



Apprentissage Statistique Automatique I

Méthodes de régularisation (RIDGE - LASSO)

Andrés F. López-Lopera
Université Polytechnique Hauts-de-France (UPHF)

1. Régression Ridge

Point de vue par SVD

Point de vue par pénalisation – Ridge / Tikhonov

Analyse du biais, de la variance et du risque

Validation Croisée (CV)

2. Régression Lasso

Sélection de variables et parcimonie

Améliorations et extensions du Lasso

- Dans un cadre d'apprentissage, on cherche à créer une règle de régression dans la classe

$$\mathcal{M}_\beta := \left\{ f(x) = \beta_0 + \sum_{j=1}^d \beta_j x_j, \beta = (\beta_0, \dots, \beta_d)^\top \in \mathbb{R}^{d+1} \right\} \quad (1)$$

avec $x = (x_1, \dots, x_d)^\top \in \mathbb{R}^d$

- On suppose qu'il existe β_* (inconnu) t.q.

$$y = X\beta_* + \varepsilon \quad (2)$$

- $y \in \mathbb{R}^n$ est le vecteur des observations
- $X \in \mathbb{R}^{n \times p}$ ($p = d + 1$) est la matrice des variables explicatives
- $\beta_* \in \mathbb{R}^p$ est les vrais paramètres du modèle que l'on veut retrouver
- $\varepsilon \in \mathbb{R}^n$ est le vecteur avec des bruits additifs

- La solution de moindres carrés (*Least Square, LS*) :

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

Limitations.

- $M = X^{\top} X$ doit être une matrice de plein rang, *i.e.*, $\text{rg}(X) = p$
- Si les plus petites valeurs propres de M s'approchent de zéro alors la solution numérique de M^{-1} n'est pas stable !

Régression Ridge

Theorem (Décomposition SVD [Golub and Van Loan, 2013])

Pour toute matrice $X \in \mathbb{R}^{n \times p}$, il existe deux matrices orthogonales $U = [u_1, \dots, u_n] \in \mathbb{R}^{n \times n}$ et $V = [v_1, \dots, v_p] \in \mathbb{R}^{p \times p}$, telles que

$$U^T X V = \Sigma \in \mathbb{R}^{n \times p},$$

avec $\Sigma = \text{diag}(s_1, \dots, s_r)$, $s_1 \geq \dots \geq s_r > 0$, et $r = \text{rang}(X)$.

$$X = U \Sigma V^T = [u_1, \dots, u_n] \left[\begin{array}{ccc|c} s_1 & \cdots & 0 & 0 \\ \vdots & \ddots & \vdots & \\ 0 & \cdots & s_r & 0 \\ \hline & & 0 & 0 \end{array} \right] \begin{bmatrix} v_1^T \\ \vdots \\ v_p^T \end{bmatrix} = \sum_{i=1}^r s_i u_i v_i^T$$

· Une solution des moindres carrés est alors :

$$\hat{\beta}^{\text{LS}} = (U \Sigma V^T)^{-1} y = (V \Sigma^{-1} U^T) y = \left(\sum_{i=1}^r \frac{1}{s_i} v_i u_i^T \right) y$$

Rem : Si s_i s'approche de zéro, la solution de la SVD n'est pas stable !

- Une solution des moindres carrés doit vérifier :

$$\begin{aligned}X^T X \beta &= X^T y \Leftrightarrow (U \Sigma V^T)^T (U \Sigma V^T) \beta = (U \Sigma V^T)^T y \\&\Leftrightarrow V \Sigma^T \Sigma V^T \beta = V \Sigma^T U^T y \\&\Leftrightarrow \Sigma^T \Sigma V^T \beta = \Sigma^T U^T y\end{aligned}$$

- $\Sigma^T \Sigma \in \mathbb{R}^{p \times p}$ est diagonale avec r éléments non nuls qui sont s_1^2, \dots, s_r^2

$$\Sigma^T \Sigma = \left[\begin{array}{ccc|c} s_1^2 & \cdots & 0 & 0 \\ \vdots & \ddots & \vdots & 0 \\ 0 & \cdots & s_r^2 & 0 \\ \hline & 0 & & 0 \end{array} \right]$$

- Il est possible de considérer une alternative régularisée
- Pour $\lambda > 0$ (paramètre de régularisation),

$$\Sigma^T \Sigma \rightarrow \Sigma^T \Sigma + \lambda I_p$$
$$\left[\begin{array}{ccc|c} s_1^2 & \cdots & 0 & 0 \\ \vdots & \ddots & \vdots & 0 \\ 0 & \cdots & s_r^2 & 0 \\ \hline & & 0 & 0 \end{array} \right] \rightarrow \left[\begin{array}{ccc|c} s_1^2 + \lambda & \cdots & 0 & 0 \\ \vdots & \ddots & \vdots & 0 \\ 0 & \cdots & s_r^2 + \lambda & 0 \\ \hline & & 0 & \lambda I_{p-r} \end{array} \right]$$

- Cela veut dire qu'on ajoute à toutes les valeurs propres de $X^T X$ un terme $\lambda > 0$ "petit" afin de rassurer la inversion de la matrice
- De cette manière, on a $(\Sigma^T \Sigma + \lambda I_p) V^T \beta = \Sigma^T U^T y$, et

$$\beta = V(\Sigma^T \Sigma + \lambda I_p)^{-1} \Sigma^T U^T y$$

- En sachant que $\Sigma = U^T X V$ (définition SVD), on trouve ensuite **[exercice]**

$$\begin{aligned}\hat{\beta}_{\lambda}^{\text{Ridge}} &= V(\Sigma^T \Sigma + \lambda I_p)^{-1} \Sigma^T U^T y \\ &= (X^T X + \lambda I_p)^{-1} X^T y\end{aligned}$$

Rappel : sous l'hypothèse de plein rang : $\hat{\beta}^{\text{LS}} = (X^T X)^{-1} X^T y$

Rem :

- $\lim_{\lambda \rightarrow 0^+} \hat{\beta}_{\lambda}^{\text{Ridge}} = \hat{\beta}^{\text{LS}}$
- $\lim_{\lambda \rightarrow +\infty} \hat{\beta}_{\lambda}^{\text{Ridge}} = \mathbf{0} \in \mathbb{R}^p$

- La prédiction associée s'obtient ainsi :

$$\hat{y} = X\hat{\beta}_{\lambda}^{\text{Ridge}} = \underbrace{X(X^{\top}X + \lambda I_p)^{-1}X^{\top}}_{H_{\lambda}} y$$

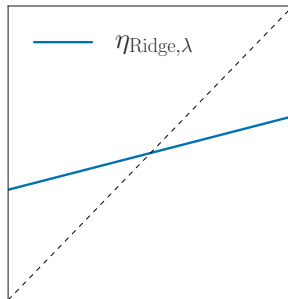
Rem :

- L'estimateur \hat{y} est linéaire en y
- L'équivalent de la matrice chapeau (*hat matrix*) est

$$H_{\lambda} = X(X^{\top}X + \lambda I_p)^{-1}X^{\top} = \sum_{j=1}^{\text{rg}(X)} \frac{s_j^2}{s_j^2 + \lambda} u_j u_j^{\top}$$

- Si X est une matrice orthogonale, i.e., $X^\top X = I_p$:

$$\begin{aligned}\hat{\beta}_{\lambda}^{\text{Ridge}} &= (X^\top X + \lambda I_p)^{-1} X^\top y \\ &= \frac{1}{1 + \lambda} X^\top y \\ \hat{y} &= \frac{1}{1 + \lambda} y \\ &= [\eta_{\text{rdg}, \lambda}(y_1), \dots, \eta_{\text{rdg}, \lambda}(y_n)]^\top\end{aligned}$$



Rem :

- Cas classique en traitement du signal (peu fréquent en statistique)
- La fonction réelle $\eta_{\text{rdg}, \lambda}$ est une contraction linéaire (*shrinkage*)

Astuce du noyau (*kernel trick*)

- En utilisant l'astuce du noyau, il est possible de démontrer **[exercice]**

$$X^T (XX^T + \lambda I_n)^{-1} y = (X^T X + \lambda I_p)^{-1} X^T y$$

Piste: utiliser l'identité matricielle de Woodbury

$$(Z + U W V^T)^{-1} = Z^{-1} - Z^{-1} U (W^{-1} + V^T Z^{-1} U)^{-1} V^T Z^{-1}$$

- On note que
 - \Leftarrow on doit résoudre un système $n \times n$
 - \Rightarrow on doit résoudre un système $p \times p$
- Selon si $n > p$ ou $n \leq p$, une méthode qui cherche à trouver une solution de Ridge par inversion peut préférer l'une des deux formulations

Rem : Cette astuce est aussi utile pour les méthodes à noyaux (e.g., SVM)

Démonstration.

Identité matricielle de Woodbury

$$(Z + U W V^T)^{-1} = Z^{-1} - Z^{-1} U (W^{-1} + V^T Z^{-1} U)^{-1} V^T Z^{-1}$$

- En utilisant l'identité matricielle de Woodbury, on a

$$(\lambda I_n + X X^T)^{-1} = \frac{1}{\lambda} I_n - \frac{1}{\lambda^2} X \left(I_p + \frac{1}{\lambda} X^T X \right)^{-1} X^T$$

- En reportant le résultat sur l'expression initial, on obtient

$$\begin{aligned} X^T (X X^T + \lambda I_n)^{-1} y &= \frac{1}{\lambda} X^T y - \frac{1}{\lambda^2} X^T X \left(I_p + \frac{1}{\lambda} X^T X \right)^{-1} X^T y \\ &= \left[\frac{1}{\lambda} I_p - \frac{1}{\lambda} X^T X \left(\lambda I_p + X^T X \right)^{-1} \frac{\lambda}{\lambda} \right] X^T y \\ &= \left(\lambda I_p + X^T X \right)^{-1} X^T y \end{aligned}$$

$$\hat{\beta}_{\lambda}^{\text{Ridge}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \underbrace{\|y - X\beta\|_2^2}_{\text{attache aux données}} + \underbrace{\lambda \|\beta\|_2^2}_{\text{régularisation}} \quad (3)$$

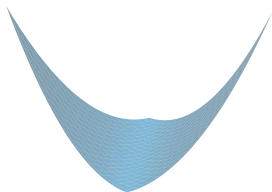
- L'estimateur Ridge est unique pour un λ fixé

Rem : On retrouve de nouveau les cas limites :

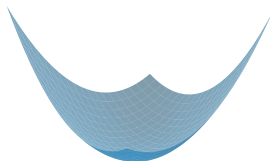
- $\lim_{\lambda \rightarrow 0^+} \hat{\beta}_{\lambda}^{\text{Ridge}} = \hat{\beta}^{\text{LS}}$ (solution de norme $\|\cdot\|_2$ minimale)
- $\lim_{\lambda \rightarrow +\infty} \hat{\beta}_{\lambda}^{\text{Ridge}} = \mathbf{0} \in \mathbb{R}^p$

- Lien avec les conditions nécessaires d'optimalité (CNO) :

$$f(\beta) = \frac{\|y - X\beta\|_2^2}{2} + \frac{\lambda \|\beta\|_2^2}{2}$$
$$\nabla f(\beta) = X^T(X\beta - y) + \lambda\beta = \mathbf{0} \quad \Leftrightarrow \quad (X^T X + \lambda I_p)\beta = X^T y$$

