

Outils d'optimisation pour les sciences des données et de la décision

21 octobre 2024

Aujourd'hui: CM

Lundi 28/10: TD avec notebook
Point sur le projet

Vendredi 08/11: TD (annales)

Au-delà de la descente de gradient

→ Descente de gradient: algorithme d'optimisation qui est applicable à toute fonction C^1 (qui admet un gradient en tout point), que cette fonction soit convexe ou non convexe

Itération

Pour un problème

$$\text{minimiser } f(w) \quad w \in \mathbb{R}^d$$

avec $f: (\mathbb{R}^d \rightarrow \mathbb{R})^{C^1}$,

l'itération k de la descente de gradient s'écrit

$$w_{k+1} = w_k - \alpha_k \nabla f(w_k)$$

\downarrow
 $\alpha_k > 0$
taille de pas

↳ Gradient de f en w_k
(vecteur de \mathbb{R}^d)

Complexité

Etant donnée une précision $\varepsilon > 0$, quel est le nombre d'itérations de descente de gradient nécessaire pour atteindre cette précision ?

(Ex) Boîne de complexité dans le cas convexe

Si $f: C^1$ convexe, alors la descente de gradient calcule w_k tel que $f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq \varepsilon$ en au plus $O(\varepsilon^{-1})$ itérations.

\uparrow
valeur optimale du problème

$C \varepsilon^{-1}$, où $C > 0$ est une constante qui ne dépend pas de ε

→ Indique à quel point le nombre d'itérations augmente lorsque la précision souhaitée diminue

- Question: . Existe-t-il de meilleurs algorithmes que la descente de gradient ?
 => Oui dans le cas convexe
- . Que faire si f n'est pas C^1 ?
 => Si f est convexe, on peut définir des algorithmes même sans gradient
 => En sachant des données, le fait que f ne soit pas C^1 provient généralement d'une régularisation. Dans ce cas, on peut définir une famille d'algorithmes qui généralise la descente de gradient : méthodes proximales.

① Accélération

Cadre de travail : minimiser $f(w)$, $w \in \mathbb{R}^d$, $f: \mathbb{R}^d \rightarrow \mathbb{R}$
 et f convexe.

$$\nabla f \text{ L-lipschitz}$$

$$\|\nabla f(v) - \nabla f(w)\| \leq L\|v-w\|$$

$$(f \text{ convexe} \Leftrightarrow \forall (v, w) \in (\mathbb{R}^d)^2, f(v) \geq f(w) + \nabla f(w)^T(v-w))$$

On suppose que f admet un minimum ($\underset{w \in \mathbb{R}^d}{\operatorname{argmin}} f(w) \neq \phi$)
 et on considère $\begin{cases} w^* & \operatorname{argmin}_{w \in \mathbb{R}^d} f(w). \\ f^* = f(w^*) & = \min_{w \in \mathbb{R}^d} f(w) \end{cases}$

↑ ensemble des solutions de minimiser $f(w)$

→ La descente de gradient garantit $f(w_h) - f^* \leq \epsilon$ en au plus $O(\epsilon^{-2})$ itérations : est-ce la meilleure complexité ?

possible ? → Non !

Réultat d'existence:
Ne dénit pas
l'algorithme en
question, mais
certifie que cet
algorithme est
optimal en termes de complexité

~1970-1980: résultat théorique "Il existe
un algorithme pour l'optimisation convexe qui
garantit $f(w_k) - f^* \leq \varepsilon$ en au plus $O(\varepsilon^{-\frac{1}{2}})$
itérations"

