar{x}) \text{ est définie positive,} \quad (8.4)$$

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

**Preuve 8.18** *Démonstration* : Raisonnons par l'absurde. On suppose que  $\bar{x}$  n'est pas un minimum local de  $f$  et que  $\nabla f(\bar{x}) = 0$  ; montrons que  $\nabla^2 f(\bar{x})$  n'est pas définie positive. Si  $\bar{x}$  n'est pas un minimum local de  $f$ , il existe une suite  $(x_k)$  convergeant vers  $\bar{x}$  dont le terme  $k$  est défini ainsi :  $x_k$  est tel que  $0 < \|\bar{x} - x_k\| \leq 1/k$  et  $f(x_k) < f(\bar{x})$ . Soit alors la direction normalisée  $d_k = \frac{x_k - \bar{x}}{\|x_k - \bar{x}\|}$  et  $t_k = \|x_k - \bar{x}\|$  (i.e.  $x_k = \bar{x} + t_k d_k$ ). Comme  $\|d_k\| = 1$ ,  $d_k$  appartient à un compact et on peut en extraire une sous-suite convergente. Soit  $d$  la limite d'une telle sous-suite notée sans perdre de généralité  $(d_k)$ . Alors  $0 > \frac{f(x_k) - f(\bar{x})}{t_k^2/2} = \frac{f(\bar{x} + t_k d_k) - f(\bar{x})}{t_k^2/2} = \langle \nabla^2 f(\bar{x}) d_k, d_k \rangle + 2\epsilon(t_k)$ , et comme  $\epsilon(t_k) \rightarrow 0$  on a par passage à la limite  $\langle \nabla^2 f(\bar{x}) d, d \rangle \leq 0$ . Cela montre que  $\nabla^2 f(\bar{x})$  n'est pas définie positive. □

**Exercice 8.19** On considère la fonction quadratique définie sur  $\mathbb{R}^n$  par  $f(x) = \frac{1}{2} x^T A x - x^T b$ , où  $A$  est carrée et symétrique montrez que si  $A$  est symétrique et définie positive, alors  $A^{-1}b$  est l'unique minimum de  $f$ . Appliquez ce résultat aux moindres carrés linéaires  $\min_{x \in \mathbb{R}^n} \|Ax - b\|_2$ , avec  $\text{rg}(A) = n$ .

**Preuve 8.19** *Démonstration* : Sans passer par les valeurs propres, le problème  $\min_{\|x\|=1} x^T A x$  admet au moins une solution  $x_{\min}$  (fonction continue sur un compact). De plus comme  $A$  est symétrique et définie positive,  $x_{\min}^T A x_{\min} = \sigma_{\min} > 0$ . Ainsi,

$$f(x) = \frac{1}{2} \|x\|^2 \frac{x^T}{\|x\|} A \frac{x}{\|x\|} - x^T b \geq \frac{1}{2} \|x\| (\sigma_{\min} \|x\| - \|b\|),$$

donc  $\lim_{\|x\| \rightarrow +\infty} f(x) = +\infty$ . La fonction  $f$  admet au moins un minimum car  $f$  est continue et coercive. D'après les exercices 8.8 et 8.15, la condition nécessaire d'optimalité s'écrit  $\nabla f(x) = Ax - b = 0$ , ce qui entraîne  $x = A^{-1}b$  car  $A$  est symétrique et définie positive, donc inversible. Il est normal que ce minimum vérifie aussi la condition nécessaire du second ordre de l'exercice 8.17 ( $\nabla^2 f(x) = A$  est semi-définie positive). De plus, il vérifie la condition suffisante du second ordre de l'exercice 8.18 ( $\nabla^2 f(x) = A$  est définie positive). On s'intéresse à présent aux moindres carrés linéaires. Alors  $f(x) = \|Ax - b\|_2^2 = x^T A^T A x - 2x^T A^T b + b^T b$ . Nous avons vu que si  $\text{rg}(A) = n$ , alors  $A^T A$  est définie positive ( $\text{Ker}(A^T A) = \text{Ker}(A)$  et formule du rang sur  $A$ ,  $n = \text{rg}(A) + \dim \text{Ker}(A)$ , d'où  $\text{Ker}(A^T A) = \{0\}$ ). La matrice  $A^T A$  est donc symétrique définie positive). L'exercice montre alors qu'une condition nécessaire et suffisante d'optimalité est  $\nabla f(x) = 2(A^T A x - A^T b) = 0$ , et on retrouve bien l'équation normale  $A^T A x = A^T b$ . □

## 8.4 Algorithmes de minimisation sans contrainte

### 8.4.1 La méthode de Newton

Cette méthode et ses variantes moins coûteuses, forment une des principales classes de méthode d'optimisation pour les problèmes sans contraintes. Cette méthode s'écrit :

*Newton's method*

---

1. Choose  $x_0$
2. For  $k=0, 2, \dots$  Do
3.     Compute if  $\nabla^2 f(x_k)$  is nonsingular
4.      $x_{k+1} = x_k - \nabla^2 f(x_k)^{-1} \nabla f(x_k)$
5. EndDo

Quelques remarques sur sa mise en œuvre :

- Cette méthode nécessite d'avoir à faire à une fonction deux fois dérivable et à ses dérivées jusqu'à l'ordre 2.
- Cette méthode nécessite aussi la résolution de systèmes linéaires (on ne calcule pas l'inverse). Cette opération peut être coûteuse pour des systèmes de grande taille.
- Cette méthode jouit de propriétés de convergence locale très intéressantes comme nous allons le voir.

Soit  $x^k \in \mathbb{R}^n$ . On considère l'approximation quadratique de  $f$ , fonction deux fois dérivable, suivante :  $m(x) = f(x_k) + \nabla f(x_k)^T(x - x_k) + \frac{1}{2}(x - x_k)^T \nabla^2 f(x_k)(x - x_k)$ . Supposons que  $\nabla^2 f(x_k)$  est symétrique et définie positive alors le minimum  $x^*$  de  $m(x)$  vérifie  $x^* - x_k = -\nabla^2 f(x_k)^{-1} \nabla f(x_k)$  d'après l'exercice 8.19. La méthode de Newton minimise donc à chaque pas où  $\nabla^2 f(x_k)$  est symétrique et définie positive l'approximation quadratique de  $f$  de l'exercice 8.10. Notez que si  $\nabla^2 f(x_k)$  a des valeurs propres négatives, l'approximation quadratique n'est pas bornée inférieurement, et le point  $x_{k+1}$  peut même dans certains cas être un maximum de  $m(x)$  (considérer  $-(x - x_k)^2$ ). Cette situation n'arrive pas si l'on est suffisamment proche de points vérifiant 8.18. Cela conduit aux conditions dites standart pour l'algorithme de Newton. Hypothèses standart en  $\bar{x} \in \mathcal{O}$ , où  $\mathcal{O}$  est un ouvert convexe de  $\mathbb{R}^n$  :

