

I. Principe du Lasso

1. Principe

Variables explicative : $X \in \mathbb{R}^{n \times p}$

Variable expliquée : $y \in \mathbb{R}^n$

Modèle de régression linéaire : Trouver β tel que $y = X\beta + \varepsilon$

$$\min_{\beta \in \mathbb{R}^p} \sum \varepsilon_i^2 = \|\varepsilon\|^2 = \|y - X\beta\|^2$$

Problème quand $p \gg n \Rightarrow$ Hypothèse : Il y a beaucoup de $\beta_j = 0 \Leftrightarrow \|\beta\|_0 \leq k^1$. Cependant, $\|\beta\|_0 \leq k$ est une contrainte non-convexe qu'il est difficile d'optimiser, on choisit donc la norme $\|\beta\|_1 \leq k$ qui est la plus petite norme entraînant une contrainte convexe.

$$\boxed{\begin{cases} \min_{\beta \in \mathbb{R}^p} & \|y - X\beta\|^2 \\ \text{s.c.} & \|\beta\|_1 \leq k \end{cases}} \Leftrightarrow \begin{cases} \min_{\beta \in \mathbb{R}^p} & \frac{1}{2} \beta^T X^T X \beta - y^T X \beta \\ \text{s.c.} & \sum_{j=1}^p |\beta_j| \leq k \end{cases} \Leftrightarrow \begin{cases} \min_{(\beta^+, \beta^-) \in \mathbb{R}^p} & \frac{1}{2} \beta^T X^T X \beta - y^T X \beta \\ \text{s.c.} & (\beta^+ + \beta^-)^T e \leq k \\ & \beta^+, \beta^- \geq 0 \\ \text{avec} & \beta = \beta^+ - \beta^- \end{cases}$$

Le problème est équivalent pour λ fixé à :

$$\boxed{\min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|y - X\beta\|^2 + \lambda \|\beta\|_1}$$

2. Front de Pareto

Tirer au hasard des valeurs de β , calculer les normes de $\|y - X\beta\|^2 + \lambda \|\beta\|_1$ (optim multi-critère)

II. Chemin de régularisation du Lasso

L'évolution de β est linéaire par morceaux, entre 2 apparitions d'une composante non nulle.

$$J(\beta) = \frac{1}{2} \|y - X\beta\|^2 + \lambda \|\beta\|_1$$

$$\underset{\beta}{\operatorname{argmin}} J(\beta) \Leftrightarrow \beta^* \text{ t. q. } \nabla J(\beta^*) = 0$$

1. Gradient de $\frac{1}{2} \|y - X\beta\|^2$

$$\nabla_{\beta} \left(\frac{1}{2} \|y - X\beta\|^2 \right) = \nabla_{\beta} (\|y\|^2 - 2\beta^T X^T y + \beta^T X^T X \beta) = X^T X \beta - X^T y$$

Si $\lambda = 0$, $\hat{\beta}_{MC}(\lambda = 0) = (X^T X)^{-1} (X^T y)$ (au sens des moindres carrés)

2. Gradient de $\|\beta\|_1$

On veut calculer un gradient d'une forme contenant des valeurs absolues, mais la dérivée de la valeur absolue n'est pas définie en 0.

On va donc utiliser la notion de sous-différentielle.

¹ $\|\cdot\|_0$ vaut le nombre de termes non nuls

COURS D'APPRENTISSAGE EN CONTEXTE

APPC – Cours

a. Sous-gradient

$$g \in \mathbb{R}^n \text{ est un sous-gradient de } J \text{ en } x_0 \text{ si } \forall x \in \mathcal{V}(x_0) J(x) \geq J(x_0) + g^\top(x - x_0)$$

Interprétation : On peut prendre la dérivée de n'importe quelle droite comprise sous la courbe initiale au voisinage de x_0 .

b. Sous-différentielle $\partial_x J(x^*)$

Ensemble des sous gradients possibles

c. Solution x^* d'une minimisation d'un coût convexe J

Si J est différentiable alors $x^* = \operatorname{argmin}_x J(x) \Leftrightarrow \nabla J(x^*) = 0$

$$\text{Si } J \text{ est non-différentiable alors } x^* \text{ est solution si } 0 \in \partial_x J(x^*) \quad (\Leftrightarrow \exists \alpha \in \partial_x J(x^*) | g_\alpha = 0)$$

d. Application à J

◆ *Sous différentielle de $J_1(x) = |x|$*

$$\partial_x J_1(x=0) = \{\gamma | J_1(x) \geq J_1(0) + \gamma x\} = \{\gamma | |x| \geq \gamma x\}$$

Si $x > 0, \gamma \leq 1$. Si $x < 0, \gamma \geq -1$. $\Rightarrow \gamma \in [-1; 1]$

◆ *Sous différentielle de J*

$$\partial_\beta J(\beta) = X^\top X \beta - X^\top y + \lambda v \text{ avec } v_i = \begin{cases} \operatorname{sign}(\beta_j) & \text{si } \beta_j \neq 0 \\ \alpha_j \in [-1; 1] & \text{si } \beta_j = 0 \end{cases}$$

Soit $I_\beta = \{j | \beta_j \neq 0\}, I_0 = \{j | \beta_j = 0\}$

On veut : $\exists \alpha_j \text{ t.q. } \partial_\beta J(\beta_\lambda^*) = 0 \Leftrightarrow 0 = X^\top X \beta - X^\top y + \lambda v$

Pour $j \in I_\beta$:

$$\left[X(:, I_\beta)^\top X(:, I_\beta) \right] \beta(I_\beta) + \lambda \operatorname{sign} \beta(I_\beta) - X(:, I_\beta)^\top y = 0$$

$$\Leftrightarrow \left[X(:, I_\beta)^\top X(:, I_\beta) \right] \beta(I_\beta) = X(:, I_\beta)^\top y - \lambda \operatorname{sign} \beta(I_\beta)$$

3. Calcul du chemin de régression

a. Initialisation

◆ $\lambda = 0$

$$J(\beta) = \frac{1}{2} \|y - X\beta\|^2 \Rightarrow X^\top X \widehat{\beta}_0 = X^\top y$$

◆ $\lambda \in V(0)$

$$X^\top X \beta_\lambda = X^\top y - \lambda \operatorname{sign}(\widehat{\beta}_0)$$

On a donc quand on fait varier λ :

$$\beta_\lambda = \widehat{\beta}_0 - \lambda \underbrace{(X^\top X)^{-1} \operatorname{sign}(\widehat{\beta}_0)}_v$$

b. Process itératif pour l'étape k

On part des MC et on cherche le λ qui va annuler un β_j

Soit $\lambda_k, \widehat{\beta}_k$

Pour $\lambda \in V(\lambda_k) > \lambda_k$,

$$\left[X(:, I_\beta)^\top X(:, I_\beta) \right] \beta_\lambda = X(:, I_\beta)^\top y - \lambda \operatorname{sign} \widehat{\beta}_k \quad (1)$$

Vrai pour λ_k

$$\left[X(:, I_\beta)^\top X(:, I_\beta) \right] \beta_k = X(:, I_\beta)^\top y - \lambda_k \operatorname{sign} \widehat{\beta}_k \quad (2)$$

(1) – (2) :

$$\boxed{\beta_\lambda = \beta_k - (\lambda - \lambda_k) \left(\left[X(:, I_\beta)^\top X(:, I_\beta) \right]^{-1} \operatorname{sign} \widehat{\beta}_k \right)} \quad \boxed{\lambda_{k+1} = \beta_{k_j} + \frac{\lambda_k v_i}{v_i} = \lambda_k + \frac{\beta_{k_j}}{v_j}}$$

