

Mapi3 Machine Learning : Optimisation pour le Machine Learning

ont de : En apprentissage statistique, il faut souvent résoudre des problèmes de la forme

$$\hat{\theta} \approx \underset{\theta \in \mathcal{H} = \mathbb{R}^d}{\operatorname{argmin}} F(\theta).$$

Parfois : Pas de solution en forme close

\Rightarrow Il faut utiliser un algorithme d'approximation, un algorithme d'optimisation!

Exemple : Minimisation du risque empirique

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

Décomposition du risque

$$\begin{aligned} R(f_{\hat{\theta}}) - R(f_{\theta^*}) &= \underbrace{R(f_{\hat{\theta}}) - \hat{R}(f_{\hat{\theta}})}_{\text{cf leçons précédentes}} + \underbrace{\hat{R}(f_{\hat{\theta}}) - \hat{R}(f_{\theta^*})}_{\text{Erreur d'optimisation}} + \underbrace{\hat{R}(f_{\theta^*}) - R(f_{\theta^*})}_{\text{cf leçons précédentes}} \\ &\leq 2 \underbrace{\sup_{\theta \in \mathcal{H}} |R(f_{\theta}) - \hat{R}(f_{\theta})|}_{\text{cf leçons précédentes}} + \underbrace{\left(\hat{R}(f_{\hat{\theta}}) - \inf_{\theta \in \mathcal{H}} \hat{R}(f_{\theta}) \right)}_{\text{Erreur d'optimisation}}. \end{aligned}$$

Idée fondamentale de l'optimisation :

Remplacer F par un dével. de Taylor, effectuer un pas en fonction de ce dével.

Exemple :

$$f(x) = \frac{1}{2} x^T A x + b^T x + c, \quad A \text{ symétrique définie positive, } b \in \mathbb{R}^d, c \in \mathbb{R}.$$

• Dynamique d'ordre 1 : $f(\theta + \delta\theta) = f(\theta) + \langle \nabla f(\theta), \delta\theta \rangle + o(\|\delta\theta\|)$

\Rightarrow La direction de plus forte descente locale est $-\nabla f(\theta)$.

\Rightarrow Algorithme de descente de gradient

Entrée : Un point de départ θ_0 .

Récurrence : $\theta_t = \theta_{t-1} - \tau_t \nabla f(\theta_{t-1})$, où $(\tau_t)_{t \geq 1}$ est la suite des "pas" ou "learning rate".

• Dynamique d'ordre 2 : $\nabla f(\theta + \delta\theta) = \underbrace{\nabla f(\theta) + \langle \nabla^2 f(\theta) \delta\theta, \delta\theta \rangle}_{T_2(\theta)} + o(\|\delta\theta\|)$

un point critique satisfait $\nabla f(\theta) = 0$, et en remplaçant ∇f par T_2 , on obtient

$$\delta\theta = -\nabla^2 f(\theta)^{-1} \nabla f(\theta).$$

\Rightarrow Algorithme de Newton

Entrée : Un point de départ θ_0 .

Récurrence : $\theta_t = \theta_{t-1} - \tau_t \nabla^2 f(\theta_{t-1})^{-1} \nabla f(\theta_{t-1})$

où (τ_t) est la suite des pas.

I. Le point de vue des systèmes dynamiques :

• Descente du gradient :

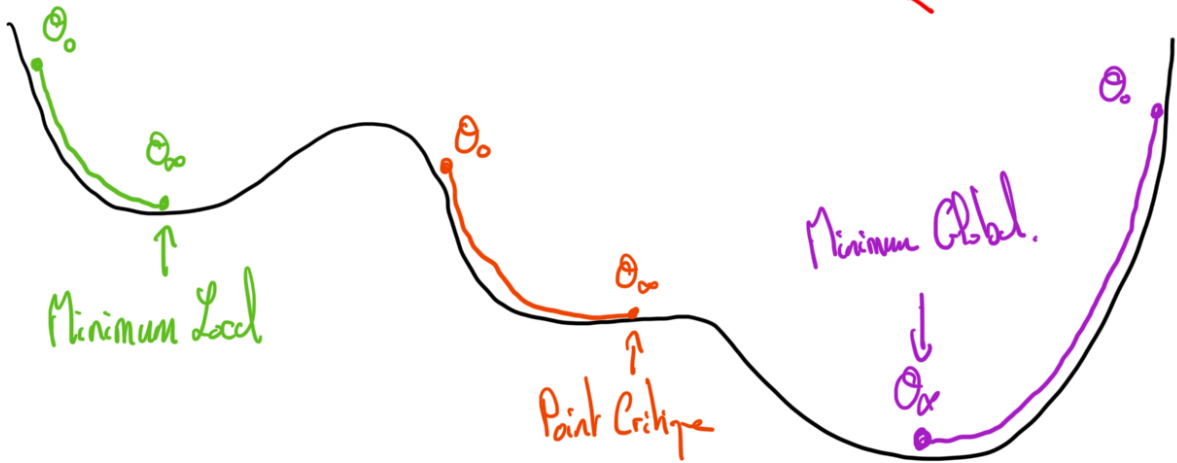
Lorsque $\sup_t \tau_t \ll 1$, on peut voir la descente du gradient comme une discrétisation d'Euler de l'ODE

