

IASD App
2025-2026

OPTIMISATION POUR L'APPRENTISSAGE AUTOMATIQUE

3 février 2026 (Après-midi)

Cette séance : Régularisation
Algorithme proximal
Quelques exemples

A venir : Feuille d'exercices mise à jour

MÉTHODES PROXIMALES ET RÉGULARISATION

Motivation:

$$(P) \quad \underset{w \in \mathbb{R}^d}{\text{minimize}} \quad \underbrace{\frac{1}{n} \sum_{i=1}^n f_i(w)}_{f(w)} \quad f_i \text{ dépend de données}$$

- Si le problème possède plusieurs solutions, laquelle choisir ?
- Si on a un a priori sur la solution, comment faire en sorte que la solution de (P) ait ces propriétés ?

(Ex) (X, y) jeu de données pour la régression

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

avec $y = Xw^* + \varepsilon$

ε : bruit
 w^* : vraie tendance linéaire

w^* possède beaucoup de coefficients nuls (99%) \leftarrow w^* est "sparse" (ou creux, ou parcimonieux)

But: Modifier le problème de départ pour prendre en compte des propriétés (réelles ou souhaitées) sur la solution

Définition: Problème régularisé

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad f(w) + \lambda \Omega(w)$$

• $f: \mathbb{R}^d \rightarrow \mathbb{R}$ terme d'attaché aux données

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

• $r: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ terme de régularisation
 (a priori indépendant des données)

• $\lambda \geq 0$: paramètre de régularisation

↳ Lorsque $\lambda = 0$, le problème est le problème d'origine
 minimiser $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$
 $w \in \mathbb{R}^d$

↳ Lorsque $\lambda \rightarrow \infty$, le problème est équivalent à
 minimiser $r(w)$ (plus d'impact des données)
 $w \in \mathbb{R}^d$

↳ Le problème avec $\lambda > 0$ correspond à pondérer l'effet de la régularisation par rapport à l'attaché aux données

① Régularisation ℓ_2 ("ridge" / écarté, Tychonov, ...)

$$r(w) = \frac{1}{2} \|w\|^2 = \frac{1}{2} \sum_{j=1}^d w_j^2$$

$$\text{minimiser } f(w) + \frac{\lambda}{2} \|w\|^2$$

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

But:

→ Réduire la sensibilité de la solution vis-à-vis des données (générique)

→ Convexifier le problème (spécifique à l'optimisation)

- $w \mapsto \frac{1}{2} \|w\|^2$ est λ -fortement convexe

(φ λ -fortement convexe)

$$\Leftrightarrow \varphi(\alpha v + (1-\alpha)w) \leq \alpha \varphi(v) + (1-\alpha) \varphi(w) - \frac{1}{2} \alpha(1-\alpha) \|v-w\|^2$$



- Si f est convexe, alors $f + \frac{\lambda}{2} \|\cdot\|^2$ est λ -fortement convexe

(unique minimum global)

Illustration

minimiser $w \in \mathbb{R}^2$ $\overbrace{\frac{1}{4} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2}_{f_\lambda(w)}$ $\lambda \geq 0$

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

$$\lambda = 0 : \operatorname{argmin}_w f_0(w) = \left\{ \begin{bmatrix} 1 \\ w_2 \end{bmatrix} \mid w_2 \in \mathbb{R} \right\}$$

$$f_0(w) = \frac{1}{4} \left((w_1 - 1)^2 + 1 \right)$$

- Supposons qu'on perturbe $X \rightarrow X_\varepsilon = \begin{bmatrix} 1 & 0 \\ 0 & \varepsilon \end{bmatrix}$ avec $0 < \varepsilon \ll 1$

Alors $\operatorname{argmin} \frac{1}{4} \|X_\varepsilon w - y\|^2 = \left\{ \begin{bmatrix} 1 \\ 1/\varepsilon \end{bmatrix} \right\}$

$$\underbrace{\begin{pmatrix} X_{\varepsilon}^T X_{\varepsilon} + 2\lambda I \end{pmatrix}}_{\begin{bmatrix} 1 & 0 \\ 0 & \varepsilon^2 \end{bmatrix} + \begin{bmatrix} 2\lambda & 0 \\ 0 & 2\lambda \end{bmatrix}} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \underbrace{X_{\varepsilon}^T y}_{\begin{bmatrix} 1 \\ \varepsilon \end{bmatrix}}$$

