

# Exercices : Gradient stochastique

Optimisation pour l'apprentissage automatique

M2 IASD Apprentissage, 2024-2025



## Exercice 1: Random reshuffling

On considère un problème de minimisation du risque empirique de la forme

$$\text{minimiser}_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}), \quad (1)$$

où chaque  $f_i$  est de classe  $\mathcal{C}^1$  et dépend uniquement du  $i$ ème point d'un jeu de données à  $n$  éléments. On suppose aussi que la fonction objectif  $f$  est convexe, et  $\mathcal{C}^{1,1}$ . On supposera enfin que le problème (1) possède une solution, et on notera  $f^* = \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ .

On supposera que  $n$  est trop large pour que le jeu de données puisse être utilisé en entier lors d'une itération d'un algorithme, et on considère donc l'algorithme du gradient stochastique

$$x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k), \quad (2)$$

avec  $\alpha_k > 0$  et  $i_k \in \{1, \dots, n\}$ . Le but de l'exercice est d'étudier les variantes basées sur le principe de *random reshuffling* (dont le principe peut être traduit par *rebattre les cartes*). Dans ces variantes, les indices  $\{i_k\}$  sont tirés selon une perturbation aléatoire de  $\{1, \dots, n\}$  qui est modifiée toutes les  $n$  itérations. Ainsi, à l'itération 0, une permutation aléatoire de  $\{1, \dots, n\}$  définit  $\{i_0, \dots, i_{n-1}\}$ , puis à l'itération  $n$ , une autre permutation aléatoire définit  $\{i_n, \dots, i_{2n-1}\}$ , et ainsi de suite.

a) Rappeler la définition d'une époque (*epoch*). Avec la stratégie du *random reshuffling*, que peut-on garantir sur les points du jeu de données qui ont été tirés au cours de la première époque ?

b) Soit un indice d'itération  $k$  correspondant à la première itération d'une époque (càd  $k = \ell n$  avec  $\ell \in \mathbb{N}$ ).

i) Montrer que

$$\mathbb{E}_{i_k} [\nabla f_{i_k}(x_k)] = \nabla f(x_k).$$

ii) Les indices  $i_k, \dots, i_{k+n-1}$  n'étant pas indépendants, justifier que la propriété de la questions b)i) n'est pas vérifiée pour les autres itérations de l'époque.

- c) Même sans l'hypothèse de la question b)i), on peut montrer des résultats de convergence pour un bon choix de taille de pas. Pour cela, on regarde la suite des itérés moyens  $\{\bar{x}_K\}_K$ , avec

$$\bar{x}_K = \frac{1}{K+1} \sum_{k=0}^K x_k \quad \forall K \in \mathbb{N}.$$

On peut ainsi montrer qu'après  $nK$  itérations, on a

$$\mathbb{E}[f(\bar{x}_{nK})] - f^* \leq \mathcal{O}\left(\frac{1}{\sqrt{nK}}\right),$$

Comparer ce taux avec celui de la descente de gradient pour le même problème.

- d) On considère une variantes par fournées du gradient stochastique, pour laquelle les indices sont tirés suivant une approche par *random reshuffling*. Si la taille de fournées est de  $n$ , quel algorithme retrouve-t-on ?
- e) Comme l'approche par *random reshuffling* ne tire pas le même indice au sein d'une même époque, on peut vouloir utiliser de l'information des itérations précédentes pour améliorer le pas. Quelle technique vue lors des séances proposeriez-vous d'utiliser pour cela ?

