

# Analyse des données prostate: quelques méthodes de régularisation $\ell_1$

[julien.chiquet@genopole.cnrs.fr](mailto:julien.chiquet@genopole.cnrs.fr)

Module MPR – option modélisation, 8 décembre 2009

## Table des matières

|   |          |
|---|----------|
| <b>1 Motivations</b>  | <b>1</b> |
| <b>2 Méthodes de régularisation s'appuyant sur la norme <math>\ell_1</math></b> | <b>2</b> |
| 2.1 Le LASSO  | 2        |
| 2.2 Quelques variantes  | 3        |
| 2.2.1 L'Elastic-Net   | 3        |
| 2.2.2 L'adaptive LASSO  | 3        |
| 2.2.3 Le group-LASSO  | 3        |
| <b>3 Code R commenté</b>  | <b>3</b> |
| 3.1 Estimation  | 4        |
| 3.1.1 LASSO   | 4        |
| 3.1.2 Elastic-Net   | 5        |
| 3.1.3 Adaptive-LASSO  | 6        |
| 3.1.4 Group-LASSO   | 7        |

## 1 Motivations

Soit le modèle linéaire défini par

$$y = \beta_0 + \sum_{i=1}^p X_i \beta_i + \varepsilon,$$

où  $y$  est la réponse à expliquer par le vecteur de variables  $X = (X_1, \dots, X_p)^\top$ , avec  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ . Notons  $\mathbf{y}$  un vecteur de taille  $n$  et  $\mathbf{X}$  la matrice de taille  $n \times p$  des  $n$  observations.

Dans la suite, afin de simplifier l'écriture, on suppose que l'intercept  $\beta_0$  est estimé par  $\bar{\mathbf{y}} = n^{-1} \sum_{i=1}^n y_i$ , que le vecteur  $\mathbf{y}$  est centré et que la matrice  $\mathbf{X}$  est centrée-réduite. Ainsi l'estimation ne porte plus que sur le vecteur  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$  de  $\mathbb{R}^p$ . Dans ce cas, le modèle s'écrit en terme matriciel sous la forme

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

Par rapport aux moindres carrés ordinaires, la régression *ridge* estime le paramètre  $\beta$  avec un léger biais mais permet de contrôler la variance. Pour un choix judicieux<sup>1</sup> du paramètre de pénalisation  $\lambda$ , l'erreur de prédiction peut être meilleure avec la régression ridge qu'avec les moindres carrés.

Nous présentons ici d'autres méthodes de régularisation qui induisent également un biais et une baisse de la variance. Elles sont fondées sur la norme  $\ell_1$ , définie, pour un vecteur  $\mathbf{x} = (x_1, \dots, x_p)^\top$ , par

$$\|\mathbf{x}\|_{\ell_1} = \sum_{i=1}^p |x_i|.$$

Ces méthodes ont un double avantage

1. elles permettent d'annuler *exactement* certaines valeurs des coefficients estimés pour des valeurs suffisamment fortes du paramètre  $\lambda$  : en *sélectionnant* automatiquement les variables, le modèle gagne en interprétabilité ;
2. le problème d'optimisation à résoudre reste convexe, c'est-à-dire relativement facile résoudre, même s'il n'existe pas de forme analytique de la solution dans le cas général.

## 2 Méthodes de régularisation s'appuyant sur la norme $\ell_1$

Dans la suite, on note (1)

$$\text{RSS}(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2,$$

la somme du carré des résidus du problème.

### 2.1 Le Lasso

La régression LASSO, très populaire, a été proposée par [Tib96], s'inspirant de méthodes de traitement du signal. Elle consiste à ajouter une pénalité en norme  $\ell_1$  du vecteur des coefficients (*i.e.*, la somme des valeurs absolues des éléments du vecteur). Ainsi, l'estimateur est défini par

$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \text{RSS}(\beta) + \lambda \|\beta\|_1 \right\}. \quad (1)$$

La LASSO (*Least Absolute Shrinkage and Selection Operator*) a la double propriété de *rétrécir* les coefficients de  $\beta$ , en limitant leur taille à la manière de la régression ridge, mais surtout de *sélectionner* ces coefficients en en annulant certains pour des valeurs suffisamment forte du paramètre de pénalité  $\lambda$ .

