

## Topic 11: Lasso And Beyond: Convex Learning

by Sai Zhang

Key points:

Disclaimer:

## 11.1 Lasso

Lasso (Least absolute Shrinkage and Selection Operator), proposed by Tibshirani (1996), aims to minimize the **SSR (sum of residual squares)** subject to the **L1-norm (sum of the absolute value)** of the coefficients being less than a constant.

### 11.1.1 Set up

For data  $(\mathbf{x}_i, y_i)_{i=1}^n$ , where

- $y_i$  is the outcome for individual  $i$
- $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})'$  is the  $p \times 1$  vector of predictors

Then the Lasso estimator  $(\hat{\alpha}, \hat{\beta})$  is defined as

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \left\{ \sum_{i=1}^n \left( y_i - \alpha - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \quad \text{s.t.} \quad \sum_{j=1}^p |\beta_j| \leq t$$

for the  $n \times 1$  response vector  $\mathbf{y} = (y_1, \dots, y_n)'$ , the  $n \times p$  design matrix  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)'$  where  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})'$  is a  $p \times 1$  vector. Here  $\hat{\alpha} = \bar{y}$ , w.l.o.g., let  $\bar{y} = 0$  and omit  $\alpha$  for simplicity.

In matrix form, we have

- constrained form:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \right\} \quad \text{s.t.} \quad \|\beta\|_1 \leq t$$

- unconstrained form:

$$\hat{\beta}(\lambda) = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

where the regularization parameter  $\lambda \geq 0$ :

- $\lambda \rightarrow \infty$ :  $\hat{\beta}_{lasso} = \mathbf{0}$
- $\lambda = 0$ :  $\hat{\beta}_{lasso} \rightarrow \hat{\beta}_{OLS}$

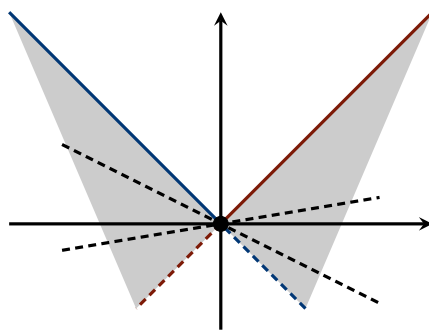
### 11.1.2 Solving Lasso

Lasso is essentially a quadratic optimization problem. Hence, the solution is given by taking the derivative (of the unconstrained question) and set it equal to 0

$$\frac{d}{d\beta} \left( \frac{1}{2n} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right) = 0$$

$\xrightarrow{\text{KKT condition}}$ 
 $\frac{1}{n} \underbrace{X'}_{p \times n} \underbrace{(y - X\beta)}_{= \epsilon, n \times 1} = \lambda \begin{cases} \text{sign}(\beta_j), & \beta_j \neq 0 \\ [-1, 1], & \beta_j = 0 \end{cases}$

this result follows the fact the L-1 norm  $\|\beta\|$  is piecewise linear (convex)<sup>1</sup>:



L1-norm (1-dimension)

For each component of the vector of the L-1 norm

$f(\beta_j) = |\beta_j|$ , we have:

-  $\beta_j > 0$ :  $f'(\beta_j) = 1$

-  $\beta_j < 0$ :  $f'(\beta_j) = -1$

-  $\beta_j = 0$ :  $df \in [-1, 1]$  (shaded area)  
which gives the results stated above.

Take another look at this result

#### Proposition 11.1.1: Lasso Parameter Selection Rule

$$\frac{1}{n} X' (y - X\beta) = \frac{1}{n} X' \epsilon = \lambda \begin{cases} \text{sign}(\beta_j), & \beta_j \neq 0 \\ [-1, 1], & \beta_j = 0 \end{cases}$$

which gives a parameter selection criterion: for  $\beta_j \neq 0$ ,  $\text{sign}(\beta_j)$  **must agree** with  $\text{Corr}(x_j, \epsilon)$ , the correlation between the  $j$ -th variable  $x_j$  and (full-model) residuals  $\epsilon = y - X\beta$ .

### 11.1.3 Algorithm: LARS

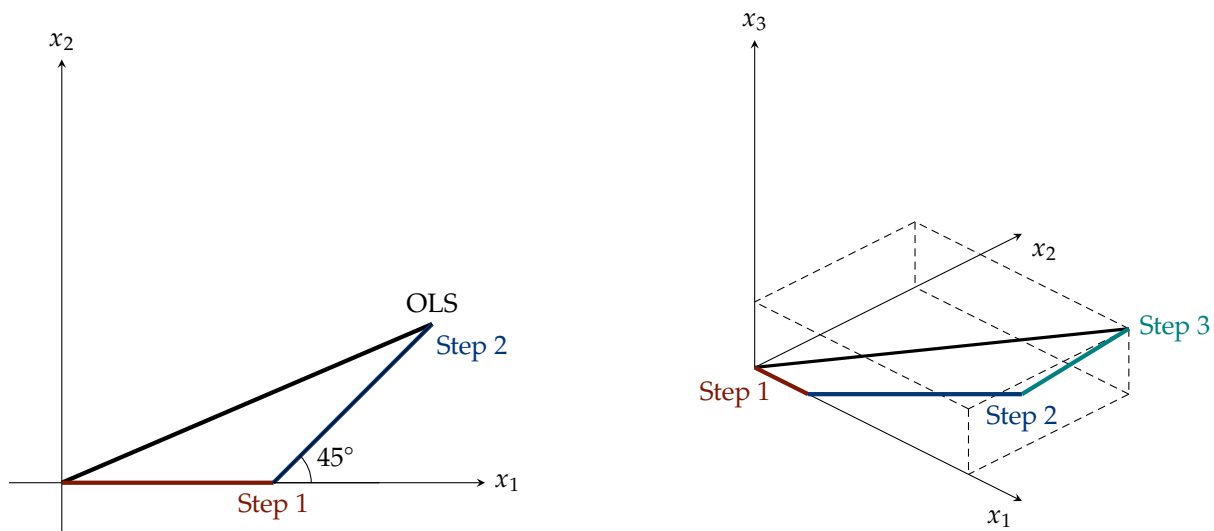
Mathematically, Lasso is quite intuitive, but computationally, it can be quite consuming. Efron et al. (2004) propose an algorithm that takes steps from a all-0 model to the biggest model (OLS), that is, **Least Angle Regression (LARS)**.