- c1  $f$  est deux fois continûment différentiable sur  $\mathcal{O}$
- c2  $x \mapsto \nabla^2 f(x)$  est Lipschitz continue sur  $\mathcal{O}$  :  $\|\nabla^2 f(y) - \nabla^2 f(x)\| \leq \gamma \|y - x\|$
- c3  $\nabla f(\bar{x}) = 0$  et  $\nabla^2 f(\bar{x})$  est définie positive

**Exercice 8.20** Sous les hypothèses standart, il existe  $\delta > 0$  et  $K > 0$ , tels que si  $\|\bar{x} - x_0\| \leq \delta$ ,  $\|\bar{x} - x_{k+1}\| \leq K \|\bar{x} - x_k\|^2$ . Si  $K\delta < 1$ ,  $(x_k)$  converge vers  $\bar{x}$ . Une telle convergence est appelée locale quadratique.

**Preuve 8.20** *Démonstration* : 1) En utilisant le Théorème de Rouché (ou un résultat de continuité des valeurs propres), il existe un voisinage de  $\bar{x}$  inclus dans  $\mathcal{O}$  où  $\nabla^2 f(x)$  est définie positive. Les fonctions  $x \mapsto \|\nabla^2 f(x)\|$  et  $x \mapsto \|\nabla^2 f(x)^{-1}\|$  sont continues dans un voisinage de  $\bar{x}$  inclus dans  $\mathcal{O}$  car  $x \mapsto \nabla^2 f(x)$  est continue, dans  $\mathcal{O}$  et  $x \mapsto \nabla^2 f(x)^{-1}$  est continue dans un voisinage de  $\bar{x}$  car  $\nabla^2 f(\bar{x})$  est inversible. Donc il existe  $\delta$  tel que

$\|\bar{x} - x\| \leq \delta$  (noté  $x \in B(\delta)$ ) entraîne

$$\|\nabla^2 f(x)\| \leq 2\|\nabla^2 f(\bar{x})\| \text{ et } \|\nabla^2 f(x)^{-1}\| \leq 2\|\nabla^2 f(\bar{x})^{-1}\|, \text{ et } \nabla^2 f(x) \text{ est définie positive.} \quad (8.5)$$

Soit  $x_k \in B(\delta)$ . Alors en utilisant l'exercice 8.12 on obtient  $\nabla f(x_k) = \int_0^1 \nabla^2 f(\bar{x} + s(x_k - \bar{x}))(x_k - \bar{x}) ds$ , qui montre que

$$\|\bar{x} - x_{k+1}\| = \|\bar{x} - x_k + \nabla^2 f(x_k)^{-1} \nabla f(x_k)\| \quad (8.6)$$

$$= \|\nabla^2 f(x_k)^{-1} \left( \nabla^2 f(x_k)(\bar{x} - x_k) + \int_0^1 \nabla^2 f(\bar{x} + s(x_k - \bar{x}))(x_k - \bar{x}) ds \right)\| \quad (8.7)$$

$$= \|\nabla^2 f(x_k)^{-1} \int_0^1 (\nabla^2 f(\bar{x} + s(x_k - \bar{x})) - \nabla^2 f(x_k))(x_k - \bar{x}) ds\| \quad (8.8)$$

$$\leq 2\gamma \|\nabla^2 f(\bar{x})^{-1}\| \int_0^1 (1-s) \|x_k - \bar{x}\|^2 ds = K \|\bar{x} - x_k\|^2. \quad (8.9)$$

Si  $K\delta < 1$ ,  $x_{k+1} \in B(\delta)$  (car  $\|\bar{x} - x_{k+1}\| \leq K\|x_k - \bar{x}\| \|x_k - \bar{x}\| \leq K\delta \|x_k - \bar{x}\|$ ) et par induction si  $x_0 \in B(\delta)$ , alors  $x_k \in B(\delta)$  pour tout  $k$ . De plus on vérifie aisément que  $\|\bar{x} - x_k\| \leq \frac{(K\delta)^{2^k}}{K}$ , ce qui montre que  $(x_k)$  converge vers  $\bar{x}$ . □

**Exercice 8.21** (Critère d'arrêt) Pour la suite  $f_n = \sum_{k=1}^n \frac{1}{k}$ , montrer que la stationnarité de  $f_n$  (i.e.  $f_{n+1} - f_n$  petit) n'indique pas la convergence. En déduire qu'arrêter une méthode d'optimisation sur  $|f(x_{k+1}) - f(x_k)| \leq \epsilon$  est dangereux. En revanche, sous les conditions standart, montrez que pour  $x_k$  suffisamment proche de  $\bar{x}$ , on a

$$\frac{\|\bar{x} - x_k\|}{4\|\bar{x} - x_0\| \text{cond}(\nabla^2 f(\bar{x}))} \leq \frac{\|\nabla f(x_k)\|}{\|\nabla f(x_0)\|} \leq \frac{4\text{cond}(\nabla^2 f(\bar{x}))\|\bar{x} - x_k\|}{\|\bar{x} - x_0\|}.$$

En déduire que la norme relative du gradient est un critère d'arrêt possible si le Hessien à l'optimum est bien conditionné.

**Preuve 8.21** Démonstration : La suite  $f_n$  diverge mais  $f_{n+1} - f_n$  tend vers 0. Par les mêmes arguments que pour la preuve de 8.20 on a pour  $x_k \in B(\delta)$ ,

$$\|\nabla f(x_k)\| = \left\| \int \nabla^2 f(\bar{x} + s(x_k - \bar{x}))(x_k - \bar{x}) ds \right\| \leq 2\|\nabla^2 f(\bar{x})\| \|\bar{x} - x_k\|.$$

Utilisant l'exercice 8.12 on obtient

$$\int (x_k - \bar{x})^T \nabla^2 f(\bar{x} + s(x_k - \bar{x}))(x_k - \bar{x}) ds = (x_k - \bar{x})^T \nabla f(x_k) \leq \|x_k - \bar{x}\| \|\nabla f(x_k)\|.$$

La matrice  $\nabla^2 f(\bar{x} + s(x_k - \bar{x}))$  étant définie positive dans  $B(\delta)$  l'inégalité matricielle pour une matrice  $A$  symétrique définie positive  $z^T z / \lambda_{\max}(A^{-1}) = \lambda_{\min}(A) z^T z \leq z^T A z$  montre alors que

$$\|x_k - \bar{x}\|^2 \int_0^1 \frac{1}{\|\nabla^2 f(\bar{x} + s(x_k - \bar{x}))^{-1}\|} ds \leq \|x_k - \bar{x}\| \|\nabla f(x_k)\|.$$

En utilisant  $\|\nabla^2 f(x_k)^{-1}\| \leq 2\|\nabla^2 f(\bar{x})^{-1}\|$ , on obtient  $\frac{\|x_k - \bar{x}\|^2}{2\|\nabla^2 f(\bar{x})^{-1}\|} \leq \|x_k - \bar{x}\| \|\nabla f(x_k)\|$ . Finalement, en rassemblant les majorations et minorations obtenues, on a pour  $x_k \in B(\delta)$   $\frac{\|x_k - \bar{x}\|}{2\|\nabla^2 f(\bar{x})^{-1}\|} \leq \|\nabla f(x_k)\| \leq 2\|\nabla^2 f(\bar{x})\| \|x_k - \bar{x}\|$  et la même inégalité pour  $x_0$  permet de conclure.

