

Optimisation pour l'apprentissage automatique

11 décembre 2024

Séance 2/8 : Théorie descente de gradient

- * Régression linéaire
- * Cas général : vitesses de convergence / complexités
- * Un résultat important dans le cas non convexe

Séance 3/8 (demain)

- * TP descente de gradient (cas convexe)
- * Exercices descente de gradient

Théorie de la descente de gradient

Cadre général

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

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ C^1

Algorithmes de descente de gradient

$$\forall k \in \mathbb{N}, \quad w_{k+1} = w_k - \alpha_k \nabla f(w_k) \quad \text{avec} \quad \alpha_k > 0$$

Deux types de résultats théoriques

- Convergence asymptotique: garanties à la limite, quand $k \rightarrow \infty$
Ex) $f(w_k) \xrightarrow{k \rightarrow \infty} \min_{w \in \mathbb{R}^d} f(w)$
↑
Valeur optimale du problème

- Convergence non-asymptotique / vitesse de convergence / complexité:
Garanties en "temps" fini
Ex) Après K itérations, que peut-on dire de $f(w_k)$? De $f(w_k) - \min_{w \in \mathbb{R}^d} f(w)$?

Cadre le plus simple: Descente de gradient pour la régression linéaire

Problème:

$$\text{minimiser}_{w \in \mathbb{R}^d} \frac{1}{2m} \|Xw - y\|^2 = \frac{1}{2m} \sum_{i=1}^m (x_i^T w - y_i)^2$$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix} \in \mathbb{R}^{m \times d}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m$$

Soit $f(w) = \frac{1}{m} \|Xw - y\|^2$ et supposons que $\underbrace{\frac{X^T X}{d \times m}}_{d \times d}$ est inversible

Sous cette hypothèse, le problème a une unique solution donnée par $w^* = (X^T X)^{-1} X^T y$.

En effet, $\nabla f(w) = 0 \iff \frac{1}{m} X^T (Xw - y) = 0$

$$\iff X^T X w - X^T y = 0$$

$$\iff X^T X w = X^T y$$

$X^T X$ inversible \rightarrow $w = (X^T X)^{-1} X^T y = w^*$

Donc w^* est l'unique point en lequel le gradient est nul.

On f est une fonction convexe, donc $\nabla f(w^*) = 0 \iff w^* \text{ est } \underset{w \in \mathbb{R}^d}{\text{argmin}} f(w)$

On applique la descente de gradient au problème avec une taille de pas constante $\alpha_k = \frac{m}{\|X^T X\|} > 0$

$$\|X^T X\| = \max_{\|v\| \neq 0} \frac{\|X^T X v\|}{\|v\|}$$

Pour tout $k \in \mathbb{N}$, on a alors

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

$$= w_k - \frac{m}{\|X^T X\|} \times \frac{1}{m} X^T (Xw_k - y)$$

$$= w_k - \frac{1}{\|X^T X\|} (X^T X w_k - X^T y)$$

$$w_{k+1} = \left(I_d - \frac{X^T X}{\|X^T X\|} \right) w_k + \frac{1}{\|X^T X\|} X^T y$$

λ valeur propre de $X^T X$
 $\iff \exists v \neq 0, X^T X v = \lambda v$

$$I_d = \begin{bmatrix} 1 & & \\ & \ddots & \\ 0 & & 1 \end{bmatrix} \in \mathbb{R}^{d \times d}$$

$$I_d w_k = w_k$$

$$\begin{aligned}
 \omega_{k+1} - \omega^* &= \left(\text{Id} - \frac{X^T X}{\|X^T X\|} \right) \omega_k + \frac{1}{\|X^T X\|} X^T y - \omega^* \\
 &= \left(\text{Id} - \frac{X^T X}{\|X^T X\|} \right) \omega_k + \left(\frac{X^T X}{\|X^T X\|} - \text{Id} \right) \omega^* \\
 &= \left(\text{Id} - \frac{X^T X}{\|X^T X\|} \right) (\omega_k - \omega^*)
 \end{aligned}$$

Pour montrer que l'algorithme converge, il faut montrer que $\omega_k \rightarrow \omega^*$, c'est-à-dire $\|\omega_k - \omega^*\| \rightarrow 0$

D'après ce qui précéde, on a

$$\|\omega_{k+1} - \omega^*\| = \left\| \left(\text{Id} - \frac{X^T X}{\|X^T X\|} \right) (\omega_k - \omega^*) \right\|$$

$$\begin{aligned}
 &\leq \left\| \text{Id} - \frac{X^T X}{\|X^T X\|} \right\| \|\omega_k - \omega^*\| \\
 &= \underbrace{\left(1 - \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)} \right)}_{\in (0, 1)} \|\omega_k - \omega^*\|
 \end{aligned}$$

plus petite valeur propre de $X^T X$
(> 0 car $X^T X$ inversible)

$$\|\omega_{k+1} - \omega^*\| < \|\omega_k - \omega^*\|$$

A chaque itération, l'algorithme se rapproche de la solution.