III. Lasso « component wise » (à la Gauss-Seidel)

1. Introduction

Pour $p = 1$

$$\min_{\beta \in \mathbb{R}} \frac{1}{2} \sum_{i=1}^n (x_i \beta - y_i)^2 + \lambda |\beta|$$

Si $\lambda = 0$, on a le problème des moindres carrés :

$$\min_{\beta \in \mathbb{R}} J_{MC} = \frac{1}{2} \sum_{i=1}^n (x_i \beta - y_i)^2$$

$$\frac{dJ_{MC}}{d\beta} = \sum_{i=1}^n (x_i \beta - y_i) x_i = 0 \Rightarrow \beta \sum_{i=1}^n x_i^2 + \sum_{i=1}^n x_i y_i = 0 \Rightarrow \hat{\beta}_{MC} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} = \frac{x^\top y}{\|x\|^2}$$

Pour tout λ , on a le Lasso a une variable :

$$\min_{\beta \in \mathbb{R}} J_{MC} = \frac{1}{2} \sum_{i=1}^n (x_i \beta - y_i)^2 + \lambda |\beta|$$

$$\partial_\beta J_{MC} = \sum_{i=1}^n (x_i \hat{\beta}_\lambda - y_i) x_i + \lambda \begin{cases} \operatorname{sign} \hat{\beta}_\lambda & \text{si } \beta \neq 0 \\ \alpha \in [-1; 1] & \text{si } \beta = 0 \end{cases} = 0$$

Si $\hat{\beta}_\lambda \neq 0$:

$$\partial_\beta J_{MC} = \sum_{i=1}^n (x_i \hat{\beta}_\lambda - y_i) x_i + \lambda \operatorname{sign} \hat{\beta}_\lambda = 0$$

$$\hat{\beta}_\lambda = \frac{x^\top y}{\|x\|^2} - \lambda \frac{\operatorname{sign} \hat{\beta}_\lambda}{\|x\|^2} = \hat{\beta}_{MC} - \lambda \frac{\operatorname{sign} \hat{\beta}_\lambda}{\|x\|^2} = \begin{cases} \hat{\beta}_{MC} - \frac{\lambda}{\|x\|^2} & \text{si } \hat{\beta}_{MC} > 0 \\ \hat{\beta}_{MC} + \frac{\lambda}{\|x\|^2} & \text{si } \hat{\beta}_{MC} < 0 \end{cases} = \operatorname{sign} \hat{\beta}_\lambda \left(|\hat{\beta}_{MC}| - \frac{\lambda}{\|x\|^2} \right)$$

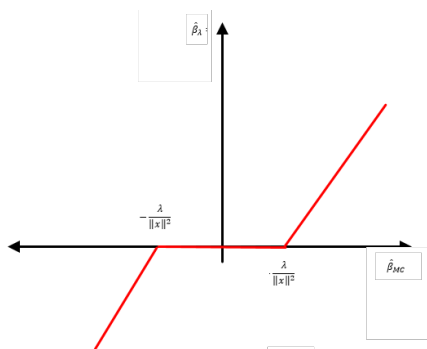
On considère $\operatorname{sign} \hat{\beta}_\lambda = \operatorname{sign} \hat{\beta}_{MC}$

Si $\hat{\beta}_\lambda = 0$:

$$\partial_\beta J_{MC} = \sum_{i=1}^n (x_i \hat{\beta}_\lambda - y_i) x_i + \lambda \alpha = 0 \quad \alpha \in [-1; 1]$$

$$\hat{\beta}_\lambda = \frac{x^\top y}{\|x\|^2} - \lambda \frac{\alpha}{\|x\|^2} = 0$$

$$-\frac{\lambda}{\|x\|^2} \leq \frac{\lambda \alpha}{\|x\|^2} \leq \frac{\lambda}{\|x\|^2} \Leftrightarrow -\frac{\lambda}{\|x\|^2} \leq \underbrace{\frac{x^\top y}{\|x\|^2}}_{\hat{\beta}_{MC}} \leq \frac{\lambda}{\|x\|^2}$$



Donc :

$$\hat{\beta}_\lambda = \text{sign } \hat{\beta}_{MC} \max\left(0, |\hat{\beta}_{MC}| - \frac{\lambda}{\|x\|^2}\right)$$

2. Algorithme component wise

a. Algorithme

Tant que (non convergé)

Choisir p

$$\beta(p) = \text{sign } \hat{\beta}_{MC} \max\left(0, |\hat{\beta}_{MC}| - \frac{\lambda}{\|x\|^2}\right)$$

Fin tant que

b. Correction du coût

Si on connaît tous les β sauf β_j

$$J = \frac{1}{2} \sum_{i=1}^n \left(x_{ij} \beta_j + \sum_{\substack{k=1 \\ k \neq j}}^p x_{ik} \beta_k - y_i \right)^2 + \lambda |\beta_j| + \sum_{\substack{k=1 \\ k \neq j}}^p |\beta_k|$$

$$\min_{\beta_j} \frac{1}{2} \sum_{i=1}^n (x_{ij} \beta_j + z_i)^2 + \lambda |\beta_j| \quad z_i = \sum_{\substack{k=1 \\ k \neq j}}^p x_{ik} \beta_k - y_i$$

3. Autres rétrécisseur

a. Hard threshold

$$\hat{\beta}_{H_j} = \begin{cases} 0 & \text{si } |\hat{\beta}_{MC_j}| < \frac{\lambda}{\|x\|^2} \\ \hat{\beta}_{MC_j} & \text{sinon} \end{cases}$$