□

Il existe des variantes inexactes de la méthode de Newton où

- le gradient  $\nabla f(x_k)$  est approximatif,
- le Hessien  $\nabla^2 f(x_k)$  est approximatif,
- la solution du système linéaire  $\nabla^2 f(x_k)s = \nabla f(x_k)$  est calculée de manière approchée,

dans le but de rendre la méthode moins coûteuse en mémoire et en temps de calcul. Pour toutes ces variantes, des théories de convergence locale existent, qui imposent un bon contrôle des approximations.

### 8.4.2 Méthodes quasi-Newton

Une façon d'approximer la Hessienne, pour éviter de calculer et de stocker les dérivées d'ordre 2 est décrite comme suit. Pour une fonction quadratique, il est aisé de montrer que  $\nabla f(x_1) - \nabla f(x_2) = \nabla^2 f(x_1)(x_1 - x_2)$ . Cela indique que la connaissance de deux vecteurs distincts  $x_1$  et  $x_2$  et de la différence de gradient associée permet d'obtenir dans le cas quadratique -ou au voisinage de la solution sous les hypothèses standard, dans les étapes ultimes de la convergence- de l'information sur la Hessienne  $\nabla^2 f(x)$ . Plus généralement, on suppose connus,  $s = x_1 - x_2$  et  $y = \nabla f(x_1) - \nabla f(x_2)$ , ainsi qu'une approximation courante  $B$  de la Hessienne. On cherche une nouvelle approximation  $\tilde{B}$  telle que  $\tilde{B}$  soit symétrique et  $\tilde{B}s = y$ . Cela ne suffit pas pour définir de manière unique  $\tilde{B}$ , et on recherche des  $\tilde{B}$  de norme minimale (pour certaines normes) pour forcer l'unicité.

**Exercice 8.22** On recherche une matrice  $\tilde{B} = B + \Delta B$ , supposée mieux approcher que  $B$  la Hessienne en  $x_2$  en considérant le problème

$$\begin{aligned} \min_{\Delta B = \Delta B^T} \quad & \|\Delta B\|_F. \\ (B + \Delta B)s = y \end{aligned}$$

La solution de ce problème est donnée par la formule Powell-symmetric-Broyden :

$$\Delta B_0 = \frac{(y - Bs)s^T + s(y - Bs)^T}{s^T s} - \frac{s^T(y - Bs)ss^T}{(s^T s)^2}.$$

**Preuve 8.22 Démonstration :** On vérifie aisément que  $\Delta B_0 s = y - Bs$  et que  $\Delta B_0$  est symétrique. Soit  $q_1 = s / \|s\|_2$ . Pour tout  $\Delta B$  qui vérifie les contraintes (et en particulier pour  $\Delta B_0$ ), on a  $\Delta B q_1 = \Delta B_0 q_1 = \frac{y - Bs}{\|s\|_2}$ . soient  $q_i$ ,  $i = 2, \dots, n$ , qui complètent  $q_1$  en une base orthonormale de  $\mathbb{R}^n$ . Alors de  $q_i^T q_1 = 0$  pour  $i > 1$ , on tire  $\Delta B_0 q_i = \frac{s(\Delta B s)^T q_i}{s^T s} = \frac{ss^T}{s^T s} \Delta B q_i$ . D'où, en notant  $Q = [q_1, \dots, q_n]$ ,  $\|\Delta B_0 Q\|_F^2 = \sum_{i=1}^n \|\Delta B_0 q_i\|_2^2 \leq$

$\sum_{i=1}^n \|\Delta B q_i\|_2^2 = \|\Delta B Q\|_F^2$ . En utilisant le fait que la norme de Frobenius est unitairement invariante, on obtient  $\|\Delta B_0\|_F \leq \|\Delta B\|_F$ , d'où le résultat.

□

**Exercice 8.23** Soit  $f$  une fonction deux fois continûment dérivable, telle que  $\nabla^2 f(x)$  est définie positive pour tout  $x$ . Soit  $G = \int_0^1 \nabla^2 f(x_1 + s(x_2 - x_1)) ds$ . La matrice  $G$  est symétrique définie positive. Soit une matrice symétrique  $W$  telle que  $W^2 = G$ . On s'intéresse au problème

$$\begin{aligned} \min \quad & \|W^{-1} \Delta B W^{-1}\|_F \\ \Delta B = \Delta B^T \quad & \\ (B + \Delta B)s = y \quad & \end{aligned}$$

La solution de ce problème est donnée par la formule de Davidon-Fletcher-Powell

$$\Delta B_0 = \frac{(y - Bs)y^T + y(y - Bs)^T}{s^T y} - \frac{s^T (y - Bs) \cdot yy^T}{(s^T y)^2}.$$

Noter qu'alors

$$B + \Delta B_0 = \left( I - \frac{ys^T}{s^T y} \right) B \left( I - \frac{sy^T}{s^T y} \right) + \frac{yy^T}{s^T y}.$$

**Preuve 8.23** *Démonstration* : On a d'après l'exercice 8.12, puisque  $s = x_1 - x_2$  et  $y = \nabla f(x_1) - \nabla f(x_2)$  que  $Gs = y$ . De plus  $G$  est définie positive (considérer  $\int_0^1 u^T \nabla^2 f(x_1 + s(x_2 - x_1))u ds$  pour tout  $u$  de norme 1, et le fait que l'intégrande est continu et strictement positif). Donc  $s^T y = s^T Gs > 0$ . Soit alors  $W$  une racine carrée positive de  $G$  (en fait elle est unique). Par changement de variable  $\Delta = W^{-T} \Delta B W^{-1}$ , le problème devient

$$\begin{aligned} \min \quad & \|\Delta\|_F \\ \Delta = \Delta^T \quad & \\ (W^{-1} B W^{-1} + \Delta)W s = W^{-1} y \quad & \end{aligned}$$

D'après l'exercice 8.22 précédent, et en notant que  $Gs = WWs = y$  et  $Ws = W^{-1}y$ , la solution s'écrit

$$\begin{aligned} \Delta_0 &= \frac{(W^{-1}y - W^{-1}B W^{-1}W s)s^T W + s(W^{-1}y - W^{-1}B W^{-1}W s)^T}{s^T W W s} \\ &\quad - \frac{s^T W (W^{-1}y - W^{-1}B W^{-1}W s)W s s^T W}{(s^T W W s)^2} \\ &= \frac{W^{-1}(y - Bs)y^T W^{-1} + W^{-1}y(y - Bs)^T W^{-1}}{s^T y} - \frac{s^T (y - B W s)W^{-1}y y^T W^{-1}}{(s^T y)^2}. \end{aligned}$$

En faisant le changement de variable  $\Delta B = W \Delta W$ , on obtient le résultat désiré.