#### Intuition

The basic intuition of LARS is quite straight-forward: covariates are considered from the **highest** correlation with  $y$  (*smallest* angle from  $y$ ) to the **least** correlated one (*largest* angle from  $y$ ) (illustrated below).

And the steps of the LARS algorithm are

<sup>1</sup>KKT condition gives the analytical optimization rule for **convex** function.



- 1 start with the null model  $\hat{\beta} = \mathbf{0}$ :  $\hat{\mu} = \mathbf{X}'\mathbf{0} = \mathbf{0}$
- 2 calculate residual vector  $\mathbf{r} = \mathbf{y} - \hat{\mu}$
- 3 determine the correlation vector between  $\mathbf{r}$  and each parameter  $\mathbf{x}_j, \forall j = 1, \dots, p$ :  $\mathbf{X}'\mathbf{r}$
- 4 pick the largest correlation  $\mathbf{x}_{\text{step1},1}^*$ , increase its  $\hat{\beta}$  to the point where its correlation with  $\mathbf{r}$  will be **equal** with that of another parameter  $\mathbf{x}_{\text{step1},2}^*$
- 5 next, increase the  $\hat{\beta}$  for both  $\mathbf{x}_{\text{step1},1}^*, \mathbf{x}_{\text{step1},2}^*$  in an **equiangular** direction between these two, until a third parameter becomes equally important

And keep looping this way, until all the predictors enter the model and eventually  $\mathbf{X}'\mathbf{r} = \mathbf{0}$

### Properties of LARS

LARS has several properties:

- geometrically travels in the direction of **equal** angle to all active covariates
- assume all covariates are independent
- computationally quick: only take  $m$  steps, where  $m$  is the number of parameters being considered

And it is in between 2 classic model-selection methods: **Forward Selection** and **Stagewise Selection**:

- **Forward Selection**

- for  $\mathbf{y}$ , select the most correlated  $\mathbf{x}_{j_1}$
- regress  $\mathbf{x}_{j_1}$  on  $\mathbf{y}$ , get the residuals
- select the most correlated  $\mathbf{x}_{j_2}$  with the residual of  $\mathbf{y}$  net of  $\mathbf{x}_{j_1}$

looping this, for a  $k$ -parameter linear model, it takes  $k$  steps. Forward Selection is an aggressive fitting technique, can be overly greedy (some important predictors may be eliminated due to correlation with already selected variables).

- **Forward Stagewise**

- also begin with  $\hat{\mu} = \mathbf{0}$
- for a current Stagewise estimate  $\hat{\mu}$ , the current residual vector is then  $\mathbf{y} - \hat{\mu}$ , its correlation with  $\mathbf{X}$  is then  $\mathbf{X}'(\mathbf{y} - \hat{\mu}) \equiv \hat{\mathbf{c}}$

- next, heavily computational, go in the direction of the greatest current correlation, but by only a **small** step

$$\hat{j} = \arg \max |\hat{c}_j|, \hat{\mu} \rightarrow \hat{\mu} + \epsilon \cdot \text{sign}(\hat{c}_{\hat{j}}) \cdot \mathbf{x}_{\hat{j}}$$

here,  $\epsilon$  is a **small** constant, hence avoiding the greediness of Forward Selection, at a cost of computational efficiency<sup>2</sup>.

LARS avoids the over-greediness of Forward Selection and computational heaviness of Forward Stagewise.

### 11.1.4 From LARS to Lasso

The Lasso algorithm is built upon LARS, with the constraint from the mathematical condition of Proposition 11.1.1:  $\text{sign}(\beta_j)$  **must agree** with  $\text{Corr}(\mathbf{x}_j, \epsilon)$ .

#### Theorem 11.1.2: Lasso Modification Condition

If  $\tilde{\gamma} < \hat{\gamma}$ , stop the ongoing LARS step at  $\gamma = \tilde{\gamma}$  and remove  $j$  from the calculation of the next equiangular direction, where

- the path at any LARS step is

$$\beta(\gamma), \beta_j(\gamma) = \hat{\beta}_j + \gamma \hat{d}_j$$

$\hat{d}_j$  specifies the **direction** to take the  $j$ -th component,  $\gamma$  is **how far** to travel in the direction of  $\hat{d}_j$  before adding in a new covariate

- $\hat{\gamma}$  represents the smallest **positive** value of  $\gamma$  s.t. some new covariate joins the active set (the set of covariates used on path)
- $\tilde{\gamma}$  represents the first time  $\beta_j(\gamma)$  **changes signs**.

The key point of 11.1.2 is that Lasso does **NOT** allow the  $\hat{\beta}_j$  to change signs, if it changes sign, it will be subtracted from the active set. Now, from this point of view, we can compare the 3 algorithms:

|           |  |
|-----------|--|
| LARS      | no sign restrictions   |
| Lasso     | $\hat{\beta}_j$ agrees in sign with $\hat{c}_j$  |
| Stagewise | successive differences of $\hat{\beta}_j$ agree in sign with the current correlation $\hat{c} = \mathbf{x}'_j(\mathbf{y} - \hat{\mu})$ |

Again, LARS requires the least steps but is most greedy, Stagewise is computationally consuming but robust. Lasso is in between.

## 11.2 Consistency of Lasso

Next, we want to establish the consistency of Lasso, by showing that Lasso selects exactly the relevant covariates asymptotically. We do this in 2 steps:

- show that Lasso at least captures all the relevant covariates
- asymptotically, under some conditions, Lasso selects exactly all the relevant covariates, not more

<sup>2</sup>Forward Selection is essentially choosing  $\epsilon = |\hat{c}_{\hat{j}}|$

### 11.2.1 Overestimation

First, Lasso tends to select a superset of the relevant covariates.