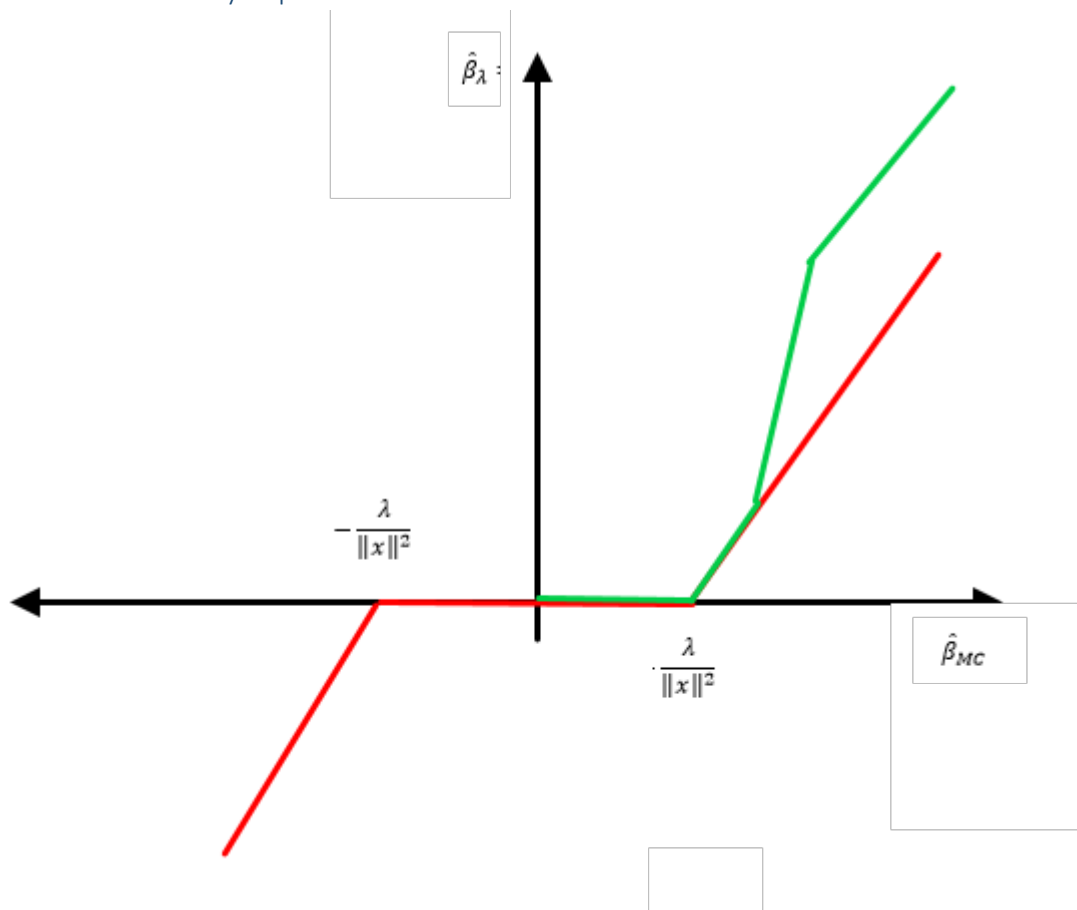
Débiaisé mais discontinu

b. Minimax Concave Penalty

$$\hat{\beta}_{H_j} = \begin{cases} 0 & \text{si } |\hat{\beta}_{MC_j}| < \frac{\lambda}{\|x\|^2} \\ \left(2|\hat{\beta}_{MC_j}| - \frac{2\lambda}{\|x\|^2}\right) \text{sign } \hat{\beta}_{MC_j} & \text{si } \frac{\lambda}{\|x\|^2} < |\hat{\beta}_{MC_j}| < \frac{2\lambda}{\|x\|^2} \\ \hat{\beta}_{MC_j} & \text{si } |\hat{\beta}_{MC_j}| > \frac{2\lambda}{\|x\|^2} \end{cases}$$

Débiaisé et continu

c. Smoothly Clipped Absolute Deviation



On peut les réécrire :

$$\min_{\beta_j} \frac{1}{2} \sum_{i=1}^n (x_i \beta - y_i)^2 + \lambda \text{pen } \beta$$

LASSO : $\text{pen } \beta = |\beta|$

$$\text{HARD : pen } \beta = \begin{cases} |\beta| & |\hat{\beta}_{MC}| < \lambda \\ cst & |\hat{\beta}_{MC}| \geq \lambda \end{cases}$$

$$\text{MCP : pen } \beta = \begin{cases} |\beta| & |\hat{\beta}_{MC}| < \lambda \\ a|\beta|^2 + b|\beta| + c & \lambda < |\beta| < 2\lambda \\ cst & |\hat{\beta}_{MC}| \geq 2\lambda \end{cases}$$

IV. Adaptive Lasso

1. Etude d'une loi binomiale (HS)

$$B \in \{-1; 1\} \sim \mathcal{B}\left(\frac{1}{2}\right) \quad \sum_{i=1}^n B_i \xrightarrow{n \rightarrow \infty} \pm \infty \quad \frac{1}{n} \sum_{i=1}^n B_i \xrightarrow{n \rightarrow \infty} 0 \quad \frac{1}{\sqrt{n}} \sum_{i=1}^n B_i \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, \sigma^2)$$

2. Modèle Lasso

$$y = X\beta^* + \varepsilon \text{ et } \mathcal{A}^* = \{j | \beta_j^* \neq 0\} \text{ (ensemble des variables actives)}$$

$$\text{avec } \varepsilon \sim \mathcal{N}(0, \Sigma)$$

3. Estimateur Lasso

$$\hat{\beta}_\lambda = \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \frac{1}{2} \|y - X\beta\|^2 + \lambda \|\beta\|_1 \quad \hat{\mathcal{A}}_\lambda = \{j | \hat{\beta}_{\lambda j} \neq 0\}$$

On veut que $(\hat{\beta}_\lambda, \hat{\mathcal{A}}_\lambda) \rightarrow (\beta^*, \mathcal{A}^*)$, c'est-à-dire :

$$\sqrt{n}(\hat{\beta}_\lambda - \beta^*) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, \Sigma) \quad \mathbb{P}(\hat{\mathcal{A}}_\lambda = \mathcal{A}^*) \xrightarrow{n \rightarrow \infty} 1$$

Théorème : $\exists c \in [0, 1[$ t. q. $\mathbb{P}(\hat{\mathcal{A}}_\lambda = \mathcal{A}^*) \leq c < 1$

On n'est donc pas sûr de converger vers la bonne valeur.

4. Problème Adaptive Lasso

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|y - X\beta\|^2 + \lambda \sum_{j=1}^p w_j |\beta_j|$$

$$w_j = \frac{1}{|\hat{\beta}_{MC j}|}$$

5. Problème Garote

$$\min_{c \in \mathbb{R}^p} \frac{1}{2} \|y - X \text{diag}(\hat{\beta}_{MC}) c\|^2 + \lambda \|c\|_1 \quad \hat{\beta}_{G_j} = \hat{\beta}_{MC_j} c_j$$

6. Etude Adaptive Lasso

a. Reformulation

$$\begin{cases} \min_{\beta \in \mathbb{R}^p} & \frac{1}{2} \|y - z\|^2 + \lambda \sum_{j=1}^p w_j |\beta_j| \\ \text{s.c.} & z = X\beta \end{cases}$$

b. Lagrangien

$$\mathcal{L}(\beta, z, \alpha) = \frac{1}{2} \|y - z\|^2 + \lambda \sum_{j=1}^p w_j |\beta_j| + \alpha^\top (z - X\beta)$$

$$\nabla_{\beta} \mathcal{L} = -X^\top \alpha + \lambda \gamma w \quad \gamma \in [-1; 1] \quad \nabla_z \mathcal{L} = z - y + \alpha \Rightarrow \underline{z = y - \alpha}$$

Dual

$$\begin{cases} \max_{\alpha} & -\frac{1}{2} \|\alpha\|^2 + \alpha^\top z \\ \text{s.c.} & |X^\top \alpha| \leq \lambda w \end{cases}$$

car $\gamma \in [-1; 1] \Rightarrow -\lambda w \leq \lambda \gamma w \leq \lambda w$ et $X^\top \alpha = \lambda \gamma w \Rightarrow -\lambda w \leq X^\top \alpha \leq \lambda w$

$$\Leftrightarrow \begin{cases} \min_{\alpha} & \frac{1}{2} \|\alpha\|^2 - \alpha^\top y \\ \text{s.c.} & |X^\top \alpha| \leq \lambda w \end{cases}$$

$$\alpha = y - z = y - X\beta$$

$$\Leftrightarrow \begin{cases} \min_{\beta} & \frac{1}{2} \|y - X\beta\|^2 - (y - X\beta)^\top y \\ \text{s.c.} & |X^\top (y - X\beta)| \leq \lambda w \end{cases} \Leftrightarrow \begin{cases} \min_{\beta} & \frac{1}{2} \|X\beta\|^2 \\ \text{s.c.} & |X^\top (y - X\beta)| \leq \lambda w \end{cases}$$

V. Elastic Net et gradient proximal

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|y - X\beta\|^2 + \lambda \|\beta\|_1 + \frac{\mu}{2} \|\beta\|^2$$

