

Optimisation pour le Machine Learning

①

Contexte: En apprentissage statistique, il est souvent utile de résoudre des problèmes de la forme

$$\hat{\theta} \approx \underset{\theta \in \Theta}{\operatorname{argmin}} F(\theta)$$


exemples:

- Minimisation du risque empirique

$$F(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i)) + \Omega(\theta)$$

- Maximum de vraisemblance

$$f(\theta) = -\log(p(x_1, \dots, x_n | \theta))$$

 Parfois, il existe des solutions en forme close (moindres carrés, Ridge, ...), mais souvent ce n'est pas le cas et il faut alors utiliser un algorithme d'optimisation (qui approche la solution). Il faut alors prendre en compte cette nouvelle source d'erreur.

Décomposition du risque:

$$M(\hat{\theta}) - M(\theta^*) = \underbrace{M(\hat{\theta}) - \hat{M}(\hat{\theta})}_{\text{deviation I}} + \underbrace{\hat{M}(\hat{\theta}) - \hat{M}(\theta^*)}_{\leq \hat{M}(\hat{\theta}) - \inf_{\theta} \hat{M}(\theta)} + \underbrace{\hat{M}(\theta^*) - M(\theta^*)}_{\text{deviation II}}$$

= erreur d'optimisation.

II. Descente de gradient

(2)

④ \mathbb{R}^d , F différentiable.

↳ Localement, $F(\theta + s\theta) = F(\theta) + \langle \nabla F(\theta), s\theta \rangle + o(\|s\theta\|)$

Donc: Localement, pour faire décroître F le plus possible, f peut partir dans la direction $-\nabla F(\theta)$.

Algorithme de descente de gradient

Entrée: Un point de départ θ_0

Mécanisme: $\theta_t = \theta_{t-1} - \gamma_t \nabla F(\theta_{t-1})$ (Descente de Gradient)

où $(\gamma_t)_{t \geq 1}$ est une suite de tailles de pas.

Le point de vue des systèmes dynamiques

Il est possible de voir la descente de gradient comme une discrétisation d'Euler explicite de l'équation

$$\dot{\theta} = -\nabla F(\theta) \quad (\text{Gradient Flow})$$

$$\text{Alors } \frac{dF(\theta)}{dt} = \langle \nabla F(\theta), \dot{\theta} \rangle = -\langle \nabla F(\theta), \nabla F(\theta) \rangle = -\|\nabla F(\theta)\|^2 \leq 0$$

↳ Tant que la dynamique n'est pas à un point d'arrêt le gradient est nul, elle décroît.

1) Fonctions lisses et fortement convexes

(3)

Définition: Une est μ -fortement convexe ($\mu > 0$) si

$$\forall \eta, \theta, F(\eta) \geq F(\theta) + \langle \nabla F(\theta)^T, \eta - \theta \rangle + \frac{\mu}{2} \|\eta - \theta\|_2^2 \quad (*)$$

exemple: Fonctions quadratiques à valeurs propres > 0 .

exemple: $\|\cdot\|_2^2$ régularisation

Propriété (Inégalité de Łojasiewicz):

Si f est différentiable et μ -fortement convexe, avec pour minimiseur unique η^* , alors

$$\|\nabla f(\theta)\|_2^2 \geq 2\mu (F(\theta) - F(\eta^*)), \quad \forall \theta$$

démonstration: ds (*), le terme de droite est minimisé par

$$\eta = \theta - \frac{1}{\mu} \nabla f(\theta) \text{ ce qui donne le résultat.}$$

Conséquence d'un point de vue continu:

$$\dot{\theta} = -\nabla F(\theta).$$

$$\Rightarrow \frac{dF(\theta)}{dt} = -\|\nabla F(\theta)\|_2^2 \leq -2\mu (F(\theta) - F(\eta^*))$$

$$\Rightarrow \frac{d(F(\theta) - F(\eta^*))}{dt} \leq -2\mu (F(\theta) - F(\eta^*))$$

$$\text{(Grönwall)} \Rightarrow F(\theta) - F(\eta_*) \leq [F(\theta_0) - F(\eta_*)] e^{-2t}$$

(4)

Pour obtenir la convergence exponentielle des itérées discrètes, il faut avoir une forme de régularité des gradients pour assurer que l'approximation discrète se rapproche de la dynamique continue.

