

OUTILS D'OPTIMISATION POUR LES SCIENCES DES DONNÉES ET DE LA DÉCISION

Programme : Gradient stochastique 2/2 + TD avec notebook

Semaine prochaine:

- Jeudi 8h30-11h45
- Vendredi 8h30-11h45

GRADIENT ET RÉDUCTION DE

① Variance dans le

Cadre: $f: \mathbb{R}^d \rightarrow \mathbb{R}$

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

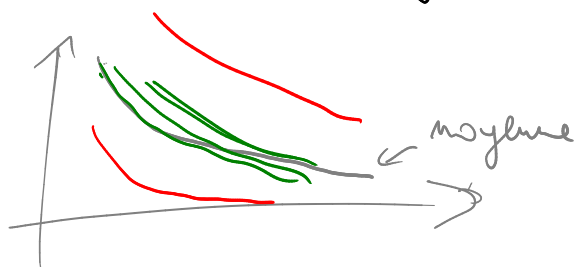
Gradient stochastique

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

↳ On suppose en général

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

⇒ Garanties sur le
du gradient stochastique
quel point on peut
comportement moyen



STOCHASTIQUE VARIANCE

gradient stochastique classique

μ -fortement convexe
 $C^{\frac{1}{2}, 1}$

f_i dépend du même
point d'un jeu de données
à n éléments

$i_k \in \{1, \dots, n\}$
tiré aléatoirement

que

(en moyenne, on retrouve
le vrai gradient)

comportement moyen
mais ne dit pas à
dériver de ce

Q: Comment analyser
l'écart entre les
réalisations de l'algorithme
et la moyenne?

Hypothèse :

$$\mathbb{E}_{i_k} [\|\nabla f_{i_k}(w_k)\|_2^2] - \|\nabla f(w_k)\|_2^2$$

↑ Moyenne par rapport à i_k ↑ Norme du gradient stochastique

$$\|\nabla f(w_k)\|_2^2 \leq \sigma^2$$

↑ Norme du vrai gradient $\sigma^2 \geq 0$

Si $\mathbb{E}_{i_k} [\nabla f_{i_k}(w_k)] = \nabla f(w_k)$, alors

$$\mathbb{E}_{i_k} [\|\nabla f_{i_k}(w_k)\|_2^2] - \|\nabla f(w_k)\|_2^2 = \mathbb{E}_{i_k} [\|\nabla f_{i_k}(w_k)\|_2^2 - \|\mathbb{E}_{i_k} [\nabla f_{i_k}(w_k)]\|_2^2]$$

"Variance du gradient stochastique"

$\mathbb{E}[a]$: moyenne
 $\text{var}[a] = \mathbb{E}[a^2] - (\mathbb{E}[a])^2$
variance

σ^2 : Paramètre de variance (ou niveau de bruit) du gradient

- Dépend des données (f_i)
- Dépend de la façon dont on construit les gradients (et tirage des i_k)

variance (ou niveau stochastique données (f_i))

façon dont on construit stochastiques (et tirage des i_k)

Théorème: Si on effectue K itérations de gradient stochastique avec $\alpha_k = \frac{1}{L}$, alors

effectue K itérations de gradient stochastique avec $\alpha_k = \frac{1}{L}$, alors

$$\mathbb{E} [f(w_K) - f^*] \leq$$

$$\frac{\sigma^2}{2L} + \left(1 - \frac{\mu}{L}\right)^K (f(w_0) - f^* - \frac{\sigma^2}{2L})$$

avec $f^* = \min_{w \in \mathbb{R}^d} f(w)$

Rappel: Sur le même problème, vérifie

la descente de gradient

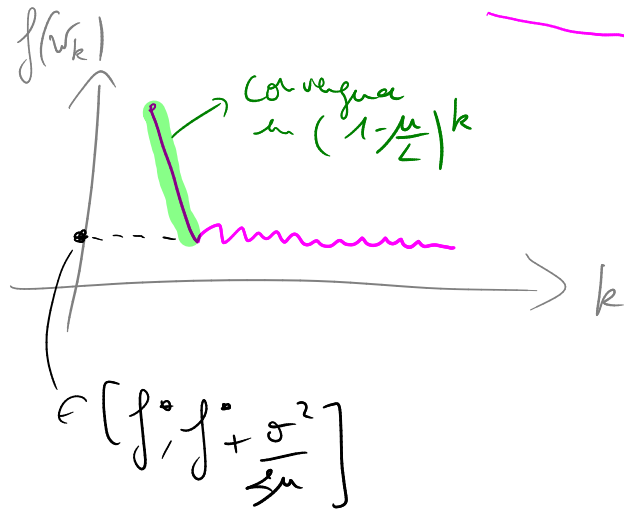
$$0 \leq f(w_k) - f^* \leq \left(1 - \frac{\mu}{L}\right)^k (f(w_0) - f^*)$$

\Rightarrow Garantit que $f(w_k) \rightarrow f^*$ à vitesse $\left(1 - \frac{\mu}{L}\right)^k$

