

# Data Science - High dimensional regression

## Summary

Linear models are popular methods for providing a regression of a response variable  $Y$ , that depends on covariates  $(X_1, \dots, X_p)$ . We introduce the problem of high dimensional regression and provide some real examples where standard linear models methods are not well suited. Then, we propose some statistical resolution through the LASSO estimator and the Boosting algorithm. A practical session is proposed in the end of this Lecture, since the knowledge of these modern methods is needed in many fields.

## 1 Back to linear models

### 1.1 Sum of squares minimization

In a standard linear model, we have at our disposal  $(X_i, Y_i)$  supposed to be linked with

$$Y_i = X_i^t \beta^* + \epsilon_i, 1 \leq i \leq n.$$

In particular, each observation  $X_i$  is described by  $p$  variables  $(X_i^1, \dots, X_i^p)$ , so that the former relation should be understood as

$$Y_i = \sum_{j=1}^p \beta_j^* X_i^j + \epsilon_i, 1 \leq i \leq n.$$

We aim to recover the unknown  $\beta^*$ .

- A classical “optimal” estimator is the MLE :

$$\hat{\beta}_{MLE} := \arg \max_{\beta \in \mathbb{R}^p} L(\beta, (X_i, Y_i)_{1 \leq i \leq n}),$$

where  $L$  denotes the likelihood of the parameter  $\beta$  given the observations  $(X_i, Y_i)_{1 \leq i \leq n}$ .

- Generically,  $(\epsilon_i)_{1 \leq i \leq n}$  is assumed to be i.i.d. replications of a centered and squared integrale noise

$$\mathbb{E}[\epsilon] = 0 \quad \mathbb{E}[\epsilon^2] < \infty.$$

A standard assumption even relies on the Gaussian structure of the errors  $\epsilon_i \sim \mathcal{N}(0, 1)$  and in this case, the log-likelihood leads to the minimization of the sum of square and

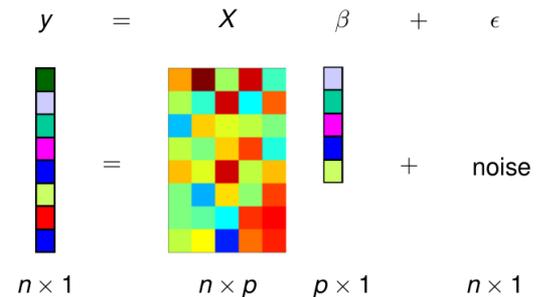
$$\hat{\beta}_{MLE} := \arg \min_{\beta \in \mathbb{R}^p} \underbrace{\sum_{i=1}^n \|Y_i - X_i^t \beta\|^2}_{:= J(\beta)}. \tag{1}$$

## 1.2 Matricial traduction & resolution

From a matricial point of view, the linear model can we written as follows :

$$Y = X \beta_0 + \epsilon, \quad Y \in \mathbb{R}^n, X \in \mathcal{M}_{n,p}(\mathbb{R}), \beta_0 \in \mathbb{R}^p$$

In this lecture, we will consider situations where  $p$  varies (typically increases) with  $n$ .



It is an easy exercice to check that (1) leads to

$$\hat{\beta}_{MLE} := (X^t X)^{-1} X^t Y.$$

This can be obtained while remarking that  $J$  is a convex function, that possesses a unique minimizer if and only if  $X^t X$  has a full rank, meaning that  $J$  is indeed strongly convex :

$$D^2 J = X^t X,$$

which is a squared  $p \times p$  symmetric and positive matrix. It is non degenerate if  $X^t X$  has full rank, meaning that necessarily  $p \leq n$ .

PROPOSITION 1. —  $\hat{\beta}_{MLE}$  is an unbiased estimator of  $\beta_0$  such that

• If  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  :  $\frac{\|X(\hat{\beta}_{MLE} - \beta^*)\|_2^2}{\sigma^2} \sim \chi_p^2$

• 
$$\mathbb{E} \left[ \frac{\|X(\hat{\beta}_{MLE} - \beta^*)\|_2^2}{n} \right] = \frac{\sigma^2 p}{n}$$

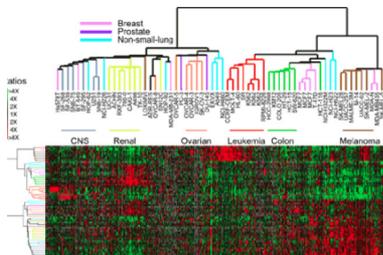
Main requirement :  $X^t X$  must be full rank (invertible) !

### 1.3 Difficulties in large dimensional case

**Example** One measures micro-array datasets built from a huge amount of profile genes expression. From a statistical point of view, we expect to find among the  $p$  variables that describe  $X$  important ones.

