

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

27 octobre 2023

Programme du jour:

- Cours gradient stochastique (1/2)
- TD gradient stochastique

GRADIENT

Contexte : Problème

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

Hypothèse : chaque f_i
dans un jeu de

$$\text{Ex) } \{ (x_i, y_i) \}_{i=1}^m \text{ avec}$$

But : Trouver un
que $x_i^T w \approx y_i$

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad f(w) = \frac{1}{2m} \sum_{i=1}^m (x_i^T w - y_i)^2$$
$$= \frac{1}{m} \sum_{i=1}^m f_i(w)$$

Forme générale en
apprentissage (supervisé)

$$\{ (x_i, y_i) \}_{i=1}^m$$

But : Trouver une
 $h(x_i) \approx y_i$

\Rightarrow En général, on
paramétrise par un

STOCHASTIQUE

en somme finie

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

dépend du i ème point
données à m éléments

$$x_i \in \mathbb{R}^d \text{ et } y_i \in \mathbb{R}$$

modèle linéaire tel

$$\text{avec } f_i(w) = \frac{1}{2} (x_i^T w - y_i)^2$$

fonction h telle que

considère que h est
vecteur $w \in \mathbb{R}^d$

\Rightarrow On quantifie
et y_i via une
 $l: (h, y) \mapsto l(h, y)$

l'écart entre $h(x_i)$
fonction de perte

On obtient le problème

d'optimisation

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

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

avec

$$f_i(w) =$$

$$l(h(x_i; w), y_i)$$

Rede
défini par w

\hookrightarrow h_i peut aussi bien
linéaire qu'un réseau

être un modèle
de neurones

$$x_i = x^0 \mapsto x^1 \mapsto x^2 \mapsto \dots \mapsto x^L = h(x_i, w)$$

$$\forall l=0..L-1, x^{l+1} = \sigma(W^l x^l + b^l)$$

$$W^l \in \mathbb{R}^{d_{l+1} \times d_l}$$

$$b^l \in \mathbb{R}^{d_{l+1}}$$

$$\sigma(v) = \begin{bmatrix} \sigma(v_i) \end{bmatrix}$$

activation

w : concaténation des
 W^l et b^l

$$d = d_0 d_1 + d_1 d_2 + d_2 d_3 + \dots$$

$$\sigma(t) = \max(t, 0)$$

$$\sigma(t) = \tanh(t)$$

BUT: Exploiter la
forme dans le
et notamment
de données
où l'accès aux
coûteux

structure en somme
processus d'optimisation
dans un contexte
matrices ($n \gg 1$)
données peut être

1) Algorithme du

↳ On considère un
fonc. et on suppose que
sont de classe C^1 .

⇒ $f = \frac{1}{n} \sum f_i$ est
et $\forall w \in \mathbb{R}^d$, $\nabla f(w)$

↳ Si on applique le
à ce problème, on obtient

$$w_{k+1} = w_k -$$
$$= w_k -$$

1 itération de
doit accéder à tout le
calculer le prochain

Itération de

$$w_{k+1} = w_k$$

gradient stochastique

problème on somme
les fonctions $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$

aussi de classe C^1

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

descente de gradient
l'itération suivante:

$$- \alpha_k \nabla f(w_k)$$

avec $\alpha_k > 0$

$$\frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_i(w_k)$$

descente de gradient
jeu de données pour
itérations

gradient stochastique (Robbins & Monro 1951)

$$- \alpha_k \nabla f_{i_k}(w_k)$$

i_k tiré aléatoirement
dans $\{1, 2, \dots, n\}$

(+) 1 itération de n accès qu'à 1 donnée

$\Rightarrow n$ fois moins de descente de d'accès aux

gradient stochastique

point du jeu de

coûts qu'une itération gradient en ramenant données

(?) Est-ce que cette l'objectif f et

\hookrightarrow En général (les f_i), cette

\hookrightarrow MAIS sur (venant d'une méthode converge en la descente de

itération améliore convergence ?

(sans hypothèses sur méthode ne converge pas

des données réelles distribution), la

moyenne mieux que gradient !

2) Comparaison entre et descente de gradient

\hookrightarrow Rappel: La descente méthode déterministe peut montrer des

Ex) Si f est convexe, itérations de descente $f(w_k) - \min_{w \in \mathbb{R}^d} f(w)$

gradient stochastique

de gradient est une et pour laquelle on vitesses de convergence

alors après $k \geq 1$ de gradient, $\leq O\left(\frac{1}{k}\right)$

↳ On peut prouver un résultat similaire sur le gradient stochastique

résultat similaire sous certaines hypothèses, notamment

$$\forall k \in \mathbb{N}, \mathbb{E}_{i_k} [\nabla f_{i_k}(w_k)] = \nabla f(w_k)$$

