Optimisation convexe - Méthodes itératives

Descentes de gradient

Bashar Dudin

April 24, 2019

EPITA

Quand on vous dit la vérité.

Principe des méthodes de descente

Prinicpe généra

Calcule du pas de la descente

Convexité forte et conditionnement de la hessienne

La classique

Les descentes de plus fortes pentes

La méthode de Newton

On s'intéresse dans ces slides à la résolution par des techniques itératives des problèmes d'optimisation 1 .

¹En premier lieu convexes.

On s'intéresse dans ces slides à la résolution par des techniques itératives des problèmes d'optimisation 1 .

La résolution analytique d'équations réelles est rarement aisée. Dans le cas où de telles expressions analytiques existent, elles peuvent être coûteuses à calculer, comme c'est déjà le cas pour la résolution de l'équation normale dans le cas de la régression linéaire.

¹En premier lieu convexes.

On s'intéresse dans ces slides à la résolution par des techniques itératives des problèmes d'optimisation 1 .

La résolution analytique d'équations réelles est rarement aisée. Dans le cas où de telles expressions analytiques existent, elles peuvent être coûteuses à calculer, comme c'est déjà le cas pour la résolution de l'équation normale dans le cas de la régression linéaire. Les solutions sont parfois même impossibles à expliciter à l'aide des fonctions usuelles et des données du problème.

¹En premier lieu convexes.

On s'intéresse dans ces slides à la résolution par des techniques itératives des problèmes d'optimisation 1 .

La résolution analytique d'équations réelles est rarement aisée. Dans le cas où de telles expressions analytiques existent, elles peuvent être coûteuses à calculer, comme c'est déjà le cas pour la résolution de l'équation normale dans le cas de la régression linéaire. Les solutions sont parfois même impossibles à expliciter à l'aide des fonctions usuelles et des données du problème.

Question

On se place dans le cadre de la régression logistique. Pourriez-vous calculer analytiquement une expression d'un point optimal?

¹En premier lieu convexes.

Quand on vous dit la vérité.

Principe des méthodes de descente

Prinicpe général

Calcule du pas de la descente

Convexité forte et conditionnement de la hessienne

La classique

Les descentes de plus fortes pentes

La méthode de Newton



Cadre

Hypothèse

Dans la suite, et à moins de faire explicitement mention du contraire, toutes nos fonctions sont supposées 2-fois différentiables de hessienne continue.

Cadre

Hypothèse

Dans la suite, et à moins de faire explicitement mention du contraire, toutes nos fonctions sont supposées 2-fois différentiables de hessienne continue.

On s'intéresse à un problème d'optimisation de la forme

Le problème (P)

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \tag{P}$$

où sera supposée convexe et vérifiant l'hypothèse de régularité ci-dessus.

Algorithm 1 Principe des descentes de gradient

Input: f : a function, x_0 : an initial point in the domain of f **Output:** x^* : an optimal solution of (P) if bounded from below

- 1: **function** gradient_descent(f, x_0)
- 2: $x \leftarrow x_0$
- 3: while not stopping condition do
- 4: compute a direction Δx to update x
- 5: compute step t > 0 of descent
- 6: $x \leftarrow x + t\Delta x$
- 7: **end while**
- 8: return x
- 9: end function

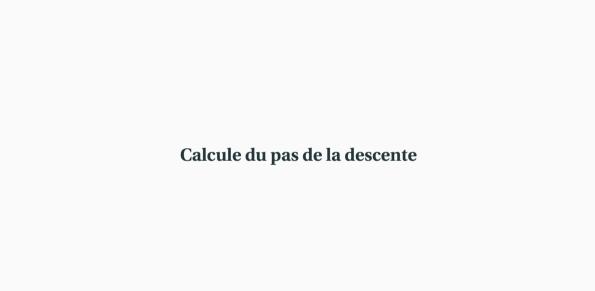
Le principe général précédent nécessite quelques remarques, qu'ils nous faudra garder en mémoire par la suite.

• Chaque choix dans un algorithme de descente impacte la convergence de celui-ci ; l'initialisation, le choix de la direction ou celle du pas.