Définition: F est L -lisse si

$$\forall \theta, \eta, |F(\eta) - F(\theta) - \langle \nabla f(\theta), \eta - \theta \rangle| \leq \frac{L}{2} \|\theta - \eta\|_2^2,$$

ou alors, de manière équivalente, si ∇f est L -lipschitz pour $\|\cdot\|_2$.

Théorème: Si F est L -lisse et μ -fortement convexe, GD avec $\tau_t = 1/L$ donne une suite $(\theta_t)_t$ satisfaisant

~~$$F(\theta_t) - F(\eta_*) \leq \left(1 - \frac{1}{\kappa}\right)^t (F(\theta_0) - F(\eta_*))$$~~

$$F(\theta_t) - F(\eta_*) \leq \underbrace{\left(1 - \frac{1}{\kappa}\right)^t}_{\leq e^{-t/\kappa}} (F(\theta_0) - F(\eta_*))$$

où $\kappa = \frac{L}{\mu}$ est le conditionnement.

démonstration:

$$F(\sigma_t) = F\left(\sigma_{t-1} - \frac{1}{L} \nabla F(\sigma_{t-1})\right)$$

$$\stackrel{\text{L-Pisze}}{\leq} F(\sigma_{t-1}) + \langle \nabla F(\sigma_{t-1}), -\nabla F(\sigma_{t-1})/L \rangle + \frac{L}{2} \|\nabla F(\sigma_{t-1})/L\|^2$$

$$= F(\sigma_{t-1}) - \frac{1}{L} \|\nabla F(\sigma_{t-1})\|_2^2 + \frac{1}{2L} \|\nabla F(\sigma_{t-1})\|_2^2$$

donc

$$(**) \quad \underline{F(\sigma_t) - F(\eta_*) \leq (F(\sigma_{t-1}) - F(\eta_*)) - \frac{1}{2L} \|\nabla F(\sigma_{t-1})\|_2^2}$$

Par μ -forte convexité et inégalité de Lojasiewicz :

$$F(\sigma_t) - F(\eta_*) \leq \left(1 - \frac{\mu}{L}\right) (F(\sigma_{t-1}) - F(\eta_*))$$

□

2) Le cas Pisze et (simplemet) convexe

Dans ce cas, l'analyse est légèrement plus complexe.

Proposition (Co-coercivité):

Si F est L -Pisze et convexe,

$$\frac{1}{2} \|\nabla F(\sigma) - \nabla F(\eta)\|_2^2 \leq \langle \nabla F(\sigma) - \nabla F(\eta), \sigma - \eta \rangle$$

$$\text{et } F(\sigma) \geq F(\eta) + \langle \nabla F(\eta), \sigma - \eta \rangle + \frac{1}{2L} \|\nabla F(\sigma) - \nabla F(\eta)\|_2^2$$

démonstration:

Nous allons prouver la deuxième inégalité, la première n'est que la somme de la deuxième et de la deuxième des ξ que nous avons échangé θ et γ .

Par convexité et ~~car~~ comme F est L -lisse,

$$F(\gamma) + \langle \nabla F(\gamma), \xi - \gamma \rangle \leq F(\xi) \leq F(\theta) + \langle \nabla F(\theta), \xi - \theta \rangle + \frac{L}{2} \|\theta - \xi\|_2^2$$

et en injecter dans cette inégalité le ξ qui minimise l'écart entre le terme de gauche et le terme de droite donne le résultat.

Corollaire: Si F est convexe et L -lisse, des les itérés de GD satisfait

$$\|\theta_t - \gamma_*\|_2^2 \leq \|\theta_{t-1} - \gamma_*\|_2^2 - \frac{1}{L} \langle \nabla F(\theta_{t-1}), \theta_{t-1} - \gamma_* \rangle$$

démonstration:

$$\begin{aligned} \|\theta_t - \gamma_*\|_2^2 &= \|\theta_{t-1} - \frac{1}{L} \nabla F(\theta_{t-1}) - \gamma_*\|_2^2 \\ &= \|\theta_{t-1} - \gamma_*\|_2^2 - \frac{2}{L} \langle \nabla F(\theta_{t-1}), \theta_{t-1} - \gamma_* \rangle \\ &\quad + \frac{1}{L^2} \|\nabla F(\theta_{t-1})\|_2^2 \\ &= \dots \\ &\quad + \frac{1}{L^2} \|\nabla F(\theta_{t-1}) - \nabla F(\gamma_*)\|_2^2 \end{aligned}$$

$$\leq \quad \text{"} \quad + \frac{1}{L} \langle \nabla F(\theta_{t-1}) - \nabla F(\gamma_0), \theta_{t-1} - \gamma_0 \rangle \quad \textcircled{2}$$

