

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

29 novembre 2024

Aujourd'hui : Séance-bilan

- Examen 2023-2024 (adapté)
- Questions diverses
- Si temps : petit sujet bonus

Examen : Vendredi 13/12, 9h-11h

Documents : 1 feuille A4 recto-verso de
nos manuscrits ou imprimés

Exercice 1: Un problème non convexe

(Examen 1D apprentissage 2023-2024)

Jeu de données $\{(x_i, y_i)\}_{i=1..n}$ $x_i \in \mathbb{R}^d$ et $y_i \in \{0, 1\}$ + i
(classification)

But: Trouver un modèle linéaire des données

Problème d'optimisation

minimiser
 $w \in \mathbb{R}^d$

$$\phi(w) = \frac{1}{n} \sum_{i=1}^n \left(y_i - \frac{1}{1+e^{-x_i^T w}} \right)^2$$

Réaction: $x_i^T w \gg 1 \Rightarrow \frac{1}{1+e^{-x_i^T w}} \approx 1$

$$x_i^T w \ll -1 \Rightarrow \frac{1}{1+e^{-x_i^T w}} \approx 0$$

La fonction ϕ est de classe C^2 et est non convexe.

- a) Justifier que 0 est un minorant de la fonction ϕ .
Est-ce nécessairement sa valeur optimale ?

$$\forall w \in \mathbb{R}^d, \quad \phi(w) = \frac{1}{n} \sum_{i=1}^n \underbrace{\left(y_i - \frac{1}{1+e^{-x_i^T w}} \right)^2}_{\geq 0 \ \forall w} \geq 0$$

(Résultat de carres, donc de termes positifs)



$$e^{-x} \geq 0 \quad \forall x \quad \text{mais} \quad e^{-x} > 0$$

O n'est pas nécessairement la valeur optimale du problème.
 Ce sera le cas si il existe $\bar{w} \in \mathbb{R}^d$ tel que $\phi(\bar{w}) = 0$

NB: La fonction $t \mapsto e^t$ est un exemple de fonction telle que $e^t \geq 0 \forall t$ mais e^t n'est jamais égal à 0.

- b) On souhaite appliquer la descente de gradient au problème.
- Ecrire l'itération de cet algorithme en utilisant une taille de pas quelconque

$$w_{k+1} = w_k - \alpha_k \nabla \phi(w_k) \quad \text{avec } \alpha_k > 0$$

- Donner deux choix possibles pour la taille de pas
 (plus de 2 possibilités)

Taille de pas constante: $\alpha_k = \alpha > 0 \forall k$ Ex) $\alpha = 0.1, \alpha = 0.01$

Taille de pas décroissante: $\alpha_k \xrightarrow{k \rightarrow \infty} 0$ Ex) $\alpha_k = \frac{1}{\sqrt{k+1}}, \alpha_0 = \frac{\alpha_0}{k+1}$ avec $\alpha_0 > 0$

Taille de pas adaptative:

α_k choisie en fonction du point courant,
 c'ds en fonction de $x_k, \nabla \phi(x_k), \phi(x_k)$

Ex) Recherche linéaire avec "backtracking"

$\alpha_k = 1$.
 Tant que $(\phi(x_k - \alpha_k \nabla \phi(x_k)) \geq \phi(x_k) - c \alpha_k \|\nabla \phi(x_k)\|^2)$
 $\alpha_k \leftarrow \alpha_k / 2$

Remarque bonus: Utiliser une recherche linéaire est plus coûteux qu'utiliser une taille de pas constante ou décroissante car cela requiert au moins un appel à $\phi(\cdot)$ par itération.

En revanche, les valeurs x_k obtenues par la recherche linéaire sont mieux adaptées à l'itéré correspondant que des valeurs constantes ou décroissantes choisies a priori.

