

Optimisation pour l'apprentissage automatique

09/12/2024

Liens utiles

clement.royer@lamsade.dauphine.fr

<https://www.lamsade.dauphine.fr/~croyer/coursOAA.html>

Aujourd'hui

Introduction à l'optimisation pour l'apprentissage

Convexité

Principes de base de descente de gradient

Logistique

- 3 semaines de cours
 - Cette semaine : 3 séances, 2 CTD + 1 pratique
 - SGD
 - La semaine prochaine : 2 séances, 1 cours + 1 pratique
 - Adam, Adagrad, etc.
 - Dernière semaine (22/01-23/01) : 2 CTD + 1 pratique
 - AutoDiff, régularisator, L-BFGS, ...
- Intervenants : Clément ROYER (+ Paul CAILLON ?)
- Notation :
 - Examen : 30/01, 10h-12h, 60% de la note
 - Projet : Deadline ~ Fin février, prolongement des TP, 40% de la note
- Emploi de temps :
 - 9h-12h15

1) Introduction

→ But du cours: Étudier les algorithmes d'optimisation utilisés en apprentissage (ML, Deep Learning, etc)

→ Problème type: Minimisation du risque empirique (ERM)

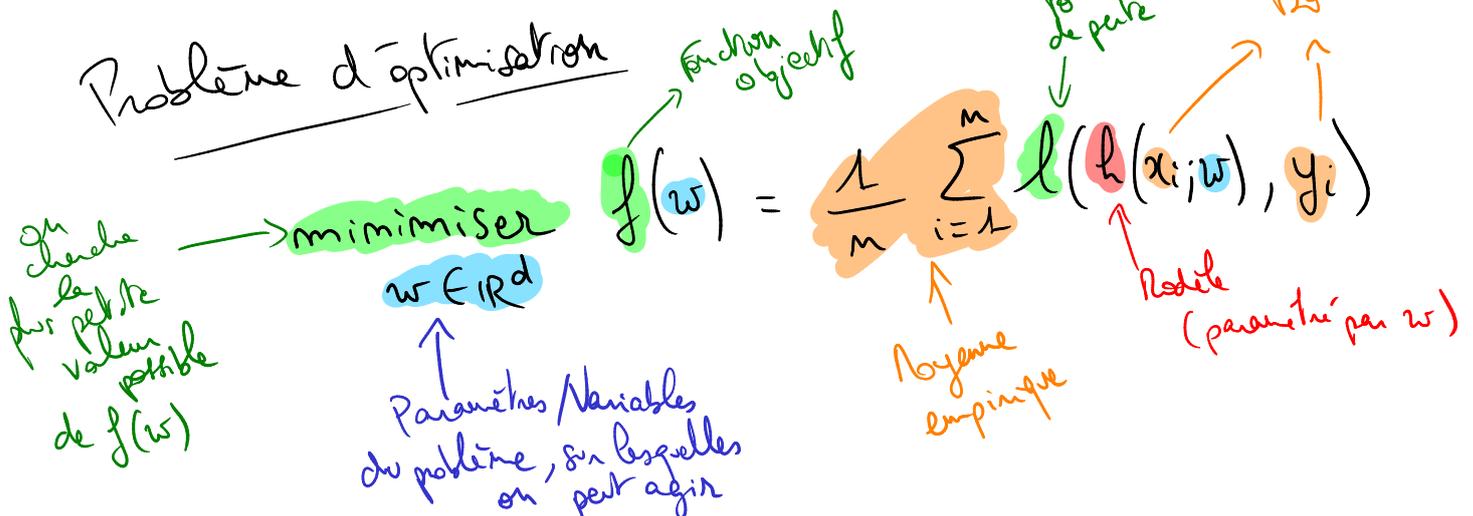
Données: $\{(x_i, y_i)\}_{i=1..n}$ x_i : entrées (vecteurs, matrices, ...)
 y_i : sorties (vecteurs, réels, classes)

Objectif: Trouver un modèle qui colle aux données, c'est-à-dire qui associe chaque x_i au y_i correspondant.

