

## Some reminders about the Lasso estimator.

Marius Kwemou <sup>1,2</sup>

Supervised by: M-L. Taupin <sup>1</sup> And A.K. Diongue <sup>2</sup>

<sup>1</sup>Université d'Evry val d'Essonne (France)  
Laboratoire de Statistique et génome.

<sup>2</sup>Université Gaston Berger de Saint-Louis (Sénégal)  
LERSTAD.

February, 2014

# Outline

- 1 Model
- 2 High dimensional data situations ( $d \gg n$ )
- 3 Oracle inequality

# Outline

- 1 Model
- 2 High dimensional data situations ( $d \gg n$ )
- 3 Oracle inequality

# linear regression model

- One observes  $n$  independent and identically distributed (i.i.d) couples  $(z_1, Y_1), \dots, (z_n, Y_n)$  from the joint distribution of  $(Z, Y) \in \mathbb{R}^d \times \mathbb{R}$  such that:

$$y_i = z_i^T \beta_0 + \epsilon_i \quad (1)$$

- $\beta_0$  is the unknown parameter to estimate.
- Matrix notation, let  $(y, X) \in \mathbb{R}^n \times \mathcal{M}_{n,d}$

$$y = X\beta_0 + \epsilon \quad (2)$$

# linear regression model

- One observes  $n$  independent and identically distributed (i.i.d) couples  $(z_1, Y_1), \dots, (z_n, Y_n)$  from the joint distribution of  $(Z, Y) \in \mathbb{R}^d \times \mathbb{R}$  such that:

$$y_i = z_i^T \beta_0 + \epsilon_i \quad (1)$$

- $\beta_0$  is the unknown parameter to estimate.
- Matrix notation, let  $(y, X) \in \mathbb{R}^n \times \mathcal{M}_{n,d}$

$$y = X\beta_0 + \epsilon \quad (2)$$

# Estimation

- If  $n > d$ , the estimator  $\hat{\beta}$  of  $\beta_0$  is defined as follows

$$\hat{\beta} := \operatorname{argmax}_{\beta \in \mathbb{R}^d} \ell(\beta),$$

where  $\ell$  is the likelihood function.

→ if  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$ ,  $\hat{\beta}$  is the well known OLS estimator.

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

# Outline

- 1 Model
- 2 High dimensional data situations ( $d \gg n$ )
- 3 Oracle inequality

# Penalization

If  $d \gg n$ , direct maximization of likelihood can lead to:

- **Overfitting**: the classifier can only behave well in training set, and can be bad in test set.
- **Unstable**: since empirical risk is data dependent, hence random, small change in the data can lead to very different estimators.

→ Penalization

$$\underbrace{\hat{L}(\beta)}_{\text{empirical risk}} + \underbrace{\text{pen}(\beta)}_{\text{penalty}}$$



# Penalization

If  $d \gg n$ , direct maximization of likelihood can lead to:

- **Overfitting**: the classifier can only behave well in training set, and can be bad in test set.
- **Unstable**: since empirical risk is data dependent, hence random, small change in the data can lead to very different estimators.

→ Penalization

$$\underbrace{\hat{L}(\beta)}_{\text{empirical risk}} + \underbrace{\text{pen}(\beta)}_{\text{penalty}}$$

## $\ell_0$ -penalization

- Akaike information criterion (AIC)

$$AIC(\beta) = 2\hat{L}(\beta) + 2\|\beta\|_0,$$

- Bayesian information criterion (BIC)

$$BIC(\beta) = 2\hat{L}(\beta) + 2 \log(n)\|\beta\|_0,$$

where  $\|\beta\|_0 = \text{card}\{j \in 1, \dots, d, \beta_j \neq 0\}$ .

★ Produce interpretable models.

↪ Non-convex optimization.

- ▶  $\text{card}(\mathcal{P}(\{1, \dots, d\})) = 2^d$  models,
- ▶ stepwise selection.