- iii) Sous les bonnes hypothèses, on peut obtenir une borne de complexité pour la descente de gradient. Quelle est l'ordre de grandeur de cette borne, et à quelle quantité s'applique-t-elle ?

$$\min_{0 \leq k \leq K-1} \|\nabla \phi(w_k)\| \leq \varepsilon \text{ en au plus } O\left(\frac{1}{\varepsilon^2}\right) \text{ itérations}$$

\Leftrightarrow L'algorithme calcule un itéré w_h tel que $\|\nabla \phi(w_h)\| \leq \varepsilon$ en au plus $O(1/\varepsilon^2)$ itérations (mais ce n'est pas forcément le dernier itéré)

- c) Supposons que la descente de gradient calcule un itéré w_h tel que $\|\nabla \phi(w_h)\| = 0$. w_h est-il nécessairement un minimum du problème ?

Non, pas nécessairement car la fonction ϕ est non convexe.

Contre-exemple: $d=1, f: w \mapsto w^3$

$Df(0)=0$ mais 0 est un point selle, pas un minimum.

d) Rappeler la condition d'optimalité nécessaire à l'ordre deux pour le problème. Un point qui vérifie cette condition est-il un minimum ?

Condition nécessaire à l'ordre 2:

$[\bar{w} \in \mathbb{R}^d \text{ minimum (local) de } \phi] \Rightarrow$

$$\left\{ \begin{array}{l} \nabla \phi(\bar{w}) = 0 \\ \nabla^2 \phi(\bar{w}) \succeq 0 \end{array} \right.$$

$$v^T \nabla^2 \phi(\bar{w}) v \geq 0$$

$v \in \mathbb{R}^d$

Un point $\bar{w} \in \mathbb{R}^d$ tel que $\nabla \phi(\bar{w}) = 0$ et $\nabla^2 \phi(\bar{w}) \succeq 0$ n'est pas nécessairement un minimum.

Exemple : $f: w \mapsto w^3$ et $\bar{w} = 0$
 $\nabla f(0) = 0, \nabla^2 f(0) = 0 \succeq 0$

Remarques bonus : Il existe des fonctions non convexes φ pour lesquelles on a

$$[\bar{w} \text{ minimum local}] \Leftrightarrow [\bar{w} \text{ vérifie } \nabla \varphi(\bar{w}) = 0 \text{ et } \nabla^2 \varphi(\bar{w}) \succeq 0]$$

Voici même

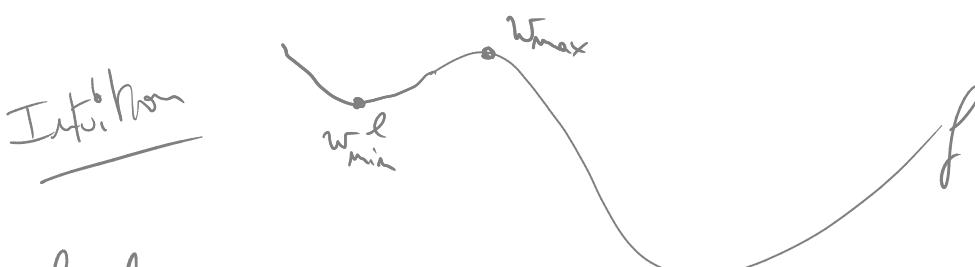
$$[\bar{w} \text{ minimum global}] \Leftrightarrow [\nabla \varphi(\bar{w}) = 0 \text{ et } \nabla^2 \varphi(\bar{w}) \succeq 0]$$

Si $\nabla \varphi(\bar{w}) = 0$ et $\nabla^2 \varphi(\bar{w}) \succ 0$ (définition positive)
 $v^T \nabla^2 \varphi(\bar{w}) v > 0 \Leftrightarrow v \neq 0 \in \mathbb{R}^d$

alors \bar{w} est un minimum local de φ .

c) Supposons que l'on applique la descente de gradient au problème de départ (minimiser $\phi(w)$) en partant d'un point w_0 tiré aléatoirement. On observe que la méthode converge vers $\bar{w} \in \mathbb{R}^d$ tel que $\nabla \phi(\bar{w}) = 0$ et $\nabla^2 \phi(\bar{w}) \succeq 0$. Comment expliquer ce comportement (étais-ce attendu) ?

