

ÉCOLE NATIONALE DE LA STATISTIQUE ET DE  
L'ADMINISTRATION ÉCONOMIQUE

## RAPPORT DU PROJET DE STATISTIQUE

JANVIER 2017

---

# Introduction à la régression pénalisée et à l'estimateur lasso

---

Auteurs :  
Mehdi ABBANA BENNANI  
Julien MATTEI

Superviseur :  
Edwin GRAPPIN

28 décembre 2016

## 1 Régression OLS

**Question 1 :**

$$\begin{aligned} \min_{\beta} \|y - X\beta\|^2 & \quad (\text{On minimise l'erreur quadratique}) \\ \nabla_{\beta} \|y - X\beta\|^2 (\beta^*) &= 0 \quad (\text{Le gradient est nul en } \beta^*) \\ \nabla_{\beta} (y - X\beta)^T (y - X\beta) (\beta^*) &= 0 \\ \boxed{\beta^* = (X^T X)^{-1} X^T y} & \quad (1) \end{aligned}$$

Dans la suite on notera  $\beta^*$  l'estimateur des moindres carrés.

**Question 2 :** Voir code.

On obtient une RMSE de 3.85, sachant que la variance de y est 90 à peu près.

## 2 Régression pénalisée : le lasso

**Question 3 :**

*Démonstration.* On ne peut pas utiliser la régression par moindres carrés si le nombre de variables est supérieur au nombre d'observations

X est une matrice de dimensions (n,p)

Donc  $X^T X$  est une matrice carrée de dimension p

On a  $rg(X^T X) = rg(X)$

Donc  $rg(X^T X) \leq \min(n, p)$

Par hypothèse, le nombre de variables est plus grand que le nombre d'observations, soit  $p > n$

Donc  $rg(X^T X) < p$

Or  $X^T X$  est une matrice de taille p

Donc  $X^T X$  n'est pas inversible

Donc on ne peut pas utiliser la méthode précédente.

La matrice  $X^T X$  n'est plus inversible car son rang est inférieur au minimum des dimensions de X.  $\square$

**Question 4 :**

Sous l'hypothèse de  $p=1$ , on résoud le problème d'optimisation :

$$\min_{\beta} (y - x\beta)^2 + \lambda|\beta|$$

A partir de l'expression (1), on a  $\hat{\beta} = \frac{y}{x}$

On injecte  $\hat{\beta}$  dans le problème d'optimisation et après simplification :

$$\min_{\beta} \frac{x^2 \beta^2}{2} - x^2 \beta \hat{\beta} + \lambda |\beta|$$

On remarque que si  $\hat{\beta} < 0$  Alors forcément  $\beta^{lasso} < 0$

De même si  $\hat{\beta} > 0$  Alors forcément  $\beta^{lasso} > 0$

On va traiter les deux cas

Si  $\hat{\beta} > 0$ , on résoud :

$$\min_{\beta} \frac{x^2 \beta^2}{2} - x^2 \beta \hat{\beta} + \lambda \beta$$

On obtient  $\beta^{lasso} = \hat{\beta} - \frac{\lambda}{x^2}$

Or  $\beta^{lasso}$  forcément positif, donc  $\beta^{lasso} = \max(0, \hat{\beta} - \frac{\lambda}{x^2})$

Si  $\hat{\beta} < 0$ , à l'aide d'un raisonnement similaire  $\beta^{lasso} = \max(0, \hat{\beta} + \frac{\lambda}{x^2})$

On en déduit que  $\beta^{lasso} = sg(\hat{\beta})(|\hat{\beta}| - t)_+$  avec  $t = \frac{\lambda}{x^2}$

(2)

Le seuillage doux permet de régler l'intervalle sur lequel on veut estimer  $\beta^*$ . Ainsi, on choisit  $\lambda$ , si notre estimateur  $\beta^*$  est compris dans l'intervalle  $[-\lambda, \lambda]$  il est ramené à 0. On ne retient que les valeurs de  $\beta^*$  en dehors de cet intervalle (même principe qu'un filtre).

L'estimateur  $\hat{\beta}^{lasso}$  est nul pour tout les  $\beta^*$  tels que  $\beta^* - t < 0$ . Donc cela induit une estimation nulle plus souvent que l'estimateur des moindres carrés (qui se produit seulement lorsque  $\beta^* = 0$ )

**Question 5 :** Dans la régression pénalisée, on pénalise les valeurs de  $\beta$  pour limiter leur valeurs possibles, dans le cas du lasso par exemple, toutes les estimations par moindres carrés inférieures à  $t$  sont ramenées zeros, l'intervalle d'estimation est  $\{0\} \cup [t, \infty[$ . Dans le cas de la régression Ridge par exemple, on pénalise les coefficients trop grands.

L'estimateur Lasso est adapté à notre problème car nous allons pouvoir réduire la taille de l'espace des paramètres jusqu'à ce que le nombre de coefficients de beta non nuls soit égal à  $n$  le nombre d'observations, en changeant la valeur de  $\lambda$ . Plus celui-ci est grand, plus le nombre de paramètres nuls sera grand.

**Question 6 :** Quelles sont les avantages du lasso sur la pénalisation ' 0 ? Même question avec la pénalisation Ridge

**Question 7 :** A implémenter

**Question 8 :** Le  $\lambda$  permet de régler la "force" de pénalisation du paramètre  $\beta$ .

**Question 9 :** Dans la méthode de cross validation, on divise le dataset en trois parties : Train Set, Validation Set, et Test Set. On entraîne le modèle sur le train set puis on estime la qualité de la prédiction sur le validation set pour plusieurs valeurs de  $\lambda$ , on choisit le lambda le plus performant sur le validation set pour entraîner le modèle sur le nouveau train set comprenant le train set précédent et le validation set.

**Question 10 :** A implémenter

Combien d'éléments de  $\beta$  estimez vous non-nuls ? Évaluez l'erreur de votre modèle. Trouvez de bonnes mesures pour évaluer votre méthode