1. Régression Ridge

$$\begin{aligned} \min_{\beta \in \mathbb{R}^p} J_R &= \frac{1}{2} \|y - X\beta\|^2 + \frac{\mu}{2} \|\beta\|^2 = \frac{1}{2} (y - X\beta)^\top (y - X\beta) + \frac{\mu}{2} \beta^\top \beta \\ &= \frac{1}{2} y^\top y - y^\top X\beta + \frac{1}{2} \beta^\top X^\top X\beta + \frac{\mu}{2} \beta^\top \beta \\ \nabla_{\beta} J_R &= -X^\top y + (X^\top X + \mu I)\beta = 0 \\ \beta &= (X^\top X + \mu I)^{-1} X^\top y \end{aligned}$$

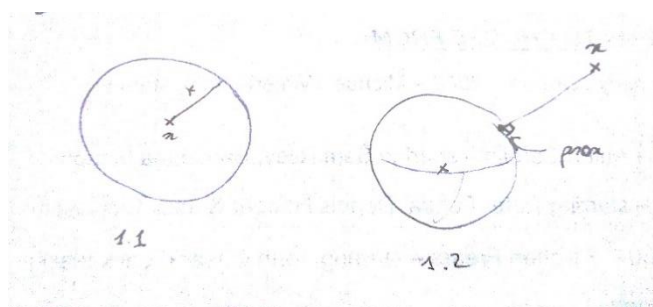
2. Opérateur proximal

a. Définition

$$\begin{array}{ccc} \Omega \text{ convexe} : \mathbb{R}^p & \rightarrow & \mathbb{R} \\ x & \rightarrow & \Omega(x) \end{array} \qquad \begin{array}{ccc} \text{prox } x : \mathbb{R}^p & \rightarrow & \mathbb{R}^p \\ \Omega & x & \rightarrow \text{prox } x \\ & & \Omega \end{array}$$

$$\text{prox}_{\Omega} x = \underset{u \in \mathbb{R}^p}{\text{argmin}} \left(\Omega(u) + \frac{1}{2} \|x - u\|^2 \right) = \underbrace{\left\{ \begin{array}{l} \text{argmin}_u \Omega(u) \\ \text{s.c. } \|x - u\|^2 \leq k \end{array} \right\}}_{\substack{\text{recherche dans une} \\ \text{région de confiance 1.1}}} \Leftrightarrow \underbrace{\left\{ \begin{array}{l} \text{argmin}_u \frac{1}{2} \|x - u\|^2 \\ \text{s.c. } \Omega(u) \leq k \end{array} \right\}}_{1.2}$$

$$\nabla_u J = \nabla_u \Omega(u) - x + u$$



b. Exemples

◆ $\Omega_1(x) = 0$

$$\text{prox}_{\Omega_1} x = u^* = \underset{u \in \mathbb{R}^p}{\text{argmin}} \left(0 + \frac{1}{2} \|x - u\|^2 \right)$$

$$\begin{aligned} \nabla_u J_1 &= -x + u^* = 0 \\ \Rightarrow \text{prox } x &= x \\ &\Omega_1 \end{aligned}$$

◆ $\Omega_2(x) = \frac{\mu}{2} \|x\|^2$

$$\text{prox}_{\Omega_2} x = u^* = \underset{u \in \mathbb{R}^p}{\text{argmin}} \left(\frac{\mu}{2} \|u\|^2 + \frac{1}{2} \|x - u\|^2 \right)$$

$$\begin{aligned} \nabla_u J_2 &= \mu u^* - x + u^* = 0 \\ \Rightarrow \text{prox } x &= \frac{1}{1 + \mu} x \\ &\Omega_2 \end{aligned}$$

◆ $\Omega_3(x) = \lambda \|x\|_1$

$$\text{prox}_{\Omega_3} x = u^* = \underset{u \in \mathbb{R}^p}{\text{argmin}} \left(\lambda \|u\|_1 + \frac{1}{2} \|x - u\|^2 \right)$$

$$\partial_u J_3 = \begin{cases} -\lambda & \text{si } u < 0 \\ \alpha \lambda & \text{si } u = 0 \\ \lambda & \text{si } u > 0 \end{cases} + u^* - x = 0 \quad \alpha \in [-1, 1]$$

$$\Leftrightarrow \begin{cases} -\lambda + u - x = 0 & \text{si } u < 0 & \Rightarrow u = x + \lambda \\ \alpha \lambda + u - x = 0 & \text{si } u = 0 & \Rightarrow x \leq |\lambda| \\ \lambda + u - x = 0 & \text{si } u > 0 & \Rightarrow u = x - \lambda \end{cases}$$

$$\Rightarrow \text{prox}_{\Omega_3} x = \begin{cases} x - \lambda & \text{si } x > \lambda \\ 0 & \text{si } x \leq |\lambda| \\ x + \lambda & \text{si } x < -\lambda \end{cases} \Leftrightarrow \text{sign } x \max(0, |x| - \lambda)$$

◆ $\Omega_4(x) = \begin{cases} 0 & \text{si } x \in C \\ \infty & \text{sinon} \end{cases}$

$$\text{prox}_{\Omega_4} x = \underset{u \in \mathbb{R}^p}{\text{argmin}} \left(1_{\{u \in C\}} + \frac{1}{2} \|x - u\|^2 \right) \quad \text{avec } 1_{\{u \in C\}} = \begin{cases} 0 & \text{si } x \in C \\ \infty & \text{sinon} \end{cases}$$

$$\text{prox}_{\Omega_4} x = \underset{u \in C}{\text{argmin}} \frac{1}{2} \|x - u\|^2 = \text{projection orthogonale de } x \text{ sur } C$$

3. Méthode du gradient proximal

a. Principe

On veut $\min_{x \in \mathbb{R}^p} C(x) + \lambda \Omega(x)$, C et Ω sont convexe. C est différentiable.

b. Algo général

On construit 2 suites $x^{(k)}$ et $\tilde{x}^{(k)}$

$$\tilde{x}^{(k+1)} = x^{(k)} - \rho^{(k)} \nabla_x C(x^{(k)})$$

$$x^{(k+1)} = \underset{\rho^{(k)} \lambda \Omega}{\text{prox}} \tilde{x}^{(k+1)} = \underset{\rho^{(k)} \lambda \Omega}{\text{prox}} \left(x^{(k)} - \rho^{(k)} \nabla_x C(x^{(k)}) \right)$$

c. Application à Elastic Net

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{2} \underbrace{\|y - X\beta\|^2}_{C(\beta)} + \frac{\mu}{2} \|\beta\|^2 + \lambda \underbrace{\|\beta\|_1}_{\Omega(\beta)}$$

$$\tilde{\beta}^{(k+1)} = \beta^{(k)} - \rho^{(k)} \left((X^T X + \mu I) \beta^{(k)} - X^T y \right)$$

$$\beta^{(k+1)} = \underset{\rho^{(k)} \lambda \Omega}{\text{prox}} \tilde{\beta}^{(k+1)}$$

$$\beta^{(k+1)} = \text{sign } \tilde{\beta}^{(k+1)} .* \max(0, |\tilde{\beta}^{(k+1)}| - \rho^{(k)} \lambda)$$

$$\underset{\rho^{(k)} \lambda \Omega}{\text{prox}} \left(x^{(k)} - \rho^{(k)} \nabla_x C(x^{(k)}) \right)$$

$$= \underset{u \in \mathbb{R}^p}{\text{argmin}} \left(\rho^{(k)} \lambda \Omega(u) + \frac{1}{2} \|x^{(k)} - \rho^{(k)} \nabla_x C(x^{(k)}) - u\|^2 \right)$$

$$= \underset{u \in \mathbb{R}^p}{\text{argmin}} \left(\lambda \Omega(u) + \frac{1}{2\rho^{(k)}} \|x^{(k)} - u\|^2 + \rho^{(k)2} \|\nabla_x C(x^{(k)})\|^2 - 2\rho^{(k)} (x^{(k)} - u)^T \nabla_x C(x^{(k)}) \right)$$

$$= \underset{u \in \mathbb{R}^p}{\text{argmin}} \left(\lambda \Omega(u) + \underbrace{C(x^{(k)}) + (u - x^{(k)})^T \nabla_x C(x^{(k)}) + \frac{1}{2\rho^{(k)}} \|x^{(k)} - u\|^2}_{\text{approximation locale de } C(u)} \right)$$

VI. Gradient proximal

1. Principe

On cherche $\min_w L(w) + \Omega(w)$

Hypothèses : L est convexe et différentiable, Ω est convexe et différentiable ou non.

