

# Étude de quelques méthodes régularisées pour la régression linéaire

[julien.chiquet@gmail.com](mailto:julien.chiquet@gmail.com)

séminaire professionnel ENSAI - 2016

## 1 Objectifs

Ces séances visent à étudier les propriétés empiriques de quelques méthodes régularisées à partir de données simulées et sur un exemple de données réelles. Un objectif connexe est la mise en place d'une méthodologie de simulations numériques, démarche très utile dans le cadre de travaux de recherche.

- utilisation de la sélection stepwise, de la régression ridge, du Lasso et de ses variantes
- mise au point d'un protocole complet de simulations dans le cadre du modèle linéaire
- acquisition de nouvelles compétences R (**glmnet**, **parallel**, **ggplot2**)
- application à un jeu de données en génomique

*Remarques* : les étudiants peuvent travailler en binôme

## 2 Mise au point d'un protocole de simulation

La simulation de données est (ou *devrait être*) un passage obligé de tout travail de recherche en statistiques – théoriques ou appliquées – lors de l'étude d'une nouvelle stratégie d'estimation. En effet, les simulations sont un moyen d'étude des propriétés des méthodes dans un cadre parfaitement *contrôlé*. Afin de préciser les *propriétés d'intérêt*, il est nécessaire de ce donner un cadre d'étude. Nous choisissons ici le problème de la sélection de variables prédictives dans le modèle linéaire.

### 2.1 Cadre d'étude: sélection de variables dans le modèle linéaire

Pour une méthode de sélection de variables, les *propriétés d'intérêt* sont

- la capacité à retrouver les variables pertinentes,
- la capacité à prédire une nouvelle observation,
- la précision de l'estimation.

La notion de *cadre contrôlé* se réfère aux hypothèses qui sont faites lors de la définition du modèle statistique. Un intérêt important des simulations numériques est de se placer dans le “bon cas” (*i.e.* celui prévu par la théorie) afin d’étudier les propriétés de la méthode dans ce cadre: on s’assure que les performances sont conformes à celles attendues. On peut ensuite se permettre de s’éloigner du cadre bien contrôlé par la théorie afin d’étudier la robustesse de la méthode à des écarts aux hypothèses.

On s’attache ici au modèle linéaire gaussien homoscedastique:

$$y_i = x_i^\top \beta^* + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n \quad (1)$$

soit matriciellement

$$\mathbf{y} = \mathbf{X}\beta^* + \boldsymbol{\varepsilon},$$

avec  $\mathbf{y} = (y_1, \dots, y_n)^\top$  un vecteur de  $\mathbb{R}^n$ ,  $\mathbf{X}$  une matrice de  $\mathcal{M}_{n,p}(\mathbb{R})$  et  $\beta^* = (\beta_1, \dots, \beta_p)^\top$  un vecteur de  $\mathbb{R}^p$  dont  $p_0$  éléments sont non-nuls. On note  $\mathcal{A}^* = \{j : \beta_j^* \neq 0\}$  l’ensemble de ces éléments. On suppose que les prédicteurs sont gaussiens multivariés de matrice de variance-covariance  $\boldsymbol{\Sigma}$ .

## 2.2 Contrôle de la difficulté

Lors de l’analyse d’une procédure statistique, il est important de pouvoir jauger précisément la difficulté du problème pour déterminer le champs d’applicabilité de cette méthode. Ceci permet en particulier de mesurer sa robustesse aux écarts aux hypothèses requit par l’analyse théorique.

Dans le cadre de la sélection de variables pour le modèle linéaire gaussien, divers phénomènes sont associés à la difficulté du problème. Dans le cas des méthodes de sélection de variables comme le LASSO, ces phénomènes sont liés par l’analyse théorique aux paramètres du modèle linéaire qui permettent de déterminer plus ou moins précisément les conditions d’applicabilité de cette méthode. On pense en particulier

- aux phénomène de grande dimension ( $n < p$ ), contrôlé par le rapport  $n/p$ ,
- au niveau de “sparsité” du problème, contrôlé par le cardinal de  $\mathcal{A}^* = p_0^*$ ,
- au niveau de bruit, qu’on peut contrôler par le rapport signal sur bruit ou le coefficient de détermination,
- à la structure de dépendance entre les prédicteurs, *i.e.*, la covariance  $\boldsymbol{\Sigma}$  du vecteur  $X = (X_1, \dots, X_p)$ .