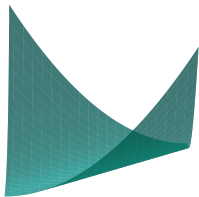


LS

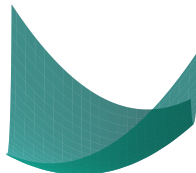


Ridge

Régulariser : simplifie le problème quand il est mal conditionné

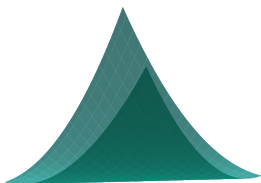


LS

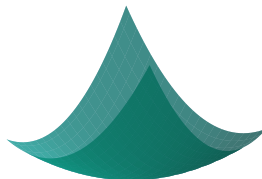


Ridge

Régulariser : simplifie le problème quand il est mal conditionné

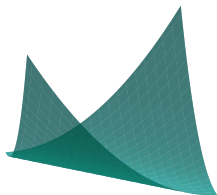


LS

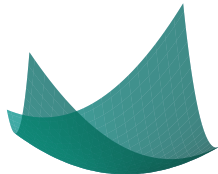


Ridge

Régulariser : simplifie le problème quand il est mal conditionné

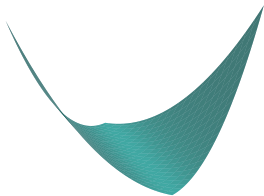


LS

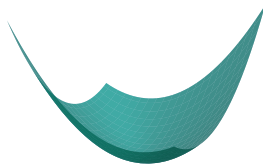


Ridge

Régulariser : simplifie le problème quand il est mal conditionné



LS



Ridge

Régulariser : simplifie le problème quand il est mal conditionné

- Un problème de la forme “Lagrangienne” suivante

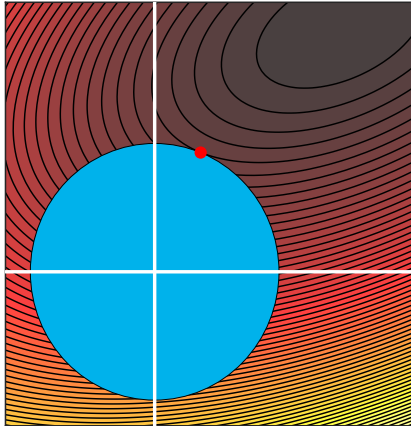
$$\operatorname{argmin}_{\beta \in \mathbb{R}^p} \underbrace{\frac{1}{2} \|y - X\beta\|_2^2}_{\text{attache aux données}} + \underbrace{\frac{\lambda}{2} \|\beta\|_2^2}_{\text{régularisation}} \quad (4)$$

admet pour un certain $T > 0$ la même solution que :

$$\begin{cases} \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 \\ \text{t.q. } \|\beta\|_2^2 \leq T \end{cases}$$

Rem : Le lien $T \leftrightarrow \lambda$ n'est pas explicite !

- Si $T \rightarrow 0^+$, on retrouve le vecteur nul $\hat{\beta}_\lambda^{\text{Ridge}} = 0 \in \mathbb{R}^p$ (i.e., $\lambda = +\infty$)
- Si $T \rightarrow +\infty$, on retrouve $\hat{\beta}_\lambda^{\text{Ridge}} = \hat{\beta}^{\text{LS}}$ (non contraint, i.e., $\lambda = 0^+$)



Optimisation sous contraintes ℓ_2

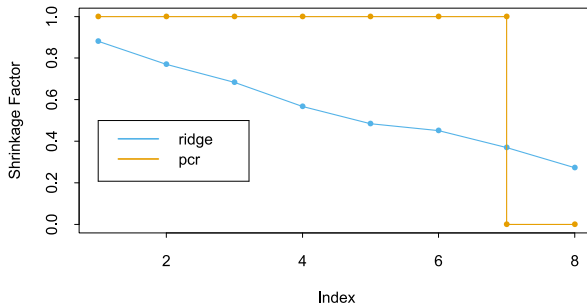


FIGURE 3.17. Ridge regression shrinks the regression coefficients of the principal components, using shrinkage factors $d_j^2/(d_j^2 + \lambda)$ as in (3.47). Principal component regression truncates them. Shown are the shrinkage and truncation patterns corresponding to Figure 3.7, as a function of the principal component index.

[Hastie et al., 2009]

Point normalisation et centrage

- Pour que la pénalisation contraigne de manière similaire toutes les variables, on peut normaliser les p variables
 - Centrer l'observation et les variables explicatives \Rightarrow pas de coefficient pour la variable constante (donc pas de contrainte)

$$\hat{\beta}_{\lambda}^{\text{Ridge}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \quad \frac{1}{2} \|y - X\beta - \beta_0 \mathbf{1}_n\|_2^2 + \lambda \sum_{j=1}^p \beta_j^2 \quad (5)$$

- Ne pas centrer les variables explicatives \Rightarrow ne pas mettre de contrainte sur la variable constante

Alternative. Si l'on ne normalise pas, on peut changer la pénalité en

$$\hat{\beta}_{\lambda}^{\text{Ridge}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \quad \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \sum_{j=1}^p \alpha_j \beta_j^2 \quad (\text{e.g., } \alpha_j = \|x_j\|^2) \quad (6)$$

Rem : Pour la validation croisée on utilisera plus naturellement

$\frac{1}{2n} \|y - X\beta\|_2^2$ que $\frac{1}{2} \|y - X\beta\|_2^2$ pour conserver l'amplitude de λ

Analyse du biais par la SVD

- Hypothèse : bruit “blanc” $\mathbf{y} = \mathbf{X}\beta^* + \varepsilon$ avec $\mathbb{E}(\varepsilon) = \mathbf{0}$:

$$\mathbb{E}(\hat{\beta}_{\lambda}^{\text{Ridge}}) = \mathbb{E} \left((\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^{\top} \mathbf{y} \right)$$

Analyse du biais par la SVD

- Hypothèse : bruit “blanc” $\mathbf{y} = \mathbf{X}\beta^* + \varepsilon$ avec $\mathbb{E}(\varepsilon) = \mathbf{0}$:

$$\begin{aligned}\mathbb{E}(\hat{\beta}_{\lambda}^{\text{Ridge}}) &= \mathbb{E}\left((\mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}^{\top}\mathbf{y}\right) \\ &= \mathbb{E}\left((\mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}^{\top}[\mathbf{X}\beta^* + \varepsilon]\right) \\ &= (\mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}^{\top}\mathbf{X}\beta^* \\ &= \sum_{i=1}^{\text{rg}(\mathbf{X})} \frac{s_i^2}{s_i^2 + \lambda} \mathbf{v}_i \mathbf{v}_i^{\top} \beta^*\end{aligned}$$

Rem :

- On retrouve $\mathbb{E}(\hat{\beta}^{\text{LS}}) = \sum_{i=1}^{\text{rg}(\mathbf{X})} \mathbf{v}_i \mathbf{v}_i^{\top} \beta^*$ quand $\lambda \rightarrow 0$
- Le biais vaut $\mathbf{b} = -\lambda(\mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\beta^*$:

$$\mathbb{E}(\hat{\beta}_{\lambda}^{\text{Ridge}}) - \mathbf{b} = (\mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}^{\top}\mathbf{X}\beta^* + \lambda(\mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\beta^* = \beta^*$$

Analyse de la variance par la SVD

- Hypothèse : bruit blanc et homoscedastique (i.e., $\mathbb{E}(\varepsilon) = \mathbf{0}$, $\mathbb{E}(\varepsilon\varepsilon^\top) = \sigma^2 I_n$)

$$V_\lambda^{\text{Ridge}} = \text{cov}(\hat{\beta}_\lambda^{\text{Ridge}}) = \text{cov}((X^\top X + \lambda I_p)^{-1} X^\top [X\beta^* + \varepsilon])$$

Analyse de la variance par la SVD

- Hypothèse : bruit blanc et homoscedastique (i.e., $\mathbb{E}(\varepsilon) = 0$, $\mathbb{E}(\varepsilon\varepsilon^\top) = \sigma^2 I_n$)

$$\begin{aligned} V_\lambda^{\text{Ridge}} &= \text{cov}(\widehat{\beta}_\lambda^{\text{Ridge}}) = \text{cov}((X^\top X + \lambda I_p)^{-1} X^\top [X\beta^* + \varepsilon]) \\ &= \text{cov}((X^\top X + \lambda I_p)^{-1} X^\top \varepsilon) \\ &= (X^\top X + \lambda I_p)^{-1} X^\top \mathbb{E}(\varepsilon \varepsilon^\top) X (X^\top X + \lambda I_p)^{-1} \\ &= \sigma^2 (X^\top X + \lambda I_p)^{-1} X^\top X (X^\top X + \lambda I_p)^{-1} \\ &= \sum_{i=1}^{\text{rg}(X)} \frac{s_i^2 \sigma^2}{(s_i^2 + \lambda)^2} v_i v_i^\top \end{aligned}$$

Rem :

- On retrouve $V^{\text{LS}} = \sum_{i=1}^{\text{rg}(X)} \frac{\sigma^2}{s_i^2} v_i v_i^\top$ quand $\lambda \rightarrow 0$
- On retrouve une variance nulle quand $\lambda \rightarrow \infty$

Analyse du risque par la SVD

- Hypothèse : modèle homoscedastique (i.e., $\mathbb{E}(\varepsilon\varepsilon^\top) = \sigma^2 I_n$)

$$R_{pred}(y, \hat{y}) = \mathbb{E} \|X\beta^\star - X\hat{\beta}_\lambda^{\text{Ridge}}\|_2^2$$

Analyse du risque par la SVD

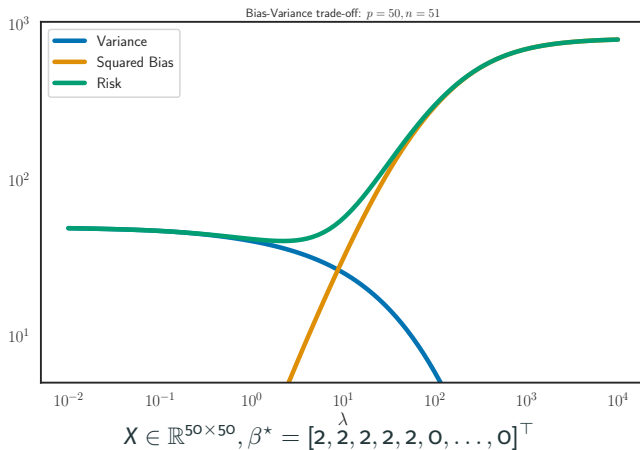
· Hypothèse : modèle homoscédastique (i.e., $\mathbb{E}(\varepsilon\varepsilon^\top) = \sigma^2 I_n$)

$$\begin{aligned} R_{pred}(y, \hat{y}) &= \mathbb{E} \|X\beta^* - X\hat{\beta}_\lambda^{\text{Ridge}}\|_2^2 \\ &= \mathbb{E} \left((\hat{\beta}_\lambda^{\text{Ridge}} - \beta^*)^\top X^\top X (\hat{\beta}_\lambda^{\text{Ridge}} - \beta^*) \right) \\ &= \mathbb{E} \left(((X^\top X + \lambda I_p)^{-1} X^\top \varepsilon)^\top X^\top X (X^\top X + \lambda I_p)^{-1} X^\top \varepsilon \right) \\ &\quad + \lambda^2 \beta^{*\top} (X^\top X) (X^\top X + \lambda I_p)^{-2} \beta^* \\ &= \sum_{i=1}^{\text{rg}(X)} \frac{s_i^4 \sigma^2}{(s_i^2 + \lambda)^2} + \lambda^2 \beta^{*\top} (X^\top X) (X^\top X + \lambda I_p)^{-2} \beta^* \end{aligned}$$