## Exercice 2 : Perte de Huber

On considère un jeu de données  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , où  $n \geq 1$ ,  $\mathbf{x}_i \in \mathbb{R}^d$  avec  $d \geq 1$  et  $y_i \in \mathbb{R}$ . On cherche un modèle linéaire qui prédise au mieux chaque  $y_i$  à partir du  $\mathbf{x}_i$  correspondant. On définit donc une famille de modèles paramétrée par  $\mathbf{w} \in \mathbb{R}^d$  comme suit :

$$\begin{aligned} h_{\mathbf{w}} : \mathbb{R}^d &\rightarrow \mathbb{R} \\ \mathbf{x} &\mapsto \mathbf{x}^T \mathbf{w} = \sum_{i=1}^d [\mathbf{x}]_i [\mathbf{w}]_i. \end{aligned}$$

Pour un modèle  $h_{\mathbf{w}}$ , on considèrera que ce modèle prédit parfaitement  $y_i$  à partir de  $\mathbf{x}_i$  si on a  $\ell(h_{\mathbf{w}}(\mathbf{x}_i) - y_i) = \ell(\mathbf{x}_i^T \mathbf{w} - y_i) = 0$ , où  $\ell : \mathbb{R} \rightarrow \mathbb{R}$  est la fonction de **perte de Huber** définie par :

$$\ell(t) = \begin{cases} \frac{1}{2}t^2 & \text{if } |t| < 1 \\ |t| - \frac{1}{2} & \text{otherwise.} \end{cases} \quad (3)$$

Cette fonction se comporte comme  $t \mapsto \frac{t^2}{2}$  pour  $|t| < 1$  et comme  $t \mapsto |t|$  lorsque  $|t|$  est très grand. Contrairement à ce que son expression peut suggérer,  $\ell$  est continûment dérivable (ou de classe  $\mathcal{C}^1$ )

L'expression  $\ell(h_{\mathbf{w}}(\mathbf{x}_i) - y_i)$  représente l'erreur du modèle en  $(\mathbf{x}_i, y_i)$ , et on cherche un modèle (c'est-à-dire un vecteur  $\mathbf{w} \in \mathbb{R}^d$ ) tel que la somme de ces erreurs soit minimale. On considère donc :

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i^T \mathbf{w} - y_i). \quad (4)$$

- a) Justifier que 0 est un minorant de (4). Est-ce sa valeur minimale ?

b) Le gradient de  $f$  en  $\mathbf{w} \in \mathbb{R}^d$  est donné par

$$\nabla f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell'(\mathbf{x}_i^T \mathbf{w} - y_i) \mathbf{x}_i, \quad (5)$$

avec

$$\ell'(t) = \begin{cases} 1 & \text{si } t > 1 \\ t & \text{si } |t| \leq 1 \\ -1 & \text{si } t < -1. \end{cases}$$

Écrire (en pseudo-code) l'itération de descente de gradient avec une taille de pas constante  $\alpha$  et en utilisant la formule (5). Que devient cette itération si le point courant est un minimum local ?

c) Une constante de Lipschitz pour  $\nabla f$  est  $L = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2$ . Comment utiliser cette constante pour définir la longueur de pas ? Lorsque  $L$  est inconnue, donner deux choix possibles pour la taille de pas.

d) La fonction  $f$  s'écrit  $f = \frac{1}{n} \sum_{i=1}^n f_i$ , où  $f_i(\mathbf{w}) = \ell(\mathbf{x}_i^T \mathbf{w} - y_i)$ . Le gradient  $f_i$  en  $\mathbf{w}$  est

$$\nabla f_i(\mathbf{w}) = \ell'(\mathbf{x}_i^T \mathbf{w} - y_i) \mathbf{x}_i.$$

Écrire (en pseudo-code) l'itération du gradient stochastique pour ce problème sans choix particulier de taille de pas.

e) On considère ici que notre unité de coût est un accès à un  $\mathbf{x}_i$ . Quel est le coût d'une itération de descente de gradient, et celui d'une itération de gradient stochastique ?

f) Quand on applique le gradient stochastique avec une longueur de pas fixe, on peut parfois observer que la méthode génère des itérés de norme de plus en plus grande, ce qui conduit à un dépassement de mémoire pour l'algorithme. Fournir une justification à ce phénomène.

g) On considère une variante par fournées (*batch*) du gradient stochastique, dans laquelle on choisit un sous-ensemble de  $n_b$  composantes dans la somme finie de (4).

i) Écrire l'itération correspondante (en pseudo-code).