On propose dans la suite d’intégrer ces divers paramètres à une fonction `rlm` permettant de générer des données issues du modèle (1).

### 2.2.1 V1: modèle linéaire simple

Écrire une fonction `rlm(n,p,p0,sigma)` qui renvoie une liste contenant un vecteur  $\mathbf{y}$ , une matrice  $\mathbf{x}$ , un vecteur `beta` avec  $p_0$  entrées non nulles de magnitude choisie selon une loi uniforme entre 1 et 2 et de signe positif ou négatif. Les prédicteurs

sont tirés selon des distribution gaussiennes univariés indépendantes de variance unitaire.

### 2.2.2 V2: intégration du coefficient de détermination

Plutôt que de fixer la variance du bruit  $\sigma^2$ , on préfère fixer la puissance du signal par rapport à la puissance du bruit, par exemple en contrôlant le niveau du coefficient de détermination, qui mesure la part de variabilité expliquée par le modèle:

$$R^2 = \frac{SCM}{SCT} = 1 - \frac{SCR}{SCT}.$$

Déterminer dans le cadre du modèle (1) quelle valeur de  $\sigma^2$  choisir pour obtenir des données ayant un coefficient de détermination donné. Amender votre fonction en conséquence: elle prend désormais en argument `rlm(n,p,p0,r2)` et renvoie une liste `list(y,x,beta,sigma)`.

### 2.2.3 V3: structure de dépendance des prédicteurs

Les prédicteurs ont pour l'instant été considérés comme indépendants. On propose de modéliser une forme de dépendance entre prédicteurs à l'aide d'une loi gaussienne multivariée telle que  $X \sim \mathcal{N}(\mathbf{0}, \Sigma)$ . Outre le cas indépendant ( $\Sigma = \mathbf{I}_p$ ), on considèrera 2 autres scénarios:

- une dépendance de type longitudinale:  $\Sigma_{ij} = \rho^{|i-j|}$ .
- une dépendance par bloc: soit une partition en  $K$  groupes, alors

$$\Sigma_{ij} = \begin{cases} 1 & \text{si } i = j, \\ \rho & \text{si } i \text{ et } j \text{ sont dans le même groupe,} \\ 0 & \text{sinon.} \end{cases}$$

Écrire deux fonctions `rlm.long(n, p, p0, r2, rho)` et `rlm.bloc(n, p, K0, r2, rho, K)` adaptées à ces deux scénarios. On remarquera que dans le second cas, la sparsité est défini en terme de nombre de bloc non nuls  $K_0$ . Pour la génération d'un vecteur gaussien multivarié, on utilisera une décomposition de Cholesky de la matrice  $\Sigma$ .

### 2.2.4 V4: ensemble test, ensemble d'apprentissage

Afin d'évaluer les performances des estimateurs, il est indispensable de générer des données de test. Amender les fonctions précédentes de sorte à prendre en argument `n.test` et renvoyer en plus des variables précédentes des ensembles `train` et `test` indiquant les observations. On pourra affecter une valeur par défaut à `n.test` dépendant de `n.train`, par exemple `10*n.train`.

---

On dispose finalement de 3 fonctions `rlm.indep`, `rlm.long` et `rlm.bloc` renvoyant une liste contenant les variables `y,x,beta,Sigma,sigma,train,test` qui serviront à la génération de données pour les simulations.

## 2.3 Implémentation des méthodes régularisées

On se propose d'étudier les procédures suivantes:

- la régression “stepwise” avec critère AIC et BIC,
- la régression “ridge”,
- la régression “lasso”,
- la variante “elastic-net” (ridge + lasso),

On comparera ces procédures dites régularisées aux méthodes de référence suivantes:

- les moindres carrés ordinaires,
- les moindres carrés oracle (où l'on suppose connaître  $\mathcal{A}^*$ ).

Écrire une fonction par estimateur, du type `getStepwiseAIC`, `getStepwiseBIC`, `getLasso`, etc. qui récupère la valeur de  $\hat{\beta}$ . Toutes les méthodes seront implémentés à l'aide de la fonction `glmnet` du package du même nom sauf la régression stepwise (fonction `step` du package `MASS`).