- Chaque choix dans un algorithme de descente impacte la convergence de celui-ci ; l'initialisation, le choix de la direction ou celle du pas.
- En pratique la boucle while apparaît à l'intérieur d'une boucle finie.

- Chaque choix dans un algorithme de descente impacte la convergence de celui-ci ; l'initialisation, le choix de la direction ou celle du pas.
- En pratique la boucle while apparaît à l'intérieur d'une boucle finie.
- Le pas est souvent (du moins dans un premier temps) pris constant. Cela peut poser des problèmes de convergence.

- Chaque choix dans un algorithme de descente impacte la convergence de celui-ci ; l'initialisation, le choix de la direction ou celle du pas.
- En pratique la boucle while apparaît à l'intérieur d'une boucle finie.
- Le pas est souvent (du moins dans un premier temps) pris constant. Cela peut poser des problèmes de convergence.
- La condition d'arrêt s'exprime souvent par le fait que la mise à jour n'est plus significative.



Il arrive parfois qu'on ne se préoccupe que peu du pas de la descente. Notamment quand on lui garde une valeur constante. C'est un choix qui comporte des dangers, un pas constant peut donner lieu à une descente divergente.

Il arrive parfois qu'on ne se préoccupe que peu du pas de la descente. Notamment quand on lui garde une valeur constante. C'est un choix qui comporte des dangers, un pas constant peut donner lieu à une descente divergente.

Pour éviter cet écueil les matheux s'arment de deux approches:

Il arrive parfois qu'on ne se préoccupe que peu du pas de la descente. Notamment quand on lui garde une valeur constante. C'est un choix qui comporte des dangers, un pas constant peut donner lieu à une descente divergente.

Pour éviter cet écueil les matheux s'arment de deux approches:

• Le calcul du pas optimal : pour une direction choisie on calcule t minimisant $f(x+t\Delta x)$.

Il arrive parfois qu'on ne se préoccupe que peu du pas de la descente. Notamment quand on lui garde une valeur constante. C'est un choix qui comporte des dangers, un pas constant peut donner lieu à une descente divergente.

Pour éviter cet écueil les matheux s'arment de deux approches:

- Le calcul du pas optimal : pour une direction choisie on calcule t minimisant $f(x+t\Delta x)$.
- Le *backtracking* : une heuristique qui mime le précédent point tout en étant moins coûteuse.

Calculer le pas | Backtracking

Algorithm 2 Backtracking

```
Input: f: a function, x: a point in the domain of f
```

Input: Δx a descent direction **Input:** $\alpha \in]0,0.5[, \beta \in]0,1[$

Output: t^* : an optimal point minimizing $f(x + t\Delta x)$ if bounded from below

- 1: **function** backtracking(f, x, $\alpha = 0.1$, $\beta = 0.8$)
- 2: $t \leftarrow 1$
- 3: **while** $f(x + t\Delta x) > f(x) + \alpha t \nabla f(x)^T \Delta x$ **do**
- 4: $t \leftarrow \beta t$
- 5: **end while**
- 6: **return** *t*
- 7: end function



À la fin de l'execution de backtracking on se retrouve dans l'une des deux situations suivantes:

Calculer le pas | Backtracking

À la fin de l'execution de backtracking on se retrouve dans l'une des deux situations suivantes:

• *t* = 1

Calculer le pas | Backtracking

À la fin de l'execution de backtracking on se retrouve dans l'une des deux situations suivantes:

- *t* = 1
- $t \in]\beta t_0, t_0]$ où t_0 est le plus grand réel satisfaisant la condtion de boucle.



On désgine par S l'ensemble $\{x \in \text{Dom}(f) \mid f(x) \leq f(x_0)\}$ où x_0 est sous-entendu être un point initial de descente de gradient. C'est un fermé de \mathbb{R}^n ; toute suite de S convergente dans \mathbb{R}^n a une limite dans S.

On désgine par S l'ensemble $\{x \in \text{Dom}(f) \mid f(x) \leq f(x_0)\}$ où x_0 est sous-entendu être un point initial de descente de gradient. C'est un fermé de \mathbb{R}^n ; toute suite de S convergente dans \mathbb{R}^n a une limite dans S.