Number of genes  $p$  (of order thousands). Number of samples  $n$  (of order hundred).

- $Y_i$  expression level of one gene on sample  $i$
- $X_i = (X_{i,1}, \dots, X_{i,p})$  biological signal (DNA micro-arrays)



Diagnostic help : healthy or ill ?

- Select among the genes meaningful elements : discover a cognitive link between DNA and the gene expression level.
- Find an algorithm with good prediction of the response ?

**Linear model ?** Difficult to imagine :  $p > n$  !

- $X^t X$  is an  $p \times p$  matrix, but its rank is lower than  $n$ . If  $n \ll p$ , then

$$rk(X^t X) \leq n \ll p.$$

- Consequence : the Gram matrix  $X^t X$  is not invertible and even very ill-conditioned (most of the eigenvalues are 0 !)
- The linear model  $\hat{\beta}_{MLE}$  completely fails.
- One standard "improvement" : use the ridge regression with an additional penalty :

$$\hat{\beta}_n^{Ridge} = \arg \min_{\beta \in \mathbb{R}^p} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$$

The ridge regression is a particular case of *penalized* regression. The penalization is still convex w.r.t.  $\beta$  and can be easily solved.

- We will attempt to describe a better suited penalized regression for high dimensional regression.
- Our goal : find a method that permits to find  $\hat{\beta}_n$  such that :
  - Select features among the  $p$  variables.
  - Can be easily computed with numerical softs.
  - Possess some statistical guarantees.

### 1.4 Goals

Important and nowadays questions :

- What is a good framework for high dimensional regression ? **A good model is required.**
- How can we estimate ? **An efficient algorithm is necessary.**
- How can we measure the performances : prediction of  $Y$  ? Feature selection in  $\beta$  ? **What are we looking for ?**
- Statistical guarantees ? **Some mathematical theorems ?**

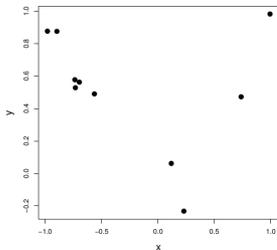
## 2 Penalized regression

### 2.1 Important balance : bias-variance tradeoff

A classical result in statistics states that a good estimator should achieve a balance between bias and variance.

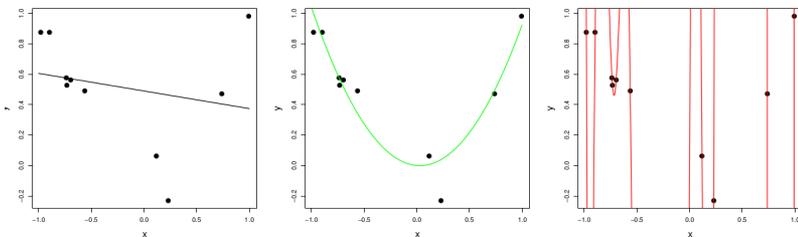
**Example :** In high dimension :

- Optimize the fit to the observed data ?
- Reduce the variability ?

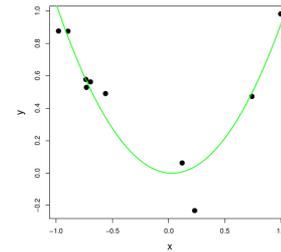


Standard question : find the best curve... In what sense ? Several regressions :

- Left : fit the best line (1-D regression)
- Middle : fit the best quadratic polynomial
- Right : fit the best 10-degree polynomial



Now I am interested in the **prediction** at point  $x = 0.5$ . What is the best ?

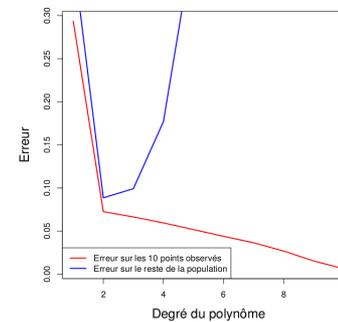


If we are looking for the best possible fit, a high dimensional regressor will be convenient.

Nevertheless, our goal is to generally to **predict**  $y$  for new points  $x$  and a standard matching criterion is

$$C(\hat{f}) := \mathbb{E}_{(X,Y)}[Y - \hat{f}(X)]^2.$$

It is a quadratic loss here, and should be replaced by other criteria (in classification for example).



- When the degree increases, the fit to the observed data (red curve) is always decreasing.
- Over the rest of the population, the generalization error starts decreasing, and after increases.
- Too simple sets of functions cannot contain the good function, and optimization over simple sets **introduces a bias**.

- Too complex sets of functions may contain the good function **but** are too rich and **generates high variance**.

