

# OUTILS D'OPTIMISATION POUR LES SDD

7 octobre 2024

Aujourd'hui :

C7 Méthode de gradient (+ illustration)

Vendredi : TD Exercices d'annales

# NÉTHODES DE GRADIENT

↳ Classe d'algorithmes d'optimisation qui forme la base des méthodes utilisées aujourd'hui en sciences des données

⇒ Algorithme phare: la descente de gradient

## ① Descente de gradient

Cadre de travail: minimiser  $f(w)$ , avec  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  de classe  $C^1$

( $\forall w \in \mathbb{R}^d$ , la dérivée de  $f$  en  $w$  existe et est représentée par le gradient  $\nabla f(w) \in \mathbb{R}^d$ )

On a vu que

$$\left[ \bar{w} \in \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} f(w) \right] \Rightarrow \text{ensemble des solutions du problème (minima globaux)}$$

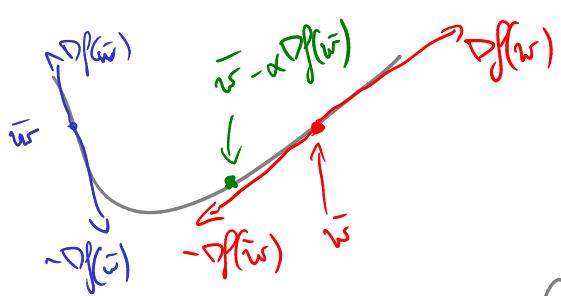
$$\begin{aligned} \nabla f(\bar{w}) &= 0_{\mathbb{R}^d} \\ \Leftrightarrow \|\nabla f(\bar{w})\| &= 0 \end{aligned}$$

Cette implication s'écrit de manière équivalente

$$\left[ \|\nabla f(\bar{w})\| \neq 0 \right] \Rightarrow \left[ \bar{w} \notin \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} f(w) \right]$$
$$\Leftrightarrow \left[ \exists w^+ \in \mathbb{R}^d, f(w^+) < f(\bar{w}) \right]$$

On peut montrer que si  $\|\nabla f(\bar{w})\| \neq 0$ , alors il existe  $w^+ \in \mathbb{R}^d$  de la forme  $w^+ = \bar{w} - \alpha \nabla f(\bar{w})$  avec  $\alpha > 0$  tel que  $f(w^+) < f(\bar{w})$

⇒ On peut diminuer la valeur de  $f$  en bougeant dans la direction opposée au gradient!

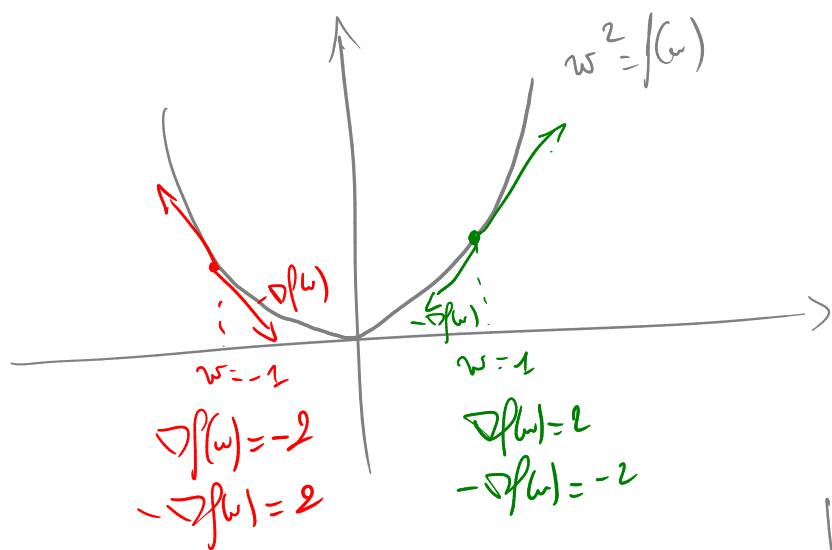


Remarques

- L'observation ci-dessus n'est valable qu'au voisinage de  $\bar{w}$  lorsque  $\|w - \bar{w}\|$  est suffisamment petite

$$(f(w) \approx f(\bar{w}) + Df(\bar{w})^T(w - \bar{w}))$$

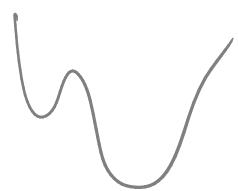
- On dit que la direction opposée au gradient est la direction de plus forte pente



$$Df(w) = 2w$$

$$Df(w) > 0 \text{ si } w > 0$$

$$Df(w) < 0 \text{ si } w < 0$$



## Algorithme de descente de gradient (pseudo-code)

Initialisation: Choisir  $w_0 \in \mathbb{R}^d$  et  $\alpha_0 > 0$ .