Define the true relevant set Lasso selection estimation  $\hat{S}_0$  aim to select as

$$S_0 = \{j : \beta_j^0 \neq 0, j = 1, \dots, p\}$$

and for some  $C > 0$ , define the relevant set w.r.t.  $C$  as

$$S_0^{\text{relevant}(C)} = \{j : |\beta_j^0| \geq C, j = 1, \dots, p\}$$

then we have

#### Theorem 11.2.1: Lasso Overestimation Condition

$\forall 0 < C < \infty$

$$\mathbb{P} \left[ \hat{S}_0(\lambda) \supset S_0^{\text{relevant}(C)} \right] \xrightarrow{n \rightarrow \infty} 1$$

### Consistency

The consistency of Lasso is established by [Meinshausen and Bühlmann \(2006\)](#) as

#### Theorem 11.2.2: Consistency of Lasso

For a suitable  $\lambda = \lambda_n \gg \sqrt{s_0 \log(p)/n}$ , Lasso is consistent, i.e.

$$\mathbb{P} \left[ \hat{S}(\lambda) = S_0 \right] \xrightarrow{n \rightarrow \infty} 1$$

if and only if it satisfies the 2 properties:

- $\beta$ -min condition (unselected coefficients non-trivial):  $\inf_{j \in S_0^c} |\beta_j^0| \gg \sqrt{s_0 \log(p)/n}$
- **irrepresentable condition**:  $\mathbf{X}$  should **NOT** exhibit too strong a degree of linear dependence w.r.t. the selected covariates, that is, for some  $0 < \theta < 1$

$$\left\| \hat{\Sigma}_{2,1} \hat{\Sigma}_{1,1}^{-1} \text{sign}(\beta_1^0, \dots, \beta_{s_0}^0) \right\|_{\infty} \leq \theta$$

**discussion on the irrepresentable condition** denote  $\hat{\Sigma} = n^{-1} \mathbf{X} \mathbf{X}'$ , and let the active set  $S_0 = \{j : \beta_j^0 \neq 0\} = \{1, \dots, s_0\}$  consists of the first  $s_0$  variables, let

$$\hat{\Sigma} = \begin{pmatrix} \hat{\Sigma}_{1,1} & \hat{\Sigma}_{1,2} \\ \hat{\Sigma}_{2,1} & \hat{\Sigma}_{2,2} \end{pmatrix}$$

where  $\hat{\Sigma}_{1,1}$  is a  $s_0 \times s_0$  var-cov matrix of the active variables  $\mathbf{X}_1$ ,  $\hat{\Sigma}_{2,2}$  is a  $(p - s_0) \times (p - s_0)$  cov-var matrix of the other variables  $\mathbf{X}_2$ <sup>3</sup>, then for a Lasso estimation  $\hat{\beta}$  that assign non-zero coefficients **only** to  $\mathbf{X}_1$ , we can,

<sup>3</sup>Here,  $\mathbf{X}_1$  is  $n \times s_0$ ,  $\mathbf{X}_2$  is  $n \times p - s_0$

following the Lasso result in Proposition 11.1.1, have

$$\begin{aligned} \frac{1}{n} \mathbf{X}' (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}) &= \lambda \cdot \text{sign}(\hat{\boldsymbol{\beta}}_1) & \Rightarrow \frac{1}{n} \mathbf{X}'_1 (\mathbf{y} - \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1) &= \lambda \cdot \text{sign}(\hat{\boldsymbol{\beta}}_1) \\ \left\| \frac{1}{n} \mathbf{X}'_2 (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}) \right\|_\infty &\leq \lambda & \Rightarrow \left\| \frac{1}{n} \mathbf{X}'_2 (\mathbf{y} - \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1) \right\|_\infty &\leq \lambda \end{aligned} \quad (11.1)$$

Now, let's assume  $\text{supp}(\hat{\boldsymbol{\beta}}) = \text{supp}(\boldsymbol{\beta}) = S_0$ , then the true model is

$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon} = \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 \end{bmatrix} \cdot \begin{bmatrix} \boldsymbol{\beta}_1 \\ \mathbf{0} \end{bmatrix} + \boldsymbol{\epsilon} = \mathbf{X}_1 \boldsymbol{\beta}_1 + \boldsymbol{\epsilon}$$

Put this back to the results in Equation 11.1, we have

- for the selected covariates  $\mathbf{X}_1$

$$\begin{aligned} \frac{1}{n} \mathbf{X}'_1 (\mathbf{y} - \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1) &= \lambda \cdot \text{sign}(\hat{\boldsymbol{\beta}}_1) \\ \Rightarrow \frac{1}{n} \mathbf{X}'_1 (\mathbf{X}_1 \boldsymbol{\beta}_1 + \boldsymbol{\epsilon} - \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1) &= \lambda \cdot \text{sign}(\hat{\boldsymbol{\beta}}_1) \\ \Rightarrow \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 (\hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_1) &= \frac{1}{n} \mathbf{X}'_1 \boldsymbol{\epsilon} - \lambda \cdot \text{sign}(\hat{\boldsymbol{\beta}}_1) \\ \Rightarrow \hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_1 &= \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \left[ \frac{1}{n} \mathbf{X}'_1 \boldsymbol{\epsilon} - \lambda \cdot \text{sign}(\hat{\boldsymbol{\beta}}_1) \right] \\ \Rightarrow \hat{\boldsymbol{\beta}}_1 &= \boldsymbol{\beta}_1 + \underbrace{\left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \frac{1}{n} \mathbf{X}'_1 \boldsymbol{\epsilon} - \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \lambda \cdot \text{sign}(\hat{\boldsymbol{\beta}}_1)}_{\text{L1-norm regularization}} \end{aligned}$$

- for the non-selected covariates  $\mathbf{X}_2$ <sup>4</sup>

$$\begin{aligned} \left\| \frac{1}{n} \mathbf{X}'_2 (\mathbf{y} - \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1) \right\|_\infty &\leq \lambda \\ \Rightarrow \left\| \frac{1}{n} \mathbf{X}'_2 \left[ \mathbf{X}_1 (\hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_1) - \boldsymbol{\epsilon} \right] \right\|_\infty &\leq \lambda \\ \Rightarrow \left\| \frac{1}{n} \mathbf{X}'_2 \left\{ \mathbf{X}_1 \left[ \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \frac{1}{n} \mathbf{X}'_1 \boldsymbol{\epsilon} - \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \lambda \cdot \text{sign}(\hat{\boldsymbol{\beta}}_1) \right] - \boldsymbol{\epsilon} \right\} \right\|_\infty &\leq \lambda \\ \xRightarrow{\text{assume sign consistency}} \left\| \frac{1}{n} \mathbf{X}'_2 \mathbf{X}_1 \left[ \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \frac{1}{n} \mathbf{X}'_1 \boldsymbol{\epsilon} - \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \lambda \cdot \text{sign}(\boldsymbol{\beta}_1) \right] - \frac{1}{n} \mathbf{X}'_2 \boldsymbol{\epsilon} \right\|_\infty &\leq \lambda \end{aligned}$$

if we assume  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ , i.e.,  $\epsilon_i$  i.i.d.  $\sim \mathcal{N}(0, \sigma^2)$ , and for each variable in the design matrix, we also assume standard normal  $\mathbf{x}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$  (think this as normalizing each variable), then we have

$$\frac{1}{n} \mathbf{x}'_j \boldsymbol{\epsilon} \sim \mathcal{N}\left(0, \frac{\sigma^2}{n}\right), \forall j = 1, \dots, p$$