The former balance is illustrated by a very simple theorem.

$$Y = f(X) + \epsilon \quad \text{with} \quad \mathbb{E}[\epsilon] = 0.$$

THÉORÈME 2. — For any estimator  $\hat{f}$ , one has

$$\begin{aligned} C(\hat{f}) = \mathbb{E}[Y - \hat{f}(X)]^2 &= \mathbb{E}[Y - \mathbb{E}[\hat{f}(X)]]^2 \\ &\quad + \mathbb{E}[\mathbb{E}[\hat{f}(X)] - \hat{f}(X)]^2 \\ &\quad + \mathbb{E}[Y - f(X)]^2 \end{aligned}$$

- The **blue** term is a **bias** term.
- The **red** term is a **variance** term.
- The **green** term is the **Bayes risk** and is independent on the estimator  $\hat{f}$ .

Statistical principle :

The empirical squared loss  $\|Y - \hat{f}(X)\|_{2,n}^2$  mimics the **bias**. It is the sum of squares in (1). Important needs to introduce something to quantify the variance of estimation : this is provided by a **penalty term**.

## 2.2 Ridge regression as a preliminary (insufficient) response

**Ridge** Ridge regression is like ordinary linear regression, but it shrinks the estimated coefficients towards zero. The ridge coefficients are defined by solving

$$\hat{\beta}_{Ridge} := \arg \min_{\beta \in \mathbb{R}^p} \|Y - X^t \beta\|_2^2 + \lambda \|\beta\|_2^2$$

Here  $\lambda \geq 0$  is a tuning parameter, which controls the strength of the penalty term. Write  $\hat{\beta}_{Ridge}$  as the ridge solution. Note that :

- When  $\lambda = 0$ , we get the linear regression estimate

- When  $\lambda = +\infty$ , we get  $\hat{\beta}_{Ridge} = 0$
- For  $\lambda$  in between, we are balancing two ideas : fitting a linear model of  $Y$  on  $X$ , and shrinking the coefficients.

**Ridge with intercept** When including an intercept term in the regression, we usually leave this coefficient unpenalized. Otherwise we could add some constant amount  $c$  to the vector  $Y$ , and this would not result in the same solution. Hence ridge regression with intercept solves

$$\hat{\beta}_{Ridge} := \arg \min_{c \in \mathbb{R}, \beta \in \mathbb{R}^p} \|Y - c - X^t \beta\|_2^2 + \lambda \|\beta\|_2^2$$

If we center the columns of  $X$ , then the intercept estimate ends up just being  $\hat{c} = \bar{Y}$ , so we usually just assume that  $Y$  and  $X$  have been centered and don't include an intercept.

Also, the penalty term  $\|\beta\|_2^2$  is unfair is the predictor variables are not on the same scale. (Why ?) Therefore, if we know that the variables are not measured in the same units, we typically scale the columns of  $X$  (to have sample variance 1), and then we perform ridge regression.

**Bias and variance of the ridge regression** The bias and variance are not quite as simple to write down for ridge regression as they were for linear regression (see Proposition 1) but closed-form expressions are still possible. The general trend is :

- The bias increases as  $\lambda$  (amount of shrinkage) increases
- The variance decreases as  $\lambda$  (amount of shrinkage) increases

Think : what is the bias at  $\lambda = 0$  ? The variance at  $\lambda = +\infty$  ?



- Assume that the effective support of  $\beta_0$  is known, then

$$y = X\beta + \epsilon \implies y = X_S \beta_S + \epsilon$$

- If  $S$  is the support of  $\beta_0$ , maybe  $X_S^t X_S$  is full rank, and linear model can be applied.

Major issue : How could we find  $S$  ?

### 3.2 Lasso relaxation

Ideally, we would like to find  $\beta$  such that

$$\hat{\beta}_n = \arg \min_{\beta: \|\beta\|_0 \leq s} \|Y - X\beta\|_2^2,$$

meaning that the minimization is embedded in a  $\ell_0$  ball.

In the previous lecture, we have seen that it is a constrained minimization problem of a convex function ... A dual formulation is

$$\arg \min_{\beta: \|Y - X\beta\|_2 \leq \epsilon} \{\|\beta\|_0\}$$

But : The  $\ell_0$  balls are not convex and not smooth !

- First (illusiv) idea : explore all  $\ell_0$  subsets and minimize ! Bullshit since :

$$C_p^s \text{ subsets and } p \text{ is large !}$$

- Second idea (existing methods) : run some heuristic and greedy methods to explore  $\ell_0$  balls and compute an approximation of  $\hat{\beta}_n$ . (See below)
- Good idea : use a convexification of the  $\|\cdot\|_0$  norm (also referred to as a convex relaxation method). How ?

Idea of the convex relaxation : instead of considering a variable  $z \in \{0, 1\}$ , imagine that  $z \in [0, 1]$ .