Rem :

- $\lim_{\lambda \rightarrow 0} R_{pred}(\beta^*, \hat{\beta}_\lambda^{\text{Ridge}}) = \text{rg}(X) \sigma^2$
- $\lim_{\lambda \rightarrow \infty} R_{pred}(\beta^*, \hat{\beta}_\lambda^{\text{Ridge}}) = \|X\beta^*\|_2^2$

Biais / Variance / Risque : exemple de simulation



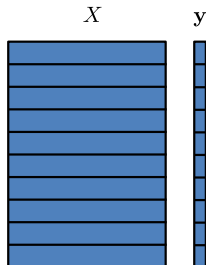
- 'svd' : méthode la plus stable, avantageuse pour calculer plusieurs λ car on ne "paye" la SVD qu'une fois
- 'cholesky' : décomposition matricielle proposant une formule fermée `scipy.linalg.solve`
- 'sparse_cg' : gradient conjugué utile dans les cas creux (*sparse*) et de grande dimension (baisser `tol` / `max_iter`)
- approche de type gradient stochastique si n est très grand

cf. le code des fonctions *Ridge*, *ridge_path*, *RidgeCV* dans le module `linear_model` de `sklearn`

Rem :

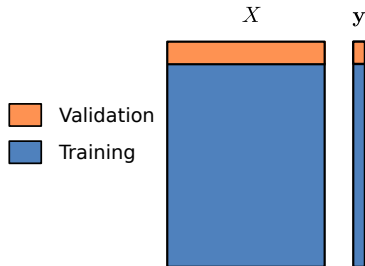
- On calcule rarement l'estimateur *Ridge* pour un λ , en général on en calcule plusieurs (10, 100, ...) et on cherche le meilleur
- Enjeu crucial de calculer des SVD de grandes tailles

- Choisir une grille de taille r de λ à tester : $\lambda_1, \dots, \lambda_r$
- Diviser (X, y) selon les observations en K blocs (*fold*) :



Validation croisée K-fold ($K = 10$)

- Choisir une grille de taille r de λ à tester : $\lambda_1, \dots, \lambda_r$
- Diviser (X, y) selon les observations en K blocs (*fold*) :



$k = 1$

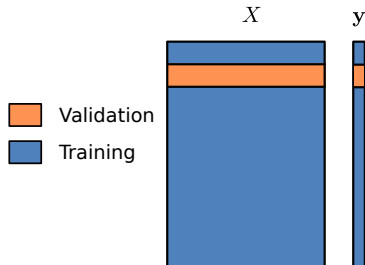
1. A partir des données d'entraînement, calculer les estimateurs pour $\lambda_1, \dots, \lambda_r$:

$$\hat{\beta}_{\lambda_1}, \dots, \hat{\beta}_{\lambda_r}$$

2. Calculer les erreurs de prédiction $\text{Erreur}_1^{(k)}, \dots, \text{Erreur}_r^{(k)}$ sur l'ensemble de validation

Validation croisée K-fold ($K = 10$)

- Choisir une grille de taille r de λ à tester : $\lambda_1, \dots, \lambda_r$
- Diviser (X, y) selon les observations en K blocs (*fold*) :



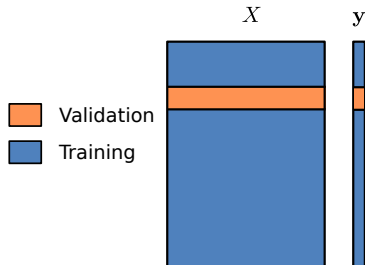
1. A partir des données d'entraînement, calculer les estimateurs pour $\lambda_1, \dots, \lambda_r$:

$$\hat{\beta}_{\lambda_1}, \dots, \hat{\beta}_{\lambda_r}$$

2. Calculer les erreurs de prédiction $\text{Erreur}_1^{(k)}, \dots, \text{Erreur}_r^{(k)}$ sur l'ensemble de validation

Validation croisée K -fold ($K = 10$)

- Choisir une grille de taille r de λ à tester : $\lambda_1, \dots, \lambda_r$
- Diviser (X, y) selon les observations en K blocs (*fold*) :



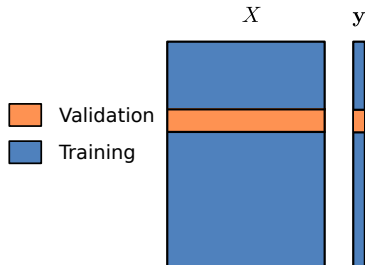
1. A partir des données d'entraînement, calculer les estimateurs pour $\lambda_1, \dots, \lambda_r$:

$$\hat{\beta}_{\lambda_1}, \dots, \hat{\beta}_{\lambda_r}$$

2. Calculer les erreurs de prédiction $\text{Erreur}_1^{(k)}, \dots, \text{Erreur}_r^{(k)}$ sur l'ensemble de validation

Validation croisée K-fold ($K = 10$)

- Choisir une grille de taille r de λ à tester : $\lambda_1, \dots, \lambda_r$
- Diviser (X, y) selon les observations en K blocs (*fold*) :



$k = 4$

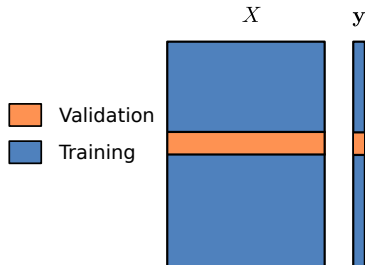
1. A partir des données d'entraînement, calculer les estimateurs pour $\lambda_1, \dots, \lambda_r$:

$$\hat{\beta}_{\lambda_1}, \dots, \hat{\beta}_{\lambda_r}$$

2. Calculer les erreurs de prédiction $\text{Erreur}_1^{(k)}, \dots, \text{Erreur}_r^{(k)}$ sur l'ensemble de validation

Validation croisée K-fold ($K = 10$)

- Choisir une grille de taille r de λ à tester : $\lambda_1, \dots, \lambda_r$
- Diviser (X, y) selon les observations en K blocs (*fold*) :



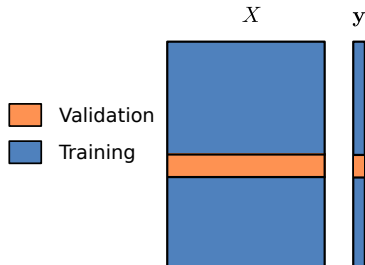
1. A partir des données d'entraînement, calculer les estimateurs pour $\lambda_1, \dots, \lambda_r$:

$$\hat{\beta}_{\lambda_1}, \dots, \hat{\beta}_{\lambda_r}$$

2. Calculer les erreurs de prédiction $\text{Erreur}_1^{(k)}, \dots, \text{Erreur}_r^{(k)}$ sur l'ensemble de validation

Validation croisée K-fold ($K = 10$)

- Choisir une grille de taille r de λ à tester : $\lambda_1, \dots, \lambda_r$
- Diviser (X, y) selon les observations en K blocs (*fold*) :



$k = 6$

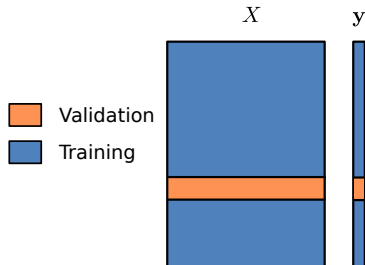
1. A partir des données d'entraînement, calculer les estimateurs pour $\lambda_1, \dots, \lambda_r$:

$$\hat{\beta}_{\lambda_1}, \dots, \hat{\beta}_{\lambda_r}$$

2. Calculer les erreurs de prédiction $\text{Erreur}_1^{(k)}, \dots, \text{Erreur}_r^{(k)}$ sur l'ensemble de validation

Validation croisée K -fold ($K = 10$)

- Choisir une grille de taille r de λ à tester : $\lambda_1, \dots, \lambda_r$
- Diviser (X, y) selon les observations en K blocs (*fold*) :



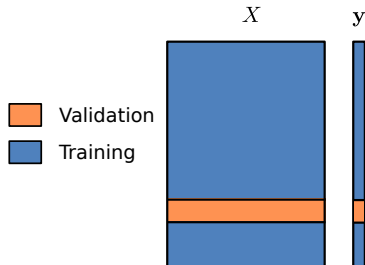
1. A partir des données d'entraînement, calculer les estimateurs pour $\lambda_1, \dots, \lambda_r$:

$$\hat{\beta}_{\lambda_1}, \dots, \hat{\beta}_{\lambda_r}$$

2. Calculer les erreurs de prédiction $\text{Erreur}_1^{(k)}, \dots, \text{Erreur}_r^{(k)}$ sur l'ensemble de validation

Validation croisée K -fold ($K = 10$)

- Choisir une grille de taille r de λ à tester : $\lambda_1, \dots, \lambda_r$
- Diviser (X, y) selon les observations en K blocs (*fold*) :



$k = 8$

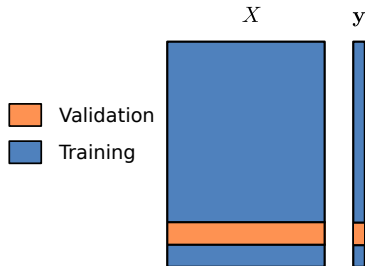
1. A partir des données d'entraînement, calculer les estimateurs pour $\lambda_1, \dots, \lambda_r$:

$$\hat{\beta}_{\lambda_1}, \dots, \hat{\beta}_{\lambda_r}$$

2. Calculer les erreurs de prédiction $\text{Erreur}_1^{(k)}, \dots, \text{Erreur}_r^{(k)}$ sur l'ensemble de validation

Validation croisée K-fold ($K = 10$)

- Choisir une grille de taille r de λ à tester : $\lambda_1, \dots, \lambda_r$
- Diviser (X, y) selon les observations en K blocs (*fold*) :



$k = 9$

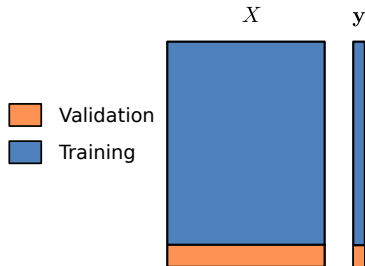
1. A partir des données d'entraînement, calculer les estimateurs pour $\lambda_1, \dots, \lambda_r$:

$$\hat{\beta}_{\lambda_1}, \dots, \hat{\beta}_{\lambda_r}$$

2. Calculer les erreurs de prédiction $\text{Erreur}_1^{(k)}, \dots, \text{Erreur}_r^{(k)}$ sur l'ensemble de validation

Validation croisée K -fold ($K = 10$)