*Quelques remarques*

- Pour les expériences en grande dimension ( $n < p$ ), l'estimateur des moindres carrés n'est pas défini de manière unique: on utilisera un estimateur avec une faible pénalité ridge pour le régulariser.
- Pour les méthodes régularisées, on choisira le paramètre  $\lambda$  par validation croisée à l'aide des fonctions `cv.glmnet`.

### 2.3.1 Lasso adaptatif

Le Lasso adaptatif est une version modifiée du lasso proposée pour palier notamment au problème de biais. L'idée est de procéder en deux temps. Dans un premier temps, on estime des paramètres  $\hat{\beta}_{\text{init}}$  à l'aide du lasso. Ce premier estimateur va être utilisé comme poids pour une deuxième étape de lasso, en fixant  $w_j = \hat{\beta}_{j,\text{init}}$  tel que:

$$\hat{\beta}_{\lambda}^{\text{ada}} = \arg \min_{\beta \in \mathbb{R}^p} \left( \frac{1}{n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \sum_{j=1}^p \frac{|\beta_j|}{|w_j|} \right).$$

L'idée est que si  $\hat{\beta}_{j,\text{init}} = 0$  alors  $\hat{\beta}_{j,\text{ada}} = 0$  de telle sorte que la première étape de lasso sert de pré-sélection. De plus, si  $\hat{\beta}_{j,\text{init}}$  est grand, le lasso adaptatif utilisera une pénalisation plus petite, donc un shrinkage plus petit pour le coefficient  $j$ , ce qui est sensé diminuer le biais pour ce coefficient.

Implémentez le lasso adaptatif à l'aide de l'option `penalty.factor` dans `glmnet` qui permet d'introduire des  $\lambda$ s différents pour chaque élément de  $\beta$ . Créer une fonction `getadalasso(X,Y,beta.init)` qui calcule l'estimateur lasso adaptatif à partir d'un premier vecteur  $\hat{\beta}_{\text{init}}$ .

### 2.3.2 Group-Lasso

Une autre modification populaire du Lasso est une version groupée du lasso – ou group-Lasso – qui suppose la connaissance d'une partition *a priori* des prédicteurs

notées  $\mathcal{G} = \{\mathcal{G}_1, \dots, \mathcal{G}_K\}$ . La sélection s'opère donc par groupe:

$$\hat{\beta}_\lambda^{\text{grp}} = \arg \min_{\beta \in \mathbb{R}^p} \left( \frac{1}{n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \sum_{k=1}^K \|\beta_{\mathcal{G}_k}\| \right).$$

Implémentez une fonction `getgrplasso` à l'aide du package `scoop`<sup>1</sup>.

## 3 Comparaison des méthodes

### 3.1 Évaluation des performances

On s'intéresse aux performances des méthodes à la fois en terme de capacité prédictive et en terme de sélection de variables. On considère à cet égard les grandeurs suivantes:

- l'erreur quadratique moyenne de  $\hat{\beta}$  (**mse**)
- l'erreur moyenne de prédiction calculée sur l'ensemble test, (**err**)
- la précision du support estimé  $\hat{A}$  (**acc**,  $(TN + TP)/p$ )
- la sensibilité de  $\hat{A}$  (**spe**,  $TP/(FN + TP)$ )
- la spécificité de  $\hat{A}$  (**sen**,  $TN/(TN + FP)$ )

On a noté TP, FP, TN, FN respectivement pour *true positive*, *false positive*, *true negative* et *false negative*.

Écrire une fonction `getPerformance` qui calcule tous ces indices pour un estimateur  $\hat{\beta}$  donné.

### 3.2 Planning de simulations

Créer un script de simulation pour chaque scénario de matrice de covariance des prédicteurs  $\Sigma$ , en commençant par exemple par le cas orthogonal.

Chaque simulation doit renvoyer un `data.frame` de la forme suivante, afin de faciliter le tracé des résultats à l'aide du package `ggplot2`.

```
##      method mse  err acc sen spe n.p  r2 simu
## 11   lasso  1.2 -1.85 0.92 0.9 0.75 0.5 0.75  1
## 12   ridge 0.06 -2.37 0.92 0.9 0.75 0.5 0.75  1
## 13 adalasso 1.54  1.03 0.92 0.9 0.75 0.5 0.75  1
```

Écrire une fonction `one.simu(i)` permettant d'effectuer la simulation numéro  $i$  pour toutes les méthodes et pour toutes les valeurs des paramètres de simulation que vous aurez choisis (commencer doucement...). Cette fonction sera ensuite facilement parallélisable (par exemple avec le package `parallel`<sup>2</sup>).