## $\ell_0$ -penalization

- Akaike information criterion (AIC)

$$AIC(\beta) = 2\hat{L}(\beta) + 2\|\beta\|_0,$$

- Bayesian information criterion (BIC)

$$BIC(\beta) = 2\hat{L}(\beta) + 2 \log(n)\|\beta\|_0,$$

where  $\|\beta\|_0 = \text{card}\{j \in 1, \dots, d, \beta_j \neq 0\}$ .

- ★ Produce interpretable models.
- ↪ Non-convex optimization.
  - ▶  $\text{card}(\mathcal{P}(\{1, \dots, d\})) = 2^d$  models,
  - ▶ stepwise selection.

# $\ell_1$ -penalization

Lasso [Tibshirani, 1996]

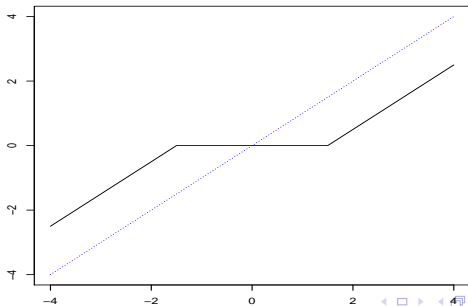
$$\hat{\beta}(\lambda) = \arg \min_{\beta \in \mathbb{R}^d} \left\{ \hat{L}(\beta) + \lambda \sum_{j=1}^d |\beta_j| \right\}$$

- ★ Regularization technique for simultaneous estimation and selection,  $\hat{\beta}_{Lj}(\lambda) = 0$  for  $j \notin K(\hat{\beta}_L) \subset \{1, \dots, d\}$ .
- ★ High dimension  $d \gg n$ .
- ★ Convex optimization.

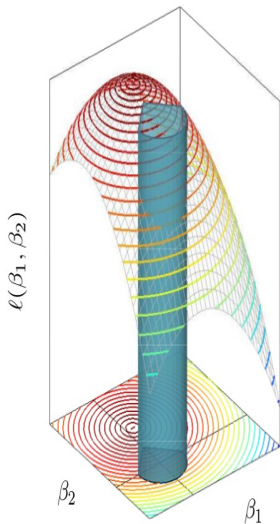
# $\ell_1$ -penalization

★ If  $\epsilon \sim \mathcal{N}(0, I)$  and  $X^T X = I$ , for  $j \in \{1, \dots, d\}$ ,

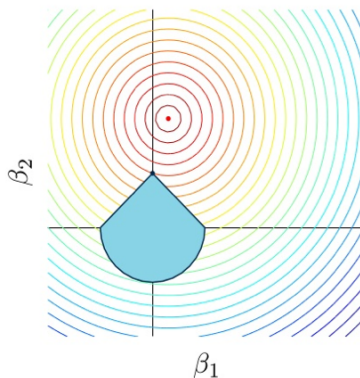
$$\beta_j(\lambda) = \begin{cases} \hat{\beta}_j^{OLS} - \lambda/2 & \text{if } \hat{\beta}_j^{OLS} > \lambda/2 \\ \hat{\beta}_j^{OLS} + \lambda/2 & \text{if } \hat{\beta}_j^{OLS} < -\lambda/2 \\ 0 & \text{otherwise} \end{cases}$$