Exemple :

$$\ell(w) = \frac{1}{2} \|y - Xw\|_2^2 \quad \Omega(w) = \frac{\lambda}{2} \|w\|_1$$

2. L'algorithme du gradient proximal est

$$w_{k+1} = \text{prox}_{\gamma\Omega} w_k - \gamma \nabla \ell(w_k)$$

3. Avec le prox

$$\text{prox}_{\Omega}(x) = \underset{u}{\operatorname{argmin}} \frac{1}{2} \|x - u\|_2^2 + \Omega(u)$$

a. Caractérisation du problème

On a $u^* = \text{prox}_{\Omega}(x)$. u^* est optimal si $0 \in \{\partial\Omega(u^*) - (x - u^*)\} \Leftrightarrow \boxed{(x - u^*) \in \partial\Omega(u^*)}$

b. Corollaire

Si $u^* = \text{prox}_{\Omega}(u^*)$ alors u^* minimise $\Omega(\cdot)$

c. Algorithme du point fixe

Si on pose $u_{k+1} = \text{prox}_{\Omega}(u_k)$, cet algorithme converge vers le min de $\text{prox}_{\Omega} \forall \Omega$ car l'opérateur proximal est contractant.

4. Théorème

$$w^* = \min_w \ell(w) + \Omega(w) \Leftrightarrow w^* = \text{prox}_{\nu\Omega}(w^* - \nu \nabla \ell(w^*)) \quad \nu > 0$$

a. Preuve

On veut minimiser $\ell(w) + \Omega(w)$. w^* est minimum si :

$$\begin{aligned} 0 &\in \{\nabla \ell(w^*) + \partial\Omega(w^*)\} \\ \Leftrightarrow -\nabla \ell(w^*) &\in \partial\Omega(w^*) \\ \Leftrightarrow -\nu \nabla \ell(w^*) &\in \partial\nu\Omega(w^*) \\ \Leftrightarrow w^* - \nu \nabla \ell(w^*) - w^* &\in \partial\nu\Omega(w^*) \end{aligned}$$

Or w^* minimise $\Omega(\cdot) \Leftrightarrow (x - w^*) \in \partial\Omega(w^*)$

On a donc $w^* = \text{prox}_{\nu\Omega}(w^* - \nu \nabla \ell(w^*))$

5. Approximation quadratique du coût

L'approx quadratique est : $Q(v, w) = \ell(w) + \nabla \ell(w)^T (v - w) + \frac{L}{2} \|v - w\|^2 + \Omega(v)$

Soit à w_k fixé :

$$Q(v, w_k) = \ell(w_k) + \nabla \ell(w_k)^T (v - w_k) + \frac{L}{2} \|v - w_k\|^2 + \Omega(v)$$

a. Propriété

$$v^* = \underset{v}{\operatorname{argmin}} Q(v, w_k)$$

On a $v^* = \operatorname{prox}_{\frac{1}{L}\Omega} \left(w_k - \frac{1}{L} \nabla \ell(w_k) \right)$

A chaque itération du proximal, on réalise en fait une approximation quadratique du coût que l'on minimise.

◆ Preuve

On cherche

$$\begin{aligned} \operatorname{argmin} Q &= \operatorname{argmin} \ell(w_k) + \nabla \ell(w_k)^\top v - (\nabla \ell(w_k)^\top w_k) + \frac{L}{2} v^\top v - \frac{L}{2} v^\top w_k + \frac{L}{2} w_k^\top w_k + \Omega(v) \\ &= \operatorname{argmin} \nabla \ell(w_k)^\top v + \frac{L}{2} v^\top v - \frac{L}{2} v^\top w_k + \Omega(v) + \ell \\ &= \operatorname{argmin} L \left[\frac{1}{2} v^\top v - v^\top \left(w_k - \frac{1}{2} \nabla \ell(w_k) \right) \right] + \Omega(v) + \ell \\ &= \operatorname{argmin} \frac{L}{2} \left\| v - \left(w_k - \frac{1}{2} \nabla \ell(w_k) \right) \right\|^2 + \Omega(v) + \ell \quad (*) \\ &= \operatorname{argmin} \frac{1}{2} \left\| v - \left(w_k - \frac{1}{L} \nabla \ell(w_k) \right) \right\|^2 + \frac{1}{L} \Omega(v) \end{aligned}$$

(*) On rajoute un terme constant en v pour pouvoir ajouter la norme.

b. Propriété

$$w_{k+1} = \operatorname{prox}_{\frac{1}{L_k}\Omega} \left(w_k - \frac{1}{L_k} \nabla \ell(w_k) \right)$$

Si $\forall k, L_k$ est t.q. $\forall v \ell(v) + \Omega(v) \leq Q(v, w_k)$

Alors on a : $\ell(w_{k+1}) + \Omega(w_{k+1}) \leq \ell(w_k) + \Omega(w_k)$

◆ Preuve

$$\ell(w_k) + \Omega(w_k) = Q(w_k, w_k) \geq Q(w_{k+1}, w_k) \geq \ell(w_{k+1}) + \Omega(w_{k+1})$$

Car $w_{k+1} = \min Q(v, w_k)$

c. Majorant de ℓ

Soit $\ell(v)$ une fonction différentiable telle que $\exists L_f \quad \forall v, w \|\nabla \ell(w) - \nabla \ell(v)\| \leq L_f \|w - v\|$

(Càd que la fonction est régulière)

Alors $\forall L \geq L_f \quad \ell(v) \leq \underbrace{\ell(w_k) + \nabla \ell(w_k)^\top (v - w_k) + \frac{L}{2} \|v - w_k\|^2}_{Q_\ell}$

Q_ℓ est majorante de ℓ .

d. Majorant de la Hessiane

Si $\ell(v)$ est doublement différentiable et $\forall w_k, \exists L_f, \|H(w_k)\| \leq L_f$

Alors $\|\nabla\ell(w) - \nabla\ell(v)\| \leq L_f \|w - v\|$

6. TP

$$l = \frac{1}{2} (1 - y(Xw + w_0))_+^T (1 - y(Xw + w_0))_+$$

$$\Omega = \|w\|_1$$

$$(\cdot)_+ = \max(0, \cdot)$$

$$\begin{aligned} \nabla_w l &= -(YX)^T (1 - Y(Xw + w_0))_+ \\ \nabla_{w_0} l &= -y(1 - Y(Xw + w_0))_+ \end{aligned}$$

VII. Factorisation non négative (Non negative matrix factorization)

1. Introduction

Soit $X \in \mathbb{R}^{n \times p}$, on cherche $\hat{X} = UV^T$ $U \in \mathbb{R}^{n \times k}$, $V \in \mathbb{R}^{k \times p}$ tel que $X \approx \hat{X}$

$$\hat{x}_i = VU_i^T$$

On a donc V dictionnaire d'information et U représente l'information pour chaque point, la contribution de chaque composante du dictionnaire à l'observation.

On peut alors contraindre U , par exemple $U \geq 0$.