La norme de la solution régularisée est d'ordre

$$O(\varepsilon) \text{ pour } \varepsilon \text{ petit, donc } \begin{matrix} \rightarrow 0 \\ \varepsilon \rightarrow 0 \end{matrix}$$

\Rightarrow le problème avec régularisation est moins sensible aux perturbations dans les données

\hookrightarrow Algorithmiquement, la régularisation l_2 est faite à intégrer dans une méthode de type gradient (descente de gradient, gradient stochastique, voire sous-gradient)

Ex) Si f est C^1 , alors

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

Descente de gradient sur le problème sans régularisation ($\lambda=0$)

$$\begin{aligned} w_{k+1} &= w_k - \alpha_k (\nabla f(w_k) + \lambda w_k) \\ &= \underbrace{(1 - \lambda \alpha_k)}_{\downarrow} w_k - \alpha_k \nabla f(w_k) \end{aligned}$$

Descente de régularisation pour le problème régularisé $\lambda > 0$

$1 - \lambda \alpha_k \in (0, 1)$
pour α_k suffisamment petit

\hookrightarrow Réduction des composantes de w_k par un facteur multiplicatif à chaque itération

"Weight decay"

NB: SGD de PyTorch / Descente de gradient avec "weight decay"

(\Rightarrow)

sur problème avec régularisation l_2

* SGD avec momentum

• Adam / AdamW ont aussi un paramètre de weight decay, mais qui n'est pas équivalent à une régularisation l_2 à cause du momentum \Rightarrow A l'origine de AdamW

(2) Régularisation l_1 (LASSO, ...)

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

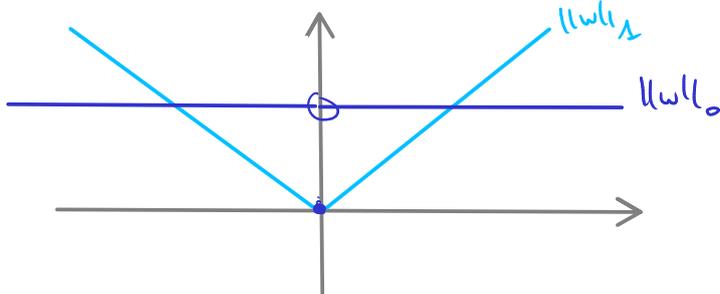
(minimiser $\|w\|_1$
 $w \in \mathbb{R}^d$
 \hookrightarrow argmin : $\begin{bmatrix} 0 \\ b \end{bmatrix}$
 \hookrightarrow min = 0

minimiser $f(w) + \lambda \|w\|_1$ pour $\lambda \geq 0$
 $w \in \mathbb{R}^d$

But: Réaliser les vecteurs qui n'ont pas beaucoup de coefficients nuls / Favoriser les vecteurs creux / "sparse"

Remarque: la régularisation "idéale" pour avoir une solution creuse serait

$\lambda \|w\|_0$, $\|w\|_0 = \begin{cases} \text{nombre de} \\ \text{coefficients non} \\ \text{nuls de } w \end{cases}$



"Norme l_0 "

\hookrightarrow Non convexe, discontinue, à valeurs discrètes (constante par morceaux)

\rightarrow Norme l_1 est convexe, continue et c'est la fonction convexe la plus proche de $\| \cdot \|_0$ autour de 0

Interprétation :

minimiser $\|Xw - y\|^2 + \lambda \|w\|_1$
 $w \in \mathbb{R}^2$

$$X = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$\lambda = 0$ minimiser $(w_1 - 1)^2 + (w_2 - 2)^2$
 $w \in \mathbb{R}^2$

$$\text{argmin}(\) = \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right\}$$