□

**Exercice 8.24** Nous avons vu que dans la méthode de Newton, il s'agit de résoudre des systèmes linéaires de la forme  $\nabla^2 f(x_k)s = \nabla f(x_k)$ . D'où l'idée d'approcher  $\nabla^2 f(x_k)^{-1}$  plutôt que  $\nabla^2 f(x_k)$ . Montrez que la formule BFGS (Broyden, Fletcher, Goldfarb, Shanno)

$$H + \Delta H_0 = \left( I - \frac{sy^T}{y^T s} \right) H \left( I - \frac{ys^T}{y^T s} \right) + \frac{ss^T}{y^T s},$$

est telle que  $\Delta H_0$  est solution de

$$\begin{aligned} \min \quad & \|\Delta H\|, \\ \Delta H = \Delta H^T \quad & \\ (H + \Delta H)y = s \quad & \end{aligned}$$

pour une norme  $\|\bullet\|$  que vous identifierez.

**Preuve 8.24** *Démonstration* : Dans la démonstration de l'exercice 8.23, on a démontré que si  $Gs = y$ , avec  $G = WW$  définie positive, alors la mise à jour DFP pour  $(B + \Delta B)s = y$  est solution du problème de mise à jour avec la norme  $\|W^{-1} \bullet W^{-1}\|_F$ . On considère maintenant l'équation  $(H + \Delta H)y = s$ . On peut appliquer DFP à ce problème en notant que  $s = G^{-1}y$  ( $G = W^{-1}W^{-1}$  est définie positive). La formule BFGS est alors la mise à jour de DFP correspondant au problème

$$\begin{aligned} \min \quad & \|W\Delta HW\|_F, \\ \Delta H = \Delta H^T \quad & \\ (H + \Delta H)y = s \quad & \end{aligned}$$

□

Deux principales difficultés sont rapportées dans la littérature sur la méthode de Newton pour la minimisation :

1. Son mauvais comportement lorsque le point de départ est loin de la solution sur des problèmes pour lesquels certains Hessiens  $\nabla^2 f(x_k)$  sont définis positifs.
2. Son mauvais comportement lorsqu'elle rencontre des Hessiens ayant des valeurs propres négatives ou nulles.

Une amélioration possible pour le problème 1) est la mise en place de stratégies de recherches linéaires. Le point 2) est souvent appréhendé en utilisant des techniques de région de confiance.

### 8.4.3 Globalisation des méthodes de Newton/quasi-Newton

**Exercice 8.25** Calculez quelques itérés de la méthode de Newton sur  $f(x) = -e^{-x^2}$ , pour  $x_0 = 10^{-1}$ ,  $x_0 = 1/2$  et  $x_0 = 1$ .

**Preuve 8.25** *Démonstration* :  $f(x) = -e^{-x^2}$ ,  $f'(x) = 2xe^{-x^2}$ ,  $f''(x) = (2 - 4x^2)e^{-x^2}$ . Alors on a  $x_{k+1} = x_k - 2x_k/(2 - 4x_k^2) = -4x_k^3/(2 - 4x_k^2)$ . Pour  $x_0 = 10^{-1}$ , on a  $x_1 \sim 2 \cdot 10^{-3}$  et  $x_2 \sim 2 \cdot 10^{-8}$ . Pour  $x_0 = 1/2$ , on a  $x_1 = -1/2$  et  $x_2 = 1/2$ . Pour  $x_0 = 1$ , on a  $x_1 \sim 2.3$  et  $x_2 \sim 2.5$ ,  $x_{23} \sim 5.4$  et  $f(x_{23}) \sim 10^{-13}$ .

□

Nous voyons dans la suite deux techniques visant à rendre la convergence moins dépendante du point de départ. Ces deux techniques sont appelées techniques de globalisation, et chercheront à approcher une convergence locale quadratique au voisinage des solutions de  $\nabla f(x) = 0$ . Ces solutions sont appelées points critiques du premier ordre.

### Recherche linéaire

Dans cette section, on suppose que la fonction  $f$  est deux fois continûment dérivable.

**Definition 8.7** Soit  $x_k \in \mathbb{R}^n$ . On dit que  $d_k$  est une direction de descente en  $x_k$  si  $\nabla f(x_k)^T d_k < 0$ .

La terminologie "direction de descente" s'explique aisément par l'exercice 8.26.

**Exercice 8.26** Si  $d_k$  est une direction de descente en  $x_k$ , alors il existe  $\eta > 0$  tel que

$$f(x_k + \alpha d_k) < f(x_k) \text{ pour tout } \alpha \in ]0, \eta].$$

**Preuve 8.26** Démonstration : Soit  $\phi(t) = f(x_k + td_k)$ . Alors  $\phi'(t) = \nabla f(x_k + td_k)^T d_k$ , donc comme  $\phi'$  est continue, et  $\phi'(0) < 0$ , il existe un intervalle  $]0, \eta]$  où  $\phi'(t) < 0$ . Alors pour  $t$  dans  $]0, \eta]$ , on a  $f(x_k + \alpha d_k) - f(x_k) = \int_{s=0}^t \phi'(s) ds < 0$ .

□

On envisage alors un premier algorithme de minimisation basé sur des directions de descente :

*Basic linesearch (bad algorithm)*

1. Choose  $x_0$
2. For  $k=0, 2, \dots$  Do
3.     Compute a descent direction such that  $\nabla f(x_k)^T d_k < 0$
4.     Compute a step such that  $f(x_k + \alpha_k d_k) < f(x_k)$ .
5.     Update  $x_{k+1} = x_k + \alpha_k d_k$ .
6. EndDo

**Exercice 8.27** L'algorithme ci-dessus ne suffit pas pour converger vers un minimum local de  $f$ . Soit  $f(x) = x^2$ ,  $x_0 = 2$ .