- Choisir une grille de taille r de λ à tester : $\lambda_1, \dots, \lambda_r$
- Diviser (X, y) selon les observations en K blocs (*fold*) :



$k = 10$

1. A partir des données d'entraînement, calculer les estimateurs pour $\lambda_1, \dots, \lambda_r$:

$$\hat{\beta}_{\lambda_1}, \dots, \hat{\beta}_{\lambda_r}$$

2. Calculer les erreurs de prédiction $\text{Erreur}_1^{(k)}, \dots, \text{Erreur}_r^{(k)}$ sur l'ensemble de validation

Choix du paramètre. $\hat{\lambda}_{cv} = \lambda_{\hat{i}_{cv}}$ avec

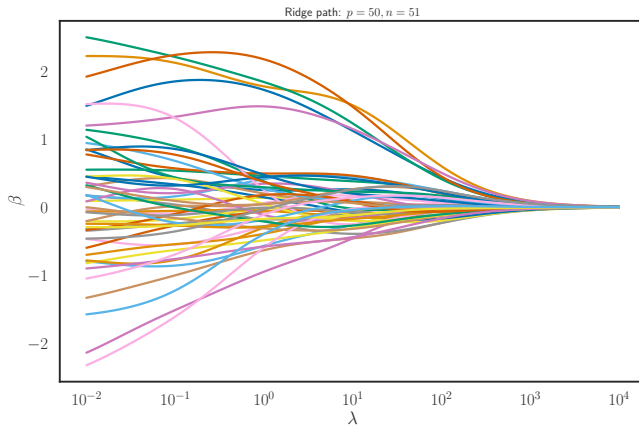
$$\hat{i}_{cv} = \underset{i \in \{1, \dots, r\}}{\operatorname{argmin}} \overline{\text{Erreur}}_i$$

Re-calibration. Calculer $\hat{\beta}_{\hat{\lambda}_{cv}}$ sur toutes les observations (X, y)

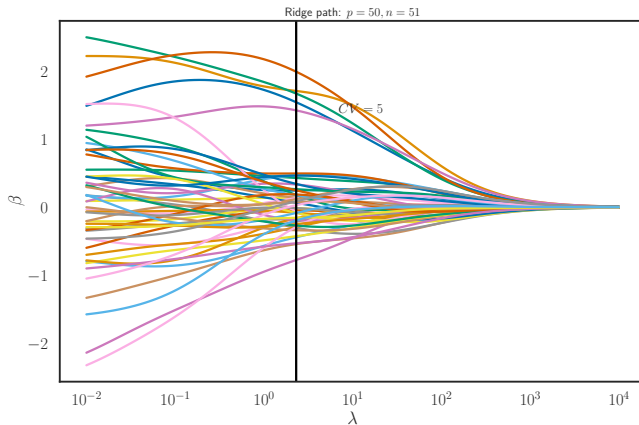
Conseils pratiques

- $K = n$ (stratégie “*leave-one-out*”) : calcul efficace pour Ridge mais assez instable
- Choix habituels : $K = 5, 10$
- “Randomiser les observations” : observations dans un ordre aléatoire afin éviter des blocs de données trop similaires (chaque sous-bloc doit être représentatif de l’ensemble)
- Prédiction : moyenner les meilleurs estimateurs obtenus plutôt que de re-calibrer sur toutes les données

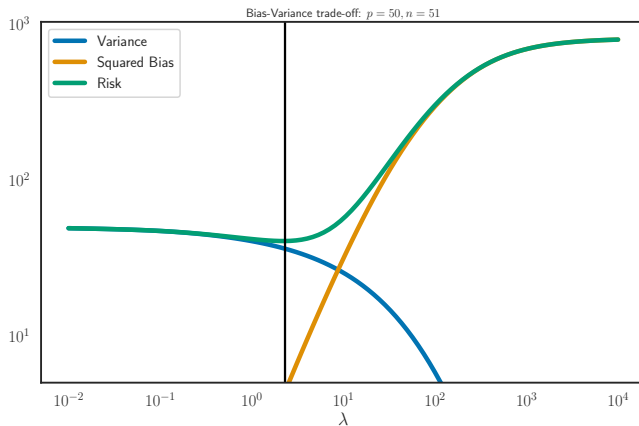
Choix de λ : exemple avec $CV = 5$



Choix de λ : exemple avec $CV = 5$



Choix de λ : exemple avec $CV = 5$



Régression Lasso

$$\hat{\beta}_{\lambda}^{\text{Ridge}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \underbrace{\|y - X\beta\|_2^2}_{\text{attache aux données}} + \underbrace{\lambda \|\beta\|_2^2}_{\text{régularisation}} \quad (7)$$

- Est-il possible de considérer d'autres normes pour la régularisation ?

- Le **support** du vecteur $\beta \in \mathbb{R}^p$ est l'ensemble des indices des coordonnées non nulles :

$$\text{supp}(\beta) = \{j \in \llbracket 1, p \rrbracket, \theta_j \neq 0\}$$

- La **pseudo-norme** ℓ_0 du vecteur β est son nombre de coordonnées non-nulles :

$$\|\beta\|_0 = \text{card}\{j \in \llbracket 1, p \rrbracket, \theta_j \neq 0\}$$

Rem :

- $\|\cdot\|_0$ n'est pas une norme, $\forall t \in \mathbb{R}^*, \|t\beta\|_0 = \|\beta\|_0$
- $\|\cdot\|_0$ n'est pas non plus convexe, $\beta_1 = (1, 0, 1, \dots, 0)$ $\beta_2 = (0, 1, 1, \dots, 0)$
et $3 = \left\| \frac{\beta_1 + \beta_2}{2} \right\|_0 \geq \frac{\|\beta_1\|_0 + \|\beta_2\|_0}{2} = 2$

- Première tentative de méthode pénalisée pour introduire de la parcimonie
- : utiliser ℓ_0 pour la pénalisation / régularisation

$$\hat{\beta}_{\lambda, \ell_0} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \underbrace{\frac{1}{2} \|y - X\beta\|_2^2}_{\text{attache aux données}} + \underbrace{\lambda \|\beta\|_0}_{\text{régularisation}} \quad (8)$$

Problème combinatoire !!! (problème “NP-dur”)

- La résolution exacte nécessite de considérer tous les sous-modèles, *i.e.*, calculer les estimateurs pour tous les supports possibles; il y en a 2^p , ce qui requiert le calcul de 2^p moindres carrés !

Exemple.

- $p = 10$ possible : $\approx 10^3$ moindres carrés
- $p = 30$ impossible : $\approx 10^{10}$ moindres carrés

Rem : avancées récentes en MIP [Bertsimas et al., 2016]

$$\hat{\beta}_{\lambda, \ell_1} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \underbrace{\frac{1}{2} \|y - X\beta\|_2^2}_{\text{attache aux données}} + \underbrace{\lambda \|\beta\|_1}_{\text{régularisation}} \quad (9)$$

où $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$

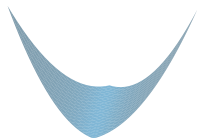
· On retrouve de nouveau les cas limites :

- $\lim_{\lambda \rightarrow 0^+} \hat{\beta}_{\lambda, \ell_1} = \hat{\beta}^{\text{LS}}$
- $\lim_{\lambda \rightarrow +\infty} \hat{\beta}_{\lambda, \ell_1} = \mathbf{0} \in \mathbb{R}^p$

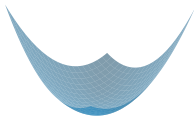
· $\hat{\beta}_{\lambda, \ell_1}$ est connu comme l'estimateur **Lasso** (*Least Absolute Shrinkage and Selection Operator*) [Tibshirani, 1996]

Rem : $\hat{\beta}_{\lambda}^{\text{Lasso}} := \hat{\beta}_{\lambda, \ell_1}$ n'est pas toujours **unique** pour un λ fixé ; prendre par exemple deux colonnes identiques

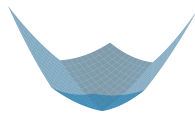
Illustration avec $p = 2$



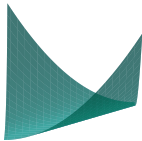
LS



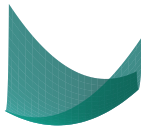
Ridge



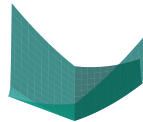
Lasso



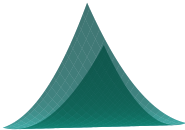
LS



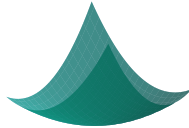
Ridge



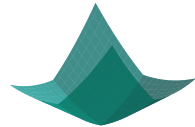
Lasso



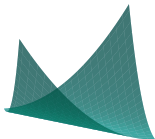
LS



Ridge



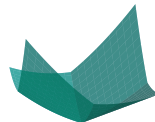
Lasso



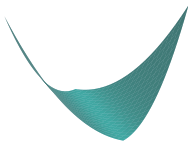
LS



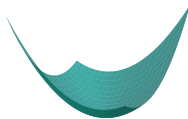
Ridge



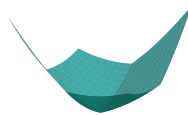
Lasso



LS



Ridge



Lasso

- Un problème de la forme “Lagrangienne” suivante

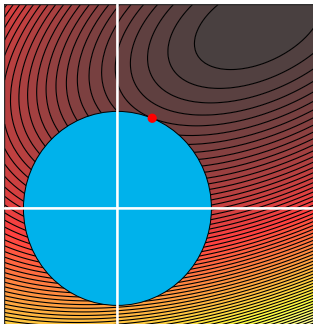
$$\hat{\beta}_{\lambda}^{\text{Lasso}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \underbrace{\frac{1}{2} \|y - X\beta\|_2^2}_{\text{attache aux données}} + \underbrace{\frac{\lambda}{2} \|\beta\|_1}_{\text{régularisation}} \quad (10)$$

admet pour un certain $T > 0$ la même solution que :

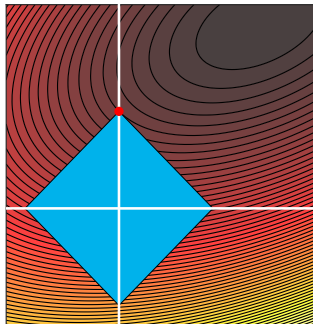
$$\begin{cases} \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \|y - X\beta\|_2^2 \\ \text{t.q. } \|\beta\|_1 \leq T \end{cases}$$

Rem : Le lien $T \leftrightarrow \lambda$ n'est pas explicite !

- Si $T \rightarrow 0^+$, on retrouve le vecteur nul $\hat{\beta}_{\lambda}^{\text{Lasso}} = 0 \in \mathbb{R}^p$ (i.e., $\lambda = +\infty$)
- Si $T \rightarrow +\infty$, on retrouve $\hat{\beta}_{\lambda}^{\text{Lasso}} = \hat{\beta}^{\text{LS}}$ (non contraint, i.e., $\lambda = 0^+$)



Optimisation sous contraintes ℓ_2



Optimisation sous contraintes ℓ_1

· Optimisation sous contrainte ℓ_1 : solution parcimonieuse !