$$\boxed{\frac{d}{dt} \theta = -\nabla f(\theta)} \quad (\text{Gradient Flow})$$

$$\frac{d}{dt}$$

Alors $\frac{dF(\theta)}{dt} = \underbrace{\langle \nabla F(\theta), \frac{d}{dt} \theta \rangle}_{\text{règle de la chaîne}} = - \underbrace{\langle \nabla F(\theta), \nabla F(\theta) \rangle}_{\text{Gradient Flow}}$

$$= - \|\nabla F(\theta)\|^2 \leq 0$$



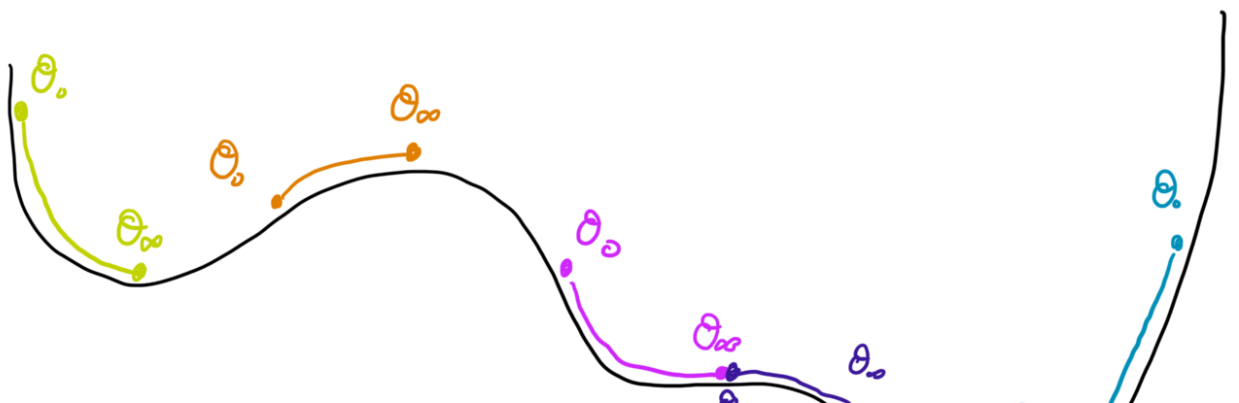
• Méthode de Newton :

$$\frac{d}{dt} \theta = - \nabla^2 F(\theta)^{-1} \nabla F(\theta)$$

Alors $\nabla^2 F(\theta) \frac{d}{dt} \theta = - \nabla F(\theta) \Rightarrow \frac{d}{dt} \nabla F(\theta) = - \nabla F(\theta)$

règle de la chaîne

donc $\nabla F(\theta) = \nabla F(\theta_0) e^{-t}$



II. Convergence du descente de gradient

1) Fonction lisse

Pour analyser le descente de gradient, il faut pouvoir "faire confiance" à l'approximation locale de la fonction par le polynôme de Taylor à l'ordre 1.

Définition: F est L -lisse si

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

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

Si F est C^2 , c'est aussi équivalent à $\|\nabla^2 F(\theta)\|_{op} \leq L \forall \theta$.

Proposition (Lemme de descente): Si F est L -lisse, GD avec $\tau_t = \frac{1}{L}$ donne

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

Preuve: $F(\theta_t) = F(\theta_{t-1} - \frac{1}{L} \nabla F(\theta_{t-1}))$

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

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

2) Le cas non-convexe et lisse

Si F est L -lisse, un télescopage dans le lemme de descente donne

$$\frac{1}{2L} \sum_{i=1}^t \|\nabla F(\theta_{i-1})\|_2^2 \leq \frac{F(\theta_0) - F(\eta_*)}{t}$$

La trajectoire converge en Moyenne de Césaro quadratique vers un point critique.