---

<sup>4</sup>Here, additionally assume **sign consistency**:

$$\text{sign}(\hat{\beta}_j) = \beta_j \neq 0, \forall j \in S_0$$

$$\text{sign}(\hat{\beta}_j) = \beta_j = 0, \forall j \in S_0^C$$

selected covariates

non-selected covariates

and also  $\max_{1 \leq j \leq p} \mathbf{x}_j \sim \sqrt{2 \log p}$ . This gives

$$\left\| \frac{1}{n} \mathbf{X}' \boldsymbol{\epsilon} \right\|_{\infty} \sim \sqrt{2 \log p} \cdot \frac{\sigma}{\sqrt{n}} = \sqrt{\frac{2 \log p}{n}} \sigma$$

hence, it can be bounded by  $\lambda = \sqrt{\frac{C \log p}{n}} \sigma$ , where the constant  $C \geq 2$ , then with large probability

$$\left\| \frac{1}{n} \mathbf{X}' \boldsymbol{\epsilon} \right\|_{\infty} \leq \frac{1}{3} \lambda \Rightarrow \begin{cases} \left\| \frac{1}{n} \mathbf{X}'_1 \boldsymbol{\epsilon} \right\|_{\infty} \leq \frac{1}{3} \lambda \\ \left\| \frac{1}{n} \mathbf{X}'_2 \boldsymbol{\epsilon} \right\|_{\infty} \leq \frac{1}{3} \lambda \end{cases}$$

now go back to the condition of non-selected covariates  $\mathbf{X}_2$

$$\begin{aligned} & \left\| \frac{1}{n} \mathbf{X}'_2 \mathbf{X}_1 \left[ \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \frac{1}{n} \mathbf{X}'_1 \boldsymbol{\epsilon} - \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \lambda \cdot \text{sign}(\boldsymbol{\beta}_1) \right] - \frac{1}{n} \mathbf{X}'_2 \boldsymbol{\epsilon} \right\|_{\infty} \\ & \leq \left\| \frac{1}{n} \mathbf{X}'_2 \mathbf{X}_1 \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \frac{1}{n} \mathbf{X}'_1 \boldsymbol{\epsilon} \right\|_{\infty} + \left\| \frac{1}{n} \mathbf{X}'_2 \mathbf{X}_1 \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \text{sign}(\boldsymbol{\beta}_1) \right\|_{\infty} \cdot \lambda + \left\| \frac{1}{n} \mathbf{X}'_2 \boldsymbol{\epsilon} \right\|_{\infty} \\ & \leq \left\| \frac{1}{n} \mathbf{X}'_2 \mathbf{X}_1 \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \right\|_{\infty} \cdot \underbrace{\left\| \frac{1}{n} \mathbf{X}'_1 \boldsymbol{\epsilon} \right\|_{\infty}}_{\leq \frac{1}{3} \lambda} + \left\| \frac{1}{n} \mathbf{X}'_2 \mathbf{X}_1 \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \right\|_{\infty} \cdot \underbrace{\left\| \text{sign}(\boldsymbol{\beta}_1) \right\|_{\infty}}_{=1} \cdot \lambda + \underbrace{\left\| \frac{1}{n} \mathbf{X}'_2 \boldsymbol{\epsilon} \right\|_{\infty}}_{\leq \frac{1}{3} \lambda} \\ & \leq \left\| \frac{1}{n} \mathbf{X}'_2 \mathbf{X}_1 \left( \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \right\|_{\infty} \cdot \frac{4}{3} \lambda + \frac{1}{3} \lambda \leq \lambda \end{aligned}$$

for the last part ( $\leq \lambda$ ) to stand, a necessary condition is  $\left\| \mathbf{X}'_2 \mathbf{X}_1 (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \right\|_{\infty} \leq \frac{1}{2}$ , or more generally, the **irrepresentable condition**:

$$\left\| \mathbf{X}'_2 \mathbf{X}_1 (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \cdot \text{sign}(\boldsymbol{\beta}_1) \right\| \leq \left\| \mathbf{X}'_2 \mathbf{X}_1 (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \right\|_{\infty} \cdot \left\| \text{sign}(\boldsymbol{\beta}_1) \right\| = \left\| \hat{\Sigma}_{2,1} \hat{\Sigma}_{1,1}^{-1} \text{sign}(\boldsymbol{\beta}_1^0, \dots, \boldsymbol{\beta}_{s_0}^0) \right\|_{\infty} \leq \theta$$

The idea of **irrepresentable** is that the **maximum** correlation between noise features and important variables needs to be **bounded**.

## 11.2.2 Oracle

Consistency is part of the requirements of an oracle procedure

### Definition 11.2.3: Oracle Property

For a fitting procedure  $\delta$ , and the estimation  $\hat{\beta}(\delta)$ , then if  $\delta$  is an oracle procedure if  $\hat{\beta}(\delta)$  asymptotically has the following properties

- **consistency** (identifying right subset model):  $\{j : \hat{\beta}_j \neq 0\} = S_0$
- **optimal estimation rate** (asymptotically normal):  $\sqrt{n} \left( \beta(\delta)_{S_0} - \beta_{S_0}^0 \right) \xrightarrow{d} \mathcal{N}(0, \Sigma_0)$ , where  $\Sigma_0$  is the true subset covariance matrix

The oracle property gives consistency, and asymptotic normality. Ideally, this is what you need.