$\forall K \in \mathbb{N}$,

$$\|\omega_K - \omega^*\| \leq \left(1 - \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)} \right)^K \|\omega_0 - \omega^*\|$$

(application de $\|\omega_{k+1} - \omega^*\| \leq \left(1 - \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)} \right) \|\omega_k - \omega^*\|$ pour $k=0 \dots K-1$)

Ce résultat s'appelle une vitesse de convergence: on dit que la descente de gradient converge linéairement (ou en vitesse linéaire) sur ce problème

$$\text{Linéaire: } \ln \left(\left(1 - \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)} \right)^K \|w_0 - w^*\| \right) = K \ln \left(1 - \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)} \right) + \ln (\|w_0 - w^*\|)$$

Fonction linéaire/affine de K

→ Pour un nombre fixé d'itérations K , on a des garanties sur

$$0 < \|w_K - w^*\| \quad \left(\leq \left(1 - \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)} \right)^K \|w_0 - w^*\| \right)$$

→ Cela implique que $\|w_K - w^*\| \xrightarrow{K \rightarrow \infty} 0$ car $\left(1 - \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)} \right)^K \rightarrow 0$

→ Cette vitesse de convergence se traduit aussi en (borne de) complexité:

Pour tout $\epsilon > 0$, on peut garantir $\|w_K - w^*\| \leq \epsilon$

· en un nombre d'itérations de l'ordre de $\lceil \ln \epsilon \rceil$

* Pour avoir $\|w_K - w^*\| \leq \epsilon$, il suffit que

$$\left(1 - \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)} \right)^K \|w_0 - w^*\| \leq \epsilon$$

$$\underbrace{K \ln \left(1 - \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)} \right) + \ln (\|w_0 - w^*\|)}_{< 0} \leq \ln \epsilon$$

$$K \geq \frac{\ln \epsilon}{\ln \left(1 - \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)} \right)} - \frac{\ln (\|w_0 - w^*\|)}{\ln \left(1 - \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)} \right)}$$

$$= \mathcal{O}(\ln \varepsilon)$$

Choix de ε

- Typiquement dans $(0, 1)$
- Peut aussi être "normalisé" par rapport au problème

$$\varepsilon = \gamma \|(\bar{w}_0 - w^*)\|$$

$$\gamma f(0, 1)$$

$\mathcal{O}(A)$: une constante $\propto A$, et la constante ne dépend pas des quantités qui apparaissent dans A

Analyse pour des classes plus générales de fonctions

↳ les résultats ci-dessus pour la régression linéaire s'étendent aux fonctions $C^{1,1}$ et fortement convexes

Def: . $f: \mathbb{R}^d \rightarrow \mathbb{R}$ est $C_L^{1,1}$ si elle est C^1 et ∇f est L -lipschitzien avec $L > 0$, c'est à dire $\forall (v, w) \in (\mathbb{R}^d)^2, \|\nabla f(v) - \nabla f(w)\| \leq L \|v - w\|$

On dit parfois que f est L -lisse ("L-smooth")

Si f est $C_L^{1,1}$, alors pour tous $(v, w) \in (\mathbb{R}^d)^2$,

$$(1) \quad f(v) \leq f(w) + \nabla f(w)^T(v - w) + \frac{L}{2} \|v - w\|^2$$

(Version exacte de $f(v) \approx f(w) + \nabla f(w)^T(v - w)$)

Conséquence: Si $w = w_k$ et $v = w_k - \alpha_k \nabla f(w_k)$

$$\cdot f(w_k - \alpha_k \nabla f(w_k)) < f(w_k) \text{ pour } \alpha_k < \frac{2}{L}$$

$$\cdot f(w_k - \frac{1}{L} \nabla f(w_k)) \leq f(w_k) - \frac{1}{2L} \|\nabla f(w_k)\|^2, \text{ et}$$

c'est la plus forte décroissance pour $\alpha_k \in (0, \frac{2}{L})$ \Rightarrow Justifie de choisir $\alpha_k = \frac{1}{L}$ + tk

Ex) Régression linéaire

$$f(w) = \frac{1}{2m} \|Xw - y\|^2 \quad \nabla f(w) = \frac{1}{m} X^T(Xw - y)$$

$$f \in C_L^{1,1} \quad \text{avec} \quad L = \frac{\|X^T X\|}{m}$$

Déf. $f: \mathbb{R}^d \rightarrow \mathbb{R}$ (¹ est μ -fondament convexe si $(\mu > 0)$)

$$(2) \forall (v, w) \in (\mathbb{R}^d)^2, \quad f(v) \geq f(w) + \nabla f(w)^T(v-w) + \frac{\mu}{2} \|v-w\|^2$$