L'analyse de convergence, sous-entendu de la vitesse de convergence, des méthodes de descentes s'effectue, outre l'hypothèse de régularité \mathscr{C}^2 , sous l'une des deux conditions suivantes:

1. La hessienne de f est majorée sur S.

On désgine par S l'ensemble $\{x \in \text{Dom}(f) \mid f(x) \leq f(x_0)\}$ où x_0 est sous-entendu être un point initial de descente de gradient. C'est un fermé de \mathbb{R}^n ; toute suite de S convergente dans \mathbb{R}^n a une limite dans S.

L'analyse de convergence, sous-entendu de la vitesse de convergence, des méthodes de descentes s'effectue, outre l'hypothèse de régularité \mathscr{C}^2 , sous l'une des deux conditions suivantes:

- 1. La hessienne de f est majorée sur S.
- **2.** f est *strictement convexe* sur S; c'est-à-dire qu'il existe m > 0 tel que pour tout $x \in S$, $\nabla^2 f(x) \ge mI$. C'est une inégalité fonctionnelle entre formes quadratiques!

On désgine par S l'ensemble $\{x \in \text{Dom}(f) \mid f(x) \leq f(x_0)\}$ où x_0 est sous-entendu être un point initial de descente de gradient. C'est un fermé de \mathbb{R}^n ; toute suite de S convergente dans \mathbb{R}^n a une limite dans S.

L'analyse de convergence, sous-entendu de la vitesse de convergence, des méthodes de descentes s'effectue, outre l'hypothèse de régularité \mathscr{C}^2 , sous l'une des deux conditions suivantes:

- 1. La hessienne de f est majorée sur S.
- **2.** f est *strictement convexe* sur S; c'est-à-dire qu'il existe m > 0 tel que pour tout $x \in S$, $\nabla^2 f(x) \ge mI$. C'est une inégalité fonctionnelle entre formes quadratiques!

La seconde condition est la plus forte des deux ; elle implique la première.

Définition

Si $f: \mathbb{R}^n \to \mathbb{R}$ est une fonction convexe \mathscr{C}^2 , le nombre de conditionnement de $\nabla^2 f(x)$ est le rapport de sa plus grande valeur propre à sa plus petite.

Définition

Si $f: \mathbb{R}^n \to \mathbb{R}$ est une fonction convexe \mathscr{C}^2 , le nombre de conditionnement de $\nabla^2 f(x)$ est le rapport de sa plus grande valeur propre à sa plus petite.

• La hessienne de f étant symétrique positive et à valeurs réelles, elle est diagonalisable sur $\mathbb R$ à valeurs propres positives.

Définition

Si $f: \mathbb{R}^n \to \mathbb{R}$ est une fonction convexe \mathscr{C}^2 , le nombre de conditionnement de $\nabla^2 f(x)$ est le rapport de sa plus grande valeur propre à sa plus petite.

- La hessienne de f étant symétrique positive et à valeurs réelles, elle est diagonalisable sur $\mathbb R$ à valeurs propres positives.
- Le nombre de conditionnement est ≥ 1 .

Définition

Si $f: \mathbb{R}^n \to \mathbb{R}$ est une fonction convexe \mathscr{C}^2 , le nombre de conditionnement de $\nabla^2 f(x)$ est le rapport de sa plus grande valeur propre à sa plus petite.

- La hessienne de f étant symétrique positive et à valeurs réelles, elle est diagonalisable sur $\mathbb R$ à valeurs propres positives.
- Le nombre de conditionnement est ≥ 1 .
- ullet Le rapport des bornes qui encadrent f dans le cas strictement convexe est un majorant des nombres de conditionnement.

Définition

Si $f: \mathbb{R}^n \to \mathbb{R}$ est une fonction convexe \mathscr{C}^2 , le nombre de conditionnement de $\nabla^2 f(x)$ est le rapport de sa plus grande valeur propre à sa plus petite.