(C-coercivité)

$$= \|\theta_{t-1} - \gamma_0\|_2^2 - \frac{1}{L} \langle \nabla F(\theta_{t-1}), \theta_{t-1} - \gamma_0 \rangle$$

□

Théorème: Si F est convexe et L -lisse, alors les itérés $(\theta_t)_t$ de GD satisfait

$$F(\theta_t) - F(\gamma_0) \leq \frac{L}{t+1} \|\theta_0 - \gamma_0\|_2^2$$

démonstration:

(**) est toujours valide, donc $(F(\theta_t) - F(\gamma_0))_t$ est décroissant.

$$\text{Donc } F(\theta_t) - F(\gamma_0) \leq \frac{1}{t+1} \sum_{i=0}^t (F(\theta_i) - F(\gamma_0)).$$

De plus, par convexité, $\forall i, F(\theta_i) - F(\gamma_0) \leq \langle \nabla F(\theta_i), \theta_i - \gamma_0 \rangle$

d'où, par le résultat précédent,

$$F(\theta_t) - F(\gamma_0) \leq \frac{L}{t+1} \sum_{i=0}^t \left(\|\theta_i - \gamma_0\|_2^2 - \|\theta_{i+1} - \gamma_0\|_2^2 \right)$$

$$\leq \frac{L}{t+1} \|\theta_0 - \gamma_0\|_2^2$$

□

3) Le cas non-converge :

Comme (**) reste valide,

$$F(\sigma_t) - F(\gamma_*) \leq F(\sigma_{t-1}) - F(\gamma_*) - \frac{1}{2L} \|\nabla F(\sigma_{t-1})\|_2^2$$

même dans le cas non-converge (mais L-lisse).

Alors, par télescopage

$$\frac{1}{2Lt} \sum_{i=1}^t \|\nabla F(\sigma_{i-1})\|_2^2 \leq \frac{F(\sigma_0) - F(\gamma_*)}{t}$$

4) Alors plus loin :

• Optimisation non lisse : Il est possible d'obtenir une convergence à vitesse $\frac{1}{\sqrt{t}}$ lorsque F n'est pas lisse (mais convexe) si celle dernière est Lipschitz (voir Bach)

• Accélération de Nesterov : Dans les cas fortement convexe, permet de passer de $\frac{1}{t}$ à $\frac{1}{\sqrt{t}}$. Dans les cas convexe, lisse, vitesse en $\frac{1}{t^2}$.

• Méthode du gradient proximal : Permet d'optimiser

$$\begin{matrix} F & + & G \\ \uparrow & & \uparrow \\ \text{Lisse} & & \text{Non-lisse (norme, indicatrice, \dots)} \end{matrix}$$

II Descente de gradient stochastique :

💡 Remplacer $\nabla F(\theta_{t-1})$ par un estimateur non biaisé $y_t(\theta_{t-1})$

$$(i.e. \mathbb{E}(y_t(\theta_{t-1}) | \mathcal{I}_{t-1}) = \nabla F(\theta_{t-1}))$$

↑
information au
temps $t-1$

Pourquoi? Ces calculs un estimateur peut être beaucoup plus rapide.

Algorithme (Stochastic Gradient Descent)

Entrée: Un point de départ θ_0

Récurrence: $\theta_t = \theta_{t-1} - \gamma_t y_t(\theta_{t-1})$

où $(\gamma_t)_t$ est une suite de tailles de pas.

Exemples (SGD en ML):

• Empirical Risk minimization

$$F(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i; f_\theta(x_i))$$

à chaque itération, un ou plusieurs indices sont choisis au hasard (mini-batch), l'estimateur du gradient est calculé sur ces indices.

• Expected Risk minimization

$$F(\theta) = \mathbb{E}(P(y, \theta(\cdot)))$$

À chaque étape, un couple (x, y) est échantillonné et l'estimateur du gradient est calculé sur ce couple.

Problème: On ne peut théoriquement voir qu'une seule fois chaque point du dataset même si en pratique c'est ce que font les gens.