2. Cas de la régression non négative

a. Le problème

$$\begin{aligned} \min_{\beta \in \mathbb{R}^k} \quad & \frac{1}{2} \|X\beta - y\|^2 \\ \text{s.c.} \quad & \beta \geq 0 \end{aligned}$$

b. Problème dual

$$\mathcal{L} = \frac{1}{2} \|X\beta - y\|^2 - \gamma^T \beta = \frac{1}{2} \beta^T X^T X \beta + \beta^T (X^T X \beta - X^T y - \gamma) + y^T y$$

$$\nabla_{\beta} \mathcal{L} = X^T (X\beta - y) - \gamma = 0$$

$$\mathcal{L} = \frac{1}{2} \beta^T X^T X \beta + \beta^T \underbrace{(X^T X \beta - X^T y - \gamma)}_{=0} + y^T y = \frac{1}{2} \beta^T X^T X \beta + y^T y$$

$$\begin{aligned} \min_{\beta \in \mathbb{R}^k} \quad & \frac{1}{2} \gamma^T (X^T X)^{-1} \gamma + y^T X (X^T X)^{-1} \gamma \\ \text{s.c.} \quad & \gamma \geq 0 \end{aligned}$$

Aussi difficile que le primal, voire plus car matrice à inverser.

c. Solution

- Problème dual (lagrangien)
- CVX
- MonQP
- Proximal

◆ *Proximal*

$$g = X^T(X\beta - y)$$

$$\hat{\beta} = \beta - \rho g$$

$$\beta = \max(0, \hat{\beta})$$

3. Cas matriciel

$$\min_U \|UV - Y\|_F^2 = \sum_{i=1}^n \sum_{j=1}^p (y_{ij} - u_{i\cdot} v_{j\cdot}^T)^2 = \sum_{i=1}^n \left\| \underbrace{y_{i\cdot}^T}_y - \underbrace{v_{\cdot}^T}_X \underbrace{u_{i\cdot}^T}_\beta \right\|^2$$

s.c. $\beta \geq 0$

Algorithme des moindres carrés alternés :

Connaissant V , $u_{i\cdot}$ est la solution de $\min_{\beta \in \mathbb{R}^k} \frac{1}{2} \|X\beta - y\|^2$ avec $X = V$, $y = Y(i, :)^T$
 s.c. $\beta \geq 0$

Connaissant U , $v_{j\cdot}^T = (U^T U)^{-1} U^T y_{\cdot j}$

$$v_{j\cdot}^T = v_{j\cdot}^T - \rho (U^T (U v_{j\cdot} - y_{\cdot j}))$$

4. Rappels de SVD

a. Rappel sur les valeurs propres

Soit $M \in \mathbb{R}^{n \times n}$, λ est valeur propre si $\underbrace{\det M - \lambda I}_{\substack{\text{polynome} \\ \text{d'ordre } n}} = 0$

- M a n valeurs propres complexes λ . Elles sont réelles positives si M est définie positive.
- v est le vecteur propre associé à λ tel que $Mv = \lambda v$

Si M est définie positive, $LL^T v = \lambda v \Rightarrow LU = \sqrt{\lambda} V \quad L^T V = \sqrt{\lambda} U$
 $\sqrt{\lambda}$ est la valeur singulière de L et U, V sont les vecteurs singuliers.

b. Reconstruction à l'ordre 1

Soit X une matrice $n \times p$, on cherche sa meilleure approximation à l'ordre 1.

$$\min_{\hat{X}} \|X - uv^T\|_F^2 \text{ avec } \text{rang } \hat{X} = 1, \|u\| = \|v\| = 1$$

$$\|X - uv^T\|_F^2 = \|X\|_F^2 - 2\langle X, uv^T \rangle + \|u\|^2 \|v\|^2$$

$$\nabla_u J = 2Xv - 2\|v\|^2 u$$

$$\nabla_v J = 2X^T u - 2\|u\|^2 v$$

La solution est :

$$Xv = \mu u$$

$$X^T u = \mu v$$

c. Théorème de décomposition

Soit X une matrice $n \times p$, (μ_k, u_k, v_k) la décomposition en valeurs singulières, on a $X = \sum_{k=1}^{\min(n,p)} \mu_k u_k v_k^T$

Le \hat{X} de rang K est $\hat{X} = \sum_{k=1}^K \mu_k u_k v_k^\top$ avec les K plus grand μ_k .

5. Séparation de sources

On veut approximer au mieux la matrice de données par le produit d'un dictionnaire D et d'une matrice d'activation A , soumis à des contraintes de régularisation tel que :

$$\min_{D,A} \|X - DA\|_F^2 + \Omega_D(D) + \Omega_A(A)$$

6. Factorisation non-négative

$$\begin{aligned} & \min_{D,A} \|X - DA\|_F^2 \\ & \text{sc. } d_{ij} \geq 0 \text{ et } a_{ij} \geq 0 \\ & \|d_{ij}\|^2 \leq 1 \end{aligned}$$

a. Algorithme

Optimisation alternée

$$\begin{aligned} A_{k+1} &= \operatorname{argmin}_A \frac{1}{2} \|X - D_k A\|_F^2 + \Omega_A(A) \\ D_{k+1} &= \operatorname{argmin}_D \frac{1}{2} \|X - D A_{k+1}\|_F^2 + \Omega_D(D) \end{aligned}$$

$t = 1$

$A_t = A_k$

$v = \frac{1}{\|D^\top D\|_*}$ (1/1^e v.p)

Répéter

$$\begin{aligned} \nabla_A J &= D_k^\top (X - D A_t) \\ A_{t+1} &= \operatorname{prox}_{v\Omega_A}(A_t - v \nabla_A J) \end{aligned}$$

Jusqu'à convergence

b. Proximaux

◆ 1^{er}

$$\begin{aligned} \operatorname{prox}(u) &= \operatorname{argmin}_x \frac{1}{2} \|x - u\|^2 = \begin{cases} 0 & \text{si } u_i \leq 0 \\ u_i & \text{sinon} \end{cases} \\ \text{s. c. } & x_i \geq 0 \end{aligned}$$

$$L = \frac{1}{2} \|x - u\|^2 - \sum \alpha_i x_i$$

$$\nabla_{x_i} L = x_i - u_i - \alpha_i = 0 \Rightarrow x_i = u_i + \alpha_i$$

a l'optimalité:

$$\alpha_i - x_i = 0$$

$$x_i \neq 0, \alpha_i = 0 \Rightarrow x_i = u_i$$

$$\alpha_i \neq 0, x_i = 0 \Rightarrow \alpha_i = -u_i$$

$$\alpha_i \geq 0 \text{ possible que si } u_i \leq 0$$

◆ 2e

$$\text{prox}_{\substack{\text{s.c.} \\ x \geq 0 \\ \|x\| \leq 1}} u = \begin{cases} 0 & \text{si } u_i < 0 \\ \frac{u_i}{\|u_+\|} & \text{si } u_i > 0 \text{ et } \|u_+\| \geq 1 \\ u_i & \text{si } u_i \geq 0 \text{ et } \|u_+\| \leq 1 \end{cases}$$

$$L = \frac{1}{2} \|x - u\|^2 - \sum \alpha_i x_i + \nu(\|x\|^2 - 1)$$

$$\nabla_{x_i} L = x_i - u_i - \alpha_i + 2\nu x_i = 0$$

$$\Leftrightarrow x_i = \frac{u_i + \alpha_i}{1 + 2\nu}$$

A l'optimalité

$$\alpha_i x_i = 0 \Rightarrow$$

$$x_i = \begin{cases} 0 & \text{si } u_i < 0 \\ \frac{u_i}{1 + 2\nu} & \text{sinon} \end{cases}$$

$$\nu(\|x\|^2 - 1) = 0 \Rightarrow$$

$$\|x\| = \frac{\|u_+\|}{1 + 2\nu} \Rightarrow \begin{cases} \nu \neq 0 \Rightarrow \|x\|^2 = 1 \Rightarrow 1 + 2\nu = \|u_+\| \\ \nu = 0 \Leftrightarrow \|x\|^2 < 1 \Leftrightarrow \|u_+\| < 1 \end{cases}$$