---

1. à télécharger sur ma page web

2. attention ne fonctionne pas sous Windows !

### 3.2.1 Interprétations des résultats

- Représentez les boxplots des indicateurs de performance en fonction de  $n/p$  et de la valeur du  $R^2$ . On utilisera le package **ggplot2**.
- Quels sont les effets de la grande dimension sur les performances des estimateurs (en estimation, en sélection) ?
- Explorer les différents scénarios. Y a-t-il des méthodes plus adaptées à certains scénarios ? Si oui, pourquoi ?

### 3.3 Pour aller plus loin (et pour ceux qui sont en avance)

*Cette partie a été rédigée par Franck Picard, merci à lui.*

Nous nous interrogerons sur les conditions sur  $n$ ,  $p$  et  $p_0$  pour que le lasso détecte bien les entrées nulles et non-nulles de  $\beta^*$ . Dans un article publié en 2009 (IEEE Transactions on Information Theory, 55:2183–2202, May 2009), M. Wainwright propose des conditions nécessaires et suffisantes pour que le lasso soit consistant en sélection pour le support signé. On notera  $\mathbb{S}_\pm(\beta)$  le vector de signes de  $\beta$  défini tel que:

$$\mathbb{S}_\pm(\beta_i) = \begin{cases} +1 & \text{si } \beta_i > 0 \\ -1 & \text{si } \beta_i < 0 \\ 0 & \text{si } \beta_i = 0 \end{cases}$$

Dans son article M. Wainwright démontre l'existence de deux constantes dépendant de  $\Sigma = \mathbb{V}(X)$ ,  $0 < \theta_\ell(\Sigma) \leq \theta_u(\Sigma) < \infty$  telles que pour une valeur

$$\lambda_n = \sqrt{\frac{2\sigma^2 \log(p_0) \log(p - p_0)}{n}}$$

du paramètre de régularisation du lasso,

- si  $n/(2p_0(\log(p - p_0))) > \theta_u(\Sigma)$  alors il est toujours possible de trouver une valeur du paramètre de régularisation  $\lambda$  telle que le lasso a une solution unique  $\hat{\beta}$  telle que  $\mathbb{P}\{\mathbb{S}_\pm(\beta^*) = \mathbb{S}_\pm(\hat{\beta})\}$  tend vers 1.
- si  $n/(2p_0(\log(p - p_0))) < \theta_\ell(\Sigma)$ , alors quelle que soit la valeur du paramètre de régularisation  $\lambda > 0$ , aucune des solutions du lasso ne spécifie correctement le support signé de  $\beta^*$ ,  $\mathbb{P}\{\mathbb{S}_\pm(\beta^*) = \mathbb{S}_\pm(\hat{\beta})\}$  tend vers 0.

Dans son article, M. Wainwright propose d'appeler la quantité  $n/(2p_0(\log(p - p_0)))$  "taille d'échantillon normalisée". C'est un indicateur qui combine les informations nécessaires à la consistance du lasso. Dans la suite, nous nous placerons dans le cas  $\Sigma = I$ , avec  $\theta_\ell(I) = \theta_u(I) = 1$ .

#### 3.3.1 Questions

- Pour les paramètres suivants, étudiez l'évolution de l'accuracy en terme de support en fonction de la taille d'échantillon normalisée, pour le lasso avec la valeur théorique de  $\lambda$  proposée ci-dessus.  $p \in \{128, 256, 512\}$ ,  $n \in \{100, \dots, 1000\}$ ,  $p_0 = \lceil 0.4 \times p \rceil$ ,  $\beta_0^* = 0.5$ ,  $\sigma = 0.5$ .

- Comparer ces performances avec celles du lasso utilisant un  $\lambda$  calibré par validation croisée, et celles du lasso adaptatif (également avec  $\lambda$  calibré par validation croisée). Discutez les différences de comportement de l'accuracy.

## 4 Analyse de jeux de données omiques

### 4.1 Jeu de données "HIV"

#### 4.1.1 Description

Jeu de données de génotypes associés au niveau du virus du VIH dans le sang. 605 individus ont été génotypés pour 300 SNPs.