DÉFINITION 3. — [Convex Envelope] The convex envelope  $f^*$  of a function  $f$  is the largest convex function below  $f$ .

THÉORÈME 4. — [Envelope of  $\beta \mapsto \|\beta\|_0$ ]

- On  $[-1, 1]^d$ , the convex envelope of  $\beta \mapsto \|\beta\|_0$  is  $\beta \mapsto \|\beta\|_1$ .
- On  $[-R, R]^d$ , the convex envelope of  $\beta \mapsto \|\beta\|_0$  is  $\beta \mapsto \frac{\|\beta\|_1}{R}$ .

Idea : Instead of solving the minimization problem :

$$\forall s \in \mathbb{N} \quad \min_{\|\beta\|_0 \leq s} \|Y - X\beta\|_2^2, \tag{2}$$

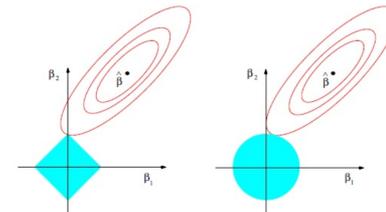
we are looking for

$$\forall C > 0 \quad \min_{\|\cdot\|_0^*(\beta) \leq C} \|Y - X\beta\|_2^2, \tag{3}$$

What's new ?

- The function  $\|\cdot\|_0^*$  is convex and thus the above problem is a convex minimization problem with convex constraints.
- Since  $\|\cdot\|_0^*(\beta) \leq \|\beta\|_0$ , it is rather reasonable to obtain sparse solutions. In fact, solutions of (3) with a given  $C$  provide a lower bound of solutions of (2) with  $s \leq C$ .
- If we are looking for good solutions of (2), then there must exists even better solution to (3).

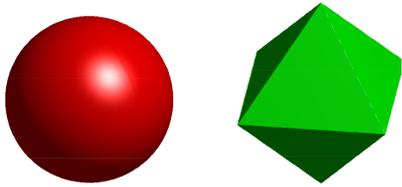
### 3.3 Geometrical interpretation (in 2 D)



Left : Level sets of  $\|Y - X\beta\|_2^2$  and intersection with  $\ell^1$  ball. Right : Same with  $\ell^2$  ball.

The left constraint problem is likely to obtain a sparse solution. Oppositely, the right constraint no !

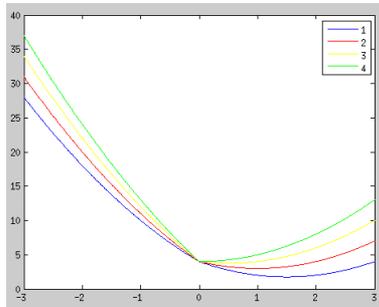
In larger dimensions the balls are even more different :



- **Analytic point of view** : why does the  $\ell^1$  norm induce sparsity ?
- From the KKT conditions (see Lecture 1), it leads to a **penalized criterion**

$$\min_{\beta \in \mathbb{R}^p: \|\beta\|_1 \leq C} \|Y - X\beta\|_2^2 \iff \min_{\beta \in \mathbb{R}^p} \underbrace{\|Y - X\beta\|_2^2}_{\text{Mimics the bias}} + \underbrace{\lambda \|\beta\|_1}_{\text{Controls the variance}}$$

- In the 1d case :  $\arg \min_{\alpha \in \mathbb{R}} \frac{1}{2} |x - \alpha|^2 + \lambda|x|$  :  
 $:= \varphi_\lambda(x)$



- The minimal value of  $\varphi_\lambda$  is reached at point  $x^*$  when  $0 \in \partial\varphi_\lambda(x^*)$ . We can check that  $x^*$  is minimal iff
    - $x^* \neq 0$  and  $(x^* - \alpha) + \lambda \text{sgn}(x^*) = 0$ .
    - $x^* = 0$  and  $d\varphi_\lambda^+(0) > 0$  and  $d\varphi_\lambda^-(0) < 0$ .
- PROPOSITION 5. — [Analytical minimization of  $\varphi_\lambda$ ]

$$x^* = \text{sgn}(\alpha)[|\alpha| - \lambda]_+ = \arg \min_{x \in \mathbb{R}} \left\{ \frac{1}{2} |x - \alpha|^2 + \lambda|x| \right\}$$

- For large values of  $\lambda$ , the minimum of  $\varphi_\lambda$  is reached at point 0.

### 3.4 Lasso estimator

We introduce the *Least Absolute Shrinkage and Selection Operator* :

$$\forall \lambda > 0 \quad \hat{\beta}_n^{Lasso} = \arg \min_{\beta \in \mathbb{R}^p} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

The above criterion is **convex** w.r.t.  $\beta$ .