Définitions

Pour $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction convexe, $u \in \mathbb{R}^d$ est un **sous-gradient** de f en x^* , si pour tout $x \in \mathbb{R}^d$ on a

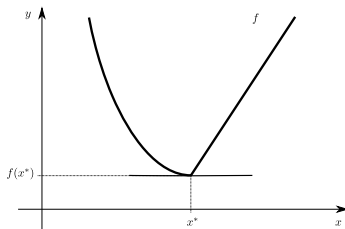
$$f(x) \geq f(x^*) + \langle u, x - x^* \rangle$$

La **sous-différentielle** est l'ensemble des sous-gradients :

$$\partial f(x^*) = \{u \in \mathbb{R}^d : \forall x \in \mathbb{R}^d, f(x) \geq f(x^*) + \langle u, x - x^* \rangle\}.$$

Rem : si le sous-gradient est unique, on retrouve le gradient

· Visuellement cela correspond à des droites tangentes sur x^*



Définitions

Pour $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction convexe, $u \in \mathbb{R}^d$ est un **sous-gradient** de f en x^* , si pour tout $x \in \mathbb{R}^d$ on a

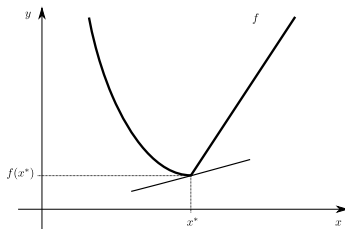
$$f(x) \geq f(x^*) + \langle u, x - x^* \rangle$$

La **sous-différentielle** est l'ensemble des sous-gradients :

$$\partial f(x^*) = \{u \in \mathbb{R}^d : \forall x \in \mathbb{R}^d, f(x) \geq f(x^*) + \langle u, x - x^* \rangle\}.$$

Rem : si le sous-gradient est unique, on retrouve le gradient

· Visuellement cela correspond à des droites tangentes sur x^*



Définitions

Pour $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction convexe, $u \in \mathbb{R}^d$ est un **sous-gradient** de f en x^* , si pour tout $x \in \mathbb{R}^d$ on a

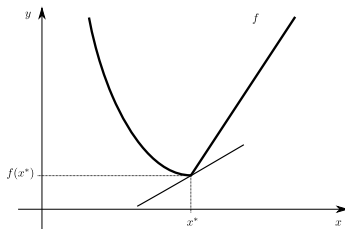
$$f(x) \geq f(x^*) + \langle u, x - x^* \rangle$$

La **sous-différentielle** est l'ensemble des sous-gradients :

$$\partial f(x^*) = \{u \in \mathbb{R}^d : \forall x \in \mathbb{R}^d, f(x) \geq f(x^*) + \langle u, x - x^* \rangle\}.$$

Rem : si le sous-gradient est unique, on retrouve le gradient

· Visuellement cela correspond à des droites tangentes sur x^*



Définitions

Pour $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction convexe, $u \in \mathbb{R}^d$ est un **sous-gradient** de f en x^* , si pour tout $x \in \mathbb{R}^d$ on a

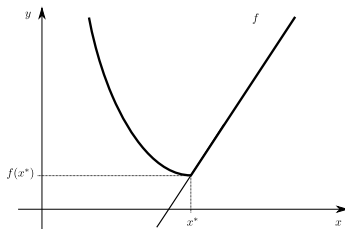
$$f(x) \geq f(x^*) + \langle u, x - x^* \rangle$$

La **sous-différentielle** est l'ensemble des sous-gradients :

$$\partial f(x^*) = \{u \in \mathbb{R}^d : \forall x \in \mathbb{R}^d, f(x) \geq f(x^*) + \langle u, x - x^* \rangle\}.$$

Rem : si le sous-gradient est unique, on retrouve le gradient

· Visuellement cela correspond à des droites tangentes sur x^*



Théorème

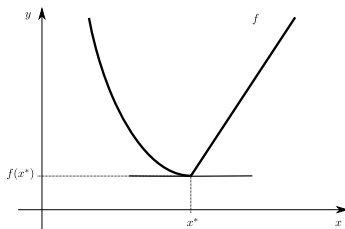
Un point x^* est un minimum d'une fonction convexe $f : \mathbb{R}^d \rightarrow \mathbb{R}$ si et seulement si $0 \in \partial f(x^*)$

Preuve : utiliser la définition des sous-gradients :

- 0 est un sous-gradient de f en x^* si et seulement si

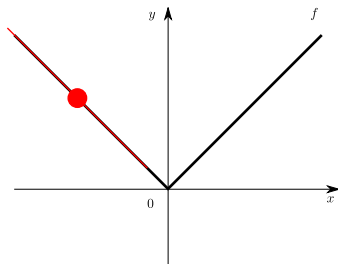
$$\forall x \in \mathbb{R}^d, f(x) \geq f(x^*) + \langle 0, x - x^* \rangle$$

- Visuellement cela correspond à une tangente horizontale



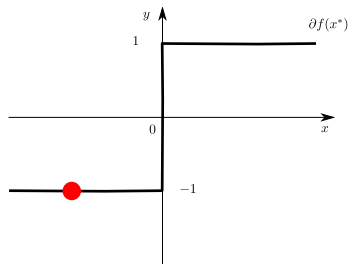
Fonction $\text{abs}(\cdot)$

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



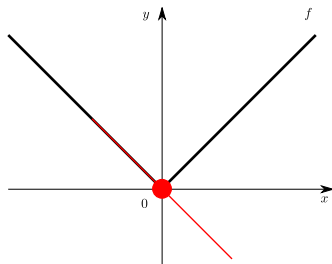
Sous-différentielle $\text{sign}(\cdot)$

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{si } x^* \in]-\infty, 0[\\ [-1, 1] & \text{si } x^* = 0 \\ \{1\} & \text{si } x^* \in]0, \infty[\end{cases}$$



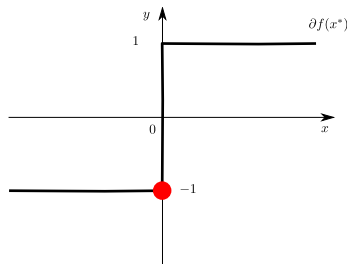
Fonction $\text{abs}(\cdot)$

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



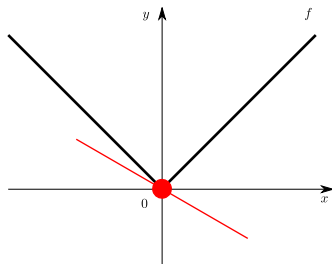
Sous-différentielle $\text{sign}(\cdot)$

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{si } x^* \in]-\infty, 0[\\ [-1, 1] & \text{si } x^* = 0 \\ \{1\} & \text{si } x^* \in]0, \infty[\end{cases}$$



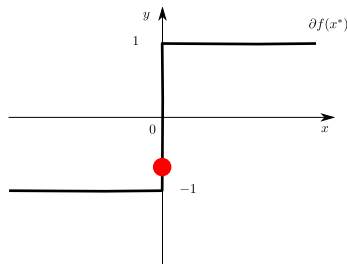
Fonction $\text{abs}(\cdot)$

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



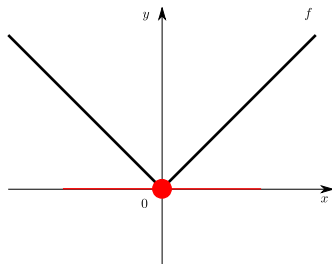
Sous-différentielle $\text{sign}(\cdot)$

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{si } x^* \in]-\infty, 0[\\ [-1, 1] & \text{si } x^* = 0 \\ \{1\} & \text{si } x^* \in]0, \infty[\end{cases}$$



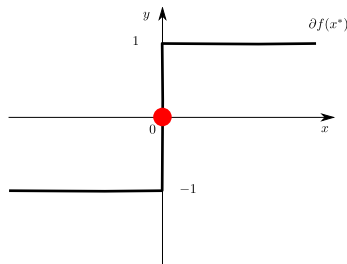
Fonction $\text{abs}(\cdot)$

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



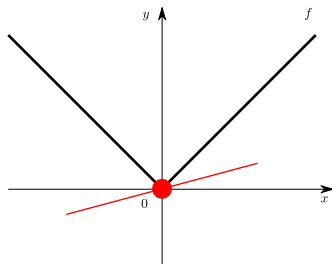
Sous-différentielle $\text{sign}(\cdot)$

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{si } x^* \in]-\infty, 0[\\ [-1, 1] & \text{si } x^* = 0 \\ \{1\} & \text{si } x^* \in]0, \infty[\end{cases}$$



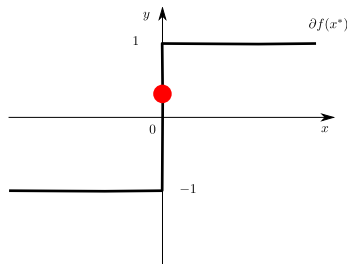
Fonction $\text{abs}(\cdot)$

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



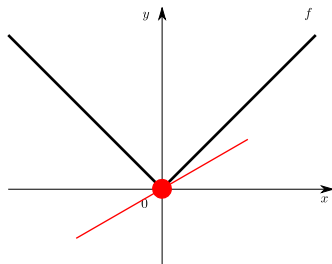
Sous-différentielle $\text{sign}(\cdot)$

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{si } x^* \in]-\infty, 0[\\ [-1, 1] & \text{si } x^* = 0 \\ \{1\} & \text{si } x^* \in]0, \infty[\end{cases}$$



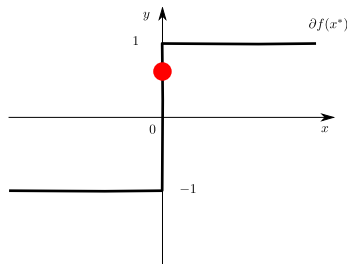
Fonction $\text{abs}(\cdot)$

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



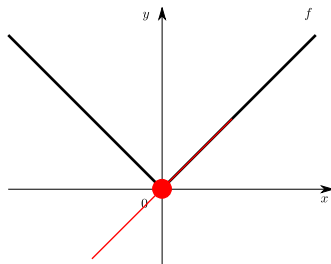
Sous-différentielle $\text{sign}(\cdot)$

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{si } x^* \in]-\infty, 0[\\ [-1, 1] & \text{si } x^* = 0 \\ \{1\} & \text{si } x^* \in]0, \infty[\end{cases}$$



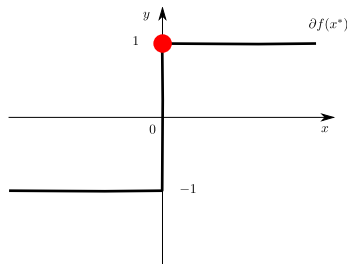
Fonction $\text{abs}(\cdot)$

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



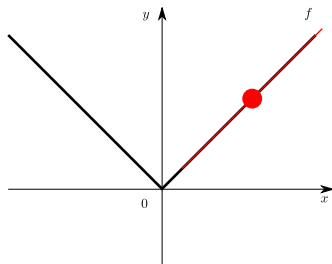
Sous-différentielle $\text{sign}(\cdot)$