⇒ Dans ce cours, un modèle sera une fonction $h(\cdot; w) : x \mapsto h(x; w)$ paramétrée par un vecteur de paramètres $w \in \mathbb{R}^d$

⇒ On cherche le **meilleur modèle** pour nos données et relativement à une certaine fonction de perte (loss) → Optimisation

Problème d'optimisation



Exemples

① Régression linéaire

$$x_i \in \mathbb{R}^d, y_i \in \mathbb{R}, h(x, w) = x^T w = \sum_{j=1}^d [x]_j [w]_j$$

$$l(h, y) = \frac{1}{2} (h - y)^2 \quad (\text{perte aux moindres carrés} / l_2)$$

Problème :

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

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

$$\text{et } \|v\|^2 = \sum_{i=1}^m [v]_i^2 \quad \forall v \in \mathbb{R}^m$$

② SVM Linéaires (SVM = Support Vector Machines)

$$\text{minimiser}_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \max(1 - y_i x_i^T w, 0)$$

$$\text{avec } x_i \in \mathbb{R}^d, h(x, w) = x^T w, y_i \in \{-1, 1\}$$

$$l(h, y) = \max(1 - yh, 0) \quad \text{"Hinge Loss"}$$

$$\text{Si } h \gg 1 \text{ et } y = +1, l(h, y) = 0$$

$$\text{Si } h \ll -1 \text{ et } y = -1, l(h, y) = 0$$

$$\text{Sinon } l(h, y) > 0$$

③ Régression logistique

$$h(\cdot; w) : x \mapsto x^T w$$

$$x_i \in \mathbb{R}^d, \quad h(x; w) = x^T w, \quad y_i \in \{-1, 1\}$$

$$l(h, y) = \log(1 + \exp(-y h))$$

Problème : minimiser $w \in \mathbb{R}^d$ $\frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w))$

Ideé : $x_i^T w \gg 1$ et $y_i = +1$
 $\log(1 + \exp(-y_i x_i^T w)) \approx 0$

④ "0/1 loss"

$$y_i \in \mathcal{Y} \quad \forall i=1..n$$

$$h(\cdot; w) : x \mapsto h(x; w) \in \mathcal{Y}$$

$$l(h, y) = \mathbb{1}(h \neq y) = \begin{cases} 1 & \text{si } h \neq y \\ 0 & \text{si } h = y \end{cases}$$

$$\text{minimiser } w \in \mathbb{R}^d \quad \frac{1}{n} \sum_{i=1}^n \mathbb{1}(h(x_i; w) \neq y_i)$$

↳ Pour une même tâche d'apprentissage, on peut définir une infinité de problèmes d'optimisation basés sur des pertes et des modèles

⇒ En pratique, on évite de choisir des fonctions discontinues (comme la "0/1 loss") et on préfère des fonctions continues, pour lesquelles on sait quantifier comment la fonction varie quand ses paramètres varient

⇒ On considère dans la grande majorité des cas des pertes

convexes

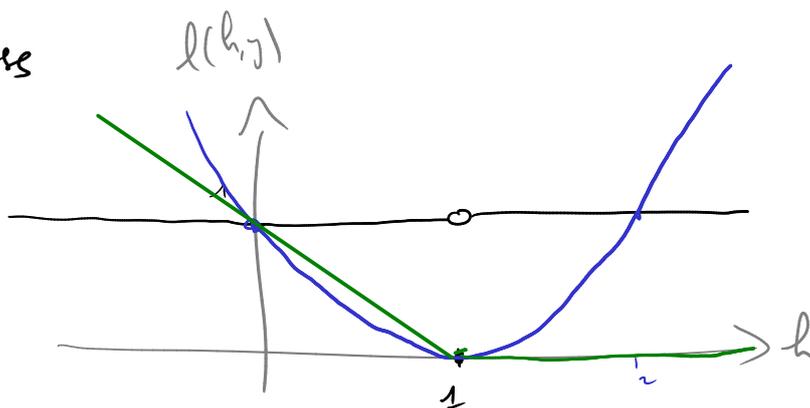
Définition: $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}$ est convexe si:

$$\forall (v, w) \in (\mathbb{R}^d)^2, \forall \alpha \in [0, 1],$$

$$\varphi(\alpha v + (1-\alpha)w) \leq \alpha \varphi(v) + (1-\alpha)\varphi(w)$$