## 11.3 Variants of Lasso

Lasso is great, but there is still space for improvement

- for high dimension  $p > n$  cases, lasso selects **at most**  $n$  variables
- for a group of variables with large **pairwise correlation**, then the lasso tends to select **only one**
- for  $n > p$  cases, if there are high correlations between predictors, lasso is dominated by ridge regression

Here are some popular variants people came up with to tackle these issues.

### 11.3.1 Adaptive Lasso

Adaptive Lasso replaces the  $L_1$ -penalty for a **re-weighted** version:

$$\hat{\beta}(\lambda) = \arg \min_{\beta} \left( \frac{1}{n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_n \sum_{j=1}^p \frac{|\beta_j|}{|\hat{\beta}_{init,j}|^\gamma} \right)$$

this is a 2-step procedure:

- Step 1: run an initial estimation (could be OLS or Lasso) as  $\hat{\beta}_{init}$ , if  $\hat{\beta}_{init,j} = 0$ , then  $\hat{\beta}_{adapt,j} = 0$
- Step 2: run the reweighted estimation, which is still an  $L_1$ -penalization, still a **convex** optimization problem (allowing the same Lasso algorithm to solve)

#### Theorem 11.3.1: Properties of Adaptive Lasso

The adaptive Lasso has the following properties

- If  $|\hat{\beta}_{init,j}|$  is large (*important* features), then the adaptive Lasso uses a small penalty
- The adaptive lasso enjoys the oracle properties

- **consistency**:  $\lim_n \mathbb{P}(\hat{S}(\lambda) = S_0) = 1$
- **asymptotic normality**:  $\sqrt{n}(\beta_{S_0} - \beta_{S_0}^*) \xrightarrow{d} \mathcal{N}(0, \sigma^2 \times C_{1,1}^{-1})$

if  $\lambda_n$  is chosen properly, i.e.,  $\frac{\lambda_n}{\sqrt{n}} \rightarrow 0$  and  $\lambda_n n^{(\gamma-1)/2} \rightarrow \infty$ . Notice: the adaptive lasso's consistency **DOES NOT** require the irrepresentable condition.

### 11.3.2 Naive Elastic Net

The idea is to do a convex combination of Lasso and Ridge:  $\forall \lambda_1, \lambda_2 \geq 0$ , the naive elastic net criterion is

$$\mathcal{L}(\lambda_1, \lambda_2, \beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \underbrace{\|\beta\|_1}_{=\sum_{j=1}^p |\beta_j|} + \lambda_2 \underbrace{\|\beta\|_2^2}_{=\sum_{j=1}^p \beta_j^2}$$

Now, notice that both  $\|\mathbf{y} - \mathbf{X}\beta\|_2^2$  and  $\|\beta\|_2^2$  are both **quadratic**, hence, we can reframe the question as a lasso-type optimization, with an augmented data  $(\mathbf{X}^*, \mathbf{y}^*)$  where

$$\mathbf{X}_{(n+p) \times p}^* = \frac{1}{\sqrt{1 + \lambda_2}} \begin{pmatrix} \mathbf{X} \\ \sqrt{\lambda_2} \mathbf{I} \end{pmatrix} \quad \mathbf{y}_{n+p}^* = \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix}$$

then we have



**Definition 11.3.2: Naive elastic net solution**

Let  $\gamma = \frac{\lambda_1}{\sqrt{1+\lambda_2}}$ , then the naive elastic net criterion can be rewritten as

$$\mathcal{L}(\gamma, \beta) = \mathcal{L}(\gamma, \beta^*) = \|\mathbf{y}^* - \mathbf{X}^* \beta^*\|_2^2 + \gamma \|\beta^*\|_1$$

which gives use the estimation for the augmented data

$$\hat{\beta}^* = \arg \min_{\beta^*} \mathcal{L}(\gamma, \beta^*)$$

and the estimation for the true model of the original data  $(\mathbf{X}, \mathbf{y})$  is

$$\hat{\beta} = \frac{1}{\sqrt{1+\lambda_2}} \hat{\beta}^*$$

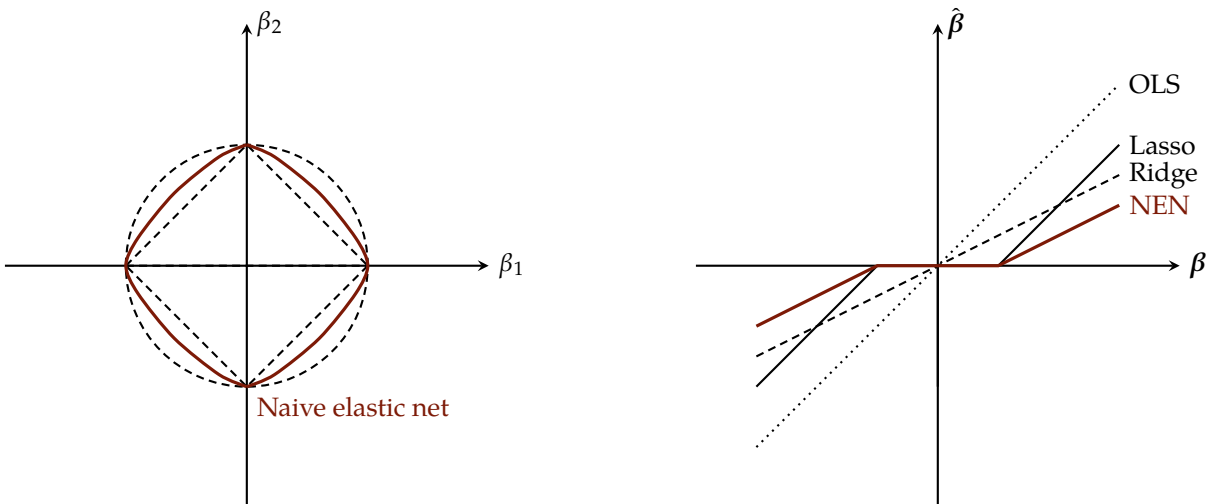
Naive elastic net enjoys 2 advantages:

- The estimation is done on the augmented data  $(\mathbf{X}^*, \mathbf{y}^*)$ , and now  $\mathbf{X}^*$  is  $(n+p) \times p$ , therefore,  $\text{rank}(\mathbf{X}^*) = p$ , hence all  $p$  predictors **can** be selected, instead of just  $n$ , like in Lasso.
- Computationally, it is similarly efficient as Lasso

If we further assume an orthogonal design  $(\mathbf{X}'\mathbf{X} = \mathbf{I}_p)$ , we can easily compare the 3 estimations

$$\begin{aligned} \hat{\beta}_j \text{ (naive elastic net)} &= \frac{\left( \left| \hat{\beta}_j^{OLS} \right| - \frac{\lambda_1}{2} \right)_+}{1 + \lambda_2} \cdot \text{sign} \left[ \hat{\beta}_j^{OLS} \right] \\ \hat{\beta}_j \text{ (ridge)} &= \frac{\hat{\beta}_j^{OLS}}{1 + \lambda_2} \\ \hat{\beta}_j \text{ (lasso)} &= \left( \left| \hat{\beta}_j^{OLS} \right| - \frac{\lambda_1}{2} \right)_+ \cdot \text{sign} \left[ \hat{\beta}_j^{OLS} \right] \end{aligned}$$

This comparison is illustrated graphically below.