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{si } x^* \in]-\infty, 0[\\ [-1, 1] & \text{si } x^* = 0 \\ \{1\} & \text{si } x^* \in]0, \infty[\end{cases}$$



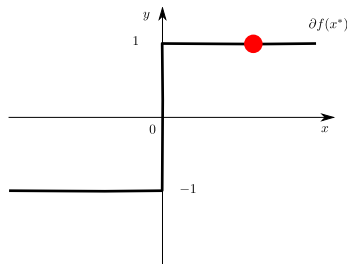
Fonction $\text{abs}(\cdot)$

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



Sous-différentielle $\text{sign}(\cdot)$

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{si } x^* \in]-\infty, 0[\\ [-1, 1] & \text{si } x^* = 0 \\ \{1\} & \text{si } x^* \in]0, \infty[\end{cases}$$



$$\hat{\beta}_{\lambda}^{\text{Lasso}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \underbrace{\frac{1}{2} \|y - X\beta\|_2^2}_{\text{attache aux données}} + \underbrace{\lambda \|\beta\|_1}_{\text{régularisation}}$$

· Conditions nécessaires et suffisantes d'optimalité (Fermat) :

$$\forall j \in \llbracket 1, p \rrbracket, x_j^\top \left(\frac{y - X\hat{\beta}_{\lambda}^{\text{Lasso}}}{\lambda} \right) \in \begin{cases} \{\operatorname{sign}(\hat{\beta}_{\lambda}^{\text{Lasso}})_j\} & \text{si } (\hat{\beta}_{\lambda}^{\text{Lasso}})_j \neq 0, \\ [-1, 1] & \text{si } (\hat{\beta}_{\lambda}^{\text{Lasso}})_j = 0. \end{cases}$$

$$\hat{\beta}_{\lambda}^{\text{Lasso}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \underbrace{\frac{1}{2} \|y - X\beta\|_2^2}_{\text{attache aux données}} + \underbrace{\lambda \|\beta\|_1}_{\text{régularisation}}$$

· Conditions nécessaires et suffisantes d'optimalité (Fermat) :

$$\forall j \in \llbracket 1, p \rrbracket, x_j^\top \left(\frac{y - X\hat{\beta}_{\lambda}^{\text{Lasso}}}{\lambda} \right) \in \begin{cases} \{\operatorname{sign}(\hat{\beta}_{\lambda}^{\text{Lasso}})_j\} & \text{si } (\hat{\beta}_{\lambda}^{\text{Lasso}})_j \neq 0, \\ [-1, 1] & \text{si } (\hat{\beta}_{\lambda}^{\text{Lasso}})_j = 0. \end{cases}$$

Démonstration.

$$\frac{\partial}{\partial \beta_j} \left(\frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right) = -x_j^\top (y - X\beta) + \lambda \frac{\partial}{\partial \beta_j} |\beta_j|$$

Si $\beta = \hat{\beta}_{\lambda}^{\text{Lasso}}$, alors $-x_j^\top (y - X\beta) + \lambda \frac{\partial}{\partial \beta_j} |\beta_j| = 0$, et on peut en déduire

$$\frac{\partial}{\partial \beta_j} |(\hat{\beta}_{\lambda}^{\text{Lasso}})_j| = x_j^\top \left(\frac{y - X\hat{\beta}_{\lambda}^{\text{Lasso}}}{\lambda} \right)$$

- Retour sur un cas simple (*design* orthogonal) : $X^\top X = I_p$

$$\|y - X\beta\|_2^2 = \|X^\top y - X^\top X\beta\|_2^2 = \|X^\top y - \beta\|_2^2$$

car X est une isométrie dans ce cas

- L'objectif du lasso devient :

$$\frac{1}{2}\|y - X\beta\|_2^2 + \lambda\|\beta\|_1 = \sum_{j=1}^p \frac{1}{2}(x_j^\top y - \beta_j)^2 + \lambda|\beta_j|$$

Problème séparable : problème qui revient à minimiser terme à terme en séparant les termes la somme. Il faut donc minimiser :

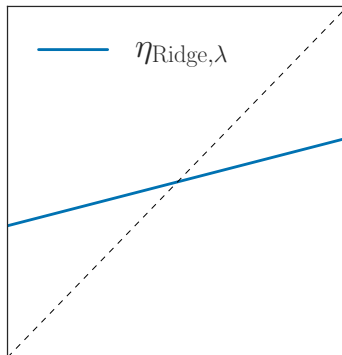
$$x \mapsto \frac{1}{2}(z - x)^2 + \lambda|x|,$$

pour $z = x_j^\top y$

Régularisation en 1D : ℓ_2 (Ridge)

Résoudre : $\eta_\lambda(z) = \operatorname{argmin}_{x \in \mathbb{R}} x \mapsto \frac{1}{2}(z - x)^2 + \frac{\lambda}{2}x^2$

$$\eta_\lambda(z) = \frac{z}{1 + \lambda}$$

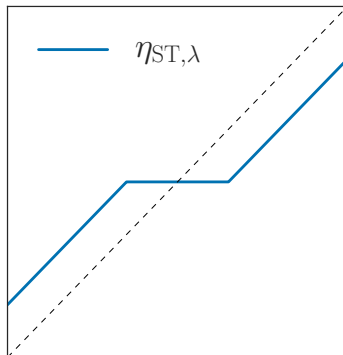


Contraction ℓ_2 : Ridge

Régularisation en 1D : ℓ_1 (Lasso)

Résoudre : $\eta_\lambda(z) = \operatorname{argmin}_{x \in \mathbb{R}} x \mapsto \frac{1}{2}(z - x)^2 + \lambda|x|$

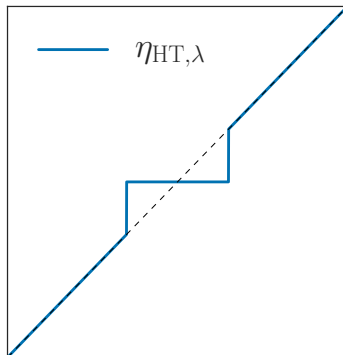
$$\eta_\lambda(z) = \operatorname{sign}(z)(|z| - \lambda)_+ \quad \textbf{[exercice]}$$



Contraction ℓ_1 : Seuillage doux (*soft thresholding*)

Résoudre : $\eta_\lambda(z) = \operatorname{argmin}_{x \in \mathbb{R}} x \mapsto \frac{1}{2}(z - x)^2 + \lambda \mathbb{1}_{x \neq 0}$

$$\eta_\lambda(z) = z \mathbb{1}_{|z| \geq \sqrt{2\lambda}}$$



Contraction ℓ_0 : Seuillage dur (*hard thresholding*)

Solution.

· Si $x < 0$:

$$\frac{d}{dx} \left(\frac{1}{2}(z-x)^2 - \lambda x \right) = x - z - \lambda = 0 \Rightarrow x = z + \lambda < 0 \quad (x < 0 \text{ si } z < -\lambda)$$

· Si $x > 0$:

$$\frac{d}{dx} \left(\frac{1}{2}(z-x)^2 + \lambda x \right) = x - z + \lambda = 0 \Rightarrow x = z - \lambda \quad (x > 0 \text{ si } z > \lambda)$$

· Si $x = 0$:

$$\partial f(x) \in [x - z - \lambda, x - z + \lambda] = [-z - \lambda, -z + \lambda]$$

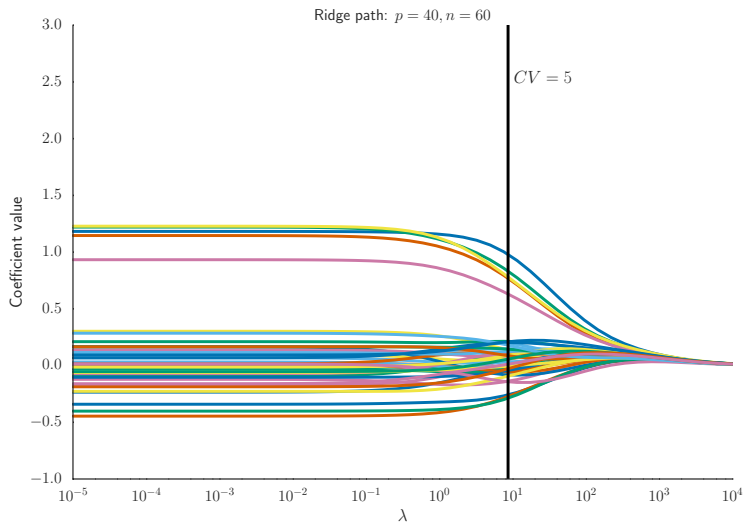
Le sous-gradient doit vérifier $\partial f(x) = 0$, alors $|z| \leq \lambda$

· On obtient finalement

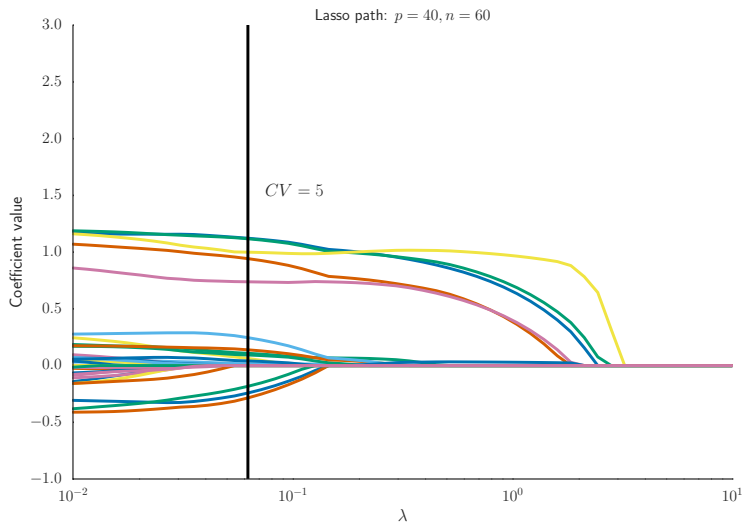
$$\eta_{\text{Lasso}, \lambda}(z) = \begin{cases} z + \lambda & \text{si } z < -\lambda \\ 0 & \text{si } |z| \leq \lambda \\ z - \lambda & \text{si } z > \lambda \end{cases}$$

- Pour cet exemple les tailles sont : $n = 60, p = 40, \sigma^2 = 1$
 - $\beta^* = (1, 1, 1, 1, 1, 0, \dots, 0) \in \mathbb{R}^p$ (5 coefficients non-nuls)
 - $X \in \mathbb{R}^{n \times p}$ a des colonnes tirées selon une loi gaussienne
 - $y = X\beta^* + \varepsilon \in \mathbb{R}^n$ avec $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$
 - On utilise une grille de 50 valeurs de λ

Choix de λ : Lasso vs Ridge (exemple avec $CV = 5$)



Choix de λ : Lasso vs Ridge (exemple avec $CV = 5$)



- Enjeu numérique : le Lasso est un problème **convexe**
- Sélection de variables / solutions parcimonieuses (sparse) : $\hat{\beta}_{\lambda}^{\text{Lasso}}$ a potentiellement de nombreux coefficients nuls
- Le paramètre λ contrôle le niveau de parcimonie : si λ est grand, les solutions sont très creuses

Exemple. Dans la simulation précédente :