3) Le cas convexe et L-lisse

Proposition (Co-coercivité) Si F est L-lisse et convexe

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

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

Preuve

$$\underbrace{F(\eta) + \langle \nabla F(\eta), \xi - \eta \rangle}_{(G)} \leq \underbrace{F(\xi)}_{L\text{-lisse}} \leq \underbrace{F(\theta) + \langle \nabla F(\theta), \xi - \theta \rangle}_{(D)} + \frac{L}{2} \|\theta - \xi\|_2^2$$

$$(D) - (G) = F(\theta) - F(\eta) + \langle \nabla F(\theta), \xi - \theta \rangle - \langle \nabla F(\eta), \xi - \eta \rangle + \frac{L}{2} \|\theta - \xi\|_2^2$$

donc $(D) - (G)$ est minimum pour $\xi = \eta$

$$\nabla F(\theta) - \nabla F(\eta) + L(\theta - \xi) = 0 \Rightarrow \xi = \theta + \frac{1}{L} (\nabla F(\theta) - \nabla F(\eta))$$

ce qui donne la dernière inégalité du résultat en réinjectant. \square

Corollaire: Si F est convexe et L-lisse, alors les itérés de GD avec $\tau_t = \frac{1}{L}$ satisfont

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

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

Preuve $\|O_t - \gamma_*\|_2 = \|O_{t-1} - \frac{1}{L} \nabla F(O_{t-1}) - \gamma_*\|_2$

$$= \|O_{t-1} - \gamma_*\|_2^2 - \frac{2}{L} \langle \nabla F(O_{t-1}), O_{t-1} - \gamma_* \rangle + \frac{1}{L^2} \|\nabla F(O_{t-1})\|_2^2$$

$$= \underbrace{\|O_{t-1} - \gamma_*\|_2^2}_{\text{C.o. Coercivité}} - \frac{2}{L} \langle \nabla F(O_{t-1}), O_{t-1} - \gamma_* \rangle + \frac{1}{L^2} \|\nabla F(O_{t-1}) - \underbrace{\nabla F(\gamma_*)}_{=0}\|_2^2$$

$$\leq \underbrace{\|O_{t-1} - \gamma_*\|_2^2}_{\text{C.o. Coercivité}} + \frac{1}{L} \langle \nabla F(O_{t-1}) - \nabla F(\gamma_*), O_{t-1} - \gamma_* \rangle$$

$$= \|O_{t-1} - \gamma_*\|_2^2 - \frac{1}{L} \langle \nabla F(O_{t-1}), O_{t-1} - \gamma_* \rangle \quad \square$$

Théorème : Si F est convexe et L -lisse, alors les itérés $(O_t)_t$ de GD avec $\gamma_t = \frac{1}{L}$ satisfont

$$F(O_t) - F(\gamma_*) \leq \frac{L}{t+1} \|O_0 - \gamma_*\|_2^2$$

Preuve : D'après le lemme de descente, $(F(O_t))_t$ est décroissant.

Donc, $F(O_t) - F(\gamma_*) \leq \frac{1}{t+1} \sum_{i=0}^t (F(O_i) - F(\gamma_*))$

$$\leq \frac{1}{t+1} \sum_{i=0}^t \underbrace{\langle \nabla F(O_i), O_i - \gamma_* \rangle}_{\text{Convexité}}$$

$$\leq \frac{L}{t+1} \sum_{i=0}^t \underbrace{\left(\|O_i - \gamma_*\|_2^2 - \|O_{i+1} - \gamma_*\|_2^2 \right)}_{\text{Corr. préc.}}$$

$$\leq \underbrace{\frac{L}{t+1} \|O_0 - \gamma_*\|_2^2}_{\text{Telescoping}}$$

4) Fonctions fortement convexes

Ⓟ Gradient Flow : $\frac{dF}{dt} = -\|\nabla F\|^2$ donc, si $z_p(F - F(O_*)) \leq \|\nabla F\|^2$,

$$\Rightarrow \frac{d(F - F(O_*))}{dt} \leq -2z_p(F - F(O_*))$$

$$\frac{d}{dt} (F(\theta_t) - F(\theta_*)) \leq -2\mu (F(\theta_t) - F(\theta_*))$$

$$\stackrel{\text{Gronwall}}{\Rightarrow} F - F(\theta_*) \leq e^{-2\mu t} (F(\theta_0) - F(\theta_*)) .$$

Définition: F est μ -fortement convexe si

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

Si F est C^2 , c'est équivalent à $\nabla^2 f(\theta) \succeq \mu I \quad \forall \theta$.

Proposition Si F est différentiable, μ -fortement convexe et que η_* est son minimiseur unique, $\|\nabla F(\theta)\|_2^2 \geq 2\mu (F(\theta) - F(\eta_*))$, $\forall \theta$ (Inégalité de Łojasiewicz).

Preuve: Minimiser en η le terme de droite dans la def de μ -fortement convexe.

Théorème: Si F est L -Lisse et μ -fortement convexe, GD avec $\gamma_t = \frac{1}{L}$ satisfait

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

$\kappa = \frac{L}{\mu}$: conditionnement.

Preuve: Inégalité de Łojasiewicz dans le lemme de descente.

III. Descente de gradient stochastique

⊙ Remplacer $\nabla F(\theta_{t-1})$ par un estimateur non-biaisé $y_t(\theta_{t-1})$
(i.e. $\mathbb{E}(y_t(\theta_{t-1}) | \theta_{t-1}) = \nabla F(\theta_{t-1})$)
↑
information au temps $t-1$.

