

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

15 novembre 2024

Programme des séances restantes

- Aujourd'hui : Cours (gradient stochastique)
- 22/11 : TD (gradient stochastique)
⚠ Exercices + Notebook Python
- 29/11 : Séance - bilan

Projet: . Sujet bientôt disponible !

. Date rendu (flexible) : 31 janvier 2025

MÉTHODES DE GRADIENT STOCHASTIQUE

① Motivation

↳ Jusqu'à maintenant, on a considéré un problème général : minimiser $f(w)$ avec $f: \mathbb{R}^d \rightarrow \mathbb{R}$ de classe C^1 et $w \in \mathbb{R}^d$

⇒ En exploitant la structure de f , on peut développer des algorithmes d'optimisation plus adaptés au problème que la descente de gradient

Exemples de problèmes structurés en sciences des données

1) Régression linéaire

$\{(x_i, y_i)\}_{i=1..m}$ jeu de données avec m éléments
points échantillons

$x_i \in \mathbb{R}^d$: vecteur d'attributs (features)

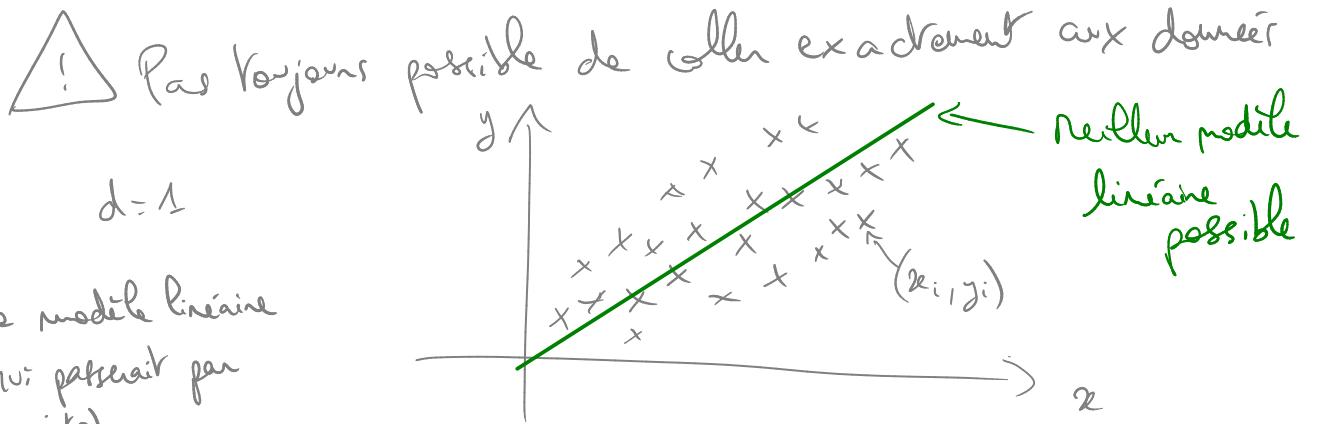
$y_i \in \mathbb{R}$: label

But: Trouver un modèle linéaire, c'est-à-dire une fonction $x \in \mathbb{R}^d \mapsto x^T w = \sum_{j=1}^d [x]_j [w]_j$,

qui colle aux données du mieux possible, c'est-à-dire

que $x_i^T w - y_i \approx 0 \quad \forall i=1..m$

"à peu près" = le plus proche possible de 0



→ Pas de modèle linéaire parfait (qui passe par tous les points)

→ On peut quand même trouver le meilleur modèle possible étant donné les points

Formulation via un problème d'optimisation

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

$$\frac{1}{m} \sum_{i=1}^m \frac{1}{2} (x_i^T w - y_i)^2$$

Rapport sur tous les points

Valeur du modèle linéaire en x_i

Ecart entre $x_i^T w$ et y_i (au carré x $\frac{1}{2}$ par convention)

- Problème convexe très fortement convexe
- Fonction objectif de classe C^2
- Si $m \gg 1$, pour être costeux de calculer la fonction et/ou son gradient
- Si de plus les données proviennent d'une même distribution, il n'est pas forcément nécessaire de considérer tous les points du jeu de données pour déterminer si la valeur de la fonction peut être améliorée

Propriétés structurales du problème

2) Réseaux de neurones

- Modèle simple de réseau "feed-forward" pour la régression
- Idee: Aller au-delà des modèles linéaires