Now consider another problem: grouping effect, that is, for a generic penalization method

$$\hat{\beta} = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda J(\beta)$$

where  $J(\cdot)$  is positive valued for  $\beta \neq 0$ , then for a group of highly correlated variables, the regression coefficients tend to be **equal**.

#### Lemma 11.3.3: Grouping effect

Assume  $\mathbf{x}_i = \mathbf{x}_j, i, j \in \{1, \dots, p\}$ , then

- If  $J(\cdot)$  is **strictly convex**, then  $\hat{\beta}_i = \hat{\beta}_j, \forall \lambda > 0$
- If  $J(\beta) = \|\beta\|_1$ , then  $\hat{\beta}_1, \hat{\beta}_2 \geq 0$  and  $\hat{\beta}^*$  is another minimizer of the generic penalization function, where

$$\hat{\beta}_k^* = \begin{cases} \hat{\beta}_k & \text{if } k \neq i \text{ and } k \neq j \\ (\hat{\beta}_i + \hat{\beta}_j) \cdot s & \text{if } k = i \\ (\hat{\beta}_i + \hat{\beta}_j) \cdot (1 - s) & \text{if } k = j \end{cases}$$

for any  $s \in [0, 1]$

As shown in Lemma 11.3.3, there is a clear distinction between **strictly convex** penalty functions and the Lasso penalty (not strictly convex). For the naive elastic net penalty,  $\lambda_2 > 0$  gives strict convexity.

And now, to take this grouping effect into consideration, we want: for two variables  $\mathbf{x}_i, \mathbf{x}_j$  that are closely correlated, we should expect the coefficient paths of them to converge, that is

#### Theorem 11.3.4: Grouping effect requirement

For the naive elastic net estimate  $\hat{\beta}(\lambda_1, \lambda_2)$ , WLOG, suppose that  $\hat{\beta}_j(\lambda_1, \lambda_2), \hat{\beta}_i(\lambda_1, \lambda_2) > 0$ , then define the **difference** between the coefficient paths of predictors  $i$  and  $j$  as

$$D_{\lambda_1, \lambda_2}(i, j) = \frac{1}{\|\mathbf{y}\|_1} |\hat{\beta}_i(\lambda_1, \lambda_2) - \hat{\beta}_j(\lambda_1, \lambda_2)|$$

then the grouping effect requires

$$D_{\lambda_1, \lambda_2}(i, j) \leq \frac{1}{\lambda_2} \sqrt{2(1 - \rho)}$$

where  $\rho = \mathbf{x}_i' \mathbf{x}_j$  is the sample correlation.

### 11.3.3 Elastic Net

One downfall of naive elastic net estimation is that it incurs **double** shrinkage (2-stage procedure, with each stage having a regularizer). To address this issue, we can rescale the naive elastic net estimation to get the **elastic net estimate**:

#### Definition 11.3.5: Elastic net solution

Again for  $\gamma = \frac{\lambda_1}{\sqrt{1 + \lambda_2}}$ , we have

$$\hat{\beta}^* = \arg \min_{\beta^*} \mathcal{L}(\gamma, \beta^*) = \arg \min_{\beta^*} \|\mathbf{y}^* - \mathbf{X}^* \beta^*\|_2^2 + \gamma \|\beta^*\|_1$$

and the elastic estimate of  $\beta$  is

$$\hat{\beta}_{\text{elastic net}} = \sqrt{1 + \lambda_2} \hat{\beta}^* = (1 + \lambda_2) \hat{\beta}_{\text{naive elastic net}}$$

The rescaling from  $\hat{\beta}_{\text{naive elastic net}}$  to  $\hat{\beta}_{\text{elastic net}}$  preserves the variable selection property (able to select all  $p$  variables), while solve the double-regularization shrinkage problem.

But how do we make sense of the scalar  $1 + \lambda_2$ ? Remember the ridge estimator is

$$\hat{\beta}_{\text{ridge}} = (\mathbf{X}'\mathbf{X} + \lambda_2 \mathbf{I})^{-1} \mathbf{X}'\mathbf{y}$$

which leads to a Lasso-looking rewriting of the elastic net estimation

#### Theorem 11.3.6: Elastic net solution: A modified version of Lasso

Given data  $(\mathbf{y}, \mathbf{X})$  and regularization parameters  $(\lambda_1, \lambda_2)$ , then the elastic net estimation is

$$\hat{\beta}_{\text{elastic net}} = \arg \min_{\beta} \beta' \left( \frac{\mathbf{X}'\mathbf{X} + \lambda_2 \mathbf{I}}{1 + \lambda_2} \right) \beta - 2\mathbf{y}'\mathbf{X}\beta + \lambda_1 \|\beta\|_1$$

comparing to Lasso estimation

$$\hat{\beta}_{\text{Lasso}} = \arg \min_{\beta} \beta' (\mathbf{X}'\mathbf{X}) \beta - 2\mathbf{y}'\mathbf{X}\beta + \lambda_1 \|\beta\|_1$$

hence, elastic net is a **stabilized version** of Lasso.

With the *stabilizing* term, elastic net dominates Lasso under collinearity.