↳ Espérance par rapport à i_k

↳ gradient stochastique i_k aléatoire dans $\{1, \dots, m\}$

$$= \nabla f(w_k)$$

↳ "vrai gradient" $\frac{1}{m} \sum_{i=1}^m \nabla f_i(w_k) \approx$ moyenne des ∇f_i

Ex) Si i_k est tiré uniformément au hasard dans $\{1, \dots, m\}$ (càd $\mathbb{P}(i_k=1) = \dots = \mathbb{P}(i_k=m) = \frac{1}{m}$) alors

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

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

Th) Si f est convexe, hypothèses sur $\{i_k\}_{k \in \mathbb{N}}$, de gradient stochastique,

alors sous les bonnes après $k \geq 1$ itérations on a:

$$\mathbb{E} [f(w_k) - \min_{w \in \mathbb{R}^d} f(w)] \leq O\left(\frac{1}{\sqrt{k}}\right)$$

↳ Espérance par rapport aux $\{i_k\}$

Interprétation: Le gradient stochastique, qui est une méthode

stochastique, qui est une aléatoire (ou "randomisée")

converge en moyenne

• Dans le cas
convergence est en $\frac{1}{\sqrt{K}}$
que la vitesse en $\frac{1}{K}$

⇒ Ces résultats
Stochastique est moins
de gradient pour un

⇒ Mais le coût d'une
stochastique est différent
de descente de gradient!

GS : 1 accès
données

DG : n accès

Métrique pertinente pour
Époques (epochs en

Def: 1 époque -
un point d'un jeu

Avec cette définition:

1 itération de DG
et 1 itération de GS

(\Leftrightarrow 1 époque est le coût

↳ La notion d'époque
les algorithmes pour un
données fixé, ce qui

convexe, la vitesse de
, ce qui est plus lent
de la descente de gradient

suggèrent que le gradient
efficace que la descente
même nombre d'itérations

itérations de gradient
de celui d'une itération

à un point du jeu de

(accès à tout le jeu de données)

comparer les deux méthodes:
anglais)

correspond à n accès à
de données à n éléments

coûte 1 époque

coûte $\frac{1}{n}$ époque

de n itérations de GS)

permet de comparer
budget d'accès aux
est une métrique

plus approprié que le nombre d'itérations

Ex: f convexe

	$K \geq 1$ itérations	$N_E \geq 1$ époques
DG	$\frac{1}{K}$	$\frac{1}{N_E}$ ← 1 époque = 1 itération
GS	$\frac{1}{\sqrt{K}}$	$\frac{1}{\sqrt{m N_E}}$ ← 1 époque m itérations

Lorsque $m \gg 1$,
général $\frac{1}{\sqrt{m N_E}} \ll \frac{1}{N_E}$

3) Gradient stochastique

"avec batch"
par fournées
le gradient stochastique
d'un point du jeu de
itération

↳ Plus général que
↳ Idee: Considérer plus
données lors d'une

Itération k

$$w_{k+1} = w_k$$

$$- \frac{\alpha_k}{|S_k|} \sum_{i \in S_k} \nabla f_i(w_k)$$

Ex). $|S_k| = 1$ $S_k = \{i_k\}$

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

⇒ Gradient
stochastique

• $|S_k| = m$ et tirage sans
remise

$$\Rightarrow S_k = \{1, \dots, m\}$$

$$w_{k+1} = w_k - \alpha_k \frac{1}{m} \sum_{i=1}^m \nabla f_i(w_k)$$

S_k est un ensemble d'indices
tirés aléatoirement dans
 $\{1, \dots, m\}$ avec ou sans remise

S_k : "fournée"

$|S_k|$: nombre d'éléments de S_k
"taille de fournée"

$$= w_k - \alpha \nabla f(w_k)$$

\Rightarrow Descente de gradient

\hookrightarrow La méthode par gradient stochastique car particuliers

\hookrightarrow Pour les cas restants, grandes classes :

* Grande fonnée :

- Comportement gradient
- Coût très celui du gradient

(Plutôt utilisé en d'autres variantes car

* Mini-fonnées :

- Coût qui reste la descente de gradient
- Bénéfice possible stochastique (meilleure stochastique)

Remarque: En pratique, classes car difficile et le choix de la problème complexe

fonnées couvre le et la descente couvre

on distingue deux

$$|S_k| \gg 1$$

proche de la descente de Supérieur (par itération) à stochastique

combinaison avec très coûteux)

$$|S_k| > 1 \text{ et } |S_k| \ll n$$

très inférieur à celui de par rapport au gradient variance du gradient

La frontière entre les deux à déterminer a priori taille de fonnée est un

↳ Un choix classique
celui du nombre de
des calculs en parallèle

⇒ Ce cadre est
par mini-journées
impossible de paralléliser
parallélise pas par

de taille de journées est
processeurs disponibles pour

(ex: $|S_k| = 32$)

favorable aux variantes

(surtout $m \gg 32$ donc

DG et GS ne se

construisent)