$\lambda > 0$ minimiser $(w_1 - 1)^2 + (w_2 - 2)^2 + \lambda \|w\|_1$
 $\lambda |w_1| + \lambda |w_2|$

$0 \in \frac{\partial}{\partial w_1} (w_1 - 1)^2 + \lambda |w_1|$
 $-\frac{\partial}{\partial w_1} (w_1 - 1)^2$

↳ Via les sous-gradients, on montre que le problème a une unique solution

$$0 \leq \lambda < 1 \quad \begin{bmatrix} 1 - \lambda \\ 2 - \lambda \end{bmatrix}$$

$$1 \leq \lambda < 2 \quad \begin{bmatrix} 0 \\ 2 - \lambda \end{bmatrix}$$

$$2 \leq \lambda \quad \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

→ la régularisation l_1 permet d'identifier les coordonnées les plus importantes en faisant varier λ . Plus une coordonnée sera importante, plus λ devra être grand pour qu'elle soit nulle en la solution

⇒ A l'origine, la régularisation l_1 était utilisée en régression linéaire pour de la sélection d'attributs (feature selection)

En termes d'algorithmes

- Si f est convexe, alors $f + \lambda \|\cdot\|_1$ aussi et on peut utiliser une méthode de sous-gradient
- Pour $f \in C^1$, il existe une méthode dérivée (ISTA) qui est une variante de l'algorithme du gradient proximal

↳ Il existe de nombreux autres termes de régularisation basés sur la norme l_1 qui encouragent la parcimonie de manière structurée

Ex) Group LASSO

$$r(w) = \sum_{g \in G} \overbrace{\|w_{g,\cdot}\|_1}^{\| \begin{bmatrix} w_{g,1} \\ \vdots \\ w_{g,m} \end{bmatrix} \|_1}$$

$$\text{ou } w = \begin{bmatrix} w_{g_1} \\ \vdots \\ w_{g_m} \end{bmatrix}$$

M groupes de variables

$$G = \{g_1, \dots, g_m\}$$

Toutes les variables d'un même groupe doivent être soit nulles, soit non nulles

③ Méthodes proximales / Gradient proximal

Idee: On cherche un algorithme qui puisse s'appliquer aux problèmes de la forme suivante

$$\underset{w \in \mathbb{R}^d}{\text{minimiser}} \quad f(w) + \lambda r(w) \quad \lambda \geq 0$$

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

Algorithme du gradient proximal pour ce problème

$$w_{k+1} \in \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \underbrace{f(w_k) + \nabla f(w_k)^\top (w - w_k)}_{\substack{\text{Approximation de} \\ f(w) \text{ autour de } w_k}} + \underbrace{\frac{1}{2\alpha_k} \|w - w_k\|^2}_{\substack{\text{Terme} \\ \text{proximal} \\ \text{pénalise les} \\ w \text{ loin de} \\ w_k}} + \underbrace{\lambda r(w)}_{\substack{\text{Terme} \\ \text{de} \\ \text{régularisation} \\ (\text{inchangé} \\ \text{par} \\ \text{rapport au} \\ \text{problème de} \\ \text{base})}} \right\}$$

avec $\alpha_k > 0$

Une itération de gradient proximal:

- calculer une solution d'un problème d'optimisation "sous-problème" où on a "linéarisé" le terme $f(w)$
- pour résoudre le sous-problème (une fois que $\nabla f(w_k)$ est calculé) il n'est pas nécessaire d'accéder à f ou à son gradient (et donc aux données)
- le sous-problème est toujours fortement convexe, et donc w_{k+1} est défini de manière unique

L'algorithme du gradient proximal n'est intéressant que si le sous-problème est facile à résoudre (plus facile au moins que le problème d'origine). Cela est souvent le cas lorsque le "prox" de r est facile à calculer

Def: $h: \mathbb{R}^d \rightarrow \mathbb{R}$ convexe

$$\forall w \in \mathbb{R}^d, \quad \operatorname{prox}_h(w) \in \underset{v \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \underbrace{\frac{1}{2} \|v - w\|^2 + h(v)}_{\text{fortement convexe}} \right\}$$

