

OPTIMISATION POUR L'APPRENTISSAGE

23 janvier 2025

Ce matin : Régularisation (Core, illustration numérique)
Exercice du TD 3 en ligne

RÉGULARISATION

Problème régularisé

mimimiser
 $w \in \mathbb{R}^d$

Terme
 d'attache aux
 données

$$(typiquement f(w) = \frac{1}{m} \sum_{i=1}^m f_i(w))$$

\Rightarrow Explique la tâche
 d'apprentissage à réaliser (régression, classification, etc)

$$f(w) + \lambda \mathcal{L}(w) \rightarrow$$

$\lambda \geq 0$: poids du terme en $\mathcal{L}(w)$
 par rapport au terme en $f(w)$

Terme de régularisation
 $\mathcal{L}: \mathbb{R}^d \rightarrow \mathbb{R}$

\Rightarrow Modélise des propriétés souhaitées pour la solution

Exemples en régression linéaire

$$X \in \mathbb{R}^{m \times d}, \quad y \in \mathbb{R}^m$$

$$\begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix}, \quad \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

Problème non régularisé:

Bal: Trouver $w \in \mathbb{R}^d$ tq $Xw \approx y$

$$\text{minimiser}_{w \in \mathbb{R}^d} \frac{1}{2m} \|Xw - y\|_2^2$$

$$\|Xw - y\|_2^2 = \sum_{i=1}^m (x_i^T w - y_i)^2$$

\rightarrow Régularisation ℓ_2 (aka "ridge", "Tikhonov", ...)

$$\text{minimiser}_{w \in \mathbb{R}^d} \frac{1}{2m} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$

$$\|w\|_2^2 = \sum_{j=1}^d w_j^2$$

$\lambda = 0 \Rightarrow$ Problème de départ

$\lambda \ll 1$: proche du problème de départ

$\lambda \gg 1$: proche du problème

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad \lambda \|w\|_2^2 \quad \Rightarrow \begin{array}{l} \text{unique} \\ \text{solution} \\ w^* = 0_{\mathbb{R}^d} \end{array}$$

Pourquoi régulariser avec la norme ℓ_2 ?

→ Améliorer la généralisation du modèle

$$\Rightarrow \text{erreur } \|X^{\text{test}} w^* - y^{\text{test}}\|_2^2$$

$$\text{avec } w^* = \underset{w \in \mathbb{R}^d}{\text{argmin}} \frac{1}{2n} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$

et $(X^{\text{test}}, y^{\text{test}})$ proviennent de la même distribution de données que (X, y)

→ Rendre le problème fortement convexe

$$w \mapsto \lambda \|w\|_2^2$$

est 2λ -fortement convexe

$$w \mapsto \frac{1}{2n} \|Xw - y\|_2^2$$

est convexe

⊕ Garantie d'unicité de la solution

⊕ Convergence plus rapide des algorithmes type descente de gradient / gradient stochastic

→ Réduire la sensibilité de la solution par rapport aux données

→ On souhaite que la solution varie peu si les données (x_i, y_i) sont légèrement perturbées

Ex) $n=d=2$

$$X = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\underset{w \in \mathbb{R}^2}{\text{minimiser}} \quad \frac{1}{2n} \|Xw - y\|_2^2 = \frac{1}{4} [(w_1 - 1)^2 + 1]$$

$$\underset{w \in \mathbb{R}^2}{\operatorname{argmin}} \quad \frac{1}{2n} \|Xw - y\|^2 = \left\{ w = \begin{bmatrix} 1 \\ w_2 \end{bmatrix} \mid w_2 \in \mathbb{R} \right\}$$

$$A > 0, \quad \underset{w \in \mathbb{R}^2}{\operatorname{minimiser}} \quad \left\{ \frac{1}{2n} \|Xw - y\|_2^2 + \lambda \|w\|_2^2 \right\}$$

$$\frac{1}{4}w_1^2 - \frac{1}{2}w_1 + \lambda w_1^2$$

$$\frac{1}{4}((w_1 - 1)^2 + 1) + \lambda w_1^2 + \lambda w_2^2$$

$$\frac{1}{2}w_1 + 2\lambda w_1 - \frac{1}{2} = 0$$

$$w_1 = \frac{1}{1+4\lambda}$$

$$\underset{w \in \mathbb{R}^2}{\operatorname{argmin}} \quad \rightarrow = \left\{ \begin{bmatrix} \frac{1}{1+4\lambda} \\ 0 \end{bmatrix} \right\}$$