Ex) Perte l_2 convexe, Hinge loss, perte logistique

Contre-ex) 0/1 loss



$y=1$
 — 0/1 loss
 — l_2 loss
 $(h-y)^2$
 — hinge loss
 $\max(1-yh, 0)$

Caractériser les fonctions convexes (tiré de "Patterns, predictions and actions", Romya Herdt Benjamin Recht 2022)

"Plus de 50% des fonctions convexes en ML sont obtenues à partir des axiomes suivants"

1. Toutes les normes sont convexes ($\|v\|^2 = \sum_{j=1}^n |v_j|^2$ est une norme, mais n'est pas la seule)
2. Si φ est convexe et que $\lambda \geq 0$, alors $\lambda\varphi$ est convexe
3. Si φ et ψ sont convexes, alors $\varphi + \psi$ est convexe
4. Si φ et ψ sont convexes, alors $\max(\varphi, \psi)$ est convexe
5. Si φ est convexe, alors pour toute matrice A et tout vecteur b de tailles appropriés, $w \mapsto \varphi(Aw + b)$ est convexe

Exercice: En utilisant les axiomes, justifier que la fonction de l'exemple (2) (SVM) est convexe

$$w \mapsto \frac{1}{n} \sum_{i=1}^n \max(1 - y_i x_i^T w, 0)$$

• $w \mapsto 0$ est convexe $0 \leq 0 + 0$

• $t \mapsto t$ est convexe $\frac{\alpha t + (1-\alpha)u}{\varphi(\alpha t + (1-\alpha)u)} \leq \frac{\alpha t}{\varphi(t)} + (1-\alpha) \frac{u}{\varphi(u)}$

• Par l'axiome 5, $w \mapsto 1 - y_i x_i^T w$ est convexe $\forall i = 1..n$

$$w \mapsto \varphi\left(\underbrace{(-y_i x_i^T)}_A w + \underbrace{1}_b\right) \text{ avec } \varphi: t \mapsto t$$

• Par l'axiome 6, $w \mapsto \max(1 - y_i x_i^T w, 0)$ est convexe $\forall i = 1..n$

• Par l'axiome 3, $w \mapsto \sum_{i=1}^n \max(1 - y_i x_i^T w, 0)$ est convexe

• Par l'axiome 2, $w \mapsto \frac{1}{n} \sum_{i=1}^n \max(1 - y_i x_i^T w, 0)$ est convexe ≥ 0

2) Descente de gradient

2-a) Descente de gradient dans le cas convexe

Problème minimiser $f(w)$ $f: \mathbb{R}^d \rightarrow \mathbb{R}$ de classe C^1
 $w \in \mathbb{R}^d$

$f \in C^1 \Rightarrow$ En tout point $w \in \mathbb{R}^d$, il existe un vecteur

appelé le gradient de f en w (et noté $\nabla f(w) \in \mathbb{R}^d$) qui représente la dérivée de f en w et donc la variation de f au voisinage de w

$$\forall v \in \mathbb{R}^d, \quad f(v) \approx f(w) + \nabla f(w)^T (v-w)$$

lorsque $\|v-w\|$ est faible

→ Avec le gradient, on approxime la fonction f par une fonction linéaire

Théorème: Si f est C^1 et convexe, alors

$$\left[\underbrace{w^* \in \operatorname{argmin}_{w \in \mathbb{R}^d} f(w)}_{\substack{\text{"argmin"} \\ \text{Ensemble des} \\ \text{solutions/des minima} \\ \text{du problème} \\ \text{minimiser } f(w) \\ w \in \mathbb{R}^d}} \right] \Leftrightarrow \left[\underbrace{\nabla f(w^*) = 0_{\mathbb{R}^d}}_{\substack{\text{Condition qui est} \\ \text{calculable/vérifiable en} \\ \text{temps fini}}} \right]$$

w^* est une solution: $f(w^*) \leq f(w) \quad \forall w \in \mathbb{R}^d$] condition qui n'est pas vérifiable en temps fini

Principe de la descente de gradient (basé sur le théorème)