1. On choisit  $d_k = (-1)^{k+1}$  et  $\alpha_k = 2 + 3 \cdot 2^{-k-1}$ . Vérifier que  $x_k = (-1)^k(1 + 2^{-k})$  et que chaque direction  $d_k$  est de descente. Vérifier aussi que la suite ne converge pas, que  $f(x_{k+1}) < f(x_k)$  et que  $\lim_{k \rightarrow +\infty} f(x_k) = 1$ . Tracer les itérés et vérifier qu'entre deux itérés successifs, la décroissance de  $f$  est très petite par rapport au pas  $|\alpha_k d_k|$ .
2. On choisit  $d_k = -1$  et  $\alpha_k = 2^{-(k+1)}$ . Vérifier que  $x_k = 1 + 2^{-k}$  et que chaque direction  $d_k$  est de descente. Vérifier aussi que la suite converge vers 1 (et pas vers 0) que  $f(x_{k+1}) < f(x_k)$  et que  $\lim_{k \rightarrow +\infty} f(x_k) = 1$ . Tracer les itérés et vérifier qu'entre deux itérés successifs, les pas  $|\alpha_k d_k|$  deviennent très petits par rapport à  $|f'(x_k) d_k|$ .

**Preuve 8.27** *Démonstration* :

1. Par récurrence,  $x_{k+1} = x_k + \alpha_k d_k = (-1)^k(1 + 2^{-k}) + (2 + 3 \cdot 2^{-k-1})(-1)^{k+1} = (-1)^{k+1}(1 + 2^{-(k+1)})$ . Direction de descente :  $f'(x_k)d_k = 2(-1)^k(1 + 2^{-k})(-1)^{k+1} < 0$ . La suite admet  $-1$  et  $1$  comme points d'accumulation et  $\lim_{k \rightarrow +\infty} f(x_k) = 1$ . De plus  $f(x_{k+1}) - f(x_k) = (1 + 2^{-k})^2 - (1 + 2^{-(k-1)})^2 < 0$ .
2. Par récurrence,  $x_{k+1} = x_k + \alpha_k d_k = 1 + 2^{-k} - 2^{-k-1} = 1 + 2^{-(k+1)}$ . Direction de descente :  $f'(x_k)d_k = 2(1 + 2^{-k})(-1) < 0$ , et  $f(x_{k+1}) - f(x_k) < 0$ .

□

**Definition 8.8** Soit  $\beta_1 \in ]0, 1[$ ,  $\beta_2 \in ]\beta_1, 1[$ , et soit  $d_k$  une direction de descente en  $x_k$ . On appelle conditions de Wolfe les deux conditions :

1.  $f(x_k + \alpha d_k) \leq f(x_k) + \beta_1 \alpha \nabla f(x_k)^T d_k$  (condition de diminution suffisante)
2.  $\nabla f(x_k + \alpha d_k)^T d_k \geq \beta_2 \nabla f(x_k)^T d_k$  (condition de progrès suffisant)

Ces deux conditions pallient respectivement les deux types de problèmes rencontrés dans l'exercice 8.27. Si  $\alpha \rightarrow f(x_k + \alpha d_k)$  admet un minimum global, celui-ci vérifie les conditions de Wolfe (mais peut être très ou trop cher à calculer à des étapes préliminaires de convergence).

**Preuve 8.28** *Démonstration* :

1. Dans le cas 1.,  $f(x_k + \alpha_k d_k) - f(x_k) = (1 + 2^{-k-1})^2 - (1 + 2^{-k})^2 = -2^{-k-1}(2 + 3 \cdot 2^{-k-1})$  et  $\nabla f(x_k)^T d_k = -2(1 + 2^{-k})$ . Donc la condition de diminution suffisante n'est pas vérifiée.
2. Dans le cas 2,  $\nabla f(x_k + \alpha_k d_k)^T d_k = -2x_{k+1}$  et  $\nabla f(x_k)^T d_k = -2x_k$ , et comme  $\{x_k\}$  tend vers 1, la condition de progrès suffisant n'est pas vérifiée.

□

**Exercice 8.28** *Validité des conditions de Wolfe.* Soient  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  une fonction différentiable, un point  $x_k \in \mathbb{R}^n$  et une direction (de descente)  $d_k \in \mathbb{R}^n$  telle que  $f$  est bornée inférieurement dans la direction  $d_k$  (c'est-à-dire il existe  $f_0$  tel que  $f(x_k + \alpha d_k) \geq f_0$  pour tout  $\alpha \geq 0$ ). Pour  $0 < \beta_1 < 1$ , il existe  $\eta$  tel que la première condition de Wolfe soit vérifiée pour tout  $\alpha_k$ ,  $0 < \alpha_k \leq \eta$ . De plus, si  $0 < \beta_1 < \beta_2 < 1$ , il existe  $\alpha > 0$  tel que les deux conditions de Wolfe soient toutes deux vérifiées.

**Preuve 8.29** *Démonstration* : On s'intéresse aux  $\alpha > 0$  tels que  $f(x_k + \alpha d_k) = f(x_k) + \beta_1 \alpha \nabla f(x_k)^T d_k$ . Cet ensemble est non vide (car sinon  $\alpha \mapsto f(x_k + \alpha d_k)$  serait en dessous de  $\alpha \mapsto f(x_k) + \beta_1 \alpha \nabla f(x_k)^T d_k$ , ce qui est impossible car  $0 < \beta_1 < 1$  et  $f$  est bornée inférieurement), fermé (image réciproque de  $\{0\}$ ) et borné inférieurement. Donc cet ensemble admet un plus petit élément  $\alpha_1$ , qui vérifie

$$f(x_k + \alpha_1 d_k) = f(x_k) + \beta_1 \alpha_1 \nabla f(x_k)^T d_k,$$

donc qui vérifie la première condition de Wolfe.

D'après Taylor-Lagrange, appliqué à  $\alpha \mapsto f(x_k + \alpha d_k)$ , entre 0 et  $\alpha_1$ , il existe  $\alpha_2$  tel que

$$f(x_k + \alpha_1 d_k) = f(x_k) + \alpha_1 \nabla f(x_k + \alpha_2 d_k)^T d_k,$$

En rassemblant les deux résultats, on obtient

$$\nabla f(x_k + \alpha_2 d_k)^T d_k = \beta_1 \nabla f(x_k)^T d_k > \beta_2 \nabla f(x_k)^T d_k,$$

donc  $\alpha_2$  vérifie la seconde condition de Wolfe. Comme  $\alpha_2 < \alpha_1$ , on a  $f(x_k + \alpha_2 d_k) < f(x_k) + \beta_1 \alpha_2 \nabla f(x_k)^T d_k$  et donc  $\alpha_2$  vérifie la première condition de Wolfe est vérifiée.