↑
mieux que $O(\varepsilon^{-1})$

C'est moins vite
quand ε diminue

↳ cet algorithme a été établi en 1983

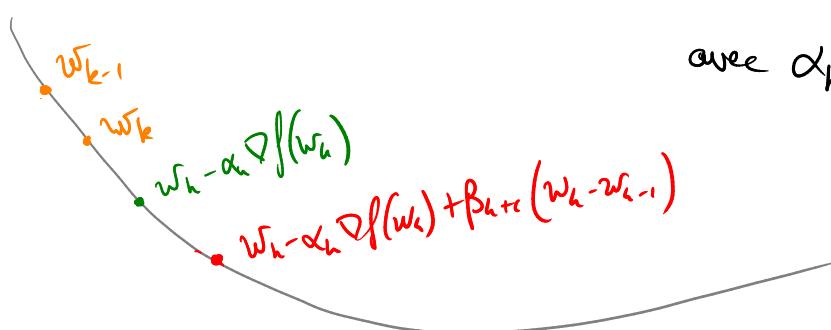
Histoire

1964 : Polyak propose la méthode "heavy ball", ou
descente de gradient avec momentum

Algorithm

- . Partir de $w_0 \in \mathbb{R}^d$, poser $w_{-1} = w_0$.
- . Iteration k Iteration de la
descende de la
dérivée de gradient

$$w_{k+1} = w_k - \alpha_k \nabla f(w_k) + \beta_{k+1} (w_k - w_{k-1})$$
↳ Paramètre de
momentum



→ Polyak montre que cette méthode est optimale pour les
fonctions quadratiques fortement convexes (sous-ensemble des
fonctions convexes)

→ Mais la méthode peut ne pas converger sur des fonctions
qui ne sont pas quadratiques (travaux de 1970-1980)

1983 : Nesterov : gradient accéléré / algorithme de Nesterov

Idee : la méthode de Polyak peut être vue comme une double récurrence (2 suites w_h, z_h)

$$z_{h+1} = w_h - \alpha_h \nabla f(w_h) \rightarrow \text{pas de gradient}$$

$$w_{h+1} = z_{h+1} + \beta_{h+1}(w_h - w_{h-1}) \rightarrow \text{pas de momentum}$$

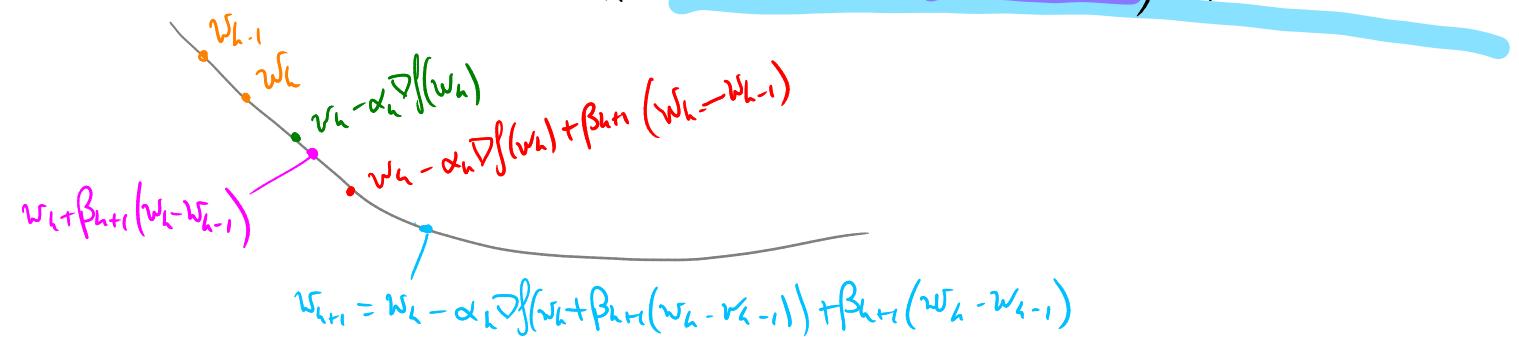
Nesterov inverse les deux équations

$$z_{h+1} = w_h + \beta_{h+1}(w_h - w_{h-1}) \rightarrow \text{pas de momentum}$$

$$w_{h+1} = z_{h+1} - \alpha_h \nabla f(z_{h+1}) \rightarrow \text{pas de gradient sur le nouveau point}$$

ce qui se réécrit

$$w_{h+1} = w_h - \alpha_h \nabla f(w_h + \beta_{h+1}(w_h - w_{h-1})) + \beta_{h+1}(w_h - w_{h-1})$$