**Chemin de régularisation.** Si l'on construit le chemin de solutions du vecteur  $\beta$  en fonction du paramètre  $\lambda$ , on s'aperçoit qu'il est *linéaire par morceaux*. Ainsi, on peut construire une suite de  $p + 1$  valeurs croissantes de  $\lambda_0 = 0 < \lambda_1 < \dots < \lambda_p$  telles

---

<sup>1</sup>typiquement, une valeur qui minimise l'erreur de prédiction calculée sur un échantillon test

que  $\beta_{\lambda_0}^{\text{lasso}}$  soit l'estimateur des moindres carrés,  $\beta_{\lambda_1}^{\text{lasso}}$  ait exactement un coefficient nul,  $\beta_{\lambda_2}^{\text{lasso}}$  deux coefficients nuls, etc. jusqu'à  $\beta_{\lambda_p}^{\text{lasso}}$  égal au vecteur nul.

La détermination de cette suite de valeurs de  $\lambda$  et des  $\beta$  correspondant se fait par l'algorithme LARS dont il existe un package R.

## 2.2 Quelques variantes

De nombreuses variantes du LASSO ont été proposées, impliquant la plupart du temps de nouveaux algorithmes de calcul.

### 2.2.1 L'Elastic-Net

L'elastic-net de [ZH05] pénalise la somme des carrés résiduelles par un mélange de pénalité ridge et de pénalité LASSO :

$$\hat{\beta}^{\text{e-net}} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \text{RSS}(\beta) + \lambda (\alpha \|\beta\|_1 + (1 - \alpha) \|\beta\|_2^2) \right\}, \quad (2)$$

### 2.2.2 L'adaptive Lasso

L'adaptive-LASSO permet de pénaliser différemment les différents coefficients du vecteur  $\beta$ .

$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta \in \mathbb{R}^{p+1}} \left\{ \frac{1}{2} \text{RSS}(\beta) + \lambda \|\mathbf{w} \circ \beta\|_1 \right\}. \quad (3)$$

où  $\mathbf{w}$  est un vecteur de poids permettant de doser le niveau de pénalité de chaque coefficient de  $\beta$ . La notation  $\circ$  est utilisée pour désigner le produit terme à terme entre deux vecteurs ou deux matrices.

### 2.2.3 Le group-Lasso

Le group-LASSO de [YL06] pénalise les variables par groupes, ainsi les variables sont activées (sélectionnées) par groupes :

$$\hat{\beta}^{\text{grp}} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \text{RSS}(\beta) + \lambda \sum_{g=1}^G \|\beta^g\|_1 \right\}, \text{ où } \beta = (\beta^1, \dots, \beta^g, \dots, \beta^G)^\top, \quad (4)$$

avec  $\text{card}(\beta^g) = p_g$ ,  $\sum_g p_g = p$ . Les groupes sont connus selon ce que représentent les variables explicatives.

## 3 Code R commenté

Les données prostate contiennent une matrice  $\mathbf{x}$ , un vecteur  $\mathbf{y}$  et un vecteur de booléens `set` indiquant les données de l'ensemble d'apprentissage. On charge les données et les fonctions qui sont stockées dans le fichier `functions.R`, puis on centre et on réduit la matrice  $\mathbf{X}$ . On construit les données de test et d'apprentissage comme pour l'étude de la régression ridge.

```

> source("functions.R")
> load("prostate.rda")
> x <- scale(as.matrix(x))
> n <- sum(set)
> n.test <- sum(!set)
> x.test <- x[!set, ]
> y.test <- y[!set]
> x <- x[set, ]
> y <- y[set]

```

Pour calculer les estimateurs, nous chargeons les packages suivants<sup>2</sup> :

```

> library(glmnet)
> library(grplasso)

```

```

-----
Please note that this is an early test release of package 'grplasso'.
It should only be used for experimental reasons. Use at your own risk!
-----

```

Le package `glmnet` a été conçu pour calculer la solution de l'elastic-net, mais il permet également de calculer la solution du LASSO. Le package `grplasso` est utile pour le group-LASSO. Un point fort de ces packages est de gérer les modèles de régressions linéaires généralisés (ainsi la régression logistique) : lorsque la réponse  $y$  à expliquer n'est pas une variable continue, le modèle linéaire simple n'est pas toujours approprié, puisqu'il suppose que  $y$  est gaussien. Ainsi, dans le cas d'une variable catégorielle (comme dans le cas des données sport –  $y$  représente le sport pratiqué), il peut être judicieux d'utiliser un modèle logistique. Ces packages permettent de le faire<sup>3</sup>.

Ici la variable explicative est continue, donc le modèle linéaire usuel est bien adapté.