◆ 3e

$$L = \frac{1}{2} (x_i - u_i)^2 + \lambda \nu |x_i| - \alpha_i x_i$$

$$\nabla_{x_i} L = x_i - u_i + \lambda \nu g - \alpha_i = 0 \quad g = \begin{cases} 1 & \text{si } x_i > 0 \\ [0; 1] & \text{si } x_i = 0 \end{cases} \text{ sous gradient}$$

$$\alpha_i g \Rightarrow x^*$$

$$x_i > 0$$


$\text{prox}_{r(x \geq 0)}(u) = \underset{x}{\text{argmin}} \frac{1}{2} \|x - u\|_2^2$
 $= \underset{x}{\text{argmin}} \frac{1}{2} \|x - u\|_2^2$

$= \begin{cases} 0 & \text{si } u < 0 \\ u_i & \text{sinon} \end{cases} = \begin{cases} 0 & \text{si } u_i < 0 \\ \frac{u_i}{\|u\|} & \text{si } u_i \geq 0 \text{ et } \|u\| > 1 \\ u_i & \text{si } u_i \geq 0 \text{ et } \|u\| \leq 1 \end{cases}$

$\mathcal{L}(x, d_i) = \frac{1}{2} \|x - u\|_2^2 - \sum d_i x_i$

$\frac{\partial \mathcal{L}}{\partial x_i} = x_i - u_i - d_i = 0$
 $\Leftrightarrow x_i = u_i + d_i$

à l'optimalité $x_i \cdot x_i = 0$
 • $x_i \neq 0, d_i = 0 \Rightarrow x_i = u_i$
 • $d_i \neq 0 \Rightarrow x_i = 0$ donc $d_i = -u_i$
 Or $d_i \geq 0$ possible que si $u_i \leq 0$



$\text{prox}_{\{x \geq 0, \|x\| \leq 1\}}(u) = \underset{x}{\text{argmin}} \frac{1}{2} \|x - u\|_2^2$
 $\mathcal{L}(x, d_i, \nu) = \frac{1}{2} \|x - u\|_2^2 - \sum d_i x_i + \nu (\|x\|_2^2 - 1)$

$\frac{\partial \mathcal{L}}{\partial x_i} = 0 \Leftrightarrow x_i - u_i - d_i + 2\nu x_i = 0$
 $\Leftrightarrow x_i = \frac{u_i + d_i}{1 + 2\nu}$

$\tilde{x}_i = x_i(1 + 2\nu) = u_i + d_i$

à l'optimalité $d_i x_i = 0$ et $\nu (\|x\|_2^2 - 1) = 0$

• $x_i = \begin{cases} 0 & \text{si } u_i < 0 \\ \frac{u_i}{1 + 2\nu} & \text{sinon} \end{cases}$

$\nu \neq 0 \Rightarrow \|x\|_2^2 = 1 \Rightarrow 1 + 2\nu = \|u\|$
 $\Rightarrow \|x\|_2^2 < 1 \Rightarrow \nu = 0$
 possible si $\|u\| < 1$

VIII. Sélection de modèles : choix des hyper-paramètres

1. Exemple du Lasso

On se donne une grille de M valeurs ordonnées λ_m

Pour chaque λ_m on résout le Lasso et on obtient $\hat{\beta}_m$

Pour chaque $\hat{\beta}_m$ on évalue sa qualité \Rightarrow SURE

2. SURE : Stein Unbiased Risk Estimator

a. Risque de l'estimateur

◆ Erreur de prédiction

$$x_i^{new}, y_i^{new}$$

$$EP(\hat{\beta}_m) = \lim_{l \rightarrow \infty} \frac{1}{l} \sum_{i=1}^l (x_i^{new} \hat{\beta}_m - y_i^{new})^2 = \mathbb{E}(x_i^{new} \hat{\beta}_m - y_i^{new})^2 \text{ avec } (x_i^{new}, y_i^{new}) \sim \mathbb{P}(x, y)$$

◆ Hypothèses

- X est déterministe
- On pose un modèle $y = X\beta + \varepsilon$ $\varepsilon_i \sim \mathcal{N}(0, \sigma_i^2)$

$$EP_{\text{simpl.}}(\hat{\beta}_m) = \mathbb{E}_y(\|X\hat{\beta}_m - y\|) \text{ espérance sur la VA } y$$

$$\text{Risque}(\hat{\beta}) = \mathbb{E}_\beta(EP(\hat{\beta}_m)) \text{ espérance sur } \beta$$

Simplification de l'erreur de prédiction

COURS D'APPRENTISSAGE EN CONTEXTE

APPC – Cours

$$\begin{aligned}
 EP(\hat{\beta}_m) &= \mathbb{E}(\|X\hat{\beta}_m - y^{new}\|^2) = \mathbb{E}(\|X\hat{\beta}_m - X\beta + X\beta - y\|^2) \\
 &= \mathbb{E}(\|X\hat{\beta}_m - X\beta\|^2) + \underbrace{2 \mathbb{E}\left(\underbrace{(X\hat{\beta}_m - X\beta)^T}_{=0} \underbrace{(X\beta - y)}_{\varepsilon}\right)}_{=0} + \mathbb{E}\left(\left\|\underbrace{(X\beta - y)}_{\varepsilon}\right\|^2\right) \\
 &= \|X\hat{\beta}_m - X\beta^*\|^2 + n\sigma^2
 \end{aligned}$$

b. Réécriture

$$\begin{aligned}
 \|X\hat{\beta}_m - X\beta^*\|^2 &= \|X\hat{\beta}_m - y + y - X\beta^*\|^2 \quad y \text{ observation} \\
 &= \underbrace{\|X\hat{\beta}_m - y\|^2}_{\text{connu}} + 2 \underbrace{(X\hat{\beta}_m - y)^T}_{r^T} \underbrace{(y - X\beta^*)}_{\varepsilon} + \underbrace{\|y - X\beta^*\|^2}_{\text{indépendant de } m} \\
 &= \underbrace{\|X\hat{\beta}_m - y\|^2}_{\text{connu}} + 2(X\hat{\beta}_m)^T (y - X\beta^*) + cst(m)
 \end{aligned}$$

c. SURE

$$SURE(\hat{\beta}_m) := \|X\hat{\beta}_m - y\|^2 + 2 \operatorname{div}(X\hat{\beta}_m) \sigma^2 - n\sigma^2$$

Identité de Stein

$$\mathbb{E}(SURE(\hat{\beta}_m)) = \mathbb{E}(\|X\hat{\beta}_m - X\beta^*\|^2)$$

Preuve

$$\mathbb{E}(\|X\hat{\beta}_m - X\beta^*\|^2) = \mathbb{E}(\|X\hat{\beta}_m - y\|^2) + 2\mathbb{E}((X\hat{\beta}_m)^T \varepsilon) - \underbrace{2\mathbb{E}((X\beta^*)^T \varepsilon)}_0 - \underbrace{\mathbb{E}\|\varepsilon_i\|^2}_{\sigma^2}$$