$w_0$  ↑  
point initial

$\alpha_0$  ↑  
longueur de pas initiale

Pour  $k=0, 1, 2, \dots$

Calculer  $w_{k+1} = w_k - \alpha_k Df(w_k)$  ← Itération de la descente de gradient

Choisir  $\alpha_{k+1} > 0$  ← longueur de pas de l'itération  $k+1$

Pour implémenter cet algorithme, on considère deux aspects fondamentaux :

- critère d'arrêt : souvent une combinaison de 2 critères :

### 1) Critère de convergence :

Ex) s'arrêter lorsque  $\|\nabla f(\omega_k)\| < \epsilon$  avec  $\epsilon > 0$

Ex) s'arrêter lorsque  $f(\omega_k) - \min_{\omega \in \Omega} f(\omega) < \epsilon$  avec  $\epsilon > 0$

Remarque

$\leq \epsilon$  ou  $\leq \epsilon$

sont tous les deux utilisés en pratique

⚠ Demande de connaître la valeur optimale

↑ valeur optimale du problème

Ex) s'arrêter lorsque  $\|\omega_k - \omega^*\| < \epsilon$  avec  $\epsilon > 0$

où  $\omega^* = \arg \min_{\omega \in \Omega} f(\omega)$

⚠ Un tel  $\omega^*$  n'est pas connu en général

### 2) Critère de budget

Ex) s'arrêter au bout de  $K$  itérations ( $K \geq 1$ )

Ex) s'arrêter au bout de  $N$  évaluations de  $f$  et/ou son gradient (peut différer du nombre d'itérations en fonction de la méthode utilisée)

Ex) s'arrêter au bout d'un certain temps  $\left\{ \begin{array}{l} \text{CPU / GPU} \\ \text{tout court} \end{array} \right.$

### Stratégie de choix de longueur de pas

(En anglais, longeur de pas = step size / learning rate )  
en ML

$$\text{Pas : } \omega_{k+1} - \omega_k = -\alpha_k \nabla f(\omega_k)$$

Grandes familles de stratégies :

1) Longeur de pas constante :  $\alpha_k = \alpha_0 = \alpha > 0 \quad \forall k \in \mathbb{N}$

- ⊕ Simple à calculer
- ⊕ Il existe des valeurs pour lesquelles l'algorithme va converger
- ⊕ Choix facile à calibrer pour un problème spécifique

- ⊖ Si  $\alpha$  est mal choisi, l'algorithme peut converger très lentement, ou même diverger
- ⊖ Les valeurs théoriques pour lesquelles l'algorithme converge peuvent être impossibles à calculer en pratique

2) Longeur de pas décroissante :  $\{\alpha_k\}_k \quad \alpha_k > 0 \quad \alpha_k \searrow 0$

- ⊕ Garantir que  $f(w_{k+1}) < f(w_k)$  pour  $k$  suffisamment grand
- ⊕ Calculable avant de lancer l'algorithme (pas besoin de calculer  $\alpha_{k+1}$  à chaque itération)

- ⊖ Si  $\{\alpha_k\}_k$  décroît trop rapidement, alors la méthode converge trop lentement et proche de pas numériquement égaux à 0
- ⊖ Difficile de choisir la vitesse de décroissance adaptée à un problème donné a priori

3) Longeur de pas adaptative

↳ A chaque itération, on choisit  $\alpha_k$  en fonction de  $w_k, f(w_k), \nabla f(w_k)$ , c'est à dire en fonction de l'information disponible au point courant  $w_k$

- ⊕ Produit en général de meilleures valeurs de  $\alpha_k$  (en termes de décroissance de  $f$ ) que des stratégies a priori
- ⊖ Coûte en général plus cher en appels à  $f$  et/ou à  $\nabla f$

Ex) Recherche linéaire avec retour arrière (backtracking)  
d'Armijo

$\alpha = 1$  ① Définir  $\alpha > 0$ .

par ex ② Partant de  $w_k$ , on calcule  $f(w_k - \alpha \nabla f(w_k))$

$c = 10^{-4}$  ③ Si  $f(w_k - \alpha \nabla f(w_k)) < f(w_k) - c \alpha \|\nabla f(w_k)\|^2$  (où  $c \in (0, \frac{1}{2})$ )

alors poser  $w_{k+1} = w_k - \alpha \nabla f(w_k)$

Condition d'Armijo dite  
de décroissance suffisante

$\theta = \frac{1}{2}$  Sinon mettre à jour  $\alpha \leftarrow \theta \alpha$  avec  $\theta \in (0, 1)$   
et retourner en ②