## 3.1 Estimation

### 3.1.1 Lasso

Par défaut, `glmnet` calcule la solution du LASSO pour 100 valeur de  $\lambda$  depuis une valeur annulant tous les coefficients jusqu'à une valeur proche de zéro, où le vecteur estimé est proche des moindres carrés. L'algorithme ne descend pas nécessairement jusqu'à de pénalités  $\lambda = 0$  pour des problèmes de stabilité. Nous calculons l'erreur au passage :

```

> out.lasso <- glmnet(x, y)
> err.lasso <- 1/n.test * colSums((y.test - out.lasso$a0[1] - x.test %*%
+   out.lasso$beta)^2)

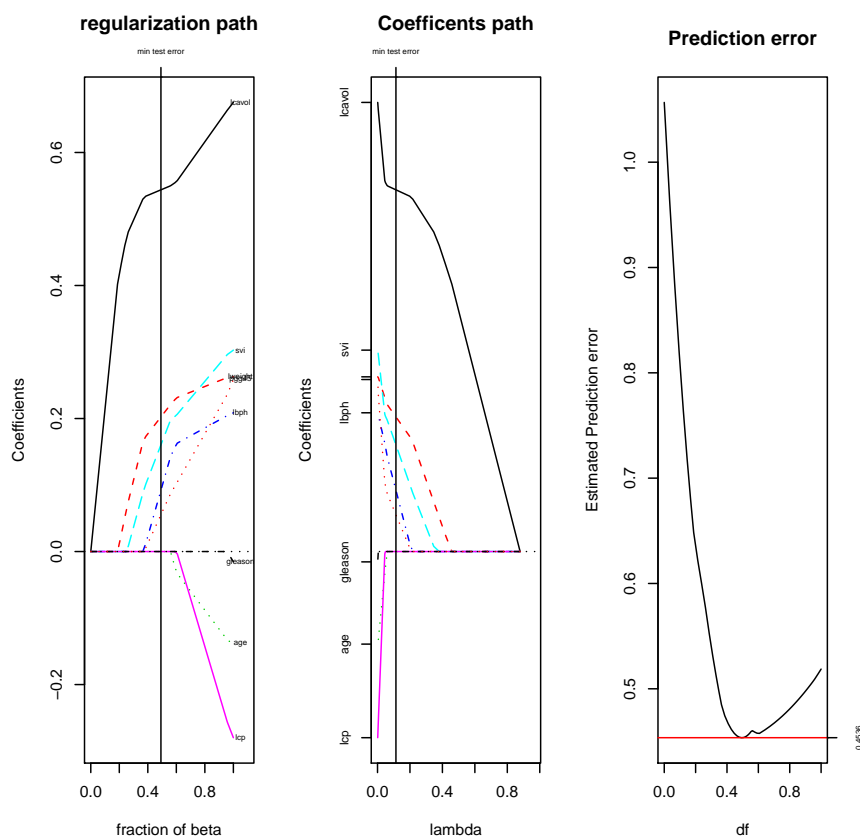
```

<sup>2</sup>N'ayez pas peur du message d'avertissement du package `grplasso`.

<sup>3</sup>à bon entendeur...

Dans l'objet `out.lasso` se trouvent `out.lasso$beta`, valeurs estimées des coefficients en fonction des valeurs `out.lasso$lambda` de la pénalité. Voici le tracé des résultats (chemin de régularisation, coefficients et erreur de prédiction) à l'aide des fonctions contenues dans le fichier `fonctions.R`.

```
> par(mfrow = c(1, 3))
> frac.lasso <- plot.path(t(out.lasso$beta), err = err.lasso)
> plot.coef(t(out.lasso$beta), out.lasso$lambda, err.lasso)
> plot.err(err.lasso, frac.lasso)
```