La descente de gradient avec initialisation aléatoire converge presque sûrement (càd en probabilité 1) vers un point stationnaire à l'ordre 2 (càd un point \bar{w} tel que $\nabla \phi(\bar{w}) = 0$ et $\nabla^2 \phi(\bar{w}) \succeq 0$). On pouvait donc s'attendre à ce résultat.



Intuition

w_{\max} : max local
 w_{\min}^l : min local
 w_{\min}^g : min global

les trois points pour lesquels le gradient est nul

$$\nabla f(w_{\max}) = 0, \quad \nabla^2 f(w_{\max}) \not\succeq 0$$

$$\nabla f(w_{\min}^l) = 0, \quad \nabla^2 f(w_{\min}^l) \succeq 0$$

$$\nabla f(w_{\min}^g) = 0, \quad \nabla^2 f(w_{\min}^g) \succeq 0$$

$w_0 \neq w_{\max} \Rightarrow$ la descente de gradient converge vers w_{\min}^l ou w_{\max}

$w_0 = w_{\max} \Rightarrow$ la descente de gradient revient en w_{\max}

$$w_k = w_{\max}$$

Si w_0 est tiré aléatoirement dans \mathbb{R} , alors $\mathbb{P}(w_0 = w_{\max}) = 0$

Exercice 2: Gradient proximal

(Examen 2023-2024)

Jeux de données: $\{(x_i, y_i)\}_{i=1..n}$

$$x_i \in \mathbb{R}^d$$

$$y_i \in \mathbb{R}$$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

attaché aux données

$$\underbrace{\frac{1}{2n} \|Xw - y\|^2}_\text{fonction d'erreur} + \underbrace{\frac{\lambda_2}{2} \|w\|^2 + \lambda_1 \|w\|_1}_\text{termes de régularisation}$$

Problème

minimiser
 $w \in \mathbb{R}^d$

$$\forall v \in \mathbb{R}^d, \quad \|v\|^2 = \sum_{j=1}^d v_j^2 \quad \text{et} \quad \|v\|_1 = \sum_{j=1}^d |v_j|$$

\Rightarrow Problème de régression linéaire avec régularisation
"elastic net" avec $\lambda_1 \geq 0$ et $\lambda_2 \geq 0$

a) Quelle est l'efficacité générale de la régularisation?

Utiliser un terme de régularisation permet de favoriser les solutions qui possèdent des propriétés souhaitées. Ces propriétés dépendent de la nature du terme de régularisation.

b) Si $\lambda_1=0$ et $\lambda_2>0$, quel est le rôle de l'usage de la régularisation ?

$\lambda_1=0, \lambda_2>0 \Rightarrow$ Régularisation par $\frac{\lambda_2}{2} \|w\|^2$
"Ridge / l_2 / énervé / -"

La régularisation l_2 :

- réduire la sensibilité de la solution vis-à-vis des données sans régularisation: une petite perturbation de (X, y) peut entraîner un grand changement dans la valeur de solution
- avec régularisation: une petite perturbation de (X, y) ne change que très peu la valeur de la solution

N.B.: Si $\lambda_2 \gg 1$, alors la solution est presque nulle

- transforme un problème convexe en problème fortement convexe (intervenant ici car $w \mapsto \frac{1}{2m} \|Xw - y\|^2$ est convexe)
 - Fonction convexe: peut avoir une infinité de minima globaux
 - Fonction fortement convexe: Unique minimum global
- réduire la norme l_2 de la solution

c) Si $\lambda_1 > 0$ et $\lambda_2 = 0$, quel est le rôle du terme de régularisation ?