□

*Descent algorithm with Wolfe linesearch*

1. Choose  $x_0$
2. For  $k=0,2, \dots$  Do
3.     Compute a descent direction such that  $\nabla f(x_k)^T d_k < 0$
4.     Compute a step such that the Wolfe conditions hold.
5.     Update  $x_{k+1} = x_k + \alpha_k d_k$ .
6. EndDo

**Théorème 8.9** Supposons de  $f$  soit continûment différentiable, bornée inférieurement, et que son gradient vérifie  $\|\nabla f(x) - \nabla f(y)\|_2 \leq \gamma \|x - y\|_2$ . supposons qu'un algorithme de descente soit employé tel que chaque pas vérifie les conditions de Wolfe. Alors soit  $\lim_{k \rightarrow +\infty} \nabla f(x_k) = 0$ , soit  $\lim_{k \rightarrow +\infty} \frac{\nabla f(x_k)^T d_k}{\|d_k\|_2} = 0$ .

**Preuve 8.30** Démonstration : Admise, voir Denis et Schnabel 1996, p.121.

□

Le théorème ci-dessus indique que si l'angle entre  $d_k$  et  $\nabla f(x_k)$  ne converge pas vers l'angle droit, la limite du gradient de l'itéré est 0 (on vérifie asymptotiquement la condition nécessaire du premier ordre) quel que soit  $x_0$ . C'est donc un résultat de convergence globale. Malheureusement cet algorithme peut avoir une convergence très lente si  $d_k$  n'est pas choisi avec soin. Par exemple, le choix  $d_k = -\nabla f(x_k)$  s'avère un très mauvais choix si l'algorithme converge vers un point  $x^*$  tel que  $\text{cond}(\nabla^2 f(x_k))$  est grand : la convergence est linéaire, avec une vitesse de convergence modeste.

Dans le cas d'une convergence vers un point  $x^*$  tel que  $\nabla^2 f(x^*)$  est défini positif (condition suffisante du second ordre), l'idée consiste alors à préconditionner la recherche linéaire et à la combiner avec la méthode de Newton qui est localement quadratiquement convergente, comme le fait l'algorithme ci-dessous. Il est possible de montrer que lorsque les itérés s'approchent d'une solution qui vérifie les conditions suffisantes d'optimalité au second ordre, le pas de Newton est accepté et la convergence est quadratique.

*Newton with linesearch*

1. Choose  $x_0$
2. For  $k=0,2, \dots$  Do
3. If  $\nabla^2 f(x_k)$  is SPD, compute the Newton step  $s^N = -\nabla^2 f(x_k)^{-1} \nabla f(x_k)$ .  
If  $s^N$  is acceptable (Wolfe) accept it. If not, perform a line search (Wolfe) in direction  $s^N$
4. If  $\nabla^2 f(x_k)$  is not SPD, add a perturbation  $E$  so that  $\nabla^2 f(x_k) + E$  is SPD,  
and perform a line search (Wolfe) in direction  $-(\nabla^2 f(x_k)^{-1} + E) \nabla f(x_k)$
5. Update  $x_{k+1} = x_k + \alpha_k d_k$ .
6. EndDo

**Région de confiance**

**Definition 8.10** *Modèle quadratique.* On appelle modèle quadratique de  $f$  en  $x_k$  une fonction quadratique  $m_k(x_k + s)$  telle que  $m_k(x_k) = f(x_k)$  et  $\nabla m_k(x_k) = \nabla f_k(x_k)$ . Il existe alors une matrice  $H_k \in \mathbb{R}^{n \times n}$  telle que

$$m_k(x_k + s) = f(x_k) + \nabla f_k(x_k)^T s + \frac{1}{2} s^T H_k s.$$

**Definition 8.11** *Région de confiance.* On appelle région de confiance Euclidienne centrée en  $x_k$ , de rayon  $\Delta_k > 0$  la sphère  $\mathcal{B}_k = x_k + \{s, \|s\|_2 \leq \Delta_k\}$ .

L'idée de l'algorithme de région de confiance et de résoudre approximativement le problème

$$\min_{x_k + s \in \mathcal{B}_k} m_k(x_k + s).$$

On note  $x_{k+1} = x_k + s_k$  le point ainsi obtenu. La condition technique portant sur  $x_{k+1}$  demandée pour les résultats de convergence est la condition dite de décroissance suffisante :

$$m_k(x_k) - m_k(x_k + s_k) \geq \kappa_{mdc} \|\nabla m_k(x_k)\|_2 \min \left( \frac{\|\nabla m_k(x_k)\|_2}{\beta_k}, \Delta_k \right), \quad (8.10)$$

où  $\kappa_{mdc} \in ]0, 1[$  et  $\beta_k = \|H_k(x)\|_2 + 1$ .

**Exercice 8.29** Le point de Cauchy  $x_k^C$  qui est, par définition, solution de

$$\begin{cases} \min & m_k(x) \\ t > 0 \\ x = x_k - t \nabla m(x_k) \in \mathcal{B}_k \end{cases}$$

vérifie

$$m_k(x_k) - m_k(x_k^C) \geq \frac{1}{2} \|\nabla m_k(x_k)\|_2 \min \left( \frac{\|\nabla m_k(x_k)\|_2}{\beta_k}, \Delta_k \right).$$

**Preuve 8.31** *Démonstration :* Posons  $g_k = \nabla_x m_k(x_k)$ . On a  $m_k(x_k - t g_k) = m_k(x_k) - t \|g_k\|^2 + \frac{1}{2} t^2 g_k^T H_k g_k$ .

1. Supposons  $g_k^T H_k g_k > 0$ . Alors le minimum de  $m_k(x_k - t g_k)$  pour  $t \in \mathbb{R}$  est atteint en  $t^* = \frac{\|g_k\|^2}{g_k^T H_k g_k} \geq 0$ .

Premier cas. Supposons d'abord que  $t^* \|g_k\| = \frac{\|g_k\|^3}{g_k^T H_k g_k} \leq \Delta_k$ , donc  $x_k - t^* g_k$  est dans la région de confiance et c'est  $x_k^C$ . Comme  $g_k^T H_k g_k \leq \beta_k \|g_k\|$ , on a alors

$$\begin{aligned} m_k(x_k) - m_k(x_k^C) &= t^* \|g_k\|^2 - \frac{1}{2} t^{*2} g_k^T H_k g_k \geq \frac{\|g_k\|^4}{g_k^T H_k g_k} - \frac{1}{2} \frac{\|g_k\|^4}{(g_k^T H_k g_k)^2} g_k^T H_k g_k \\ &= \frac{1}{2} \frac{\|g_k\|^4}{g_k^T H_k g_k} \geq \frac{1}{2} \frac{\|g_k\|^2}{\beta_k}. \end{aligned}$$