Pourquoi ? Car $y_t(\theta_{t-1})$ peut être bien plus rapide à calculer que $\nabla F(\theta_{t-1})$.

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

On choisit $B_t \subseteq \{1, \dots, n\}$ au hasard, puis

$$y_t(\theta) = \frac{1}{|B_t|} \sum_{i \in B_t} \ell(y_i, f_{\theta}(x_i)).$$

Complexité : $O(n) \longrightarrow O(|B_t|)$.

Algorithme (Stochastic Gradient Descent)

Entrée : Un point de départ θ_0 .

Mécanisme : $\theta_t = \theta_{t-1} - \sigma_t y_t(\theta_{t-1})$. (σ_t) : suite de "pas"

H1 : $\mathbb{E}(y_t(\theta_{t-1}) | \theta_{t-1}) = \nabla F(\theta_{t-1})$ (Estimateur sans biais)

H2 : $\|y_t(\theta_{t-1})\|_2^2 \leq B^2$ p.s. (Variance contrôlée).

1) Le cas général lisse

Lemme (de descente) Si F est L -lisse et sous H_1 :

$$\mathbb{E}(F(\theta_t) | \theta_{t-1}) \leq F(\theta_{t-1}) - \sigma_t \|\nabla F(\theta_{t-1})\|_2^2 + \frac{L}{2} \sigma_t^2 \mathbb{E}(\|y_t(\theta_{t-1})\|_2^2 | \theta_{t-1})$$

Preuve : Par L-lissitude

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

et l'espérance conditionnelle donne le résultat. \square

Donc sous H_2 , une somme télescopique donne ($\sigma_t = \sigma$)

$$\frac{1}{t} \sum_{i=0}^{t-1} \mathbb{E}(\|\nabla F(\theta_i)\|_2^2) \leq \frac{F(\theta_0) - F(\theta_*)}{\sigma t} + \frac{L}{2} \sigma B^2$$

et $\sigma = \frac{1}{\sqrt{t}}$ donne la vitesse finale.

2) Le cas convexe

Théorème : Si F est convexe et B -Lipschitz, admet un minimiseur θ_* tel que $\|\theta_0 - \theta_*\|_2 \leq D$ et si H_1 et H_2 sont satisfaites, si $\sigma_t = \left(\frac{D}{B}\right) \frac{1}{\sqrt{t}}$, les itérées de SGD satisfont

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

$$\text{où } \bar{\theta}_t = \left(\sum_{i=1}^t \sigma_i \theta_{i-1} \right) / \left(\sum_{i=1}^t \sigma_i \right).$$

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

de plus,
$$\mathbb{E}(\langle y_t(\theta_{t-1}), \theta_{t-1} - \theta_* \rangle) = \mathbb{E}(\mathbb{E}(\langle y_t(\theta_{t-1}), \theta_{t-1} - \theta_* \rangle | \theta_{t-1}))$$

$$\begin{aligned}
&= \mathbb{E} \left(\langle \mathbb{E}(g(\theta_{t-1}) | \theta_{t-1}), \theta_{t-1} - \theta_* \rangle \right) \\
&= \mathbb{E} \left(\langle \nabla f(\theta_{t-1}), \theta_{t-1} - \theta_* \rangle \right) \geq \langle \nabla f(\theta_{t-1}), \theta_{t-1} - \theta_* \rangle
\end{aligned}$$

$$\text{donc } \mathbb{E}(\|\theta_t - \theta_*\|_2^2) \leq \mathbb{E}(\|\theta_{t-1} - \theta_*\|_2^2) - 2\sigma_t \mathbb{E}(\langle \nabla f(\theta_{t-1}), \theta_{t-1} - \theta_* \rangle) + \sigma_t^2 B^2$$

$$\text{donc : } \sigma_t \mathbb{E}(f(\theta_{t-1}) - f(\theta_*)) \leq \frac{1}{2} \left(\mathbb{E}(\|\theta_{t-1} - \theta_*\|_2^2) - \mathbb{E}(\|\theta_t - \theta_*\|_2^2) \right) + \frac{1}{2} \sigma_t^2 B^2$$

$$\Rightarrow \frac{1}{\sum_{i=1}^t \sigma_i} \sum_{i=1}^t \sigma_i \mathbb{E}(f(\theta_{i-1}) - f(\theta_*)) \leq \frac{\|\theta_0 - \theta_*\|_2^2}{2 \sum_{i=1}^t \sigma_i} + B^2 \frac{\sum_{i=1}^t \sigma_i^2}{\sum_{i=1}^t \sigma_i}$$

Puis on peut minorer le terme de Gauche par Jensen et majorer le terme de droite par comparaison sans intégrale.

□