- Efficient algorithms to solve the LASSO, even for very large  $p$ .
- The minimizer may not be unique since the above criterion is not strongly convex.
- Predictions  $X\hat{\beta}_n^{Lasso}$  are always unique.
- $\lambda$  is a **penalty constant that must be carefully chosen**.
- A **large value of  $\lambda$**  leads to a **very sparse** solution, with an important bias.
- A low value of  $\lambda$  yields overfitting with no penalization (too much important variance).
- We will see that a careful balance between  $s$ ,  $n$  and  $p$  exists. These parameters as well as the variance of the noise  $\sigma^2$  influence a “good” choice of  $\lambda$ .

Alternative formulation :

$$\hat{\beta}_n^{Lasso} = \arg \min_{\beta \in \mathbb{R}^p: \|\beta\|_1 \leq C} \|Y - X\beta\|_2^2$$

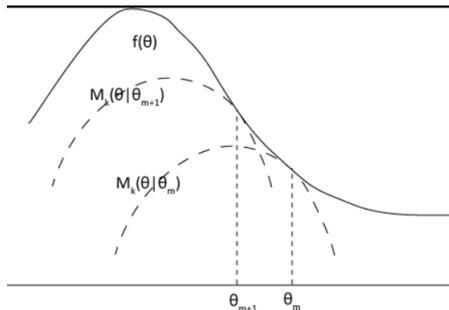
### 3.5 Principle of the MM algorithm

Algorithm is needed to solve the minimization problem

$$\arg \min_{\beta \in \mathbb{R}^p} \underbrace{\|Y - X\beta\|_2^2 + \lambda \|\beta\|_1}_{:= \varphi_\lambda(\beta)}$$

An efficient method follows the method of "**Minimize Majorization**" and is referred to as MM method.

- MM are useful for the minimization of a convex function/maximization of a concave one.
- Geometric illustration



- Idea : Build a sequence  $(\beta_k)_{k \geq 0}$  that converges to the minimum of  $\varphi_\lambda$ .
- A particular case of such a method is encountered with the E.M. algorithm useful for clustering and mixture models.
- MM algorithms are powerful, especially they can convert non-differentiable problems to smooth ones.

1. A function  $g(\beta, \beta_k)$  is said to *majorize*  $f$  at point  $\beta_k$  if

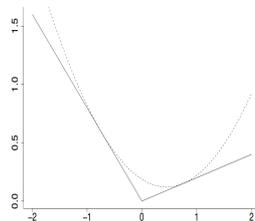
$$g(\beta_k | \beta_k) = f(\beta_k) \quad \text{and} \quad g(\beta | \beta_k) \geq f(\beta), \forall \beta \in \mathbb{R}^p.$$

2. Then, we define

$$\beta_{k+1} = \arg \min_{\beta \in \mathbb{R}^p} g(\beta | \beta_k)$$

3. We wish to find each time a function  $g(\cdot, \beta_k)$  whose minimization is easy.

4. An example with a quadratic majorizer of a non-smooth function :



5. **Important remark** : The MM is a descent algorithm :

$$\begin{aligned} f(\beta_{k+1}) &= g(\beta_{k+1} | \beta_k) + f(\beta_{k+1}) - g(\beta_{k+1} | \beta_k) \\ &\leq g(\beta_k | \beta_k) = f(\beta_k) \end{aligned}$$

(4)

### 3.6 MM algorithm for the LASSO

We can deduce for the LASSO the [coordinate descent algorithm](#)

1. Define a sequence  $(\beta_k)_{k \geq 0} \iff$  find a suitable majorization.
2.  $g : \beta \mapsto \|Y - X\beta\|^2$  is convex, whose Hessian matrix is  $X^t X$ . A Taylor's expansion leads to

$$\forall y \in \mathbb{R}^p \quad g(y) \leq g(x) + \langle \nabla g(x), y - x \rangle + \rho(X) \|y - x\|^2,$$

where  $\rho(X)$  is the spectral radius of  $X$ .

3. We are naturally driven to upper bound  $\varphi_\lambda$  as

$$\begin{aligned} \varphi_\lambda(\beta) &\leq \varphi_\lambda(\beta_k) + \langle \nabla g(\beta_k), \beta - \beta_k \rangle + \rho(X) \|\beta - \beta_k\|_2^2 + \lambda \|\beta\|_1 \\ &= \psi(\beta_k) + \rho(X) \left\| \beta - \left( \beta_k - \frac{\nabla g(\beta_k)}{\rho(X)} \right) \right\|_2^2 + \lambda \|\beta\|_1 := \varphi_k(\beta) \end{aligned}$$

The important point with this majorization is that it is “tensorized” : each coordinates acts separately on  $\varphi_k(\beta)$ .