~~Analyse des cas~~

Hypothèse 1: $\mathbb{E}(g_t(\theta_{t-1}) | \theta_{t-1}) = \nabla F(\theta_{t-1})$
(Estimateur sans biais)

Hypothèse 2: $\|g_t(\theta_{t-1})\|_2^2 \leq B^2$ presque sûrement (v.t.).

(Peut être remplacé par une hypothèse de type Lipschitz sur les gradients).

Théorème: Si F est convexe et B -Lipschitz, admet un minimiseur θ_* et $\|\theta_0 - \theta_*\|_2 \leq D$, et si H1 et H2 sont satisfaites, alors si $\gamma_t = \left(\frac{D}{B}\right) \frac{1}{\sqrt{t}}$, les itérations $(\theta_t)_{t \geq 0}$ de SGD satisfait

$$\mathbb{E}(F(\bar{\theta}_t) - F(\theta_*)) \leq DB \frac{2 + \log(t)}{2\sqrt{t}}$$

où $\bar{\theta}_t = \left(\sum_{i=0}^{t-1} \gamma_i \theta_{i+1} \right) / \left(\sum_{i=0}^{t-1} \gamma_i \right)$

démonstration:

$$\begin{aligned} \mathbb{E}(\|\sigma_t - \theta_*\|_2^2) &= \mathbb{E}(\|\sigma_{t-1} - \gamma_t y_t(\sigma_{t-1}) - \theta_*\|_2^2) \\ &= \mathbb{E}(\|\sigma_{t-1} - \theta_*\|_2^2) - 2\gamma_t \mathbb{E}(\langle y_t(\sigma_{t-1}), \sigma_{t-1} - \theta_* \rangle) \\ &\quad + \gamma_t^2 \mathbb{E}(\|y_t(\sigma_{t-1})\|_2^2) \end{aligned}$$

de plus,

$$\begin{aligned} \mathbb{E}(\langle y_t(\sigma_{t-1}), \sigma_{t-1} - \theta_* \rangle) &= \mathbb{E}(\mathbb{E}(\langle \cdot, \cdot \rangle \mid \sigma_{t-1})) \\ &= \mathbb{E}(\langle \mathbb{E}(y_t(\sigma_{t-1}) \mid \sigma_{t-1}), \sigma_{t-1} - \theta_* \rangle) \\ &= \mathbb{E}(\langle \nabla F(\sigma_{t-1}), \sigma_{t-1} - \theta_* \rangle) \end{aligned}$$

donc

$$\mathbb{E}(\|\sigma_t - \theta_*\|_2^2) \leq \mathbb{E}(\|\sigma_{t-1} - \theta_*\|_2^2) - 2\gamma_t \mathbb{E}(\langle \nabla F(\sigma_{t-1}), \sigma_{t-1} - \theta_* \rangle) + \gamma_t^2 B^2$$

de plus, par convexité, $F(\sigma_{t-1}) - F(\theta_*) \leq \langle \nabla F(\sigma_{t-1}), \sigma_{t-1} - \theta_* \rangle$,

donc

$$\gamma_t \mathbb{E}(F(\sigma_{t-1}) - F(\theta_*)) \leq \frac{1}{2} \left(\mathbb{E}(\|\sigma_{t-1} - \theta_*\|_2^2) - \mathbb{E}(\|\sigma_t - \theta_*\|_2^2) \right) + \frac{1}{2} \gamma_t^2 B^2$$

Donc, en sommant :

$$\frac{1}{\sum_{i=1}^t \gamma_i} \sum_{i=1}^t \gamma_i \mathbb{E}(F(\sigma_{i-1}) - F(\theta_*)) \leq \frac{\|\sigma_0 - \theta_*\|_2^2}{2 \sum_{i=1}^t \gamma_i} + B^2 \frac{\sum_{i=1}^t \gamma_i^2}{\sum_{i=1}^t \gamma_i}$$

Ensuite, le terme de gauche provient de Jensen, et le terme de droite résulte de calculs simples (comparaison séries-intégrales) \square (12)

Remarque: Le cas non-stochastique et convexe découle de cette preuve.

Aller plus loin

• Le cas fortement convexe

Il est possible d'obtenir une vitesse en $\frac{1}{t}$ dans le cas fortement convexe

• La réduction de variance

Garder en mémoire des gradients peut permettre de réduire la variance et de gagner en vitesse de

convergence

• Méthodes adaptatives

En pratique, des méthodes adaptatives (Adagrad, RMS Prop, Adam, ...) sont utilisées pour trouver un bon pré-conditionnement au problème.