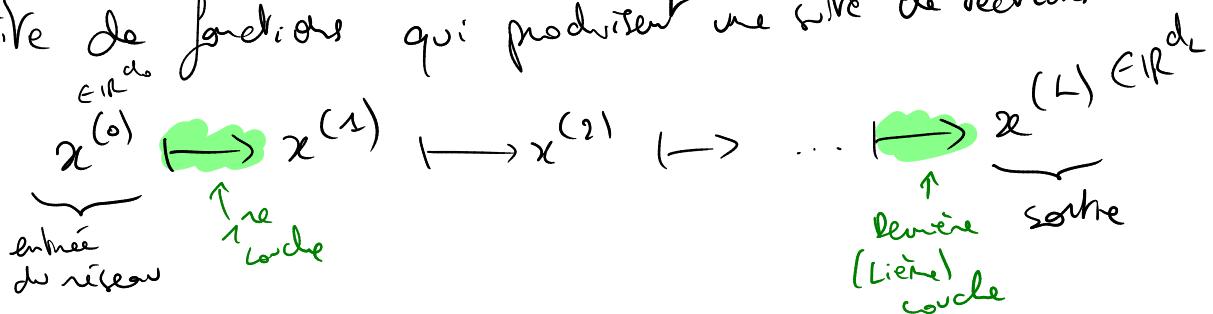
Données: $\{(x_i, y_i)\}_{i=1..n}$

$x_i \in \mathbb{R}^{d_o}$ $d_o \geq 1$

$y_i \in \mathbb{R}^{d_L}$ $d_L \geq 1$

Réseau de neurones à L couches

Suite de fonctions qui produisent une suite de vecteurs



$$\forall l=1..L, \quad x^{(l)} = \sigma(W^{(l)} x^{(l-1)} + b^{(l)})$$

L, Transformation linéaire de $x^{(l-1)}$

où $W^{(l)} \in \mathbb{R}^{d_{l-1} \times d_l}$ matrice de poids (weights)
 Paramètres de la couche l → $b^{(l)} \in \mathbb{R}^{d_l}$ vecteur de biais

σ fonction d'activation non linéaire qui est appliquée à chaque coordonnée de son vecteur argument

$$\forall v \in \mathbb{R}^{d_L}, \quad \sigma(v) = \begin{bmatrix} \sigma([v]_1) \\ \vdots \\ \sigma([v]_{d_L}) \end{bmatrix} \in \mathbb{R}^{d_L}$$

(Ex) $\sigma(t) = \max(t, 0)$: Activation ReLU (la plus fréquente en pratique)

$\sigma(t) = \tanh(t)$ (populaire en traitement naturel des langages)

But: • Préduire y_i via $x_i^{(L)}$, où $x_i^{(L)}$ est la sortie du réseau quand x_i est donné en entrée
 \Rightarrow On souhaite que $\|x_i^{(L)} - y_i\|^2$ soit la plus faible possible ($y_i \in \mathbb{R}^{d_L} \Rightarrow$ on utilise une norme pour mesurer un écart)

, Variables du problème : poids $\{W^{(l)}\}_{l=1..L}$ et $\{b^{(l)}\}_{l=1..L}$

\Rightarrow On concatène toutes ces variables sous la forme d'un vecteur $w \in \mathbb{R}^d$

$$w = \begin{bmatrix} w_{1,1}^{(1)} \\ w_{d_1,1}^{(1)} \\ b_{1,1}^{(1)} \\ \vdots \end{bmatrix}$$

$$d = d_1 d_0 + d_1 + d_2 d_1 + d_2 + \dots + d_L d_{L-1} + d_L$$

$$\underbrace{w^{(1)} \quad b^{(1)}}_{w^{(L)} \quad b^{(L)}}$$

\Rightarrow On note $x_i^{(L)}(w)$ la sortie du réseau de neurones avec x_i en entrée et w en paramètres

Problème d'optimisation

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

$$\frac{1}{n} \sum_{i=1}^n \|x_i^{(L)}(w) - y_i\|^2$$

Prédiction à partir de x_i

- Si $L > 1$, problème non convexe
- Si $\sigma = \max(\cdot, 0)$, fonction objectif non lisse
- Si $L \gg 1$ et/ou $d \gg 1$, très coûteux de calculer la fonction objectif ou son gradient
- Si $n \gg 1$, coût supplémentaire qui n'est pas forcément nécessaire si les données proviennent d'une même distribution

↳ Dans les deux exemples, la fonction objectif du problème d'optimisation possède une structure de somme finie ($\frac{1}{m} \sum_{i=1}^m \dots$), où chaque terme de la somme dépend d'un point distinct du jeu de données.