Remarque : Dans le cas général

$$\underset{w \in \mathbb{R}^d}{\operatorname{minimiser}} \quad f(w) + \lambda \|w\|_2^2$$

avec $f \in C^1$, l'algorithme de descente de gradient sur ce problème donne

$$w_{k+1} = w_k - \alpha_n [\nabla f(w_k) + 2\lambda w_k]$$

$$\nabla (\lambda \|w\|_2^2)(w_k) = 2\lambda w_k$$

$$= (1 - 2\lambda \alpha_n) w_k - \alpha_n \nabla f(w_k)$$

Pour $\alpha_n < \frac{1}{2\lambda}$, "weight decay": on réduit l'amplitude des coefficients de w_k en plus de faire un pas de gradient

→ Régularisation ℓ_1 / LASSO

Pour la régression linéaire,

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad \frac{1}{2n} \|Xw - y\|_2^2 + \lambda \|w\|_1$$

$$\|w\|_1 = \sum_{j=1}^d |w_j|$$

Pourquoi la norme ℓ_1 ?

→ Favorise les solutions avec des coordonnées nulles

$$(\text{cas extrême } \lambda \gg 1 : \underset{w \in \mathbb{R}^d}{\text{minimiser}} \lambda \|w\|_1)$$

$$\downarrow$$

$$w^* = 0_{\mathbb{R}^d}$$

→ lorsque λ augmente, la solution du problème a de plus en plus de coordonnées nulles, et les coordonnées sont "mises à 0" par ordre croissant

d'importance

(Ex)

$$\underset{w \in \mathbb{R}^2}{\text{minimiser}}$$

$$\left[(w_1 - 1)^2 + (w_2 - 2)^2 \right] \Rightarrow w^* = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$\underset{w \in \mathbb{R}^2}{\text{minimiser}}$$

$$\left[(w_1 - 1)^2 + (w_2 - 2)^2 \right] + \lambda \|w\|_2^2$$

$$2w_1 - 2 + 2\lambda w_1 = 0$$

$$w_1 = \frac{1}{1+\lambda}$$

$$\Rightarrow w^* = \begin{bmatrix} \frac{1}{1+\lambda} \\ \frac{2}{1+\lambda} \end{bmatrix} \xrightarrow{\lambda \rightarrow +\infty} \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\underset{\mathbf{w} \in \mathbb{R}^2}{\text{minimiser}} \quad \left[(\omega_1 - 1)^2 + (\omega_2 - 2)^2 \right] + \lambda \|\mathbf{w}\|_1$$

$$(\omega_1 - 1)^2 + \lambda |\omega_1|$$

$$2(\omega_1 - 1) \in \lambda \partial |\omega_1|$$

$$0 \leq \lambda < 2 \quad \mathbf{w}^* = \begin{bmatrix} 1 - \frac{\lambda}{2} \\ 2 - \frac{\lambda}{2} \end{bmatrix}$$

$$2 \leq \lambda < 4 \quad \mathbf{w}^* = \begin{bmatrix} 0 \\ 2 - \lambda/2 \end{bmatrix}$$

$$4 \leq \lambda \quad \mathbf{w}^* = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Q) Existe-t-il un algorithme qui s'applique à tout problème de la forme $\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \quad f(\mathbf{w}) + \lambda \varphi(\mathbf{w})$?