↳ Avec les bons choix de α_h et β_{h+1} , l'algorithme de Nesterov possède des garanties de complexité optimales pour les fonctions convexes ainsi que pour les fonctions fortement convexes

⚠ les choix de $\{\alpha_h\}$ et $\{\beta_{h+1}\}$ sont différents selon que la fonction est convexe ou fortement convexe : en ce sens, l'algorithme de Nesterov est moins versatile que la descente de gradient

Illustration : Pour f convexe (non fortement convexe), on prend

$$\alpha_h = \frac{1}{L} \quad (\text{si } f \in \mathcal{C}^1 \text{ et } \frac{L}{L})$$

et $\beta_{h+1} = \frac{t_h - 1}{t_{h+1}}$ avec $t_0 = 0$

$$t_{h+1} = \frac{1}{2} (1 + \sqrt{1 + 4t_h^2})$$

indépendant du problème

- Pour $f \in C_L^{1,1}$ et μ -lentement convexe ($\mu > 0$)

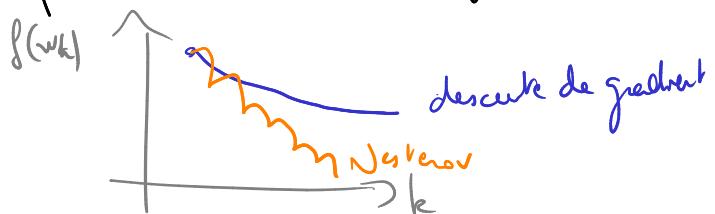
$$\alpha_h = \frac{1}{L} \quad \text{et} \quad \beta_{h+1} = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}$$

Table de complexité

	$f \in C_L^{1,1}$ convexe	$f \in C_L^{1,1}$ μ -lentement convexe
Descente de gradient $\alpha_h = 1/L$	$f(w_h) - f^* \leq \varepsilon$ en au plus $O(\varepsilon^{-1})$ itérations	$f(w_h) - f^* \leq \varepsilon$ en au plus $O\left(\frac{L}{\mu} \ln(\varepsilon^{-1})\right)$ itérations $L/\mu \geq 1$ par définition
Heavy-ball (Polyak)	X	(Si) f est quadratique, alors $f(w_h) - f^* \leq \varepsilon$ en au plus $O\left(\sqrt{\frac{L}{\mu}} \ln(\varepsilon^{-1})\right)$ itérations
Newton	$f(w_h) - f^* \leq \varepsilon$ en au plus $O(\varepsilon^{-1/2})$ itérations $\beta_{h+1} = \frac{t_h - 1}{t_{h+1}}$...	$f(w_h) - f^* \leq \varepsilon$ en au plus $O\left(\sqrt{\frac{L}{\mu}} \ln(\varepsilon^{-1})\right)$ itérations avec $\beta_{h+1} = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}$

A venir : Notebook

- Les méthodes avec accélération peuvent conduire à augmenter la valeur de f !



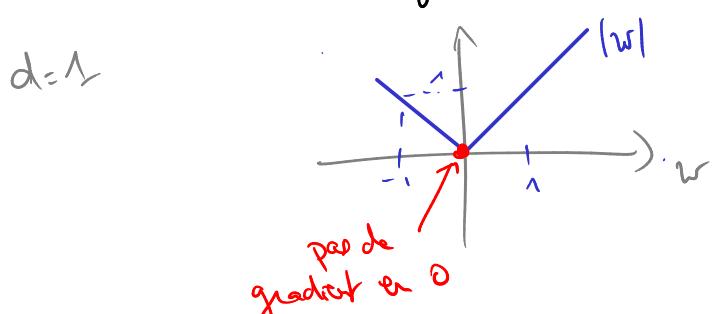
- La méthode de Nesterov pour fonctions convexes n'est pas la meilleure sur les problèmes fortement convexes, idem pour la variante pour fonctions fortement convexes (pas la meilleure sur les problèmes convexes)
- Une implémentation efficace des algorithmes de Polyak et Nesterov requiert 3 vecteurs de stockage (contre 1 ou 2 pour la descente de gradient)

Réponse: Dans le cas non convexe, le momentum/l'accélération ne marche pas en général (mais on peut l'implémenter !)
 \Rightarrow Par contre, lorsque l'on fait du gradient stochastique, utiliser du momentum donne de bons résultats pratiques.