L'itération du gradient proximal se réécrit

$$w_{k+1} \in \underset{w}{\operatorname{argmin}} \left\{ f(w_k) + \nabla f(w_k)^\top (w - w_k) + \frac{1}{2\alpha_k} \|w - w_k\|^2 + \lambda r(w) \right\}$$

$$\Leftrightarrow w_{k+1} = \text{prox}_{\lambda \alpha_k \pi} \left(w_k - \alpha_k \nabla f(w_k) \right)$$

↑
itération de la descente de gradient

$$\text{prox}_{\lambda \alpha_k \pi}(w) \in \underset{v}{\text{argmin}} \left\{ \frac{1}{2} \|v-w\|^2 + \lambda \alpha_k \pi(v) \right\}$$

Exemples

• $\pi(w) = 0$ (ou $\lambda = 0$)

$$\text{prox}_0(w) \in \underset{v}{\text{argmin}} \left\{ \frac{1}{2} \|v-w\|^2 + 0 \right\} = \{w\}$$

Sans régularisation, l'itération de gradient proximal

devient $w_{k+1} = w_k - \alpha_k \nabla f(w_k) \Rightarrow$ Descente de gradient!

• $\pi(w) = \frac{1}{2} \|w\|^2$

$$w_{k+1} = \text{prox}_{\lambda \alpha_k \pi}(w_k - \alpha_k \nabla f(w_k)) = \frac{1}{1 + \lambda \alpha_k} w_k - \frac{\alpha_k}{1 + \lambda \alpha_k} \nabla f(w_k)$$

↓
"weight decay"

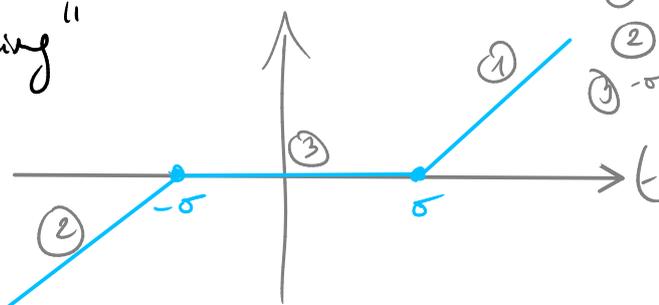
↓
"gradient decay"

≠ descente de gradient appliqué au problème avec régularisation

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

• $\pi(w) = \|w\|_1$

↳ Le "prox" de la norme l_1 est défini par une fonction dite de "soft-thresholding"



- ① $t > \sigma \Rightarrow t - \sigma$
- ② $t < -\sigma \Rightarrow t + \sigma$
- ③ $-\sigma \leq t \leq \sigma \Rightarrow 0$

Pour le gradient proximal, on obtient

$$\forall j=1..d, \quad [w_{k+1}]_j = \begin{cases} [w_k - \alpha_k \nabla f(w_k)]_j - \lambda \alpha_k & \text{si } [w_k - \alpha_k \nabla f(w_k)]_j > \lambda \alpha_k \\ [w_k - \alpha_k \nabla f(w_k)]_j + \lambda \alpha_k & \text{si } [w_k - \alpha_k \nabla f(w_k)]_j < -\lambda \alpha_k \\ 0 & \text{sinon} \end{cases}$$

→ Cette formule produit une suite d'itérés qui ont nécessairement au moins autant de coordonnées nulles que les itérés de la descente de gradient

$$\forall k, \quad \|w_{k+1}\|_0 \leq \|w_k - \alpha_k \nabla f(w_k)\|_0$$

→ La méthode du gradient proximal avec régularisation l_1 est connue (notamment en traitement du signal) comme la méthode ISTA (Iterative Soft-Thresholding Algorithm)

NB: ISTA avec momentum = FISTA

(Le gradient proximal peut être défini en versions stochastique, avec sous-gradients, etc)

→ Cas d'application principal de ISTA / régularisation l_1 :
"Sparse recovery"

$$(X, y) \in \mathbb{R}^{n \times d} \times \mathbb{R}^n$$

$d \gg n$

$$y = X w^* + \varepsilon$$

"grand bruit"
↑
bruit
(par ex. gaussien)