$\lambda_1 > 0, \lambda_2 = 0 \Rightarrow$ Régularisation $\lambda_1 \|w\|_1$
"on nomme $\ell_1/\ell_1/\text{LASSO}"$

Rôle : . Favorise les solutions parcimonieuses (avec beaucoup de coefficients nuls)

- . Réduire la norme $\ell_1 \|.\|_1$ de la solution
- . Permet d'identifier les coordonnées les plus importantes pour l'attribution aux données

d) On rappelle que $\ell: w \mapsto \frac{1}{2n} \|Xw - y\|^2$ est C^1 , et

on peut donc appliquer l'algorithme du gradient proximal au problème

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad \ell(w) + \frac{\lambda_2}{2} \|w\|^2 + \lambda_1 \|w\|_1$$

Ecrire l'édition du gradient proximal pour ce problème avec une taille de pas quelconque.

Terme proximal: Garantit
que $\|w - w_k\|$

$$w_{k+1} \in \underset{w \in \mathbb{R}^d}{\text{argmin}} \left\{ \ell(w_k) + \nabla \ell(w_k)^T (w - w_k) + \frac{1}{2\alpha_k} \|w - w_k\|^2 + \frac{\lambda_2}{2} \|w\|^2 + \lambda_1 \|w\|_1 \right\}$$

avec $\alpha_k > 0$

Approximation
de $\ell(w)$
pertinente lorsque
 $\|w - w_k\| \ll 1$

Terme de régularisation
représenté par bleu

NB: Utile en pratique si le sous-problème à chaque itération est plus facile à résoudre que le problème de départ (c'est le cas pour la régularisation corrodinée!)

e) lorsque $\lambda_2=0$ et $\lambda_1>0$, à quel algorithme le gradient proximal est-il équivalent?

ISTA (Iterative Soft-thresholding Algorithm)

→ Dans ce cas, on peut écrire explicitement w_{k+1} en fonction de w_k :

$$\forall j=1..d, \quad [w_{k+1}]_j = \begin{cases} [w_k - \alpha_k \nabla l(w_k)]_j - \alpha_k \lambda_1 & \text{si } [w_k - \alpha_k \nabla l(w_k)]_j > \alpha_k \lambda_1 \\ [w_k - \alpha_k \nabla l(w_k)]_j + \alpha_k \lambda_1 & \text{si } [w_k - \alpha_k \nabla l(w_k)]_j < -\alpha_k \lambda_1 \\ 0 & \text{si } [w_k - \alpha_k \nabla l(w_k)]_j \in [-\alpha_k \lambda_1, \alpha_k \lambda_1] \end{cases}$$

f) Si $\lambda_1=0$ et $\lambda_2>0$, quel algorithme peut-on utiliser au lieu du gradient proximal pour résoudre le problème?

$$(\lambda_1=0, \lambda_2>0 \Rightarrow \text{minimise}_{w \in \mathbb{R}^d} \underbrace{\frac{1}{2m} \|Xw-y\|^2}_C + \underbrace{\frac{\lambda_2}{2} \|w\|^2}_C)$$

→ La fonction $w \mapsto \frac{1}{2m} \|Xw-y\|^2 + \frac{\lambda_2}{2} \|w\|^2$ étant C^2 , on peut appliquer la descente de gradient à ce problème.

→ Cette fonction est même fortement convexe et quadratique: on peut donc aussi appliquer le gradient accéléré et la

méthode "heavy-ball" / de Polyak

$$\rightarrow \text{En écrivant } \frac{1}{2n} \|x - y\|^2 + \frac{\lambda_2}{2} \|w\|^2 = \frac{1}{n} \sum_{i=1}^m \left(\frac{1}{2} (x_i^T w - y_i)^2 + \frac{\lambda_2}{2} \|w\|^2 \right)$$

on voit qu'on peut aussi appliquer le gradient stochastique (et toutes ses variantes) au problème, qui s'écrit comme une somme finie de termes qui dépendent chacun d'un point différent du jeu de données

g) Si $\lambda_1 > 0$ et $\lambda_2 > 0$, alors le sous-problème vaut

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \left\{ l(w_h) + \nabla l(w_h)^T (w - w_h) + \frac{1}{2\lambda_1} \|w - w_h\|^2 + \frac{\lambda_2}{2} \|w\|^2 + \frac{1}{2} \|w\|_1^2 \right\}$$

Parmi les algorithmes vus en cours, lequel peut être appliqué à la résolution de ce problème ? Tushar.