(2) fournit un minorant de $f(v)$

\rightarrow Toute fonction fondament convexe est convexe

$$[f \text{ } C^1 \text{ convexe}] \Leftrightarrow \forall (v, w) \in (\mathbb{R}^d)^2, f(v) \geq f(w) + \nabla f(w)^T(v-w)$$

mais l'inverse n'est pas vrai (contre ex: fonction constante)

\rightarrow Toute fonction fondament convexe f possède un unique minimum, qui est l'unique vecteur w tel que $\nabla f(w) = 0$

↑
minimum
réel
contient un seul vecteur

$$\text{Ex) Régression linéaire : } f(w) = \frac{1}{2m} \|Xw - y\|^2$$

Si $X^T X$ est inversible, f est μ -fortement convexe avec $\mu = \frac{\lambda_{\min}(X^T X)}{m}$

Si non, f est simplement convexe et le problème a une infinité de solutions.

Vitesses de convergence / Complexité dans le cas fortement convexe

Soit $f \in \mathcal{C}_L^1$ et μ -fortement convexe.

On applique la descente de gradient avec $\alpha_n = \frac{1}{L}$.

Alors, si w^* est le minimum de f , on a :

$$\text{. } \forall K \in \mathbb{N}, \|w_k - w^*\| \leq \left(1 - \frac{\mu}{L}\right)^K \|w_0 - w^*\| \quad \text{avec } C > 0$$

\Rightarrow vitesse de convergence linéaire vers la solution

$$\text{et } f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq \left(1 - \frac{\mu}{L}\right)^K \left(f(w_0) - \min_{w \in \mathbb{R}^d} f(w)\right)$$

\Rightarrow vitesse de convergence linéaire vers la valeur optimale

$$\text{avec } t = 1 - \frac{\mu}{L} \in (0, 1)$$

$$\text{. Pour tout } \varepsilon > 0, \text{ on a } f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq \varepsilon$$

après au plus $O(\ln \varepsilon)$ itérations

$$(\text{idem pour } \|w_k - w^*\| \leq \varepsilon)$$

\rightarrow Se généralise au choix d'un α_n constant égal à $\alpha \in (0, \frac{2}{L})$

Vitesse de convergence de la descente de gradient dans le cas fortement convexe
linéaire en $\mathcal{O}(t^K)$

Comparaison avec le cas non fortement convexe

$f \in C_L^{1,1}$, f admet un minimum ($\arg\min_{w \in \mathbb{R}^d} f(w) + \phi$) et on note
 $f^* = \min_{w \in \mathbb{R}^d} f(w)$

Descr de gradient avec $\alpha_k = \frac{1}{L}$

f μ -fortement convexe	$f(w_k) - f^* \leq O((1 - \frac{\mu}{L})^k)$	$f(w_k) - f^* \leq \varepsilon$ en $O(\frac{1}{\mu\varepsilon})$ itérations
f convexe	$f(w_k) - f^* \leq O(\frac{1}{k})$	$f(w_k) - f^* \leq \varepsilon$ en $O(\frac{1}{\varepsilon})$ itérations
f non convexe	$\min_{0 \leq k \leq K-1} \ \nabla f(w_k)\ \leq O(\frac{1}{\sqrt{K}})$	$\ \nabla f(w_k)\ \leq \varepsilon$ en $O(\frac{1}{\varepsilon^2})$ itérations

Analys de ces résultats

- Le cas convexe est plus difficile que le cas fortement convexe, dans le sens où la descr de gradient converge moins vite dans le cas convexe (en v. v. sous-linéaire $\frac{1}{k}$)

Illustration: $\frac{\mu}{L} = \frac{1}{2}$, $K = 100$ $\frac{1}{K} = 0.01 = 10^{-2}$
 $(1 - \frac{\mu}{L})^K = \frac{1}{2^{100}} \approx 10^{-30}$

\Rightarrow Dans les deux cas cependant, on montre que $f(w_k) - f^* \rightarrow 0$

\Rightarrow Dans le cas fortement convexe, comme le minimum est unique, on peut en plus garantir que w_k converge vers ce minimum

- Le cas non convexe est plus difficile que le cas convexe, et donne des garanties plus faibles:

- En termes de vitesse de convergence

$\frac{1}{\sqrt{k}}$ tend vers 0 moins vite que $\frac{1}{k}$

$$K = 100 \quad \frac{1}{k} = 10^{-2} \quad \frac{1}{\sqrt{k}} = 10^{-1}$$

- En termes de quantité à laquelle s'applique cette vitesse

$f(w_k) - f^*$ pour f convexe/fortement convexe

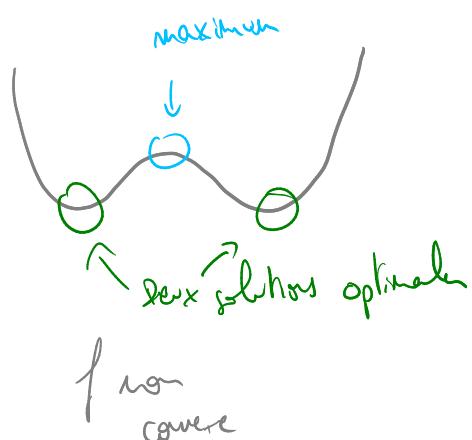
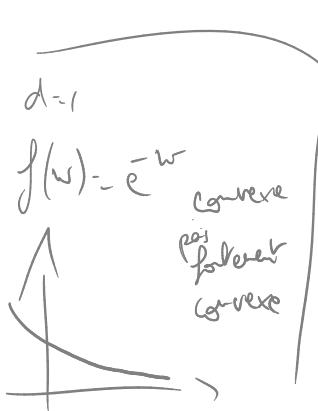
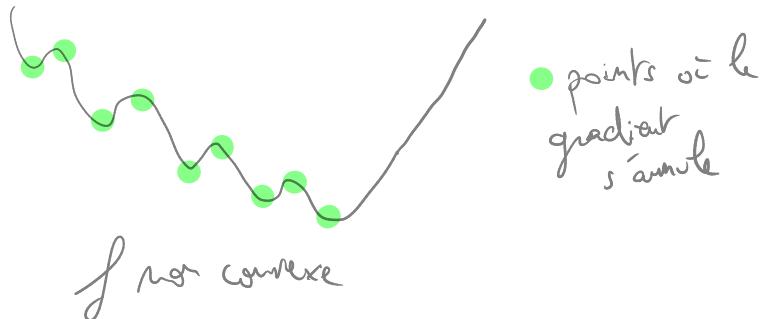
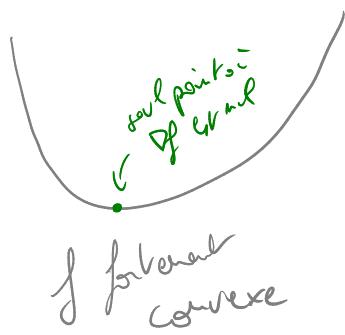
$$\min_{0 \leq k \leq K-1} \|\nabla f(w_k)\| \text{ pour } f \text{ non convexe}$$

\Rightarrow Cette différence s'explique par le fait que $\|\nabla f(\bar{w})\| = 0 \Leftrightarrow \bar{w}$ tangente $f(w)$ lorsque f

est convexe, alors que on a seulement

$$\left[\begin{array}{c} \bar{w} \text{ tangente } f(w) \\ w_{\text{final}} \end{array} \right] \Rightarrow \|\nabla f(\bar{w})\| \neq 0$$

pour f non convexe en général.

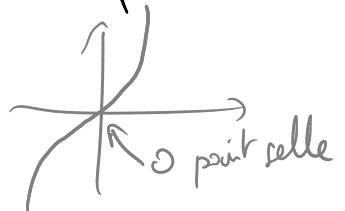


$$f\left(\begin{bmatrix} w_1 \\ w_2 \end{bmatrix}\right) = w_1^2$$

convexe
pas fortement convexe

- La descente de gradient converge vers un point où le gradient est nul.
- Si la fonction est non convexe, alors un tel point peut être :
- . Un minimum global (donc une solution de minimiser $f(w)$)
 - . Un minimum local non global ($f(\bar{w}) < f(w)$ et w tel que $\|w - \bar{w}\|$ est faible)
 - . Un maximum local/global
 - . Un point selle, c'est un point qui est un maximum local dans certaines directions et un minimum local pour d'autres

$$d=1, f(w) = w^3$$



- On peut trouver des exemples de fonctions non convexes pour lesquelles la descente de gradient converge vers un maximum ou un point selle
- Mais en pratique, et notamment sur des formulations non convexes d'apprentissage (complétiion de matrice, réseaux de neurones peu profonds, apprentissage de dictionnaire, ...), on observe que la descente de gradient converge vers un minimum local voire global.
- (cf exercices)

Théorème (informel, 2015)

Si on applique la descente de gradient avec $w_0 \in \mathbb{R}^d$ choisi aléatoirement dans \mathbb{R}^d , alors on converge presque sûrement vers un minimum local.

$$\mathbb{P}_{w_0} \left(\text{on converge vers un maximum ou un point selle} \right) = 0$$