### 11.3.4 Group Lasso

So far, all the variants are built on the model with one big set of covariates  $\mathbf{X}$  ( $n \times p$ ) and a long vector of coefficient  $\beta$  ( $p \times 1$ ), with a strong assumption of orthogonal design:  $\mathbf{X}'\mathbf{X} = \mathbf{I}_p$ , what if we can pool some covariates into a factor, and do a group estimation? This is the idea of group lasso. Consider

$$\mathbf{y} = \sum_{j=1}^J \mathbf{X}_j \beta_j + \epsilon$$

where  $\mathbf{y}$  is still a  $n \times 1$  vector,  $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ , and  $\mathbf{X}_j$  ( $j = 1 \dots J$ ) are a decomposition of  $\mathbf{X}$  into  $J$  submatrix (factors):

$$\mathbf{X} = \begin{bmatrix} \underbrace{\mathbf{X}_1}_{n \times p_1} & \cdots & \underbrace{\mathbf{X}_j}_{n \times p_j} & \cdots & \underbrace{\mathbf{X}_J}_{n \times p_J} \end{bmatrix}$$

Here, each  $\mathbf{X}_j$  is orthonormal:  $\mathbf{X}_j' \mathbf{X}_j = \mathbf{I}_{p_j}$ <sup>5</sup>. Similarly, the vector of coefficients will be decomposed into

$$\beta = [\beta_1' \cdots \beta_j' \cdots \beta_J']'$$

<sup>5</sup>This is a much weaker assumption than the full orthogonal design  $\mathbf{X}'\mathbf{X} = \mathbf{I}_p$

within which each component  $\beta_j$  ( $p_j \times 1$ ) corresponds to a factor.

Now for this problem, we have the **group Lasso** estimation

**Definition 11.3.7: Group Lasso**

The group Lasso estimate is defined as

$$\min_{\beta_j} \frac{1}{2} \|\mathbf{y} - \sum_{j=1}^J \mathbf{X}_j \beta_j\|_2^2 + \lambda \sum_{j=1}^J \|\beta_j\|_{\mathbf{K}_j}$$

where the matrix norm  $\|\cdot\|_{\mathbf{K}}$  is defined as

$$\|\eta\| = \sqrt{\eta' \mathbf{K} \eta}$$

with a symmetric  $d \times d$  positive definite matrix  $\mathbf{K}$ , for a vector  $\eta \in \mathbb{R}^d$

In Definition 11.3.7, the problem is reduced to Lasso when  $p_1 = \dots = p_J = 1$ .

**How to solve this?** Notice that this is still a convex optimization problem, hence KKT condition still holds, and assume the kernel  $\mathbf{K}_j = p_j \mathbf{I}_{p_j}$ , we have the estimation

$$\beta_j = \left(1 - \frac{\lambda \sqrt{p_j}}{\|\mathbf{S}_j\|_1}\right)_+ \mathbf{S}_j$$

where

$$\mathbf{S}_j = \mathbf{X}_j' (\mathbf{y} - \mathbf{X} \beta_{-j}), \text{ with } \beta_{-j} = [\beta'_1, \dots, \beta'_{j-1}, 0, \beta'_{j+1}, \dots, \beta'_J]'$$

if and only if

$$\begin{aligned} -\mathbf{X}_j' (\mathbf{y} - \mathbf{X} \beta) + \frac{\lambda \beta_j \sqrt{p_j}}{\|\beta_j\|_1} &= 0, & \forall \beta_j \neq 0 \\ \|\mathbf{X}_j' (\mathbf{y} - \mathbf{X} \beta)\|_1 &\leq \lambda \sqrt{p_j}, & \forall \beta_j = 0 \end{aligned}$$

**Group LARS** similar to LARS, group LARS also finds factor that has the smallest angle with  $\mathbf{y}$  and proceeds equiangularly each time a new factor is selected. But the angle (between a vector  $\mathbf{r}$  and a factor matrix  $\mathbf{X}_j$ ), is constructed as

$$\cos^2(\theta(\mathbf{r}, \mathbf{X}_j)) = \frac{\|\mathbf{X}_j' \mathbf{r}\|_2^2}{\|\mathbf{r}\|_2^2}$$

when  $p_j = \bar{p}, \forall j$ <sup>6</sup>.

### 11.3.5 Other Variants

There are also some other useful variants of Lasso

<sup>6</sup>Otherwise, use the scaled angle  $\|\mathbf{X}_j' \mathbf{r}\|_2^2 / p_j$

- **Positive Lasso:** Constrains the  $\hat{\beta}_j$  to enter the prediction equation in their **defined** directions, non-negative here

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \right\} \quad \text{s.t. } \|\beta\|_1 \leq t \text{ and } \beta_j > 0, \forall j$$

- **LARS-OLS hybrid:** Use the covariates selected by LARS, but use  $\hat{\beta}$  from the OLS model
- **Main effects first:**
  - Step 1: run LARS for a model, considering **only** main effects
  - Step 2: run LARS again, with the chosen main effects, and **all possible interactions** between them
- **Backward Lasso:** start from the **full** OLS model, and eliminate covariates **backwards** (by the order of correlation going 0 the earliest)

## 11.4 Dantzig Selector and Lasso

Now, back to the basic problem

$$\mathbf{y} = \mathbf{X}\beta + \epsilon$$

We consider another way of estimation, **Dantzig Selector**

### Definition 11.4.1: Dantzig Selector

The Dantzig selector is a solution to the  $L_1$ -regularization problem

$$\min_{\beta \in \mathbb{R}^p} \|\beta\|_1 \quad \text{s.t.} \quad \|n^{-1}\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta)\|_\infty \leq \lambda$$

similarly,  $n^{-1}\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta) = n^{-1}\mathbf{X}'\epsilon$ , the constraint aims to bound the correlation of covariates and residual.

Though not exactly, Dantzig selector can be thought as the dual problem of Lasso, hence has similar performance in certain scenarios.