- La hessienne de f étant symétrique positive et à valeurs réelles, elle est diagonalisable sur $\mathbb R$ à valeurs propres positives.
- Le nombre de conditionnement est ≥ 1 .
- Le rapport des bornes qui encadrent f dans le cas strictement convexe est un majorant des nombres de conditionnement.

Convergence

Les nombres de conditionnements de f sont corrélés à la vitesse de convergence de la descente de gradient.

Quand on vous dit la vérité.

Principe des méthodes de descente

Prinicpe général

Calcule du pas de la descente

Convexité forte et conditionnement de la hessienne

La classique

Les descentes de plus fortes pentes

La méthode de Newton

La descente de gradient à l'ancienne

Algorithm 3 Descente de gradient

Input: f : a function, x_0 : an initial point in the domain of f, ε : tolerance **Output:** x^* : an optimal solution of (P) if bounded from below

1: **function** GRADIENT_DESCENT (f, x_0, ε) 2: $x \leftarrow x_0$ 3: **while** $\|\nabla f(x)\| > \varepsilon$ **do** 4: $\Delta x \leftarrow -\nabla f(x)$ 5: compute step t > 0 of descent 6: $x \leftarrow x + t\Delta x$ 7: **end while** 8: **return** x9: **end function**

Proposition

Supposons $f \mathscr{C}^2$ ayant une hessienne majorée ; il existe $M \in \mathbb{R}_+$ tel que pour tout $x \in S$, $\nabla^2 f(x)$. La descente de gradient avec un pas constant $t \leq \frac{1}{M}$ ou via backtracking garantit

$$|f(x_k) - f(x^*)| \le \frac{\|x_0 - x^*\|_2^2}{2ck}$$

où

- x_0 est le point initial de la descente
- x_k le k-ème itéré de la descente
- c est égal à t dans le cas du pas constant et à min $\{1, \frac{\beta}{M}\}$ dans le cas du backtracking.

$$|f(x_k) - f(x^*)| \le \frac{\|x_0 - x^*\|_2^2}{2ck}$$

• La valeur absolue du membre de gauche est superflue.

$$|f(x_k) - f(x^*)| \le \frac{\|x_0 - x^*\|_2^2}{2ck}$$

- La valeur absolue du membre de gauche est superflue.
- La vitesse de convergence dépend du point initial (étonnant ...).

$$|f(x_k) - f(x^*)| \le \frac{\|x_0 - x^*\|_2^2}{2ck}$$

- La valeur absolue du membre de gauche est superflue.
- $\bullet\,$ La vitesse de convergence dépend du point initial (étonnant ...).
- Pour un point sous-optimal à ε près on est sur une complexité en $\mathcal{O}\left(\frac{1}{\varepsilon}\right)$.

$$|f(x_k) - f(x^*)| \le \frac{\|x_0 - x^*\|_2^2}{2ck}$$

- La valeur absolue du membre de gauche est superflue.
- La vitesse de convergence dépend du point initial (étonnant ...).
- Pour un point sous-optimal à ε près on est sur une complexité en $\mathcal{O}\left(\frac{1}{\varepsilon}\right)$.
- L'intérêt du *backtracking* réside dans le fait qu'on n'a pas à calculer le pas à la main. Si β est proche de 1 on ne perd pas grand chose par rapport au pas constant.

Proposition

Supposons $f \mathcal{C}^2$ fortement convexe ; il existe $m, M \in \mathbb{R}_+$ encadrant la hessienne en tout point de S. La descente de gradient à pas constant $t \leq \frac{2}{M}$ ou via backtracking garantit

$$|f(x_k) - f(x^*)| \le c^k ||x_0 - x^*||_2$$

avec $c \in]0,1[$.