Exemple d'application: Problèmes canons linéaires

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad f(w) = \frac{1}{2m} \|Xw - y\|^2 \quad X \in \mathbb{R}^{m \times d} \quad y \in \mathbb{R}^m$$

↪ Pour ce problème, si  $\nabla f(w) \neq 0_{\mathbb{R}^d}$ , alors

$$f(w - \alpha \nabla f(w)) \leq f(w) - \frac{1}{2} \alpha \|\nabla f(w)\|^2$$

$$\text{avec } \alpha = \frac{1}{L} \text{ et } L = \frac{\|X^T X\|}{m}$$

## ② Analyse théorique

↪ On veut analyser l'algorithme de gradient et déterminer s'il converge ou non vers une solution du problème  $\Rightarrow$  Sous quelles conditions la méthode converge-t-elle ?

↪ Les résultats de convergence pour la descente de gradient sont toujours valides :

- pour une classe de problèmes (de fonctions  $f$ )
- pour un choix de longueur de pas (constante, décroissante, etc.)

$\Rightarrow$  Pour aujourd'hui, on regardera les fonctions  $C^1$  à gradient Lipschitz et une logique de pas constante bien choisie.

Définition: Soit  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  une fonction de classe  $C^1$ .

On dit que  $f$  est  $C^1$  à gradient  $L$ -lipschitzien avec  $L > 0$  (ce que l'on note  $f \in C_L^{1,1}$ ) si

$$\forall (v, w) \in (\mathbb{R}^d)^2, \quad \|\nabla f(v) - \nabla f(w)\| \leq L \|v - w\|$$

Ex).  $f(w) = \underbrace{\frac{1}{2m} \|Xw - y\|^2}_{\text{régression linéaire}}$  est  $C_L^{1,1}$  avec  $L = \frac{\|X^T X\|}{m}$

. Autres exemples: régression logistique, fonctions quadratiques, certains types de réseaux de neurones avec contraintes sur les paramètres

Remarque: les résultats pour les fonctions  $C_L^{1,1}$  s'adaptent à de nombreuses classes de fonctions

Proposition: Si  $f$  est  $C_L^{1,1}$  (de  $\mathbb{R}^d$  dans  $\mathbb{R}$ ), alors

$$\forall (v, w) \in (\mathbb{R}^d)^2, \quad f(v) \leq f(w) + \nabla f(w)^T (v - w) + \frac{L}{2} \|v - w\|^2$$

Inégalité fondamentale pour les fonctions  $C_L^{1,1}$   
en optimisation

$\hookrightarrow$  L'inégalité fondamentale permet de trouver des conditions sur  $v$  qui garantissent que  $f(v) < f(w)$  lorsque  $\nabla f(w) \neq 0_{\mathbb{R}^d}$

Théorème: Soit  $f: \mathbb{R}^d \rightarrow \mathbb{R}$   $C_L^1$  avec  $L > 0$ . Soit  $w \in \mathbb{R}^d$  tel que

$\nabla f(w) \in \mathbb{R}^d$ . Alors

$$f\left(w - \frac{1}{L} \nabla f(w)\right) \leq f(w) - \frac{1}{2L} \|\nabla f(w)\|^2 < f(w)$$

Démonstration: On applique l'inégalité fondamentale avec  $v = w - \frac{1}{L} \nabla f(w)$

$$\begin{aligned} f\left(w - \frac{1}{L} \nabla f(w)\right) &\leq f(w) + \nabla f(w)^T \left( w - \frac{1}{L} \nabla f(w) - w \right) + \frac{L}{2} \left\| w - \frac{1}{L} \nabla f(w) - w \right\|^2 \\ &= f(w) + \nabla f(w)^T \left( -\frac{1}{L} \nabla f(w) \right) + \frac{L}{2} \left\| -\frac{1}{L} \nabla f(w) \right\|^2 \\ &= f(w) - \frac{1}{L} \nabla f(w)^T \nabla f(w) + \frac{L}{2} \times \left( -\frac{1}{L} \right)^2 \|\nabla f(w)\|^2 \\ &= f(w) - \frac{1}{L} \nabla f(w)^T \nabla f(w) + \frac{1}{2L} \|\nabla f(w)\|^2 \\ &= f(w) - \frac{1}{L} \|\nabla f(w)\|^2 + \frac{1}{2L} \|\nabla f(w)\|^2 \\ &= f(w) - \frac{1}{2L} \|\nabla f(w)\|^2 \end{aligned}$$