# Geometric view of penalization



$$\max_{\beta} \ell(\beta) \quad \text{s.t.} \quad \text{pen}(\beta) \leq \xi$$



# Geometry of the Lasso

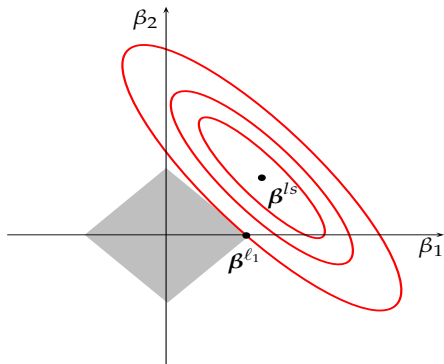


Figure : *Lasso solution.*

# Choice of $\lambda$

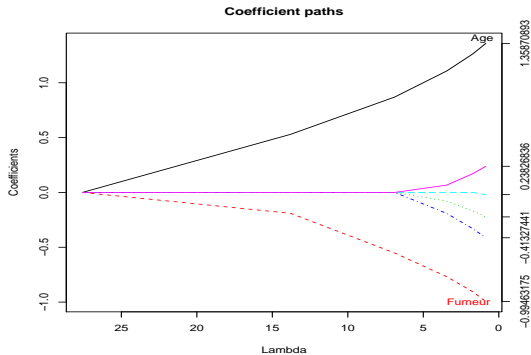


Figure : *Regularization path.*



## Choice of $\lambda$

- **Cross validation** which is recommended when the goal is to minimize the prediction error. It can be computationally slow.
- **Information criteria** such as AIC or BIC can be defined for the Lasso

$$\begin{aligned}AIC(\lambda) &= L_n(\hat{\beta}(\lambda)) + \frac{2}{n}df(\lambda), \\BIC(\lambda) &= L_n(\hat{\beta}(\lambda)) + \frac{\log n}{n}df(\lambda),\end{aligned}$$

where  $df(\lambda)$  is the degrees of freedom of the Lasso for a given parameter  $\lambda$ .

# Choice of $\lambda$

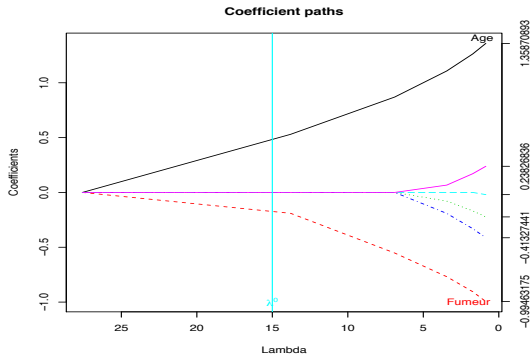


Figure : *Regularization path*

# Properties of the Lasso

- **Prediction**  
Does the Lasso provide a good approximation of  $X\beta_0$ ?
- **Estimation**  
Does the Lasso provide a good approximation of  $\beta_0$ ?
- **Selection**  
Does the Lasso select the right covariates?
- **Sign recovery**  
Does the Lasso select the right covariates and identify correctly the signs of their coefficients?

# Outline

- 1 Model
- 2 High dimensional data situations ( $d \gg n$ )
- 3 Oracle inequality

- Model :

$$y_i = f_0(x_i) + \epsilon_i$$

- 

$$\Gamma \subseteq \left\{ \beta \in \mathbb{R}^d, f_\beta(\cdot) = \sum_{j=1}^p \beta_j \phi_j(\cdot) \right\}.$$

- Oracle inequality

$$R(f_\beta, f_0) \leq C \times \inf_{\beta \in \Gamma} \left\{ R(f_\beta, f_0) + \Delta_n \right\}. \quad (3)$$

- Model :

$$y_i = f_0(x_i) + \epsilon_i$$

- 

$$\Gamma \subseteq \left\{ \beta \in \mathbb{R}^d, f_\beta(\cdot) = \sum_{j=1}^p \beta_j \phi_j(\cdot) \right\}.$$

- Oracle inequality

$$R(f_\beta, f_0) \leq C \times \inf_{\beta \in \Gamma} \left\{ R(f_\beta, f_0) + \Delta_n \right\}. \quad (3)$$

- Model :

$$y_i = f_0(x_i) + \epsilon_i$$

- 

$$\Gamma \subseteq \left\{ \beta \in \mathbb{R}^d, f_\beta(\cdot) = \sum_{j=1}^p \beta_j \phi_j(\cdot) \right\}.$$

- Oracle inequality

$$R(f_{\hat{\beta}}, f_0) \leq C \times \inf_{\beta \in \Gamma} \left\{ R(f_\beta, f_0) + \Delta_n \right\}. \quad (3)$$



Tibshirani, R. (1996).

Regression shrinkage and selection via the lasso.

*J. Roy. Statist. Soc. Ser. B*, 58(1):267–288.