ii) Si  $n_b$  correspond au nombre de processeurs disponibles pour les calculs, quel peut être l'intérêt de choisir  $n_b$  comme taille de fournée ?

iii) Donner un autre intérêt plus général des méthodes par fournées en comparaison avec l'algorithme du gradient stochastique basique.

iv) Supposons que l'on utilise plusieurs tailles de fournées et que l'on observe une amélioration en termes de convergence quand  $n_b$  augmente pour  $1 \leq n_b \leq \frac{n}{10}$ . Supposons que l'on observe aussi qu'augmenter  $n_b$  au-delà de  $n/10$  conduise à une dégradation de la performance. Comment expliquer ces observations ?

## Correction

### Correction de l'exercice 1: Random reshuffling

- a) Une époque est une unité de coût correspond à  $n$  accès à un exemple du jeu de données dans un jeu de données à  $n$  éléments. Une époque de gradient avec *random reshuffling* correspond à une passe dans le jeu de données.
- b)  $k = \ell n$  avec  $\ell \in \mathbb{N}$  (première itération d'une époque).

- i) Puisque  $i_k$  est le premier indice d'une permutation aléatoire, il suit une distribution uniforme dans  $\{1, \dots, n\}$ , indépendamment de  $\mathbf{w}_k$ . Par conséquent, on a

$$\mathbb{E}_{i_k} [\nabla f_{i_k}(\mathbf{w}_k)] = \sum_{i=1}^n \mathbb{P}(i_k = i) \times \nabla f_i(\mathbf{w}_k) = \sum_{i=1}^n \frac{1}{n} \times \nabla f_i(\mathbf{w}_k) = \nabla f(\mathbf{w}_k).$$

- ii) Soit l'itération  $k + 1$ . Comme la valeur de  $i_{k+1}$  dépend de celle de  $i_k$ , on ne peut plus appliquer le raisonnement de la question précédente. Le raisonnement s'applique également aux indices suivants.
- c) Après  $nK$  itérations de descente de gradient, on obtiendrait un itéré  $\mathbf{w}_{nK}^G$  tel que

$$f(\mathbf{w}_{nK}^G) - f^* \leq \mathcal{O}\left(\frac{1}{nK}\right),$$

Il s'agit d'une meilleure vitesse que l'approche par *random reshuffling*, à la fois parce que la vitesse de convergence est plus rapide quand  $K \rightarrow \infty$  et parce que la garantie est déterministe.

Dans le même temps, une itération de descente de gradient est  $n$  fois plus coûteuse qu'une itération de gradient stochastique en termes d'accès aux données. Pour un nombre fixe d'époques  $N_E$ , la méthode de gradient stochastique effectue  $nN_E$  itérations, tandis que celle de la descente de gradient n'en effectue que  $N_E$  itérations. Lorsque  $n$  est suffisamment grand, on a  $\sqrt{nN_E} \gg N_E$ , et donc la vitesse de convergence sera meilleure pour le gradient stochastique.

- d) Avec une taille de fournée de  $n$ , on obtient une itération de descente de gradient car les indices seront tirés sans remise.
- e) Des techniques de gradient stochastique basées sur du momentum, comme le gradient stochastique avec momentum et Adam utilisent les informations des itérations précédentes, et seraient donc appropriées.

### Correction de l'exercice 2: Perte de Huber

- a) La fonction  $\ell$  est positive sur  $\mathbb{R}$ . Pour tout  $\mathbf{w} \in \mathbb{R}^d$ , on a donc

$$f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i^\top \mathbf{w} - y_i) \geq \frac{1}{n} \sum_{i=1}^n 0 = 0.$$

Cela montre que 0 est un minorant du problème (4). Cette valeur est atteinte uniquement lorsqu'il existe un point  $\mathbf{w}$  tel que  $\mathbf{x}_i^T \mathbf{w} - y_i = 0$  pour tout  $i$ . Cela n'est pas toujours le cas (prendre par exemple  $n = 2, d = 1, \mathbf{x}_1 = 1, \mathbf{x}_2 = -1, y_1 = y_2 = 1$ ), donc 0 n'est pas nécessairement la valeur minimale du problème.