Remarque: On peut montrer plus généralement que  $\forall \alpha \in (0, \frac{2}{L})$ ,

$$f(w - \alpha \nabla f(w)) \leq f(w) - \frac{1}{2} \alpha \|\nabla f(w)\|^2 < f(w)$$

$\Rightarrow$  lorsque la valeur de  $L$  est inconnue ou impossible à calculer, la recherche linéaire d'Armijo permet de satisfaire une condition de décaissement similaire

$$(cf. plus haut): f(w - \alpha \nabla f(w)) \leq f(w) - c \alpha \|\nabla f(w)\|^2 \text{ avec } c \in (0, \frac{1}{2})$$

↳ le théorème permet de garantir  $f(\bar{w}_{k+1}) < f(w_k)$  lorsque  $\alpha_k = \frac{1}{L}$   
 $\Rightarrow$  Propriété propre à l'itération k

Pour analyser l'algorithme en entier, on se donne un critère de convergence et on borne le nombre d'itérations nécessaire pour satisfaire ce critère  $\Rightarrow$  **Borne de complexité**

## Théorème (Complexité de la descente de gradient dans le cas $C_L^{1,1}$ )

On considère minimiser  $f(w)$  avec  $f: \mathbb{R}^d \rightarrow \mathbb{R}$   $C_L^{1,1}$ .  
 $w \in \mathbb{R}^d$

On applique la descente de gradient en partant de  $w_0 \in \mathbb{R}^d$   
avec le pas  $\alpha_k = \frac{1}{L} \quad \forall k \in \mathbb{N}$ .

Alors, pour tout  $\varepsilon > 0$ , l'algorithme calcule  $\bar{w}_k$  tel que  
 $\| \nabla f(\bar{w}_k) \|_H \leq \varepsilon$  en au plus  $\underbrace{\left[ 2L \left( f(w_0) - \min_{w \in \mathbb{R}^d} f(w) \right) \varepsilon^{-2} \right]}_{= O(\varepsilon^{-2})} \text{ itérations}$

$\rightarrow$  Borne de complexité: fonction de  $\varepsilon$  (précision demandée)

On s'intéresse à comment cette fonction augmente quand  $\varepsilon$  diminue

$$O(\varepsilon^{-2}) = C \varepsilon^{-2}$$

$\varepsilon = 10^{-2} \rightarrow 10^4 \text{ C itérations}$   
 $\varepsilon = 10^{-3} \rightarrow 10^6 \text{ C itérations}$   
 divise la précision par 10  
 $\Rightarrow$  multiplie le nombre d'itérations

$\rightarrow$  On dira: "La complexité de la descente de gradient sur une fonction  $C_L^{1,1}$  est en  $O(\varepsilon^{-2})$ "  
 (s'étend au-delà du choix  $\alpha_k = \frac{1}{L}$ )

Esquisse de preuve: On a vu précédemment que

$$f(w_{k+1}) = f\left(w_k - \frac{1}{L} \nabla f(w_k)\right) \leq f(w_k) - \frac{1}{2L} \|\nabla f(w_k)\|^2$$

Supposons que  $\|\nabla f(w_k)\| \geq \varepsilon \quad \forall k=0, \dots, K-1$

Alors,  $\forall k=0, \dots, K-1$ , on a

$$f(w_{k+1}) \leq f(w_k) - \frac{1}{2L} \|\nabla f(w_k)\|^2 \leq f(w_k) - \frac{1}{2L} \varepsilon^2$$

$$\text{d'où} \quad \frac{1}{2L} \varepsilon^2 \leq f(w_k) - f(w_{k+1}) \quad \forall k=0, \dots, K-1$$

En sommant ces inégalités, on obtient

$$\sum_{k=0}^{K-1} \frac{1}{2L} \varepsilon^2 \leq \sum_{k=0}^{K-1} (f(w_k) - f(w_{k+1}))$$