Comme $w \mapsto \|w\|_1$ n'est pas C¹, on ne peut pas appliquer la descente de gradient, le gradient stochastique ou les méthodes de gradient accélérée / heavy ball.

$$\text{Par contre, en posant } g(w) = l(w_h) + \nabla l(w_h)^T (w - w_h) + \frac{1}{2\lambda_1} \|w - w_h\|^2 + \frac{\lambda_2}{2} \|w\|^2$$

le problème réécrit

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad g(w) + \lambda_1 \|w\|_1$$

avec $g \in C^1$

On peut donc appliquer ISTA.

Exercice 3 : Somme finie et optimisation distribuée

(Exam 2023-2024)

Problème : minimiser $w \in \mathbb{R}^d$ $\sum_{i=1}^m f_i(w)$

$f_i : \mathbb{R}^d \rightarrow \mathbb{R}$
de classe $C^1 \forall i$

a) Solutions du problème

i) Donner la définition d'un minimum global du problème

$w^* \in \mathbb{R}^d$ est un minimum global

$$\text{si: } \sum_{i=1}^m f_i(w^*) \leq \sum_{i=1}^m f_i(w) \quad \forall w \in \mathbb{R}^d$$

ii) On pose $f(w) = \frac{1}{m} \sum_{i=1}^m f_i(w)$. Justifie que les problèmes
 minimiser $\sum_{i=1}^m f_i(w)$ et minimiser $f(w)$ ont le même
 $w \in \mathbb{R}^d$ ensemble de solutions.

Comme $\frac{1}{m}$ est une constante strictement positive, les deux problèmes ont le même nombre de solutions

En revenant à la définition, on voit que

$$\left[\begin{array}{l} w^* \text{ est un minimum} \\ w \in \mathbb{R}^d \end{array} \right] \Leftrightarrow \left[\begin{array}{l} \sum_{i=1}^m f_i(w^*) \leq \sum_{i=1}^m f_i(w) \\ \forall w \in \mathbb{R}^d \end{array} \right]$$

$$\Leftrightarrow \left[\begin{array}{l} \frac{1}{m} \sum_{i=1}^m f_i(w^*) \leq \frac{1}{m} \sum_{i=1}^m f_i(w) \\ \forall w \in \mathbb{R}^d \end{array} \right]$$

$$\Leftrightarrow \left[\begin{array}{l} w^* \text{ est un minimum} \\ w \in \mathbb{R}^d \end{array} \right] \frac{1}{m} \sum_{i=1}^m f_i(w)$$

iii) Si $f: w \mapsto \frac{1}{m} \sum_{i=1}^m f_i(w)$ est fortement convexe, que peut-on dire des minima globaux du problème ?

Il existe un unique minimum global

b) minimiser $f(w) = \frac{1}{m} \sum_{i=1}^m f_i(w)$
 $w \in \mathbb{R}^d$

- i) Écrire l'algorithme de descente de gradient (cf phys heur)
- ii) Si toutes les f_i sont des fonctions convexes, quelle est la complexité de la descente de gradient sur ce problème (sous les bonnes hypothèses) ?

$$f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq \varepsilon \text{ en au plus } O(\varepsilon^{-1}) \text{ itérations}$$

(NB. N'oubliez que le cas non convexe)

- iii) Donner un algorithme qui possède une meilleure complexité que la descente de gradient sous les hypothèses de la question ii)
- Le gradient accéléré / la méthode de Nesterov

NB: Pour Nesterov,

$$f(w_h) - \min_{w \in \mathbb{R}^d} f(w) \leq \varepsilon \text{ en au plus } O(\varepsilon^{1/2}) \text{ itérations}$$

Nesterov:

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