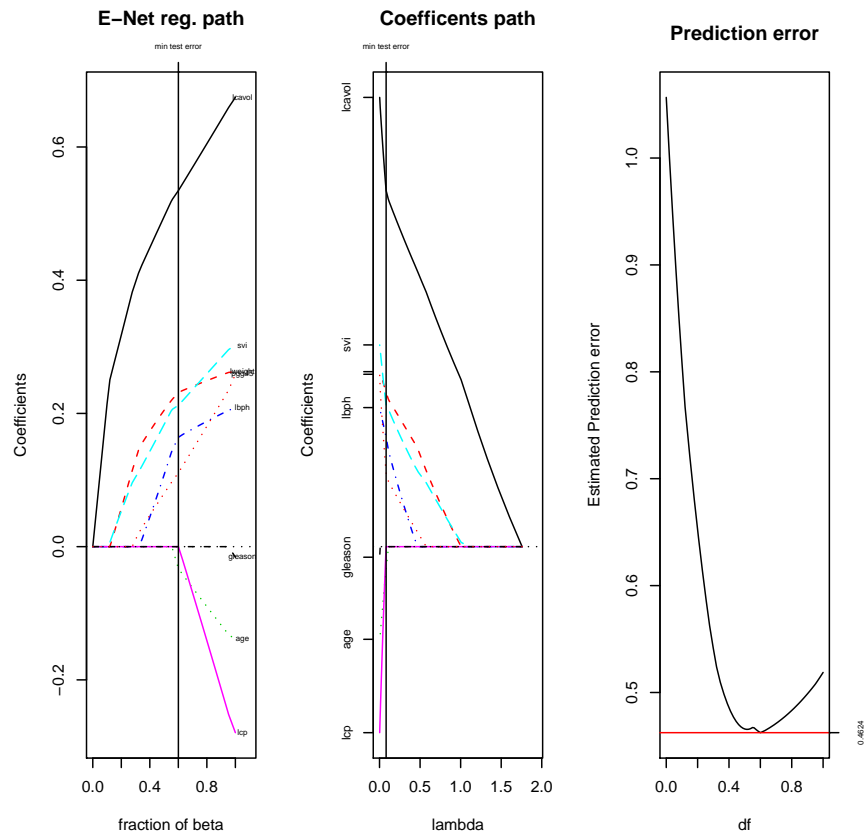


### 3.1.2 Elastic-Net

Pour l'estimateur de l'elastic-net, choisissons par exemple  $\alpha = 0.5$ , c'est-à-dire accordant autant d'importance à la pénalité en norme  $\ell_1$  qu'en norme  $\ell_2$  :

```
> out.elas <- glmnet(x, y, alpha = 0.5)
> err.elas <- 1/n.test * colSums((y.test - out.elas$a0[1] - x.test %*%
+   out.elas$beta)^2)
```

Graphiquement, les résultats sont :



### 3.1.3 Adaptive-Lasso

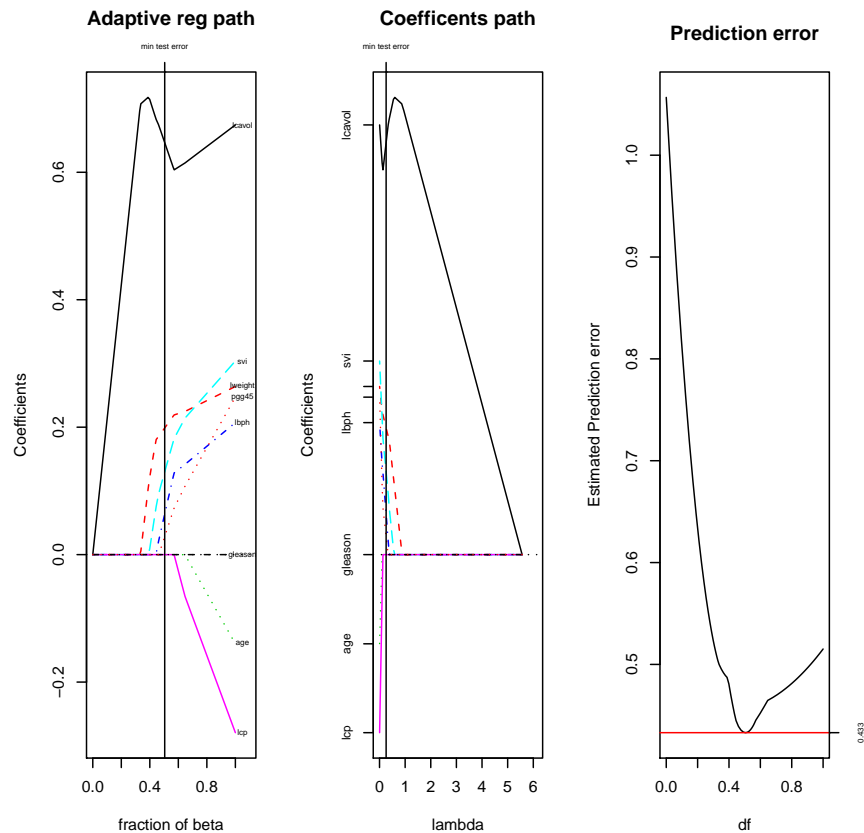
Pour l'adaptive-LASSO un choix usuelle pour les poids est le suivants :

$$\mathbf{w} = \max\left(1, 1/\beta^{\text{ols}}\right),$$

d'où le code :

```
> beta.ols <- lm(y ~ x)$coefficients[-1]
> p.fact <- pmax(1, abs(beta.ols)^(-1))
> out.lasso <- glmnet(x, y, penalty.factor = p.fact)
> err.lasso <- 1/n.test * colSums((y.test - out.lasso$a0[1] -
+   x.test %*% out.lasso$beta)^2)
```