$\uparrow$   
Somme telescopicque  
 $f(w_0) - \underbrace{f(w_1) + f(w_1) - f(w_2) + \dots + f(w_{K-1}) - f(w_K)}_{=0}$

$$\begin{aligned} f(w_K) &\geq \min_{w \in \mathbb{R}^d} f(w) \\ -f(w_K) &\leq -\min_{w \in \mathbb{R}^d} f(w) \end{aligned}$$

$$\sum_{k=0}^{K-1} (f(w_k) - f(w_{k+1})) = f(w_0) - f(w_K)$$

$$\text{D'où} \quad K \left( \frac{1}{2L} \varepsilon^2 \right) \leq f(w_0) - f(w_K) \leq f(w_0) - \min_{w \in \mathbb{R}^d} f(w)$$

$$(=) \quad K \leq 2L \left( f(w_0) - \min_{w \in \mathbb{R}^d} f(w) \right) \varepsilon^{-2}$$

On veut démontrer que si  $\|\nabla f(w_k)\| > \varepsilon \quad \forall k=0, \dots, K-1$ , alors  $K \leq 2L \left( f(w_0) - \min_{w \in \mathbb{R}^d} f(w) \right) \varepsilon^{-2}$

Pour conséquent, si  $K \geq \lceil 2L \left( f(w_0) - \min_{w \in \mathbb{R}^d} f(w) \right) \varepsilon^{-2} \rceil$

( $\Gamma A$ : plus petit entre supérieur égal à  $A$ )

il existe  $k \in \{0, \dots, k-1\}$  tel que  $\|\nabla f(w_k)\| < \varepsilon$ .

↳ lorsque la fonction  $f$  à minimiser est convexe et plus d'être  $C_L^{1,1}$ , on peut montrer des résultats plus forts sur la descente de gradient. avec longueur de pas  $\alpha_k = \frac{1}{L}$  (l'algorithme ne change pas par rapport au théorème précédent !)

Propriétés de $f$	Critère de convergence	borne de complexité
$f \in C_L^{1,1}$	$\ \nabla f(w_k)\  < \varepsilon$	$O(\varepsilon^{-2})$ itérations
$f \in C_L^{1,1}$ + convexe	$f(w_k) - \min_{w \in \mathbb{R}^d} f(w) < \varepsilon$	$O(\varepsilon^{-1})$ itérations
$f \in C_L^{1,1}$ + $\mu$ -fortement convexe	$\ w_k - w^*\  < \varepsilon$ avec $w^* \in \arg\min_{w \in \mathbb{R}^d} f(w)$	$O\left(\frac{L}{\mu} \ln(\varepsilon^{-1})\right)$ itérations $= O(\ln(\varepsilon^{-1}))$ itérations

↳ Interprétation des bornes: Pour le même algorithme, on obtient de meilleures garanties (critère + borne de complexité) lorsque la fonction est fortement convexe que lorsque elle est convexe, et de meilleures garanties lorsque la fonction est convexe que lorsque elle ne l'est pas.

⇒ pour cette raison, on dit qu'il est plus facile d'optimiser des fonctions fortement convexes que des fonctions convexes (et que des fonctions non convexes)

Def: Une fonction  $f: \mathbb{R}^d \rightarrow \mathbb{R}$   $C^1$  est dite  $\mu$ -fortement convexe avec  $\mu > 0$  si

$$\forall (v, w) \in (\mathbb{R}^d)^2, \quad f(v) \geq f(w) + \nabla f(w)^T(v-w) + \frac{\mu}{2} \|v-w\|^2$$

( $\hookrightarrow$  Autre inégalité fondamentale)

$\hookrightarrow$  Si  $f$  est  $C^1$  et  $\mu$ -fortement convexe, alors:

$$f(w) + \nabla f(w)^T(v-w) + \frac{L}{2} \|v-w\|^2 \geq f(v) \geq f(w) + \nabla f(w)^T(v-w) + \frac{\mu}{2} \|v-w\|^2$$

$\Rightarrow$  les deux inégalités sont utiles pour montrer que la descente de gradient vérifie  $\|w_k - w^*\| \leq \varepsilon$  en  $O\left(\frac{L}{\mu} \ln(\varepsilon')\right)$  itérations

Propriété importante: Si  $f$   $C^1$  et  $\mu$ -fortement convexe, alors elle possède un unique minimum global, qui est l'unique point  $w^* \in \mathbb{R}^d$  tel que  $\|\nabla f(w^*)\| = 0$

On note alors  $w^* = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} f(w)$



Difference entre convexe et fortement convexe

- Toute fonction fortement convexe est convexe.

- L'inverse n'est pas vrai (prendre  $f$  constante)

- Une fonction convexe vérifie

$$f(v) \geq f(w) + \nabla f(w)^T(v-w) \quad \forall (v, w) \quad (\approx \mu=0)$$

- Une fonction convexe peut avoir plusieurs minima globaux