Q) Si l'on dispose d'un algorithme (de type descente de gradient) pour résoudre le problème, comment l'adapter pour prendre en compte la structure en somme finie et la distribution sous-jacente des données ?

② Algorithme du gradient stochastique

→ Développé dans les années 1940-1950 par des statisticiens (Robbins & Monro)

→ regain d'intérêt à la fin des années 2000 avec l'utilisation de méthodes de gradient stochastique pour entraîner les réseaux de neurones (+ bien avec l'algorithme du perceptron)

⇒ Succès flag en 2012 : Performances de réseaux de neurones en classification d'images

Problème

(P)

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

$$f(w) = \frac{1}{m} \sum_{i=1}^m f_i(w)$$

avec $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$, de classe C^1 et qui dépend du i^e point d'un jeu de données $\approx m$ éléments

$f \in C^1$ et $\Rightarrow f$ est C^1 , et on peut donc appliquer la descente de gradient au problème (P)

Itération de la descente de gradient

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

avec $\alpha_k > 0$

gradient de f en w_k
 w_k n'est pas minimum
 $\|\nabla f(w_k)\| \neq 0$

Ran (P),

$$\nabla f(w_k) = \frac{1}{m} \sum_{i=1}^m \nabla f_i(w_k)$$

\Rightarrow Un calcul de ∇f nécessite d'accéder à tous les points du jeu de données (coûteux !)

Itération du gradient stochastique

$\nabla f_i(w)$: gradient de f_i évalué en w

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

avec $\alpha_k > 0$

et i_k tiré aléatoirement entre 1 et m
 $i_k \in \{1, \dots, m\}$

- (+) Gén: 1 seul accès à un point du jeu de données par itération
- (-) Approximation du gradient de f par celui de f_{i_k} \Rightarrow risque d'utiliser une direction qui n'améliore pas la fonction
- (?) Aléatoire: comment tirer les i_k ? Est-ce que l'aléatoire est bénéfique?

2 Théorie du gradient stochastique

(P) minimiser $w \in \mathbb{R}^d$

$$f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$$

Hypothèses :

- $f(w) \geq \bar{f}$, $\bar{f} \in \mathbb{R}$
- $f_i \in C^1$ et $i \in \{1, \dots, n\}$
- $f \in C_L^{1,1}$, c'est à dire que

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

avec $L > 0$

$$\Rightarrow \forall (v, w) \in (\mathbb{R}^d)^2,$$

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

Inégalité fondamentale pour la théorie de la descente de gradient

Pour la descente de gradient

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

et donc par l'inégalité fondamentale

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

$$< f(w_k) \text{ si } \alpha_k < \frac{2}{L}$$

- Permet de définir des valeurs pour α_k qui garantissent une diminution de f à chaque itération
- Permet aussi de montrer des garanties de complexité pour la descente de gradient

Pour $\{f_k\}$ bien choisie (ex: $\alpha_k = \frac{1}{L} + k$), la descente de gradient calcule w_k tel que $\|\nabla f(w_k)\| \leq \varepsilon$ en au plus

$\mathcal{O}(\varepsilon^{-2})$ itérations

$C\varepsilon^2$, où $C > 0$ ne dépend pas de ε

Pour le gradient stochastique

$$w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k) \quad \alpha_k > 0$$

\nwarrow indice aléatoire entre 1 et m

Inégalité fondamentale appliquée à (w_{k+1}, w_k)

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

↑
vecteur
choisi aléatoirement

→ L'inégalité ne permet plus de garantir $f(w_{k+1}) < f(w_k)$ à cause de l'aléatoire sur i_k .