Méthode	$\ \hat{\beta}_{\hat{\lambda}_{cv}}\ _0$ ($\ \hat{\beta}^*\ _0 = 5, p = 40$)
RidgeCV	40
LassoCV	17

Analyse théorique. Plus poussée que pour LS ou que pour Ridge [voir, e.g., Bühlmann and van de Geer, 2011, pour des résultats théoriques]

- Le lasso est biaisé : il contracte les grands coefficients vers 0

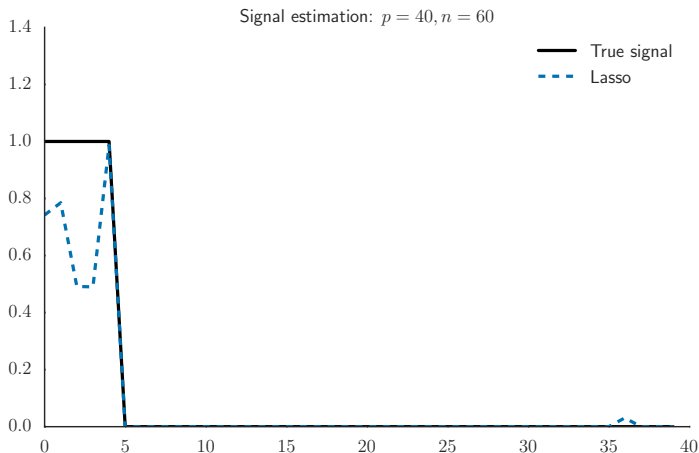


Illustration sur l'exemple

- Le lasso est biaisé : il contracte les grands coefficients vers 0

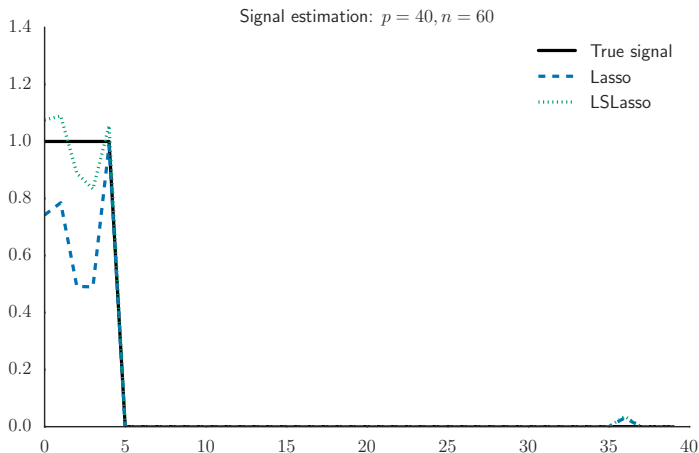


Illustration sur l'exemple

- Il est possible d'utiliser une procédure en deux étapes

Algorithm LSLasso (Least Square Lasso)

Données d'entrée : (X, y) , paramètres pour le Lasso (e.g., λ , nb d'itérations, tolérance)

1: Lasso pour obtenir $\hat{\beta}_{\lambda}^{\text{Lasso}}$ (sélection du modèle)

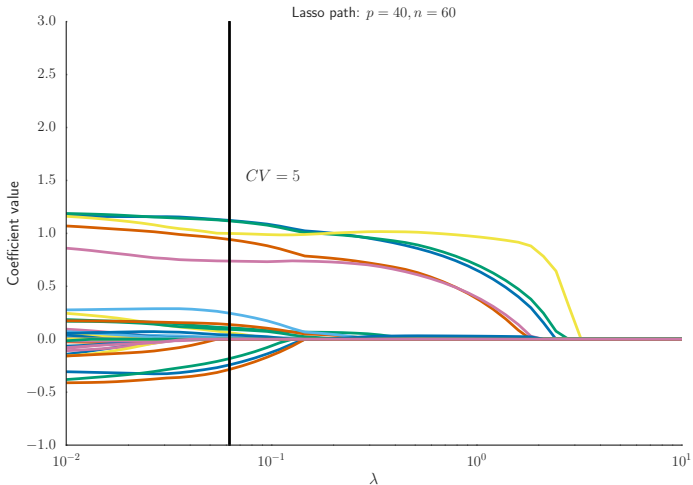
2: LS sur les variables actives (correction du biais)

$$\hat{\beta}_{\lambda}^{\text{LSLasso}} = \underset{\substack{\beta \in \mathbb{R}^p \\ \text{supp}(\beta) = \text{supp}(\hat{\beta}_{\lambda}^{\text{Lasso}})}}{\text{argmin}} \quad \frac{1}{2} \|y - X\beta\|_2^2$$

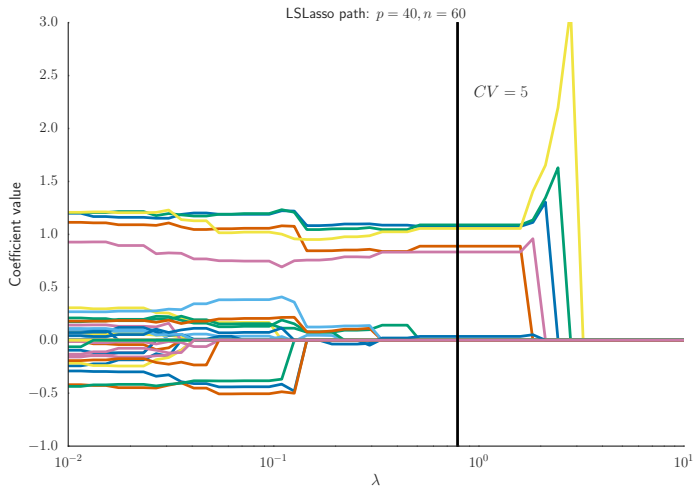
Rem :

- Il faut faire la CV sur la procédure entière
 - Choisir λ par LassoCV puis faire les LS garde trop de variables !
- LSLasso n'est pas forcément codé dans les packages usuels

Choix de λ : LS Lasso vs Lasso (exemple avec $CV = 5$)



Choix de λ : LS Lasso vs Lasso (exemple avec $CV = 5$)



Avantages.

- Les “vrais” grands coefficients sont moins atténués
- En faisant la CV on récupère moins de variables parasites (amélioration de l'interprétabilité)

Méthode	$\ \hat{\beta}_{\hat{\lambda}_{CV}}\ _0$ ($\ \hat{\beta}^*\ _0 = 5, p = 40$)
RidgeCV	40
LassoCV	17
LSLassoCV	6 (1 faux positif)

LSLasso : utile pour l'estimation

Limites.

- La différence en prédiction n'est pas toujours flagrante
- Nécessite plus de calcul : re-calculer autant de LS que de paramètres λ
- Non packagé

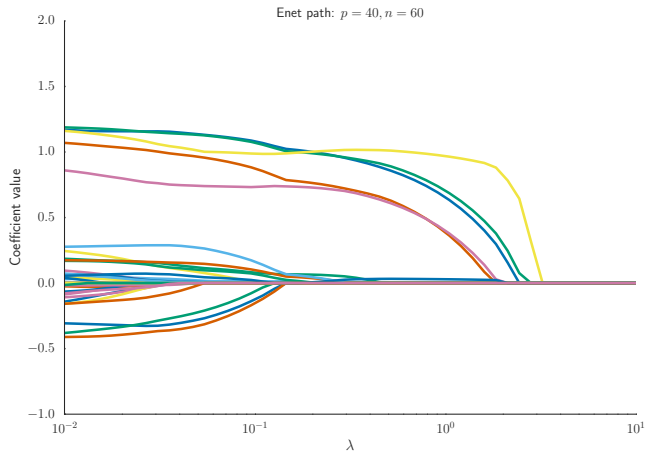
- L'Elastic Net introduit par [Zou and Hastie, 2005] est solution de

$$\widehat{\beta}_{\lambda, \gamma}^{\text{E-Net}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left[\frac{1}{2} \|y - X\beta\|_2^2 + \lambda \left(\gamma \|\beta\|_1 + (1 - \gamma) \frac{\|\beta\|_2^2}{2} \right) \right]$$

Rem :

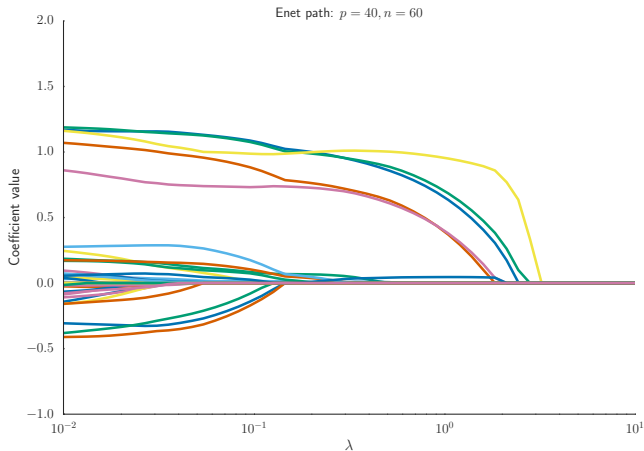
- Deux paramètres de régularisation : λ pour une régularisation globale, γ pour un contrôle de l'influence Ridge vs. Lasso
- La solution est unique et la taille du support de $\widehat{\beta}_{\lambda, \gamma}^{\text{E-Net}}$ est plus petite que $\min(n, p)$