\hookrightarrow En comparaison, le gradient stochastique

résultat sur le garantit en moyenne

que $f(w_k) \rightarrow \left[f^* - f^* + \frac{\sigma^2}{2\mu} \right]$



convergence dans un intervalle de valeurs défini par le niveau de bruit

Interprétation

* Plus les gradients vrai gradient, plus une peut varier par rapport à
 * En pratique, cela converge jusqu'à

stochastiques différent de l'exécution de l'algorithme sa moyenne signifie qu'on ne peut un certain niveau de

bruit sur les gradients

stochastiques

(2) Méthodes de

↳ Idee: Avec un
stochastique classique,
variance sur les

⇒ Comment modifier
réduire cette variance?

a) Utiliser une famille

Gradient stochastique par

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

$$\nabla f(w_k) \approx \frac{1}{m_b} \sum_{i \in S_k} \nabla f_i(w_k)$$

et si on a un niveau
 $m_b = 1$ (GS de base), alors

bruit $\frac{\sigma^2}{m_b}$ pour $m_b > 1$

réduction de variance

algorithme de gradient

on a une certaine

gradients stochastiques

l'algorithme pour

de gradients stochastiques

familles

Soit ensemble d'indices
tirés aléatoirement
dans $\{1, \dots, m\}$
(avec ou sans remise)

$$|S_k| = m_b$$

de bruit σ^2 lorsque
on a un niveau de

indices tirés avec
remise

b) Moyenne des itérés

↳ Dans le cas convexe
montrer que $\left\{ \frac{1}{k+1} \sum_{l=0}^k w_l \right\}$
variance que $\{w_k\}$
meilleure valeur de

↳ En pratique, peu
coût de stockage de
trop grand et des
possibles dans le calcul

notamment, on peut
à une plus faible
et converge vers une
fonction

utilisé à cause du
la moyenne qui peut être
erreurs numériques
de la moyenne

c) Agrégation de gradients

↳ Combinaison des
stochastiques avec des
d'itérations précédentes

⇒ Conserve la
stochastique

↳ Algorithmes: SVRG, SAGA (cf scikit-learn)

⇒ Très efficaces sur
faible dimension et

⇒ Le coût de calcul
dans certains contextes (ex:

pas de gradients
vrais gradients issus

valeur du gradient

des problèmes de
convexes

peut être trop élevé
apprentissage profond)

③ Variantes du gradient l'apprentissage profond

Forme générale

- Calculer g_k
- Calculer

$$W_{k+1} = W_k$$

m_k : direction du pas
 $m_k \in \mathbb{R}^d$

v_k : vecteur de normalisation

↳ A chaque itération k ,
coordonnée j est $\frac{\alpha}{[v_k]_j}$

Cas de base

$$m_k = g_k$$

⇒ Gradient

① Gradient stochastique

$$m_k = \beta_1 m_{k-1} + (1 - \beta_1) g_k$$

$\beta_1 \in (0, 1)$
↑ direction précédente

stochastique pour

$$J_{\text{fin}}(w_k) \quad \frac{1}{n_b} \sum_{i \in S_k} J_i(w_k)$$

↓
(gradient stochastique / fourni)

— α m_k v_k

taille ↑
de pas
constante
 $\alpha > 0$

division
composante a
composante
 $a \otimes b = \begin{bmatrix} [a]_i \\ [b]_j \end{bmatrix}$

(dépend de g_k)

de la taille de pas

le pas selon la

$$v_k = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

stochastique classique

avec momentum

$$v_k = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

↳ Adaptation de
au cas d'un gradient
↳ Méthode utilisée
réseaux lors de l'émou
en 2012

l'idée de momentum
stochastique
pour entraîner les
de l'apprentissage profond

⑥ Adagrad

$$m_k = g_k$$

$$\forall j=1..d$$

$$[v_k]_j = \sqrt{\sum_{l=0}^k [g_l]_j^2}$$

⇒ taille de pas selon
la coordonnée j

$$\frac{\alpha}{\sqrt{\sum_{l=0}^k [g_l]_j^2}}$$

→ Méthode efficace
Gradients parcimonieux
nulles), typique des

pour des problèmes à
(beaucoup de coordonnées
systèmes de recommandation

→ Peut cependant conduire
très petites ϵ trop rapidement

à des tailles de pas

→ Des variantes comme
autre formule :

RMS Prop utilisent une

$$[v_k]_j = \beta_2 \epsilon (0,1)$$

$$\sqrt{\beta_2 [v_{k-1}]_j^2 + (1-\beta_2) [g_k]_j^2}$$

à l'itération $k-1$ valeur courante

© Adam (~2015)

↳ Combinaison pondérée par les gradients stochastiques poids aux gradients calculés

$$m_k = \frac{(1-\beta_1)}{1-\beta_1^{k+1}} \sum_{l=0}^k \beta_1^{k-l} g_l$$

↳ En pratique, β_1 et β_2 (0.9, 0.99)

↳ LA méthode par la plupart des architectures

⇒ Très efficace dans les tâches de traitement naturel des langues

Remarque: La théorie de présentée pour Adam fautive en 2018.

⇒ Toujours des théories

⇒ Mais la méthode en pratique, et elle reste

de l'information fournie pour donner plus de la plus récemment

$$t_j = 1-d,$$

$$[v_k]_j = \sqrt{\frac{1-\beta_2}{1-\beta_2^{k+1}} \sum_{l=0}^k \beta_2^{k-l} [g_l]_j^2}$$

sont pris très proches de 1

déjà pour entraîner neurales actuellement

problèmes de traitement (NLP)

convergence initialement en 2015 s'est révélée

zones d'ombre dans la

est extrêmement efficace donc très utilisée