→ On peut cependant garantir une décroissance sur moyenne / espérance sous les bonnes hypothèses

Hypothèses sur les i_k

À l'itération k , on suppose que

1) i_k est tiré indépendamment de i_0, \dots, i_{k-1}

2) $E_{i_k} [\nabla f_{i_k}(w_k)] = \nabla f(w_k)$

$$3) \quad \mathbb{E}_{i_k} \left[\|\nabla f_{i_k}(w_k)\|^2 \right] \leq \sigma^2 \quad \text{avec } \sigma^2 > 0$$

Interprétation

- 1) A chaque itération, on a la même probabilité sur chaque valeur entre $\{1, \dots, n\}$
- 2) En moyenne, le gradient stochastique $\nabla f_{i_k}(w_k)$ équivaut à $\nabla f(w_k)$

A variable
qui dépend
de i_k , $\mathbb{E}_{i_k} [A(i_k)] = \sum_{i=1}^m P(i_k=i) A(i)$

$$\mathbb{E}_{i_k} [\nabla f_{i_k}(w_k)] = \sum_{i=1}^m P(i_k=i) \nabla f_i(w_k)$$

- 3) La norme des gradients stochastiques ne varie pas trop

Exemple casque (typique de l'implémentation)

Tirer i_k selon une loi uniforme

$$P(i_k=i) = \frac{1}{m}$$

$\Rightarrow 1), 2) (+ 3) \text{ en général}$



$$\begin{aligned} \mathbb{E}_{i_k} [\nabla f_{i_k}(w_k)] &= \sum_{i=1}^m P(i_k=i) \nabla f_i(w_k) \\ &= \sum_{i=1}^m \frac{1}{m} \nabla f_i(w_k) = \underbrace{\frac{1}{m} \sum_{i=1}^m \nabla f_i(w_k)}_{\text{définition}} = \nabla f(w_k) \end{aligned}$$

↪ Sous ces hypothèses sur α_k , on montre qu'à chaque itération

$$E_{\text{th}}[f(w_{k+1})] \leq f(w_k) - \alpha_k \|\nabla f(w_k)\|^2 + \frac{\frac{L}{2} \alpha_k^2}{2}$$

Ne dépend plus de i_k

→ Permet de garantir $E_{\text{th}}[f(w_{k+1})] < f(w_k)$

"Décroissance en moyenne"

pour des bons choix de α_k

$$\text{Ex)} \quad \{\alpha_k\}_k \rightarrow 0 \quad (\text{par défaut})$$

$$0 < \alpha_k \leq \frac{1}{L} \frac{\|\nabla f(w_k)\|^2}{\sigma^2} \rightarrow \text{par défaut par l'itération } k$$

→ Définir une stratégie de choix de pas constant pour la théorie

→ En pratique, on choisit un pas constant (ex: dans PyTorch, $\alpha_k = 0.01$)

NB: Les stratégies de recherche linéaire ne sont pas utilisées pour les méthodes de gradient stochastique parce qu'elles nécessitent d'évaluer f , donc d'accéder à tous les points du jeu de données.

Théorème

Sous les hypothèses précédentes, supposons que l'on choisit $\alpha_k = \frac{\varepsilon}{L\sigma^2} \gamma_k$ avec $\varepsilon > 0$.

Alors, la méthode de gradient stochastique

calculé en moyenne un point w_k tel que

$$\|\nabla f(w_k)\| \leq \varepsilon \text{ en au plus } O(\varepsilon^{-4})$$

iéritations.

Formellement, après $O(\varepsilon^4)$ itérations, $E\left[\min_{0 \leq l \leq k} \|\nabla f(w_l)\|\right] \leq \varepsilon$

Comparaison avec la descente de gradient

<u>Descente de gradient:</u>	$\alpha_k = \frac{1}{L} + k$	$\min_{0 \leq l \leq h} \ \nabla f(w_l)\ \leq \varepsilon$	$O(\varepsilon^{-2})$ itérations
<u>Graident stochastique</u>	$\alpha_k = \frac{\varepsilon^2}{L\sigma^2} + k$ $\rightarrow DG: \text{pas indépendant de } \varepsilon$ $\rightarrow GS: \text{pas fonction de } \varepsilon, \text{ qui doit être redéfini si on change la valeur de } \varepsilon$	$E\left[\min_{0 \leq l \leq h} \ \nabla f(w_l)\ \right] \leq \varepsilon$ $\rightarrow DG: \text{Convergence déterministe}$ $\rightarrow GS: \text{Convergence en moyenne}$	$O(\varepsilon^{-4})$ itérations $\varepsilon^2 \text{ augmente moins vite que } \varepsilon^{-4} \text{ quand } \varepsilon \text{ diminue}$ $\varepsilon \rightarrow \varepsilon/10$ $\varepsilon^2 = 100\varepsilon^{-2}$ $\varepsilon^{-4} = 10000\varepsilon^{-1}$