4. To minimize the majorization of  $\varphi_\lambda$ , we then use the above proposition of soft-thresholding :

- Define

$$\tilde{\beta}_k^j := \beta_k^j - \nabla g(\beta_k)^j / \rho(X).$$

- Compute

$$\beta_{k+1}^j = \text{sgn}(\tilde{\beta}_k^j) \max \left[ |\tilde{\beta}_k^j| - \frac{2\lambda}{\rho(X)}, 0 \right]_+$$

## 4 Running the Lasso

### 4.1 Choice of the regularization parameter

It is an important issue to obtain a good performance of the method, and could be almost qualified as a “tarte à la crème” issue.

We won't provide a sharp presentation of the best known results to keep the level understandable.

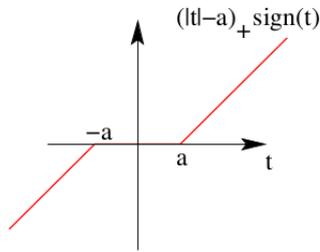
It is important to have in mind **the extremely favorable situation of an almost orthogonal design** :

$$\frac{X^t X}{n} \simeq I_p.$$

In this case solving the lasso is equivalent to

$$\min_w \frac{1}{2n} \|X^t y - w\|_2^2 + \lambda \|w\|_1$$

Solutions are given by ST (Soft-Thresholding) :



$$w_j = ST_\lambda \left( \frac{1}{n} X_j^t y \right) = ST_\lambda \left( \theta_j^0 + \frac{1}{n} X_j^t \epsilon \right)$$

We would like to **keep the useless coefficients to 0**, which requires that

$$\lambda \geq \frac{1}{n} X_j^t \epsilon, \forall j \in J_0^c.$$

The random variables  $\frac{1}{n} X_j^t \epsilon$  are i.i.d. with a variance  $\sigma^2/n$ .

PROPOSITION 6. — *The expectation of the maximum of  $p - s$  Gaussian standard variables is*

$$\mathbb{E} \left[ \max_{1 \leq i \leq p-s} X_i \right] \sim \sqrt{2 \log(p-s)}.$$

We are naturally driven to the choice

$$\lambda = A\sigma \sqrt{\frac{\log p}{n}}, \quad \text{with} \quad A > \sqrt{2}.$$

Precisely :

$$\mathbb{P} \left( \forall j \in J_0^c : |X_j^t \epsilon| \leq n\lambda \right) \geq 1 - p^{1-A^2/2}.$$

## 4.2 Theoretical consistency

An additional remark is that we expect  $ST_\lambda \mapsto Id$  to obtain a consistency result. It means that  $\lambda \mapsto 0$ , so that

$$\frac{\log p}{n} \mapsto 0$$

Hence, a good behaviour of the lasso can be expected only if we have the next dimensional settings :

$$p_n = \mathcal{O}(\exp(n^{1-\xi})).$$

THÉORÈME 7. — *Assume that  $\log p \ll n$ ,  $X$  has norm 1 and  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ , then under a coherence assumption on the design matrix  $X^t X$ , one has*

- i) *With high probability,  $J(\hat{\theta}_n) \subset J_0$ .*
- ii) *There exists  $C$  such that, with high probability,*

$$\frac{\|X(\theta_n - \theta_0)\|_2^2}{n} \leq \frac{C}{\kappa^2} \frac{\sigma^2 s_0 \log p}{n},$$

where  $\kappa^2$  is a positive constant that depends on the correlations in  $X^t X$ .

One can also find results on the exact support recovery, as well as some weaker results without any coherence assumption.

N.B. : Such a coherence is measured through the almost orthogonality of the columns of  $X$ . It can be traduced in terms of

$$|\sup_{i \neq j} \langle X_i, X_j \rangle| \leq \epsilon.$$

## 4.3 Practical calibration of $\lambda$

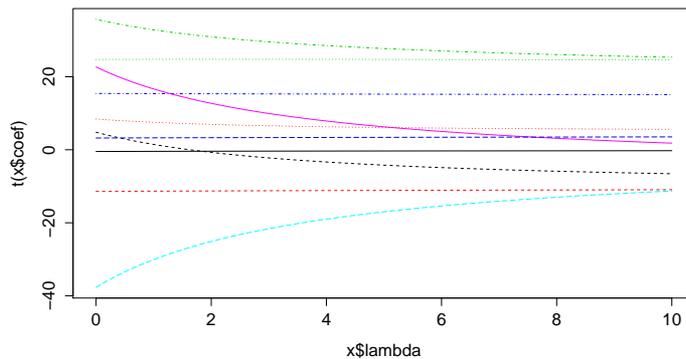
In practice,  $\lambda$  is generally chosen according to a criterion that is *data dependent*, e.g. a criterion that is calibrated on the observations through a cross-validation approach. In general, the packages implement this automatic choice of the regularization parameter with a CV option.