- x_0 est le point initial de la descente
- x_k le k-ème itéré de la descente
- c est égal à $(1-\frac{m}{M})$ dans le cas du pas constant et à $1-\min\{2m\alpha,2\beta\alpha\frac{m}{M}\}$ dans le cas du backtracking.

$$|f(x_k) - f(x^*)| \le c^k ||x_0 - x^*||_2$$

• La vitesse de convergence qu'on obtient ainsi est dite linéaire ; elle l'est si graphée contre une échelle logarithmique.

$$|f(x_k) - f(x^*)| \le c^k ||x_0 - x^*||_2$$

- La vitesse de convergence qu'on obtient ainsi est dite linéaire ; elle l'est si graphée contre une échelle logarithmique.
- Pour un point sous-optimal à ε près on est sur une complexité en $\mathcal{O}(\ln(\frac{1}{\varepsilon}))$.

$$|f(x_k) - f(x^*)| \le c^k ||x_0 - x^*||_2$$

- La vitesse de convergence qu'on obtient ainsi est dite linéaire ; elle l'est si graphée contre une échelle logarithmique.
- Pour un point sous-optimal à ε près on est sur une complexité en $\mathcal{O}(\ln(\frac{1}{\varepsilon}))$.
- La constante c dépend très fortement de $\frac{M}{m}$.

Quand on vous dit la vérité.

Principe des méthodes de descente

Prinicpe général

Calcule du pas de la descente

Convexité forte et conditionnement de la hessienne

La classique

Les descentes de plus fortes pentes

La méthode de Newton

Généraliser la descente classique

L'intuition qui mène à la descente de gradient classique répond en réalité à un problème de minimisation.

Généraliser la descente classique

L'intuition qui mène à la descente de gradient classique répond en réalité à un problème de minimisation. Étant donné un point $x \in S$ et un vecteur v de norme assez petite, on peut écrire

$$f(x + \nu) = f(x) + \nabla f(x)^T \nu + o(\nu).$$

La direction de descente qui minimise au plus la valeur objective est donnée par

$$\Delta x_{nsd} = \operatorname{argmax} \left\{ \nabla f(x)^T v \mid ||v||_2 = 1 \right\}.$$

C'est une conséquence de l'inégalité de Cauchy-Shwarz.

La descente de plus forte pente

La démarche précédente nous permet de définir différentes variantes des descentes de gradients ; en réalité l'essentielle des descentes de gradients.

La descente de plus forte pente

La démarche précédente nous permet de définir différentes variantes des descentes de gradients ; en réalité l'essentielle des descentes de gradients.

Définition

Soient $x \in S$ et $\|\cdot\|$ une norme sur \mathbb{R}^n . On désigne par Δx_{nsd} (nsd pour *nomralized steepest descent*) la quantité

$$\Delta x_{nsd} = \operatorname{argmax} \left\{ \Delta f(x)^T v \mid ||v|| = 1 \right\}$$

La direction de descente de plus forte pente (sous-entendu pour la norme $\|\cdot\|)$ est donnée par

$$\Delta_{sd} = \|\nabla f(x)\| \Delta x_{sd}.$$

La descente de plus forte pente | Géométrie

La descente de plus forte pente pour une norme $\|\cdot\|$ s'interprète comme

le vecteur de plus grande projection sur $-\nabla f(x)$.

La descente de plus forte pente | Géométrie

La descente de plus forte pente pour une norme $\|\cdot\|$ s'interprète comme

le vecteur de plus grande projection sur $-\nabla f(x)$.

Pour tout vecteur v dans la sphère unité pour la norme $\|\cdot\|$, $\nabla f(x)^T v$ a pour valeur absolue la norme de la projection orthogonale de v sur $\nabla f(x)$. Dans la mesure où l'on cherche un minimisant c'est la plus grande projection contre $-\nabla f(x)$.

La descente de gradient de plus forte pente

Algorithm 4 Descente de gradient de plus forte pente

Input: f: a function, x_0 : an initial point in the domain of f, ε : tolerance, $\|\cdot\|$ une norme sur \mathbb{R}^n . **Output:** x^* : an optimal solution of (P) if bounded from below

```
1: function STEEPEST_GRADIENT_DESCENT(f, x_0, \varepsilon, \| \cdot \|)

2: x \leftarrow x_0

3: while \|\nabla f(x)\| > \varepsilon do

4: Compute steepest descent direction \Delta x_{sd} for \| \cdot \|.

5: compute step t > 0 of descent

6: x \leftarrow x + t\Delta x_{sd}

7: end while

8: return x

9: end function
```

La descente de plus forte pente | Analyse

• La stratégie de descente de plus forte pente permet de varier les types de descentes de gradients. Les géométries sous-jacentes sont parfois plus adaptées à certains problèmes qu'à d'autres.