Deuxième cas. Supposons maintenant que  $\frac{\|g_k\|^3}{g_k^T H_k g_k} \geq \Delta_k$ . Alors  $g_k^T H_k g_k \leq \frac{\|g_k\|^3}{\Delta_k}$  et le minimum dans la région de confiance est donc atteint sur la frontière (faire un dessin). Alors  $t^* \|g_k\| = \Delta_k$  et  $x_k^C = x_k - \Delta_k g_k$  et

$$m_k(x_k) - m_k(x_k^C) = \Delta_k \|g_k\| - \frac{1}{2} \frac{\Delta_k^2}{\|g_k\|^2} g_k^T H_k g_k \geq \Delta_k \|g_k\| - \frac{1}{2} \frac{\Delta_k^2}{\|g_k\|^2} \frac{\|g_k\|^3}{\Delta_k} = \frac{1}{2} \Delta_k \|g_k\|.$$

2. Supposons  $g_k^T H_k g_k \leq 0$ . Le minimum est à nouveau atteint sur la frontière de la région de confiance et puisque  $-g_k^T H_k g_k \geq 0$

$$m_k(x_k) - m_k(x_k^C) = \Delta_k \|g_k\| - \frac{1}{2} \frac{\Delta_k^2}{\|g_k\|^2} g_k^T H_k g_k \geq \Delta_k \|g_k\|.$$

En regroupant les différents sous-cas, on obtient le résultat. □

Le calcul de  $x_{k+1}$  (donc de  $s_k$ ) est bien moins cher que la résolution du problème initial  $\min_x f(x)$  car

1.  $m_k$  est une fonction quadratique
2. la décroissance suffisante est obtenue à faible coût, en calculant le point de Cauchy, et en cherchant éventuellement à diminuer encore  $m_k$  à partir de  $x_k^C$ . La méthode des régions de confiance a donc un rapport étroit avec la recherche linéaire suivant la direction  $-\nabla f_k(x_k)$ .

On introduit le ratio de la réduction observée sur  $f$  par rapport à la réduction prédite sur  $m_k$  :

$$\rho_k = \frac{f(x_k) - f(x_{k+1})}{m(x_k) - m(x_{k+1})}.$$

Si  $\rho_k$  est suffisamment proche de 1, le modèle représente la fonction de manière fiable, on accepte le pas, et on augmente éventuellement le rayon de la région de confiance. Si  $\rho_k$  est faible, voire négatif, le modèle n'est pas assez fiable, et l'on réduit la région de confiance (notez que pour  $\Delta_k$  suffisamment petit modèle et fonction sont égaux au premier ordre). Nous sommes en mesure de présenter à présent l'algorithme des régions de confiance :

*Basic trust region algorithm (Conn, Gould, Toint (2000 p.116))*

1. Choose  $x_0$ , an initial  $\Delta_0 > 0$ , and constants  $0 < \eta_1 \leq \eta_2 < 1$  and  $0 < \gamma_1 \leq \gamma_2 < 1$
2. For  $k=0,2, \dots$  Do
3. Compute a step  $s_k$  that sufficiently reduces  $m_k$  in  $\mathcal{B}_k$  ( 8.10).
4. Define  $\rho_k = \frac{f(x_k) - f(x_k + s_k)}{m(x_k) - m(x_k + s_k)}$ .
5. If  $\rho_k \geq \eta_1$  then define  $x_{k+1} = x_k + s_k$ ; otherwise define  $x_{k+1} = x_k$
6. Trust region update. Set
  - $\Delta_{k+1} \in [\Delta_k, +\infty[$  if  $\rho_k \geq \eta_2$  or
  - $\Delta_{k+1} \in [\gamma_2 \Delta_k, \Delta_k]$  if  $\eta_1 \leq \rho_k < \eta_2$  or
  - $\Delta_{k+1} \in [\gamma_1 \Delta_k, \gamma_2 \Delta_k]$  if  $\rho_k < \eta_1$
7. If converged, exit,
8. EndDo

**Théorème 8.12** *On suppose que l'algorithme est appliqué à une fonction*

- deux fois différentiable,
- bornée inférieurement sur  $\mathbb{R}^n$ ,
- à Hessien borné ( $\|\nabla^2 f(x)\|_2 \leq \kappa_{ufh}$  pour  $x \in \mathbb{R}^n$ ),

et que les modèles  $m_k$  sont

- quadratiques,
- ont même valeur et gradient que  $f$  en  $x_k$  (cohérence au premier ordre)
- ont des Hessien bornés ( $\|\nabla^2 f(x)\|_2 \leq \kappa_{umh}$  pour  $x \in \mathcal{B}_k$ ).

alors pour tout  $x_0$ , l'algorithme des régions de confiance produit une suite d'itérés telle que  $\lim_{k \rightarrow +\infty} \nabla f(x_k) = 0$ .

**Preuve 8.32** Démonstration : Admise (Conn, Gould, Toint (2000 p.136)).

□

Le théorème 8.12 montre une manière aisée d'obtenir un algorithme globalement convergent : il suffit de choisir  $\nabla^2 m_k(x_k) = H_k = 0 \in \mathbb{R}^{n \times n}$  et de prendre pour itéré le point de Cauchy. Par contre on obtient alors un algorithme qui converge aussi peu rapidement que celui implantant systématiquement la recherche linéaire dans la direction  $-\nabla f(x_k)$ . Pour obtenir un algorithme plus performant et approcher la convergence locale de l'algorithme de Newton, il convient de choisir un pas  $s_k$  qui soit voisin du pas de Newton dans les étapes ultimes de la convergence.

Ceci est réalisé si l'on utilise pour algorithme de calcul de pas l'algorithme de gradient conjugué tronqué proposé par Steihaug et Toint et si le Hessien du modèle approche celui de la fonction. Cet algorithme commence par calculer le point de Cauchy puis poursuit la minimisation de la quadratique  $m(x_k + s)$  par la méthode des gradients conjugués, en s'arrêtant au premier itéré sortant de la région de confiance  $\mathcal{B}_k$ . On a ainsi minimisé davantage  $m(x_k + s)$  que  $m(x_k^C)$ , et donc on a, à la fin de cette procédure de gradient conjugué tronqué, la décroissance suffisante :

$$m(x_k) - m(x_k + s_k) \geq m(x_k) - m(x_k^C) \geq \frac{1}{2} \|\nabla_x m_k(x_k)\|_2 \min \left( \frac{\|\nabla_x m_k(x_k)\|_2}{\beta_k}, \Delta_k \right).$$

Dans le cas où la convergence a lieu vers un point  $x^*$  où le Hessien est défini positif et si  $\nabla^2 m_k(x_k) \sim \nabla^2 f_k(x_k)$ , le comportement typique de l'algorithme est alors le suivant :

1. les pas deviennent de plus en plus petits (on converge),
2. comme le modèle et la fonction sont cohérents au premier ordre,  $\rho_k$  devient proche de 1,
3. la région de confiance a un rayon qui augmente,
4. l'algorithme des gradients conjugués ne rencontre plus le bord de la région de confiance,
5. les gradient conjugués résolvent alors le système  $\nabla^2 f(x_k)s_k + \nabla f(x_k) = 0$  ce qui correspond bien à la méthode de Newton, qui a une convergence locale quadratique.