Étant donné un point $w \in \mathbb{R}^d$: or f C^1 convexe,

1) Soit $\nabla f(w) = 0$ et dans ce cas w est une solution

2) Soit $\|\nabla f(w)\| \neq 0$ ($\Leftrightarrow \nabla f(w) \neq 0_{\mathbb{R}^d}$) et dans ce cas

pour un $\alpha > 0$ suffisamment petit

à noter v ,
 $v=0 \Leftrightarrow \|v\|=0$

$$f(w - \alpha \nabla f(w)) \approx f(w) + \nabla f(w)^T (w - \alpha \nabla f(w) - w)$$

$$\Leftrightarrow f(w - \alpha \nabla f(w)) \approx f(w) - \alpha \underbrace{\nabla f(w)^T \nabla f(w)}_{= \|\nabla f(w)\|^2}$$

$$\Leftrightarrow f(w - \alpha \nabla f(w)) \approx f(w) - \underbrace{\alpha \|\nabla f(w)\|^2}_{> 0} < f(w)$$

Algorithme de descente de gradient

Initialisation: Choisir $w_0 \in \mathbb{R}^d$.

Boucle:

Pour $k = 0, 1, 2, \dots$

- Calculer $\nabla f(w_k)$. Si $\|\nabla f(w_k)\| = 0$, terminer.
- Poser $w_{k+1} = w_k - \alpha_k \nabla f(w_k)$ avec $\alpha_k > 0$.

→ En pratique:

- On utilise un critère d'arrêt pour l'algorithme, qui est souvent une combinaison d'un critère de convergence et d'un critère de budget

- Critère de convergence: $\|\nabla f(w_k)\| = 0$, $\|\nabla f(w_k)\| \leq \epsilon$ avec $\epsilon > 0$

$$\|w_{k+1} - w_k\| = \alpha_k \|\nabla f(w_k)\| \leq \epsilon \text{ avec } \epsilon > 0$$

$$f(w_{k+1}) - f(w_k) \leq \epsilon, \epsilon > 0$$

Requiert de calculer ces valeurs de fonction, ce que ne fait pas l'algorithme de base

- Critère de budget: $k = \underbrace{K}_{\text{nombre max d'itérations}}$, temps CPU/GPU = temps_max, nombre maximum d'appels à f (dépend de l'implémentation) ou à ∇f ,

⇒ La théorie de la descente de gradient permet de lier

- les 2 critères :
- Si on souhaite $\|\nabla f(w_k)\| \leq \epsilon$, on peut borner le nombre d'itérations nécessaires à la descente de gradient
 - Si on fait tourner la descente de gradient pendant k itérations, on peut garantir une précision sur $\|\nabla f(w_k)\|$
 - S'applique à la majorité des problèmes (critères de convergence / de budget)
- complexité (au pire cas) []
vitesses de convergence []

→ Choix de α_k (α_k : taille de pas, toujours > 0)

- Taille de pas constante: $\alpha_k = \alpha > 0 \forall k \in \mathbb{N}$
 - ⊕ Il existe des valeurs constantes pour lesquelles la convergence est garantie sur les formulations classiques (régression linéaire / régression logistique)
 - ⊖ Si $\alpha > 0$ trop grande, la méthode diverge ($\|w_k\| \rightarrow \infty$)
 - ⊖ Si $\alpha > 0$ trop petite, numériquement la méthode stagne
- Taille de pas décroissante (fixée a priori): $\alpha_k \xrightarrow[k \rightarrow \infty]{} 0$

$$(\text{Ex}) \alpha_k = \frac{\alpha_0}{k+1} \text{ avec } \alpha_0 > 0, \alpha_k = \frac{\alpha_0}{(k+1)^t} \quad t > 0$$

- ⊕ Garantit que α_k est une bonne taille de pas pour k suffisamment grand (en particulier, $f(w_{k+1}) < f(w_k)$ pour k suffisamment grand)

⊖ Si α_k décroît trop vite, numériquement la méthode va stagner

• Hybride : ("Scheduling")

• Démarrer avec un pas constant $\alpha > 0$

• Soit au bout de K itérations, soit lorsqu'on stagne (ex: $f(v_{k+1}) \approx f(v_k)$)

$\alpha \leftarrow \alpha/10$ (plus généralement $\alpha \leftarrow \theta \alpha$
 $\theta \in (0, 1)$)

Partie 1 des TP (et du projet) : Tester des tailles de pas