Les résultats sont :



### 3.1.4 Group-Lasso

Commençons par déterminer des groupes : au vu du chemin de régularisation du LASSO, on a envie d'activer simultanément les variables 2 (`lweight`) et 5 (`svi`), tout comme les variables 3 (`age`) et 6 (`scp`); de même pour les variables 4 (`lbph`) et 8 (`pgg45`). D'où le vecteur indiquant les groupes

```
> group <- c(1, 2, 3, 4, 2, 3, 5, 4)
```

La package `grplasso` est moins facile à utiliser que `glmnet`. Nous devons créer à la main un objet R contenant les valeurs de l'intercept et ainsi qu'un vecteur de pénalité. Nous prenons les mêmes valeurs de  $\lambda$  que pour le LASSO que l'on remet à échelle (effet de bord du package `grplasso`) :

```
> out.group <- list()
> out.group$beta0 <- mean(y)
> out.group$lambda <- 2 * n * out.lasso$lambda
```

Nous sommes prêt à faire appelle à la fonction `grplasso` : on notera le paramètre modèle (ici `LinReg()`, pouvant être `LogReg()`). Le paramètre `grpl.control()` est utilisé ici pour éviter les messages intempestifs.

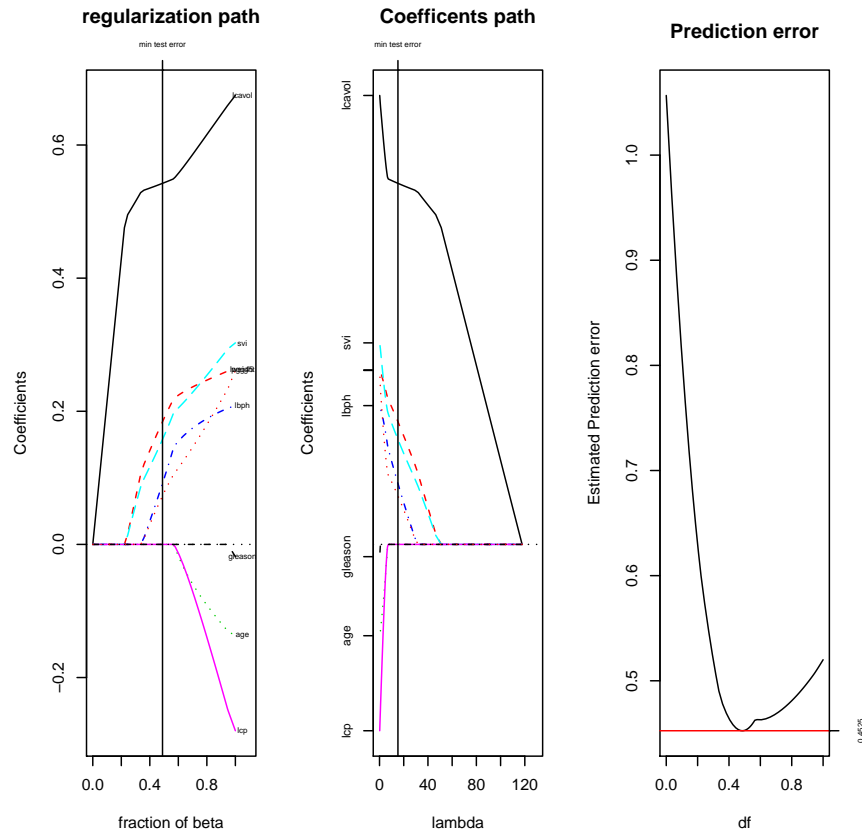
```
> out <- grplasso(x, y = y - mean(y), index = group, lambda = out.group$lambda,
+ model = LinReg(), control = grpl.control(trace = 0))
```

```

> out.group$beta <- rbind(t(out$coefficients))
> err.group <- 1/n.test * colSums((y.test - mean(y) - x.test %*%
+   t(out.group$beta))^2)

```

Les résultats donnent



## Références

- [Tib96] Tibshirani, R., Regression Shrinkage and Selection via the Lasso, *Journal of the Royal Statistical Society*, 1996.
- [ZH05] Zou, H. and Hastie T., Regularization and variable selection via the elastic net, *Journal of the Royal Statistical Society*, 2005.
- [Zou06] Zou, H. , The adaptive Lasso and its Oracle properties, *Journal of the American Statistical Society*, 2006.
- [YL06] Yuan, M., Lin, Y., Model Selection and Estimation in Regression with Grouped Variables, *Journal of the Royal Statistical Society*, 2006.