\Rightarrow Si l'on regarde uniquement ces résultats, la descente de gradient semble bien meilleure que le gradient stochastique

\Rightarrow Mais ces résultats ne tiennent pas compte du coût de calcul des gradients, qui était la raison pour laquelle on a introduit le gradient stochastique

1 itération de descente de gradient = 1 calcul de $\nabla f(w)$
= m accès à des points du jeu de données

1 itération de gradient stochastique = 1 calcul de $\nabla f_i(x)$
 = 1 accès à un point du jeu de données

⇒ En termes d'accès aux données, 1 itération de gradient stochastique est m fois moins coûteuse qu'une itération de descente de gradient

=> Pour pouvoir comparer ces deux algorithmes de manière équitable, il faut utiliser une autre quantité que le nombre d'itérations

Définition: Une époque (epoch en anglais) est définie comme le coût d'accès à n points d'un jeu de données à m éléments

NB: Si on accède m fois au même point, on compte m accès

1 itération de descente de gradient a un coût de 1 époque

1 itération de gradient stochastique a un coût de $\frac{1}{m}$ époques

Consequences sur la complexité

Descente de gradient	$O(\varepsilon^{-2})$ itérations	$O(\varepsilon^{-2})$ époques
Gradient stochastique	$O(\varepsilon^{-4})$ itérations	$O(\frac{1}{m}\varepsilon^{-4})$ époques

Lorsque $m \gg 1$, on peut avoir $m \gg \varepsilon^{-2}$ et alors la complexité du gradient stochastique est meilleure

⇒ Explique en partie pourquoi une méthode de gradient stochastique converge plus vite qu'une descente de gradient sur les premières époques.

③ Variantes de la méthode du gradient stochastique

Utiliser plus qu'un accès au jeu de données par itération

DG: m accès/itération

GS: l'accès/itération

$1 < m_b < m$ accès par itération?

Méthode de gradient stochastique par fournées (batch stochastic gradient)

$$w_{k+1} = w_k - \alpha_k \left(\frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(w_k) \right)$$

$$w_{h+1} = w_h - \frac{\alpha_h}{|S_h|} \sum_{i \in S_h} \nabla f_i(w_h)$$

$|S_h|$: taille de l'ensemble

S_h

avec $\alpha_h > 0$ et S_h un ensemble d'indices tirés aléatoirement entre 1 et m , avec ou sans remise

→ Généralise à la fois le gradient stochastique ($S_h = \{i_h\}$, $|S_h|=1$) et la descente de gradient ($S_h = \{1, \dots, n\}$, $|S_h|=n$)

tirage sans remise

→ Permet notamment de prendre en compte les possibilités de calcul parallèle des $\nabla f_i(w_h)$

\Rightarrow Si on dispose de r processeurs qui peuvent calculer des gradients en parallèle, on pose $|S_h| = r$
(en général $r \ll n$)

- \rightarrow Deux catégories de méthodes par lot (taille de lot m_b):
- Grande taille: $m_b \gg 1$, $m_b \ll n$
 - (+) Pas de variabilité entre les exécutives
 - (-) Coût comparable à celui de la descente de gradient
 - Petites/mini-tailles (mini-batch): $m_b > 1$, $m_b \ll n$
 - (+) Coût bien plus faible que celui de la descente de gradient
 - (+) Réduit la variance des gradients stochastiques
 - (mathématiquement, $E[\|\nabla f_{m_b}(w_t)\|^2] \leq \sigma^2$)
 $\rightarrow E\left[\left\|\frac{1}{|S_h|} \sum_{i \in S_h} \nabla f_i(w_t)\right\|^2\right] \leq \frac{\sigma^2}{m_b}$
 - (-) Plus cher (à chaque itération) que le gradient stochastique, et il peut arriver que le gradient stochastique classique soit quand même meilleur en pratique
 - (-) Choisir la bonne taille de lot est un problème en soi, qui dépend beaucoup des données considérées

Autres variantes

- Pour f non lisse, on peut utiliser des sous-gradients à la place des gradients ou des variantes du gradient proximal
- On peut considérer des problèmes régularisés (gradient proximal stochastique)

- On peut faire des méthodes d'ordre 2 (cf projet)
- On peut incorporer du momentum dans l'algorithme du gradient stochastique

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

↑
Terme de momentum

- (+) Marche bien en pratique, base de la méthode d'entraînement des réseaux modernes (Adam)
2011
- (-) Pas de théorie complète de convergence