② Régularisation et méthodes proximales

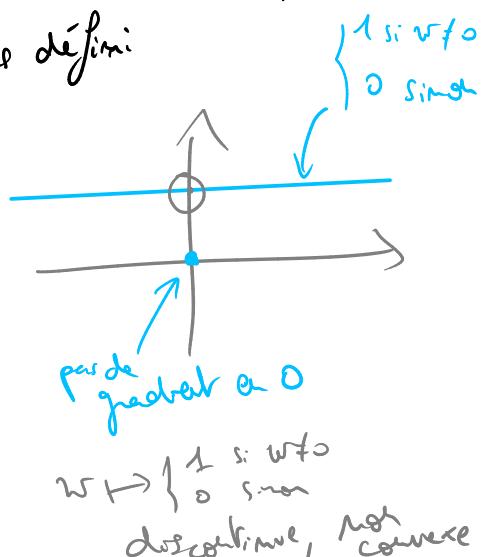
↳ Q: Que faire si f n'est pas C^1 ?

\Rightarrow On dira que f est non lisse s'il existe des points où lesquels le gradient de f n'est pas défini



$$w \mapsto \begin{cases} 1 & \text{si } w \neq 0 \\ 0 & \text{sinon} \end{cases}$$

cas où on peut appliquer des algorithmes similaires à la descente de gradient



Ces pathologies

→ En sciences des données, et notamment en traitement du signal et de l'image, les fonctions non lisses apparaissent souvent à cause de la régularisation.

Problème régularisé

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad \left\{ f(w) + \lambda \mathcal{L}(w) \right\}$$

Terme d'attache aux données, dépend des données du problème

BUT: Déterminer $w \in \mathbb{R}^d$ qui donne le meilleur modèle des données possible

Paramètre de régularisation: donne le poids de la régularisation par rapport à l'attache aux données

Terme de régularisation, ne dépend pas des données
→ Typiquement convexe
→ Souvent non lisse

↓
Favorise les $w \in \mathbb{R}^d$ qui possèdent des propriétés souhaitées

Exemple: Régression linéaire

Données: $\{(x_i, y_i)\}_{i=1..m}^{n \times d \rightarrow \mathbb{R}}$

$$X = \begin{bmatrix} x_1^\top \\ \vdots \\ x_m^\top \end{bmatrix} \quad \underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad \|Xw - y\|^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

On cherche $w \in \mathbb{R}^d$ tel que $x_i^\top w - y_i \approx 0$

$$\frac{1}{2m} \|Xw - y\|^2 = \frac{1}{2m} \sum_{i=1}^m (x_i^\top w - y_i)^2$$

Terme d'attache aux données

↳ Deux choix de régularisation classiques

• Régularisation l_2 / "élevée" (ridge)

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad \frac{1}{2m} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|_2^2$$

Terme de régularisation l_2

$$\|w\|_2^2 = \sum_{j=1}^d (w_j)^2$$

avec $\lambda > 0$

⇒ Pour $\lambda \gg 1$, ce problème équivaut à minimiser $w \in \mathbb{R}^d \quad \frac{1}{2} \|w\|^2$
dont la solution est $w = 0_{\mathbb{R}^d}$

\Rightarrow pour des valeurs modestes de λ , la solution du problème régularisé (dont on peut montrer qu'elle est unique) est moins sensible aux perturbations des données que la solution du problème d'origine ($\lambda=0$)