Soit φ fct de \mathbb{R} dans \mathbb{R}

$$\mathbb{E}(\varepsilon\varphi(\varepsilon)) = \int_{\varepsilon} \varepsilon\varphi(\varepsilon) \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\varepsilon^2}{2\sigma^2}}\right) d\varepsilon$$

$$u = \varphi(\varepsilon) \quad u' = \varphi'(\varepsilon)$$

$$v' = \varepsilon e^{-\frac{\varepsilon^2}{2\sigma^2}} \quad v = \sigma^2 e^{-\frac{\varepsilon^2}{2\sigma^2}}$$

IPP :

$$\begin{aligned}
 &= \sigma^2 \int \varphi'(\varepsilon) \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\varepsilon^2}{2\sigma^2}} d\varepsilon \\
 &= \sigma^2 \mathbb{E}(\varphi'(\varepsilon))
 \end{aligned}$$

$$\mathbb{E}(\|X\hat{\beta}_m - X\beta^*\|^2) = \mathbb{E}(\|X\hat{\beta}_m - y\|^2) + 2\sigma^2 \mathbb{E}\left(\sum_{\operatorname{div}(X\hat{\beta}_m)} \frac{\partial X\hat{\beta}_m}{\partial \varepsilon_i}\right) - \underbrace{2\mathbb{E}((X\beta^*)^T \varepsilon)}_0 - \underbrace{\mathbb{E}\|\varepsilon_i\|^2}_{\sigma^2}$$

◆ Calcul de la divergence

$$\varphi: \mathbb{R}^n \rightarrow \mathbb{R}^n \quad \varphi(v) = Mv = w$$

COURS D'APPRENTISSAGE EN CONTEXTE

APPC – Cours

$$\text{div } \varphi(v) = \frac{\partial \varphi(v)}{\partial v_1} + \dots + \frac{\partial \varphi(v)}{\partial v_n} = M_{11} + M_{22} + \dots + M_{nn} = \text{trace } M$$

3. Résumé

$$\begin{aligned} \|X\hat{\beta}_m - X\beta^*\|^2 &= \|X\hat{\beta}_m - y\|^2 + 2(X\hat{\beta}_m)^\top \varepsilon - 2(X\beta^*)^\top \varepsilon - \|\varepsilon\|^2 \\ \text{SURE}(\hat{\beta}_m) &= \|X\hat{\beta}_m - y\|^2 + 2\sigma^2 \text{div}(X\hat{\beta}_m) - n\sigma^2 \\ \text{div}(Mv) &= \text{trace } M \end{aligned}$$

4. Exemple de div

Ex: le cas des moindres carrés

$$\hat{\beta}_{\text{OC}} = (X^\top X)^{-1} X^\top y$$

$$X \hat{\beta}_{\text{OC}} = X (X^\top X)^{-1} X^\top (X\beta^* + \varepsilon)$$

$$= X\beta^* + \underbrace{X (X^\top X)^{-1} X^\top}_M \varepsilon$$

$\text{Tr}(X(X^\top X)^{-1}X^\top) = \text{Tr}\left(\begin{matrix} (X^\top X)^{-1} & \\ & I_p \end{matrix}\right) = p$

$\text{div}(X \hat{\beta}_{\text{OC}}) = \text{Tr}\left(\underbrace{X (X^\top X)^{-1} X^\top}_M\right) = p$
(matrice de Projection)

pour le LASSO on a:

$$\text{div}(X \hat{\beta}_{\text{LASSO}}) = \text{card}\{\beta_i \neq 0\}$$

le nombre de composantes non nulles de $\hat{\beta}_{\text{LASSO}}$.

IX. Optimisation non-convexe & MCP

DC (Difference of Convex) / CCCP (Concave Convex Procedure) / LLA (Local Linear Approximation) / MM (MaxMin) / Iterative reweighted

= Relaxation convexe iterative

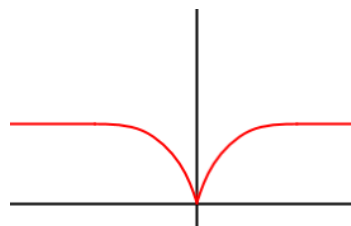
1. Problème du MCP (Minimum Concave Penalty)

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

$$\min_{\beta} \frac{1}{2} \|X\beta - y\|^2 + \lambda \sum_{j=1}^p \text{pen}(|\beta_j|)$$

Exemples de pénalités :

- $\text{pen}(t) = t \Rightarrow$ lasso
- $\text{pen}(t) = \frac{t}{w} \Rightarrow$ adaptative lasso
- $\text{pen}_{\lambda, \gamma}(t) = \begin{cases} \lambda t - \frac{t^2}{2\gamma} & \text{si } t < \gamma\lambda \\ \frac{\gamma\lambda^2}{2} & \text{sinon} \end{cases}$



Le MCP s'écrit aussi :

$$\min_{\beta} \frac{1}{2} \|X\beta - y\|^2 + \lambda \sum_{j=1}^p |\beta_j| - \sum_{j=1}^p h(|\beta_j|, \lambda, \gamma)$$

$$h(t, \lambda, \gamma) = \frac{t^2}{2\lambda\gamma} \mathbb{I}_{\{t \leq \gamma\lambda\}} + \left(t - \frac{\gamma\lambda}{2}\right) \mathbb{I}_{\{t > \gamma\lambda\}} = \begin{cases} \frac{t^2}{2\lambda\gamma} & \text{si } t \leq \gamma\lambda \\ t - \frac{\gamma\lambda}{2} & \text{sinon} \end{cases}$$

2. Algorithme

Soit $\min_x f(x) - h(x)$

Tant que on a pas convergé

$$x^{new} = \text{argmin}_x f(x) - \nabla_x h(x^{old})^\top x$$

Fin tant que

3. Application au MCP

$$h'(t) = \frac{t}{\gamma} \mathbb{I}_{\{t \leq \gamma\lambda\}} + \lambda \mathbb{I}_{\{t > \gamma\lambda\}}$$

$$\beta^{new} = \text{argmin}_x \frac{1}{2} \|X\beta - y\|^2 + \lambda \sum_{j=1}^p |\beta_j| - \sum_{j=1}^p \left(\frac{|\beta_j^{old}|}{\gamma} \mathbb{I}_{\{t \leq \gamma\lambda\}} + \lambda \mathbb{I}_{\{t > \gamma\lambda\}} \right) |\beta_j^{old}|$$

$$\Leftrightarrow \beta^{new} = \text{argmin}_x \frac{1}{2} \|X\beta - y\|^2 + \lambda \sum_{j=1}^p w_j |\beta_j|$$

$$w_j = \begin{cases} 1 - \frac{|\beta_j^{old}|}{\lambda\gamma} & \text{si } |\beta_j^{old}| < \gamma\lambda \\ 0 & \text{sinon} \end{cases}$$

$$\partial_{\beta} J = \underbrace{X^T(X\beta - y)}_g + \left\{ \begin{array}{ll} \alpha_j & \text{si } \beta_i = 0 \\ \lambda \text{sign}(\beta_j) - \frac{\beta_j}{\gamma} & \text{si } 0 < |\beta_j| < \gamma\lambda \\ 0 & |\beta_j| \geq \gamma\lambda \end{array} \right\}$$

$$0 \in \partial_{\beta} J \text{ si } \left\{ \begin{array}{ll} (g - \lambda)(g + \lambda) \leq 0 & \text{si } \beta_i = 0 \\ g + \lambda \text{sign}(\beta_j) - \frac{\beta_j}{\gamma} = 0 & \text{si } 0 < |\beta_j| < \gamma\lambda \\ g_i = 0 & |\beta_j| \geq \gamma\lambda \end{array} \right\}$$