## 5 Numerical example

### 5.1 Very brief R code

#### 5.1.1 About the use of the Ridge regression

```
library(lars)
data(diabetes)
library(MASS)
diabetes.ridge <- lm.ridge(diabetes$y ~ diabetes$x,
  lambda=seq(0,10,0.05))
plot(diabetes.ridge, lwd=3)
```



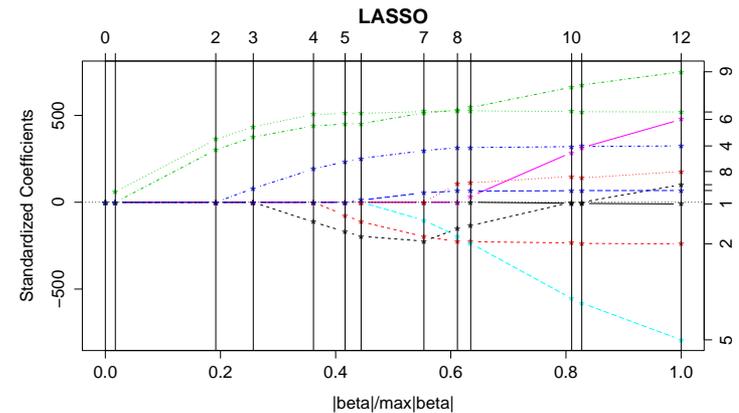
We can see that the influence of the regularization parameter  $\lambda$  of the ridge regression is important ! But a good choice of  $\lambda$  is difficult and should be data-driven. That is why a cross-validation procedure is needed. Does the ridge regression performs variable selection ?

#### 5.1.2 About the use of the Lasso regression

```
library(lars)
data(diabetes)
diabetes.lasso = lars(diabetes$x, diabetes$y,
```

```
  type='lasso')
plot(diabetes.lasso)
```

Lars algorithm : solves the Lasso less efficiently than the coordinate descent algorithm.



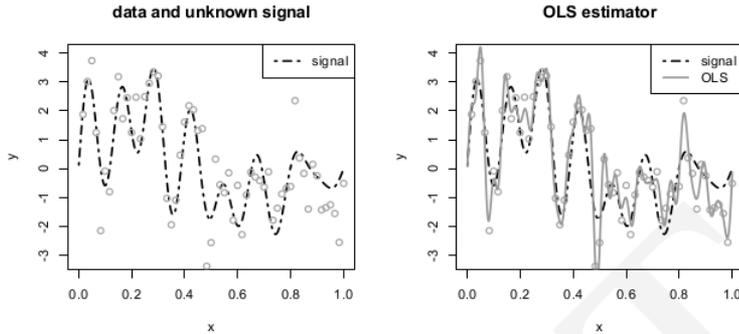
Typical output of the Lars software :

- The greater  $\ell^1$  norm, the lower  $\lambda$
- Sparse solution with small values of the  $\|\cdot\|_1$  norm.

We can see that each variable of the diabetes dataset enter the model successively as long as  $\lambda$  decreases to 0. Again, the choice of  $\lambda$  should be done carefully with a data-driven criterion.

### 5.2 Removing the bias of the Lasso

Signal processing example :



We have  $n = 60$  noisy observations  $Y(i) = f(i/n) + \epsilon_i$ .  $f$  is an unknown periodic function defined on  $[0, 1]$ , sampled at points  $(i/n)$ .  $\epsilon_i$  are independent realizations of Gaussian r.v. We use the 50 first Fourier coefficients :

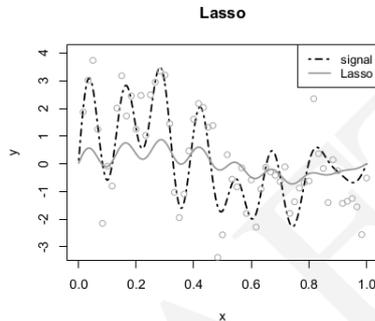
$$\varphi_0(x) = 1, \quad \varphi_{2j}(x) = \sin(2j\pi x) \quad \varphi_{2j+1}(x) = \cos(2j\pi x),$$

to approximate  $f$ . The OLS estimator is

$$\hat{f}^{OLS}(x) = \sum_{j=1}^p \hat{\beta}_j^{OLS} \varphi_j(x) \quad \text{with} \quad \hat{\beta}^{OLS} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - \sum_{j=0}^p \beta_j \varphi_j(i/n))^2.$$

The OLS does not perform well on this example.