- b) En un point  $\mathbf{w}_k \in \mathbb{R}^d$ , l'itération de la méthode de descente de gradient avec un pas constant  $\alpha$  s'écrit :

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\alpha}{n} \sum_{i=1}^n \ell'(\mathbf{x}_i^T \mathbf{w}_k - y_i) \mathbf{x}_i.$$

Si  $\mathbf{w}_k$  est un minimum local, on a  $\nabla f(\mathbf{w}_k) = 0$ , et l'itération devient  $\mathbf{w}_{k+1} = \mathbf{w}_k$ .

- c) Si on connaît une constante de Lipschitz  $L$  pour le gradient, on peut alors choisir un pas constant égal à  $\alpha = \frac{1}{L}$ .

Si on ne connaît pas cette valeur, il est possible d'utiliser un pas décroissant (par exemple  $\alpha_k = \frac{1}{k+1}$ ) ou d'effectuer une recherche linéaire pour calculer un pas approprié à l'itération.

- d) En  $\mathbf{w}_k \in \mathbb{R}^d$ , l'itération du gradient stochastique (avec pas  $\alpha_k$ ) se décompose en deux parties. On tire tout d'abord un indice  $i_k$  au hasard dans  $\{1, \dots, n\}$ ; on calcule ensuite le nouvel itéré  $\mathbf{w}_{k+1}$  via

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k) = \mathbf{w}_k - \alpha_k \ell'(\mathbf{x}_{i_k}^T \mathbf{w}_k - y_{i_k}) \mathbf{x}_{i_k}.$$

- e) Chaque itération de descente de gradient doit accéder à toute la donnée pour calculer le gradient : comme notre unité de coût correspond à un accès à un point  $\mathbf{x}_i$ , le coût d'une itération de descente de gradient est de  $n$ . Quant à l'itération du gradient stochastique, elle a un coût de 1, puisqu'elle n'accède qu'à un point du jeu de données ( $\mathbf{x}_{i_k}$ , avec  $i_k$  tiré aléatoirement).
- f) L'algorithme du gradient stochastique est aléatoire par nature, car son exécution dépend d'un tirage aléatoire d'une suite d'indices : il se peut donc qu'il ne converge pas (même s'il converge en moyenne), et c'est ce qui se produit ici.

- i) L'itération de gradient stochastique par fournées (*batch*) de taille  $n_b$  en  $\mathbf{w}_k \in \mathbb{R}^d$  consiste d'abord à tirer aléatoirement un ensemble d'indices  $S_k \subset \{1, \dots, n\}$  tel que  $|S_k| = n_b$ , puis à effectuer l'itération :

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\alpha_k}{|S_k|} \sum_{i \in S_k} \nabla f_i(\mathbf{w}_k),$$

avec  $\alpha_k > 0$  la longueur de pas.

- ii) Si  $n_b$  processeurs sont disponibles et que les gradients des  $f_i$  peuvent être calculés en parallèle, alors le coût de la fournée peut être réparti sur ces  $n_b$  processeurs.
- iii) Ces méthodes utilisent un estimateur du gradient ( $\frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(\mathbf{w}_k)$ ) dont la variance est plus faible que celle de l'estimateur utilisé par le gradient stochastique ( $\nabla f_{i_k}(\mathbf{w}_k)$ ).
- iv) Utiliser plus d'un point à chaque itération permet de réduire la variance des itérations tout en bénéficiant de plus de données, ce qui peut expliquer pourquoi une taille de fournée  $n_b = n/10$  conduit à une meilleure performance que  $n_b = 1$  (qui correspond au gradient stochastique classique). En revanche, lorsque  $n_b$  se rapproche de  $n$ , le coût de la méthode se rapproche de celui d'une itération de descente de gradient, et dans le même temps la méthode est plus sensible aux redondances dans le jeu de données. Cela explique que la convergence de la méthode se dégrade, ici lorsque  $n_b > n/10$ .