\Rightarrow Autres propriétés de la régularisation ℓ_2

- Si f est convexe, alors $w \mapsto f(w) + \frac{\lambda}{2} \|w\|^2$ est λ -fortement convexe ($\lambda > 0$) \Rightarrow la fonction régularisée a un unique minimum

- Lorsque λ augmente, les composantes de la solution diminuent de manière uniforme (on dit que la régularisation ℓ_2 "lisse" le problème)

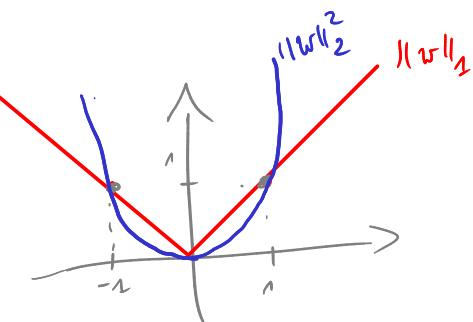
- Contrainte implicite sur la norme de la solution

Régularisation ℓ_1 / Lasso (~1996)

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad \frac{1}{2m} \|Xw - y\|^2 + \lambda \|w\|_1$$

Terme de régularisation ℓ_1

$$\|w\|_1 = \sum_{j=1}^d |w_j|$$



\hookrightarrow plus λ est important, plus la solution du problème aura de coordonnées nulles

(\neq de la régularisation ℓ_2 : les coordonnées ne sont nulles que pour $\lambda \gg 1$, et dans ce cas elles sont toutes nulles)

\hookrightarrow La régularisation ℓ_1 conserve uniquement les coordonnées les plus importantes pour l'attache aux données \Rightarrow on dit que la régularisation permet de sélectionner les attributs (c'est les dimensions de x_i les plus importantes pour l'apprentissage)

\hookrightarrow Mathématiquement, on peut montrer que la solution du problème régularisé a toujours moins (ou autant) de coordonnées nulles qu'une solution du problème non régularisé

$\Rightarrow \|w\|_1$ est convexe et le problème Lasso est fortement convexe.

↳ On parle parfois de la solution du LASSO comme d'une solution parcimonieuse (sparse en anglais)

⇒ On voudrait construire un algorithme qui s'applique à la régression linéaire, à la régression ℓ_2 et à la régression LASSO.

différence
de structure
que l'on peut vraiment exploiter

la descente de gradient
ne fait pas de distinction
entre ces deux problèmes!

$\|w\|_1$ non lissée (pas de gradient en tout point qui a une coordonnée nulle)

On ne peut pas appliquer la
descendre de gradient!

→ Les méthodes proximales font une classe d'algorithmes dédiée aux problèmes de la forme minimiser $f(w) + \lambda S_L(w)$, qui s'appliquent même quand S_L est non lissée (NB: elles s'appliquent aussi quand f est non lissée)

Algorithme principal: Le gradient proximal

Hypothèse :

minimise $w \in \mathbb{R}^d f(w) + \lambda S_L(w)$

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ C^1
(par morseanement convexe)

$S_L: \mathbb{R}^d \rightarrow \mathbb{R}$ convexe
(par morseanement C^1)

Itération

$$w_{k+1} \in \arg \min_{w \in \mathbb{R}^d} \left\{ f(w_k) + \nabla f(w_k)^T (w - w_k) + \frac{1}{2\alpha_k} \|w - w_k\|^2 + \lambda S_L(w) \right\}$$

$\approx f(w)$ pour w proche de w_k

Terme proximal:
force la solution à être
proche de w_k ($\alpha_k > 0$)

→ À chaque itération, l'algorithme du gradient proximal résout un problème d'optimisation auxiliaire (dit "sous-problème") défini par f , λ , Ω et l'hôte constant w_k .