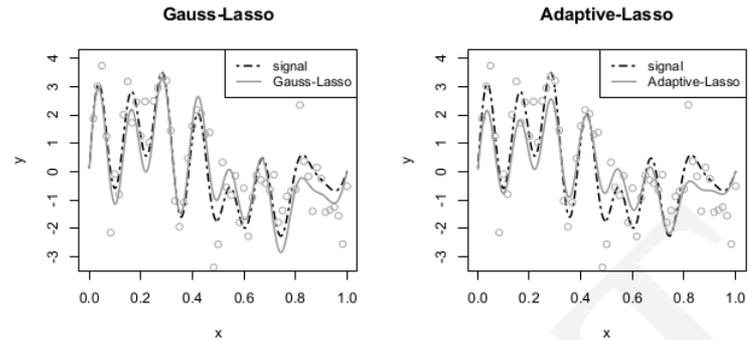
We experiment here the Lasso estimator with  $\lambda = 3\sigma\sqrt{\frac{2\log p}{n}}$  and obtain



We define

$$\hat{f}^{\text{Gauss}} = \pi_{\hat{J}_0}(Y) \quad \text{with} \quad \hat{J}_0 = \text{Supp}(\hat{\theta}^{\text{Lasso}}),$$

where  $\pi_{\hat{J}_0}$  is the  $\mathbb{L}^2$  projection of the observations on the features selected by the Lasso.



The Adaptive Lasso is almost equivalent :

$$\beta^{\text{Adaptive Lasso}} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \|Y - X\beta\|_2^2 + \mu \sum_{j=1}^p \frac{|\beta_j|}{|\hat{\beta}_j^{\text{Gauss}}|} \right\}$$

This minimization remains convex and the penalty term aims to mimic the  $\ell^0$  penalty.

The Adaptive Lasso is very popular and tends to select more accurately the variables than the Gauss-Lasso estimator.

- Lasso estimator reproduces the oscillations of  $f$  but these oscillations are shrunked to 0.
- When considering the initial minimization problem, the  $\ell^1$  penalty select nicely the good features, but introduces also a bias (introduces a shrinkage of the parameters).
- Strategy : select features with the Lasso and run an OLS estimator using the good variables.

## 6 Homework

Length limitation : 6 pages !

Deadline : 22th of February.

Group of 2 students allowed.

1. You are asked first to follow the practical session on the Cookie database, that can be found in Moodle.

You will need to install some packages with R.

2. Once you finish this practical session, please produce a short summary with a few quantity of numerical illustrations and comments. This content could form the first part of your report. A special attention should be paid to Ridge, Lasso regression and cross validation. Since this last method has not been described in this lecture, I expect a brief description of the method in the report and a discussion about its use.
3. To produce the same document, you are asked to complement your production with a gentle introduction to either

(a) Weak greedy algorithm (Boosting methods).

- V. N. Temlyakov. Weak Greedy Algorithms. *Advances in Computational Mathematics*, 12(2,3) :213-227, 2000.
- J. A. Tropp. Greed is good : algorithmic results for sparse approximation. *IEEE Trans. Inform. Theory*, 50(10) :2231-2242, 2004.
- M. Champion, C. Cierco-Ayrolles, S. Gadat, M. Vignes, Sparse regression and support recovery with L2-Boosting algorithms. *Journal of Statistical Planning and Inference* <http://dx.doi.org/10.1016/j.jspi.2014.07.006>.
- P. Buhlmann and B. Yu. Boosting. *Wiley Interdisciplinary Reviews : Computational Statistics* 2, pages 69-74, 2010.

(b) Aggregation method (Exponential Weighting Aggregation)

- A. Dalalyan, A. Tsybakov (2012). Sparse regression learning by aggregation and Langevin Monte-Carlo. *J. Comput. System Sci.*, 78(5), pp. 1423-1443.
- A. Dalalyan and A. Tsybakov (2009). Sparse Regression Learning by Aggregation and Langevin Monte-Carlo. In COLT 2009

- The 22nd Conference on Learning Theory, Montreal, Quebec, Canada, June 18-21, 2009, pp. 1-10.

- J. Salmon and A. Dalalyan (2011). Optimal aggregation of affine estimators. In *Journal of Machine Learning Research - Proceedings Track* 19, pp. 635-660.

The description of the algorithm could then form the second part of your report. You are asked to mainly describe the general behaviour of the algorithm.

4. The last part of your report should compare this new algorithm with the Lasso from a numerical point of view. You can either code the Boosting procedure (in matlab, R, python, ...) or use a package. Please, provide a reproducible simulation with sources code. This comparison should be understood from :
  - speed of the algorithm
  - statistical accuracy of the method
  - ability to handle large datasets
  - easy to use calibration of parameters