*Steihaug Toint Conjugate Gradient algorithm*

- |   |
|---|
| <ol style="list-style-type: none"> <li>0. Input parameters : <math>x_k, \nabla f(x_k), H = \nabla^2 f(x_k)</math>. Output : <math>s</math></li> <li>1. Compute <math>s_0 = 0, g_0 = \nabla f(x_k)</math></li> <li>2. For <math>k=0, 2, \dots</math> Do</li> <li>3. <math>\kappa_k = p_k^T H p_k</math></li> <li>4. If <math>\kappa_k \leq 0</math>, then           <ul style="list-style-type: none"> <li>compute <math>\sigma_k</math> as the positive root of <math>\ s_k + \sigma p_k\ _2 = \Delta_k</math></li> <li><math>s_{k+1} = s_k + \sigma_k p_k</math> and stop.</li> </ul>           End If         </li> <li>5. <math>\alpha_k = r_k^T r_k / \kappa_k</math></li> <li>6. If <math>\ s_k + \alpha_k p_k\ _2 \geq \Delta_k</math>, then           <ul style="list-style-type: none"> <li>compute <math>\sigma_k</math> as the positive root of <math>\ s_k + \sigma p_k\ _2 = \Delta_k</math></li> <li><math>s_{k+1} = s_k + \sigma_k p_k</math> and stop.</li> </ul>           End If         </li> <li>4. <math>s_{k+1} = s_k + \alpha_k p_k</math></li> <li>5. <math>g_{k+1} = g_k + \alpha_k H p_k</math></li> <li>7. <math>\beta_k = g_{k+1}^T g_{k+1} / g_k^T g_k</math></li> <li>8. <math>p_{k+1} = g_{k+1} + \beta_k p_k</math></li> <li>9. if converged then stop</li> <li>10. EndDo</li> </ol> |
|---|

#### 8.4.4 Globalisation des moindres carrés non-linéaires

**Exercice 8.30** *Fonctionnelle des moindres carrés non linéaires.* Soit  $f$  définie sur un ouvert  $\mathcal{O} \subset \mathbb{R}^n$ , deux fois différentiable, à valeurs dans  $\mathbb{R}^m$ . On définit la fonction  $F(x)$  des moindres carrés non linéaires par  $F(x) = \frac{1}{2} \|f(x)\|_2^2$ . Montrez que le gradient de  $F$  en  $x$  est  $f'(x)^T f(x) = D_f(x)^T f(x)$  et que la matrice Hessienne de  $F$  en  $x$  est  $D_f(x)^T D_f(x) + \sum_{i=1}^m f_i(x) \nabla^2 f_i(x)$ .

**Preuve 8.33** *Démonstration :* Considérons  $\phi(x) = f_i(x)^2$ . Alors, par dérivation d'une composée,  $\frac{\partial \phi(x)}{\partial x_j} = 2f_i(x) \frac{\partial f_i(x)}{\partial x_j}$ , et donc  $\frac{\partial F(x)}{\partial x_j} = \sum_{i=1}^m \frac{\partial f_i(x)}{\partial x_j} f_i(x)$ , ce qui implique

$$\nabla F(x) = \begin{pmatrix} \frac{\partial F(x)}{\partial x_1} \\ \vdots \\ \frac{\partial F(x)}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \frac{\partial f_1(x)}{\partial x_1} & \cdots & \frac{\partial f_m(x)}{\partial x_1} \\ \vdots & \vdots & \vdots \\ \frac{\partial f_1(x)}{\partial x_n} & \cdots & \frac{\partial f_m(x)}{\partial x_n} \end{pmatrix} \begin{pmatrix} f_1(x) \\ \vdots \\ f_m(x) \end{pmatrix} = f'(x)^T f(x) = D_f(x)^T f(x).$$

Pour la dérivée seconde, si on note  $\psi(x) = 2f_i(x)\frac{\partial f_i(x)}{\partial x_j}$ , on a

$$\frac{\partial^2 \phi(x)}{\partial x_k \partial x_j} = \frac{\partial \psi(x)}{\partial x_k} = 2 \frac{\partial f_i(x)}{\partial x_k} \frac{\partial f_i(x)}{\partial x_j} + 2f_i(x) \frac{\partial^2 f_i(x)}{\partial x_k \partial x_j}.$$

On a alors  $\frac{\partial^2 F(x)}{\partial x_k \partial x_j} = \sum_{i=1}^m \frac{\partial f_i(x)}{\partial x_k} \frac{\partial f_i(x)}{\partial x_j} + f_i(x) \frac{\partial^2 f_i(x)}{\partial x_k \partial x_j}$ . Ce terme est bien le terme  $(k, l)$  de la matrice  $D_f(x)^T D_f(x) + \sum_{i=1}^m f_i(x) \nabla^2 f_i(x)$ .

□

Nous avons vu dans l'exercice 8.30 que pour la fonction des moindres carrés non linéaires,  $F(x) = \frac{1}{2} \|f(x)\|_2$ , le gradient de  $F$  en  $x$  est  $f'(x)^T f(x) = D_f(x)^T f(x)$  et la matrice Hessienne de  $F$  en  $x$  est  $D_f(x)^T D_f(x) + \sum_{i=1}^m f_i(x) \nabla^2 f_i(x)$ . Il est possible donc d'utiliser des variantes de la méthode de Newton pour minimiser  $F(x)$ , en utilisant une recherche linéaire ou une région de confiance.

On remarque que  $\nabla^2 f(x)$  s'écrit sous la forme d'un terme ne faisant intervenir que des dérivations ( $D_f(x)^T D_f(x)$ ) et un terme faisant intervenir des dérivations d'ordre 2 ( $\sum_{i=1}^m f_i(x) \nabla^2 f_i(x)$ ). Il est donc tentant d'approcher  $\nabla^2 f(x)$  par le terme  $D_f(x)^T D_f(x)$  pour éviter le calcul de dérivées d'ordre 2. La variante de Newton faisant cette approximation s'appelle la méthode de Gauss-Newton

$$(GN) : x_{k+1} = x_k - (D_f(x_k)^T D_f(x_k))^{-1} D_f(x_k)^T D_f(x_k) = x_k - D_f(x_k)^+ f(x_k).$$

Cette méthode n'est même pas toujours localement convergente (il existe des points fixes répulsifs). En la globalisant par une recherche linéaire où des régions de confiance on obtient des méthodes globalement convergentes très utilisées en pratique.