#### 4.1.2 Format

Lors du chargement des données, deux objets sont créés :

1. X - une matrice 605x300 donnant les génotypes de 605 individus pour 300 SNPs.
2. y - un vecteur de taille 605 donnant le niveau du virus du VIH dans le sang, pour chaque individu.

```
load("data/HIVdata.rda")
ls(); str(X); str(y)
```

```
## [1] "X" "y"
```

```
## num [1:605, 1:300] 3 2 2 3 2 2 1 1 3 2 ...
```

```
## - attr(*, "dimnames")=List of 2
```

```
## ..$ : chr [1:605] "030101" "030102" "060101" "060102" ...
```

```
## ..$ : chr [1:300] "rs1264550" "rs9261947" "rs1079541" "rs9295871" ...
```

```
## num [1:605] 4.56 5.18 5.2 6.18 6.11 4.64 2.91 6 5.9 4.74 ...
```

#### 4.1.3 objectifs :

1. Sélectionner les SNPs qui sont associés au niveau du virus dans le sang.
2. Sélectionner des ensembles de SNPs intégrant une forme de structure de corrélation dans les données.
3. Évaluer les performances prédictives du modèles

#### 4.1.4 référence

Dalmasso, C., Carpentier, W., Meyer, L., Rouzioux, C., Goujard, C., Chaix, M. L., ... & Theodorou, I. (2008). Distinct genetic loci control plasma HIV-RNA and cellular HIV-DNA levels in HIV-1 infection: the ANRS Genome Wide Association 01 study. *PloS one*, 3(12), e3907-e3907.

## 4.2 Jeu de données "colorectal"

### 4.2.1 Description

Jeu de données de niveaux d'expression de gène associés à des tissus tumoraux ou sain dans le cancer du colon. 62 tissus ont été analysés pour 2000 gènes ou assimilés.

### 4.2.2 Format

Lors du chargement des données, trois objets sont créés :

1. X - une matrice 62x2000 donnant les niveaux d'expression (log transformés) relevé dans les tissus du colons de 62 patients.
2. y - un vecteur de taille 62 indiquant le statut du tissus (-1: tumoral, 1: sain).
3. genes.info - une liste de longueur 2000 donnant des informations sur les 2000 gènes considérés.

```
load("data/colorectal.rda")  
ls()
```

```
## [1] "genes.info" "X"          "y"
```

### 4.2.3 objectifs :

1. Sélectionner les gènes liés au cancer colo-rectal à l'aide d'un modèle gaussien
2. Sélectionner les gènes liés au cancer colo-rectal à l'aide d'un modèle logistique.
3. Prédire le statut d'un tissu.

### 4.2.4 référence

U. Alon, N. Barkai, D. A. Notterman, K. Gish, S. Ybarra, D. Mack, and A. J. Levine, "Broad patterns of gene expression revealed by clustering of tumor and normal colon tissues probed by oligonucleotide arrays", *PNAS*, vol. 96, 1999.

## 4.3 Jeu de données 'Bardet'

### 4.3.1 description

Jeu de données simplifié d'expression de gènes associées au syndrome de Bardet-Biedl. Les échantillons ont été biopsiés à partir de tissus d'oeil de 120 rats.



### 4.3.2 format

Lors du chargement des données, la liste `bardet` est créé, contenant deux objets:

1. `x` - une matrice 120 x 100, donnant les expressions associées à 120 rats pour 100 sondes associées à 20 gènes. 5 sondes consécutives correspondent au même gène.
2. `y` - un vecteur de taille 120 donnant le niveau d'expression du gène TRIM32.

```
load("data/bardet.rda")
str(bardet)
```

```
## List of 2
## $ x: num [1:120, 1:100] 0.44705 0.46684 0.01498 0 0.00281 ...
## $ y: num [1:120] 8.42 8.36 8.41 8.29 8.27 ...
```

### 4.3.3 objectifs

1. Sélectionner les sondes les plus prédictives de l'expression de TRIM32.
2. Opérer une sélection de sonde "par groupe" associée à chaque gène, sachant que 5 prédicteurs consécutifs sont des sondes associées au même gène.
3. Évaluer les performances prédictives du modèle

### 4.3.4 référence

T. Scheetz, K. Kim, R. Swiderski, A. Philp, T. Braun, K. Knudtson, A. Dorrance, G. DiBona, J. Huang, T. Casavant, V. Sheffield, E. Stone .Regulation of gene expression in the mammalian eye and its relevance to eye disease. *Proceedings of the National Academy of Sciences of the United States of America*, 2006.