Elastic-Net : $\gamma\|\beta\|_1 + (1 - \gamma)\|\beta\|_2^2/2$



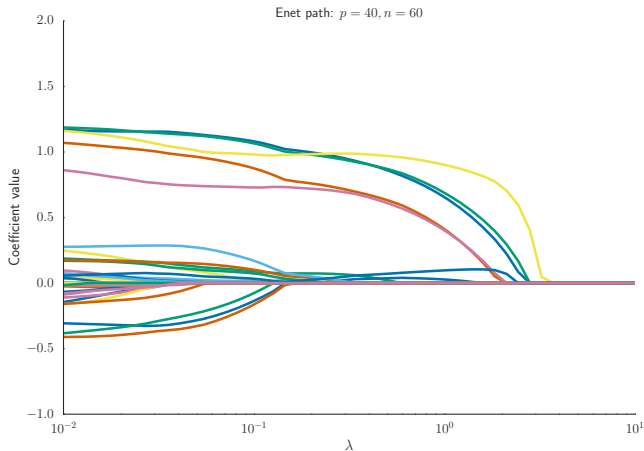
$\gamma = 1.00$

Elastic-Net : $\gamma\|\beta\|_1 + (1 - \gamma)\|\beta\|_2^2/2$



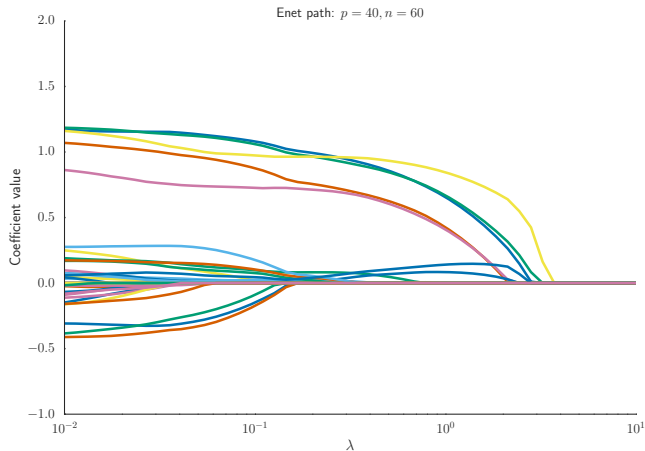
$$\gamma = 0.99$$

Elastic-Net : $\gamma\|\beta\|_1 + (1 - \gamma)\|\beta\|_2^2/2$



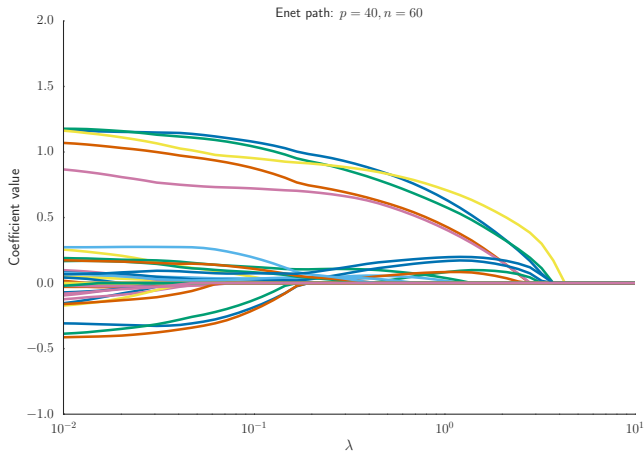
$$\gamma = 0.95$$

Elastic-Net : $\gamma\|\beta\|_1 + (1 - \gamma)\|\beta\|_2^2/2$



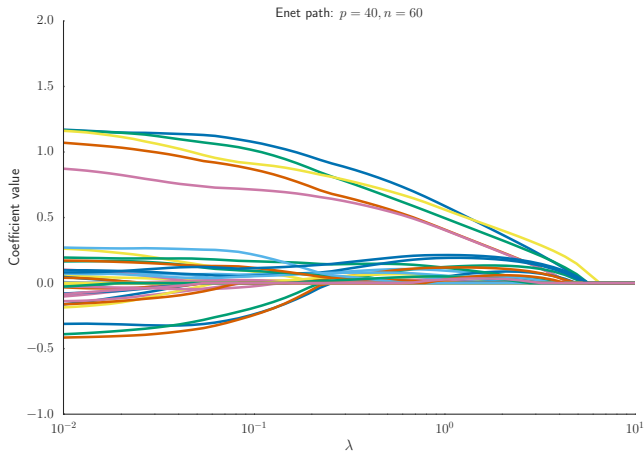
$$\gamma = 0.90$$

Elastic-Net : $\gamma\|\beta\|_1 + (1 - \gamma)\|\beta\|_2^2/2$



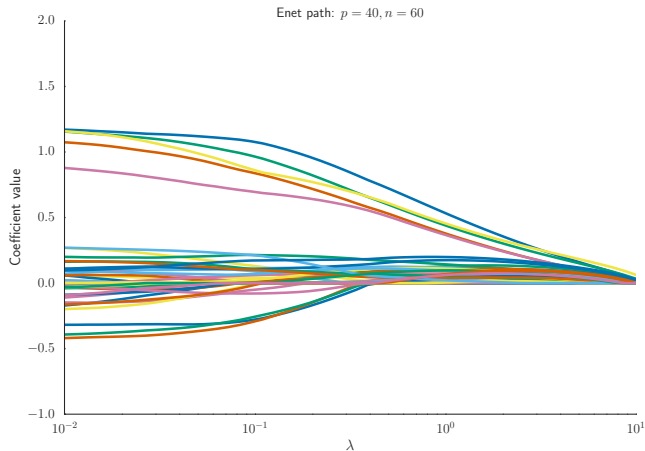
$$\gamma = 0.75$$

Elastic-Net : $\gamma\|\beta\|_1 + (1 - \gamma)\|\beta\|_2^2/2$



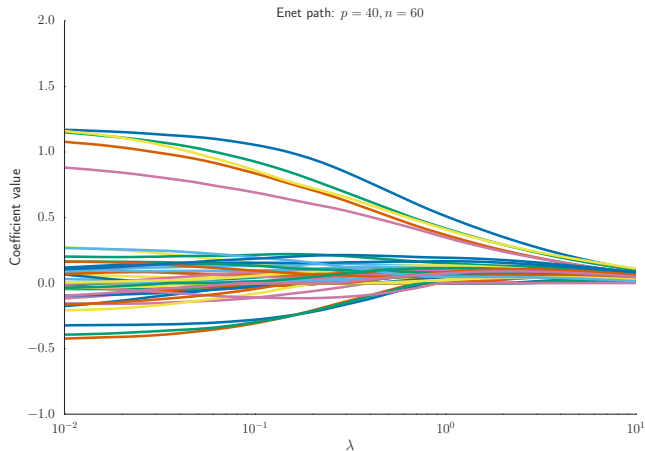
$$\gamma = 0.50$$

Elastic-Net : $\gamma\|\beta\|_1 + (1 - \gamma)\|\beta\|_2^2/2$



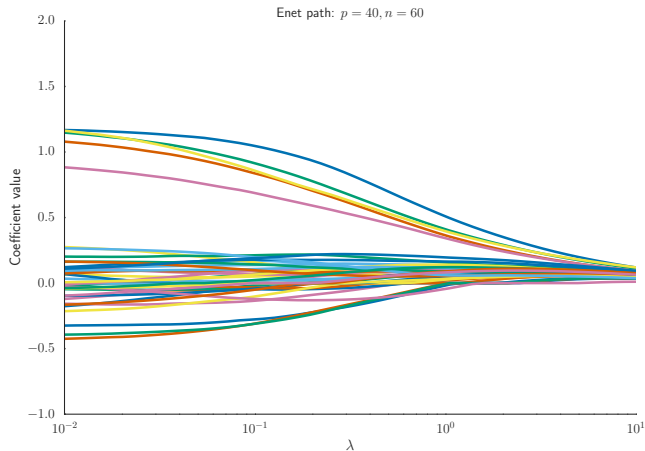
$$\gamma = 0.25$$

Elastic-Net : $\gamma\|\beta\|_1 + (1 - \gamma)\|\beta\|_2^2/2$



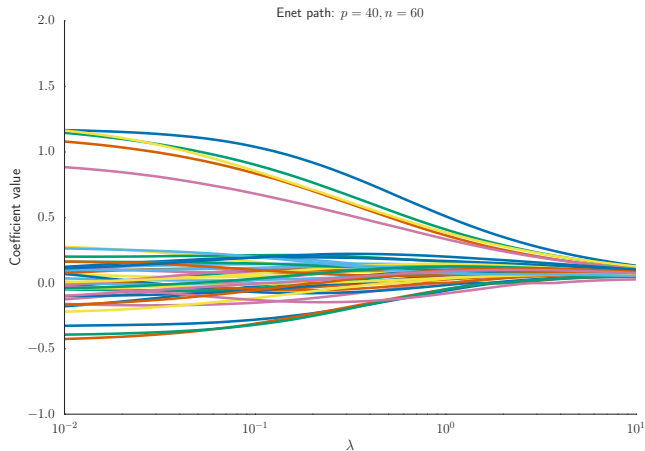
$\gamma = 0.1$

Elastic-Net : $\gamma\|\beta\|_1 + (1 - \gamma)\|\beta\|_2^2/2$



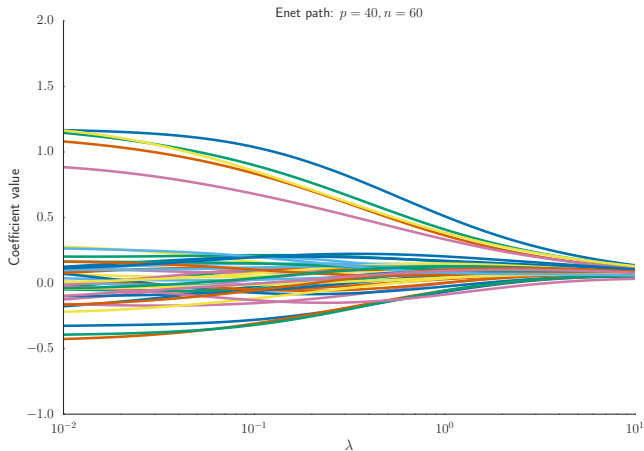
$$\gamma = 0.05$$

Elastic-Net : $\gamma\|\beta\|_1 + (1 - \gamma)\|\beta\|_2^2/2$



$$\gamma = 0.01$$

Elastic-Net : $\gamma\|\beta\|_1 + (1 - \gamma)\|\beta\|_2^2/2$



$\gamma = 0.00$

1. Pénalités non-convexes

- *Adaptive Lasso* [Candès et al., 2008, Zou, 2006] :

$$\text{pen}_\lambda(\beta) = \lambda \sum_{j=1}^p |\beta_j|^q \text{ avec } 0 < q < 1$$

2. Pénalités par groupes

- *Group Lasso* [Meier et al., 2008, Yuan and Lin, 2006] :

$$\text{pen}_\lambda(\beta) = \lambda \sum_{g=1}^G \|\beta_{\mathcal{I}_g}\|_2,$$

avec \mathcal{I}_g l'ensemble d'indices appartenant au g -ème groupe de variables

3. Cadre bayésien

- *Bayesian Lasso* [Park and Casella, 2008] : pour $\sigma > 0$ et $\tau > 0$,

$$y|\beta, \sigma^2 \sim \mathcal{N}(y|X\beta, \sigma^2 I_n),$$

$$\beta|\tau \sim \left(\frac{\tau}{2}\right)^2 \exp(-\tau \|\beta\|_1)$$

- Le mode $\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmax}} p(\beta|y, \sigma^2, \tau)$ converge vers $\hat{\beta}_\lambda^{\text{Lasso}}$ avec $\lambda = 2\tau\sigma^2$

- D. Bertsimas, A. King, and R. Mazumder. Best subset selection via a modern optimization lens. *Annals of Statistics*, 44(2):813–852, 2016.
- P. Bühlmann and S. van de Geer. *Statistics for high-dimensional data*. Springer Series in Statistics. Springer, Heidelberg, 2011. Methods, theory and applications.
- E. J. Candès, M. B. Wakin, and S. P. Boyd. Enhancing sparsity by reweighted l_1 minimization. *J. Fourier Anal. Applicat.*, 14(5-6):877–905, 2008.
- G. H. Golub and C. F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, Baltimore, 4 edition, 2013.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer, New York, second edition, 2009.
- L. Meier, S. Van De Geer, and P. Bühlmann. The group Lasso for logistic regression. *Journal of the Royal Statistical Society. Ser. B.*, 70(1):53–71, 01 2008.
- T. Park and G. Casella. The Bayesian Lasso. *Journal of the American Statistical Association*, 103(482):681–686, 2008.
- R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society. Ser. B. Methodological*, 58(1):267–288, 1996.
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society. Ser. B.*, 68(1):49–67, 2006.
- H. Zou. The adaptive Lasso and its oracle properties. 101(476):1418–1429, 2006.
- H. Zou and T. Hastie. Regularization and variable selection via the elastic net. 67(2): 301–320, 2005.