La descente de plus forte pente | Analyse

- La stratégie de descente de plus forte pente permet de varier les types de descentes de gradients. Les géométries sous-jacentes sont parfois plus adaptées à certains problèmes qu'à d'autres.
- L'analyse de la convergence dans le cas de descente classique s'étend au cas de plus forte pente. La raison en est l'équivalence des différentes normes sur \mathbb{R}^n . L'essentiel des propriétés utilisés lors des encadrements étant partagées par celles-ci.

Quand on vous dit la vérité.

Principe des méthodes de descente

Prinicpe généra

Calcule du pas de la descente

Convexité forte et conditionnement de la hessienne

La classique

Les descentes de plus fortes pentes

La méthode de Newton

Les descentes de gradient qu'on a pu voir jusqu'à présent sont des méthodes de descentes dites de premier ordre ; elle minimise sur la boule unité pour une norme donnée l'approximation au premier ordre de la fonction objectif.

Les descentes de gradient qu'on a pu voir jusqu'à présent sont des méthodes de descentes dites de premier ordre ; elle minimise sur la boule unité pour une norme donnée l'approximation au premier ordre de la fonction objectif.

La méthode de Newton est une méthode de second ordre ; on cherche à mininiser l'approximation de second ordre de la fonction objectif en un itéré.

Soit f une fonction objectif convexe et x un point du domaine de f. Le DL de f au second ordre en x s'écrit

$$f(x + v) = f(x) + \nabla f(x)^{T} v + \frac{1}{2} v^{T} \nabla^{2} f(x) v + ||v||^{2} \varepsilon(v)$$

Soit f une fonction objectif convexe et x un point du domaine de f. Le DL de f au second ordre en x s'écrit

$$f(x + v) = f(x) + \nabla f(x)^{T} v + \frac{1}{2} v^{T} \nabla^{2} f(x) v + ||v||^{2} \varepsilon(v)$$

On choisit d'approcher f(x + v) par l'expression de second ordre

$$f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v$$

C'est une fonction convexe en v qu'on sait minimiser. On obtient ici un minimisant donné par

$$\Delta x_N = -\left(\nabla f(x)\right)^{-1} \nabla f(x).$$

Méthode de Newton - Condition d'arrêt

Traditionnellement, la condition d'arrêt de la méthode de Newton est décrite par le fait que le carré du coefficient de décroissance suivant tombe sous un certain seuil de tolérance.

$$\lambda(x) = \left(\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x)\right)^{1/2}.$$

Méthode de Newton - Condition d'arrêt

Traditionnellement, la condition d'arrêt de la méthode de Newton est décrite par le fait que le carré du coefficient de décroissance suivant tombe sous un certain seuil de tolérance.

$$\lambda(x) = \left(\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x)\right)^{1/2}.$$

Celui-ci estime l'écart entre le minimisant de l'approximation du second ordre de f en un point et la valeur de f en ce même point.

Algorithme de Newton

Algorithm 5 Méthode de Newton

Input: f: a function, x_0 : an initial point in the domain of f, ε : tolerance, $\|\cdot\|$ une norme sur \mathbb{R}^n . **Output:** x^* : an optimal solution of (P) if bounded from below

- 1: **function** Newton_Method(f, x_0 , ε , $\|\cdot\|$)
- $2: \quad x \leftarrow x_0$
- 3: $\Delta x_N \leftarrow -(\nabla^2 f(x))^{-1} \nabla f(x)$
- 4: $\lambda^2(x) = -\nabla f(x)^T \Delta x_N$ 5: **while** $\frac{\lambda^2(x)}{2} > \varepsilon$ **do**
- 6: $\Delta x_N \leftarrow -(\nabla^2 f(x))^{-1} \nabla f(x)$
- 7: $\lambda^2(x) = -\nabla f(x)^T \Delta x_N$
- 8: compute step t > 0 of descent
- 9: $x \leftarrow x + t\Delta x_N$
- 10: end while
- 11: **return** x