→ L'idée est de remplacer le problème de départ par une suite de sous-problèmes pour lesquels on remplace f par une approximation en utilisant son gradient \Rightarrow N'a d'intérêt que lorsque le sous-problème est plus simple à résoudre que le problème de départ

→ Si on considère $\Omega(w) = \frac{\lambda}{2} \|w\|_2^2$, alors l'itération de gradient proximal donne

$$w_{k+1} = w_k - \alpha_k \nabla f(w_k)$$

: c'est l'itération de la descente de gradient

→ Pour Ω convexe, on peut montrer que le gradient proximal converge avec une complexité du même ordre que celle de la descente de gradient

(Ex) Si f convexe, on obtient une valeur de fonction à ε de l'optimum

$$(f + \lambda \Omega)(w_k) - \min_w (f + \lambda \Omega)(w) \leq \varepsilon$$

en au plus $\mathcal{O}(\varepsilon^{-1})$ itérations

⚠ Une itération de gradient proximal a un coût différent (en général plus grand) de celui d'une itération de descente de gradient

→ Pour des régularisations ℓ_2 et ℓ_1 , l'itération du gradient proximal s'écrit de manière explicite

Gradient proximal pour minimiser $\min_{w \in \mathbb{R}^d} f(w) + \frac{\lambda}{2} \|w\|^2$ $f(\cdot)$

$$w_{k+1} = \frac{1}{1+\lambda\alpha_k} w_k - \frac{\alpha_k}{1+\lambda\alpha_k} \nabla f(w_k)$$

$\rightarrow \lambda = 0$: on retrouve la descente de gradient

$\rightarrow \lambda \gg 1 \Rightarrow w_{\text{fin}} \approx 0$

\rightarrow le facteur $\frac{1}{1+\lambda \alpha_k}$ réduit les composantes de w_k de manière uniforme à chaque itération

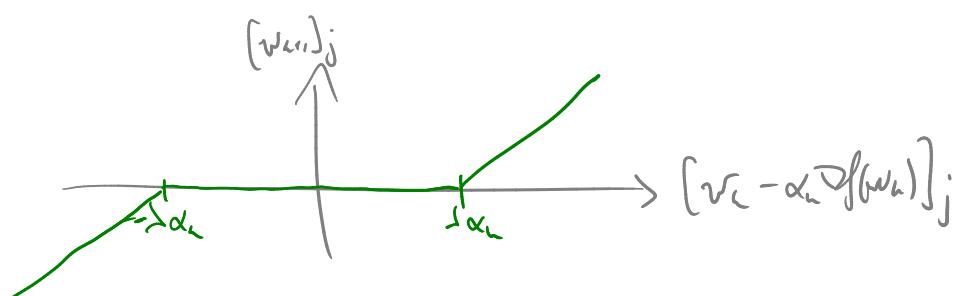
\Rightarrow liens avec la "Batch normalization" en apprendrage profond

Gradient proximal pour minimiser $f(w) + \lambda \|w\|_2$ $f \in C^1$
 $w \in \mathbb{R}^d$

aussi appelé ISTA

(Iterative Soft-thresholding Algorithm)

$$\forall j=1..d, \quad [w_{k+1}]_j = \begin{cases} [w_k - \alpha_k \nabla f(w_k)]_j - \lambda \alpha_k & \text{si } [w_k - \alpha_k \nabla f(w_k)]_j > \lambda \alpha_k \\ [w_k - \alpha_k \nabla f(w_k)]_j + \lambda \alpha_k & \text{si } [w_k - \alpha_k \nabla f(w_k)]_j < -\lambda \alpha_k \\ 0 & \text{sinon} \end{cases}$$



$\rightarrow \lambda = 0$: on retrouve l'itération de descente de gradient

\rightarrow plus λ est grand, plus $[-\lambda \alpha_k, \lambda \alpha_k]$ sera un grand intervalle et plus $[w_{k+1}]_j$ aura de chances d'être mis à zéro

\rightarrow Par rapport au pas de descente de gradient $w_k - \alpha_k \nabla f(w_k)$, w_{k+1} aura nécessairement plus ou autant de coordonnées nulles.

NB: Si f est convexe, on peut accélérer le gradient proximal
(Pour régularisé l1: FISTA, 2009 13999 iterations
du Google Scholar)