Hypothèses : $f \in C^1$ (mais pas forcément convexe)

φ convexe (mais pas forcément dérivable)

Itération du gradient proximal :

$$w_{k+1} \in \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ f(w_k) + \nabla f(w_k)^T (w - w_k) + \frac{1}{2\alpha_h} \|w - w_k\|_2^2 + \eta L(w) \right\}$$

↑
 $\approx f(w)$
 autour de w_k

Tame proximal : garantir que
 w est "proche" de w_k
↑
 avec $\alpha_h > 0$
↑
 Régularisation
 (inchangée)

Définition implicite : w_{k+1} est défini comme la solution d'un problème d'optimisation

Remarques importantes : Si $L \equiv 0$ ($L(w) = 0 \forall w$)

alors

$$\underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ f(w_k) + \nabla f(w_k)^T (w - w_k) + \frac{1}{2\alpha_h} \|w - w_k\|_2^2 \right\}$$

$$= w_k - \alpha_h \nabla f(w_k)$$

- Dans le cas général, l'algorithme du gradient proximal n'est intéressant que si le sous-problème résolu à chaque itération est plus facile à résoudre que le problème de départ !
 - ⇒ Pour les régularisations classiques (ℓ_2, ℓ_1, \dots), on peut résoudre le sous-problème de manière explicite
 - ⇒ Cas $L(w) = \|w\|_1$: l'algorithme du gradient proximal correspond à l'algorithme ISTA (Iterative Soft-Thresholding Algorithm)

Exo 2 TD 3

a) La régularisation serv à favoriser les vecteurs/modèles qui vérifient une propriété souhaitée (permanence, par exemple).

\Leftrightarrow pénaliser les vecteurs/modèles qui ne vérifient pas une propriété souhaitée.

d) $w_{k+1} \in \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \phi(w_h) + \nabla \phi(w_h)^T (w - w_h) + \frac{1}{2\alpha_h} \|w - w_h\|_2^2 + \frac{\lambda_2}{2} \|w\|_2^2 + \lambda_1 \|w\|_1 \right\}$

$$\phi(w) = \frac{1}{2m} \|Xw - y\|_2^2$$

$$\Leftrightarrow w_{k+1} \in \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \frac{1}{2m} \|Xw_h - y\|_2^2 + \frac{1}{m} (Xw_h - y)^T X (w - w_h) + \frac{1}{2\alpha_h} \|w - w_h\|_2^2 + \frac{\lambda_2}{2} \|w\|_2^2 + \lambda_1 \|w\|_1 \right\}$$

e) $\lambda_2 = 0$

$$w_{k+1} \in \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \phi(w_h) + \nabla \phi(w_h)^T (w - w_h) + \frac{1}{2\alpha_h} \|w - w_h\|_2^2 + \lambda_1 \|w\|_1 \right\}$$

ISTA = gradient proximal pour ce problème
avec formule explicite pour w_{k+1}

g) minimieren $\phi'(w_k) + \nabla \phi(w_k)^T (w - w_k) + \frac{1}{2\alpha_k} \|w - w_k\|_2^2 + \frac{\lambda_2}{2} \|w\|_2^2 + \lambda_1 \|w\|_1$

$f(w)$

$\lambda_2 > 0 \quad \lambda_1 > 0$
 $w_k \in \mathbb{R}^d, \alpha_k > 0$

→ Sors-gradient ✓

→ ISTA ✓

↪ Problem : minimieren $w \in \mathbb{R}^d \quad f(w) + \lambda_1 \|w\|_1$

ISTA

$$w^{(j)}, \dots, w^{(j)}, w^{(j+1)}, \dots$$

↪ Fors-problem:

$$\begin{aligned} w^{(j+1)} = \underset{\bar{w} \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ f(w^{(j)}) + \nabla f(w^{(j)})^T (\bar{w} - w^{(j)}) \right. \\ \left. + \frac{1}{2\beta^{(j)}} \|\bar{w} - w^{(j)}\|_2^2 + \lambda_1 \|\bar{w}\|_1 \right\} \end{aligned}$$

$$\text{mit } \beta^{(j)} > 0$$