$\|w^*\|_0 \ll d$

But: Parmi tous les modèles possibles (w), trouver celui le plus parcimonieux qui colle aux données

\Rightarrow minimiser $w \in \mathbb{R}^d$ $\frac{1}{2} \|Xw - y\|^2 + \lambda \|w\|_1$, résolu par ISTA/gradient proximal

Autres exemples en traitement du signal et des images

• Débruitage

z : image bruitée ($z^* + \varepsilon$)

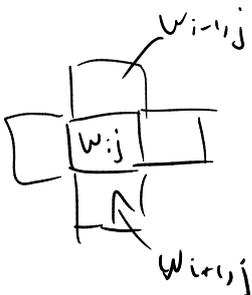
Postulat: Pour une image bruitée, la différence entre des pixels voisins doit être plus grande que pour une image non bruitée

\rightarrow Pour débruiter l'image, on peut alors calculer le

"prox" suivant $\text{prox}_{\lambda TV}(z) = \arg \min_w \left\{ \frac{1}{2} \|w - z\|^2 + \lambda \|w\|_{TV} \right\}$

$\|w\|_{TV}$: Variation totale

$$\|w\|_{TV} = \sum_{(i,j)} \|(\nabla w)_{ij}\| \quad (i,j) \text{ indice pixel}$$



$$\begin{aligned} (\nabla w)_{ij} = & (w_{ij} - w_{i,j+1})^2 \\ & + (w_{ij} - w_{i,j-1})^2 \\ & + (w_{ij} - w_{i+1,j})^2 \\ & + (w_{ij} - w_{i-1,j})^2 \end{aligned}$$

High-dimensional data analysis with low-dimensional models
J. Wright, Y. Ma (2009)

Approximation de matrice

$X \in \mathbb{R}^{m \times d}$: image (donnée)

But: Récupérer une structure sous-jacente à l'image

1. $X = L + E$, L matrice de rang faible $r \ll \min(m, d)$
 et E une perturbation de rang plein
 $\text{rang}(E) = \min(m, d)$

Pb: Retrouver L à partir de X

Formulation:

minimiser $W \in \mathbb{R}^{m \times d}$

$$\frac{1}{2} \|W - X\|_F^2 + \lambda \|W\|_*$$

$\sum_{i,j} (w_{ij} - x_{ij})^2$
 "Norme de Frobenius"

Norme nucléaire
 (pénalise les matrices de rang élevé)

$\|W\|_* =$ somme des valeurs singulières
 $\sigma_1 \geq \dots \geq \text{rang}(W)$

$\forall A \in \mathbb{R}^{m \times d}$,

$$\text{A rang } r \Leftrightarrow A = U \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_r \end{bmatrix} V^T \quad \left(\begin{smallmatrix} \text{SVD} \\ \text{de} \\ A \end{smallmatrix} \right)$$

avec $\sigma_1 \geq \dots \geq \sigma_r > 0$

$$U \in \mathbb{R}^{m \times m}, U^T U = I_m$$

$$V \in \mathbb{R}^{d \times d}, V^T V = I_d$$

→ On peut calculer explicitement le prox de $\|\cdot\|_*$
 → donc on peut faire du gradient proximal !

$$I_m = \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix}$$

$$\|I_m\|_1 = \sum_{i,j} |(I_m)_{ij}| = m \leq m^2$$

$\frac{m}{m^2}$ coeffs normés

$$\|I_m\|_* = m \text{) rang maximal}$$

2) Approximation de rang faible + parcimonieuse

Hypothèse: $X = \underbrace{L}_{\text{rang}(L) \ll \min(m,d)} + \underbrace{S}_{\text{parcimonieuse}} \quad \|S\|_0 \ll nd$

minimiser
 $W \in \mathbb{R}^{m \times d}$
 $U \in \mathbb{R}^{m \times d}$

$$\frac{1}{2} \|X - (W+U)\|_F^2 + \lambda \underbrace{\|W\|_1}_{\text{rang faible}} + \mu \underbrace{\|U\|_1}_{\text{parcimonie}}$$