**Why constrain**  $\|n^{-1}\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta)\|_\infty$  There are 2 justifications for constraining  $\|n^{-1}\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta)\|_\infty$  instead of just the size of the residuals  $\|\mathbf{y} - \mathbf{X}\beta\|_\infty$

- Stability: Dantzig constraint is **invariant** over any **orthonormal** transformation

$$n^{-1}(\mathbf{U}\mathbf{X})'(\mathbf{U}\mathbf{y} - \mathbf{U}\mathbf{X}\beta) = n^{-1}\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta)$$

- No accidental neglect: some important covariates might be noisy, could be mis-eliminated if using  $\|\mathbf{y} - \mathbf{X}\beta\|_\infty \leq \lambda$

**How to solve?** Dantzig selector problem can be recast as a linear programming (LP) problem

$$\min_{\mathbf{u} \in \mathbb{R}^p} \sum_{i=1}^p u_i$$

s.t.

$$\begin{aligned} -\mathbf{u} &\leq \boldsymbol{\beta} \leq \mathbf{u} \\ -\lambda_p \sigma \mathbf{1} &\leq n^{-1} \mathbf{X}'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \leq \lambda_p \sigma \mathbf{1} \end{aligned}$$

In their elegant work, **Candes and Tao (2007)** illustrated some properties of Dantzig selector

**Theorem 11.4.2: Properties of the Dantzig Selector (Candes and Tao, 2007)**

- Suppose  $\boldsymbol{\beta}_0 \in \mathbb{R}^p$  is any  $S$ -sparse vector of parameter obeying  $\delta_{2S} + \theta_{S,2S} < 1$ . Choose  $\lambda_p = \sqrt{2 \log p}$ . Then with large probability,  $\hat{\boldsymbol{\beta}}$  obeys

$$\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\|_2^2 \leq C_1^2 \cdot (2 \log p) \cdot S \cdot \sigma^2$$

with  $C_1 = \frac{4}{1 - \delta_S - \theta_{S,2S}}$ . Hence, for small value of  $\delta_S + \theta_{S,2S}$ ,  $C_1 \approx 4^a$ . For concreteness, if one chooses

$\lambda_p := \sqrt{2(1+a) \log p}$ ,  $\forall a \geq 0$ , the bound holds with probability exceeding  $1 - \left(\sqrt{\pi \log p} \cdot p^a\right)^{-1}$  with the proviso that  $\lambda_p^2$  substitutes  $2 \log p$ .

- $\theta_{S,2S}$  refers to the restricted orthogonality constant for  $S + S' \leq p$  to be the smallest quantities s.t.

$$|\langle \mathbf{X}_T \mathbf{c}, \mathbf{X}_{T'} \mathbf{c}' \rangle| \leq \theta_{S,S'} \cdot \|\mathbf{c}\|_2 \|\mathbf{c}'\|_2$$

<sup>a</sup>  $\delta_S$  refers to the smallest quantity ( $S$ -restricted isometry constant) s.t.

$$(1 - \delta_S) \|\mathbf{c}\|_2^2 \leq \|\mathbf{X}_J \mathbf{c}\|_2^2 \leq (1 + \delta_S) \|\mathbf{c}\|_2^2$$

for all subsets  $J$  with  $|J| \leq S$  and coefficient sequences  $(c_j)_{j \in J}$ . Notice that  $[1 - \delta_S, 1 + \delta_S]$  is the smallest interval that is symmetric around 1 and includes all eigenvalues of  $\mathbf{X}_J' \mathbf{X}_J$ .

Theorem 11.4.2 gives that the **oracle LS estimator**,  $\hat{\boldsymbol{\beta}}_S$ , with the location of the  $S$  non-zero entries of the parameter vector

$$\mathbb{E} \|\hat{\boldsymbol{\beta}}_S - \boldsymbol{\beta}_0\|_2^2 \geq \frac{1}{1 + \delta_S} \mathbb{E} \|\mathbf{X}_S (\hat{\boldsymbol{\beta}}_S - \boldsymbol{\beta}_0)\|_2^2 = \frac{1}{1 + \delta_S} \cdot |S| \cdot \sigma^2$$

which gives that

- the mean squared error is simply proportional, up to a logarithmic factor, to the **true number** of unknowns times the noise level  $\sigma^2$
- reliable estimate of  $\boldsymbol{\beta}$  can be achieved only by a **simply linear program**

But what about the idea risk,  $\min_{I \subset \{1, \dots, p\}} \mathbb{E} \|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}_I\|^2$ ?

**Theorem 11.4.3: Bound of Risk of the Dantzig Selector**

Choose  $t > 0$ , and set  $\lambda_p := (1 + t^{-1}) \sqrt{2 \log p}$ . Then if  $\boldsymbol{\beta}$  is  $S$ -sparse with  $\delta_{2S} + \theta_{S,2S} < 1 - t$ , then Dantzig selector obeys

$$\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2^2 \leq C_2^2 \cdot \lambda_p^2 \cdot \left( \sigma^2 + \sum_{i=1}^p \min(\beta_i^2, \sigma^2) \right)$$

with large probability. Here,  $C_2$  may only depend on  $\delta_{2S}$  and  $\theta_{S,2S}$  and the additional  $C_2^2 \lambda_p^2 \sigma^2$  reflects bias due to regularization.

### 11.4.1 Connecting Dantzig Selector and Lasso

The equivalence of Lasso and Dantzig was established by [James et al. \(2009\)](#), by introducing the following theorem:

#### Theorem 11.4.4: Equivalence of Lasso and Dantzig

Denote the support of Lasso estimate  $\hat{\beta}_{\text{Lasso}}$  as  $I_L$ , and construct a  $n \times |I_L|$  matrix  $\mathbf{X}_L$  by taking  $\mathbf{X}_{I_L}$  and multiplying its columns by the signs of the corresponding coefficients in  $\hat{\beta}_L$ .

Assume that  $\mathbf{X}_L$  has **full rank**  $|I_L|$  and  $\lambda_D = \lambda_L$ , then  $\hat{\beta}_D = \hat{\beta}_L$  if

$$\begin{aligned} \mathbf{u} &= (\mathbf{X}_L' \mathbf{X}_L)^{-1} \mathbf{1} \geq 0 \\ \|\mathbf{X}_L' \mathbf{X}_L \mathbf{u}\|_{\infty} &\leq 1 \end{aligned}$$

where  $\mathbf{1} = (1, \dots, 1)'$ , with the length  $|I_L|$ .

Naturally from this theorem, if  $\mathbf{X}'\mathbf{X} = \mathbf{I}_p$  (orthonormal design matrix), then Lasso and Dantzig coefficient paths are identical.

## 11.5 Penalized Least Square Estimation

Lasso is one special class of Penalized Least Square (PLS) Estimation. For the linear regression model  $\mathbf{y} = \mathbf{X}\beta + \epsilon$ , if  $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ , we have PLS as

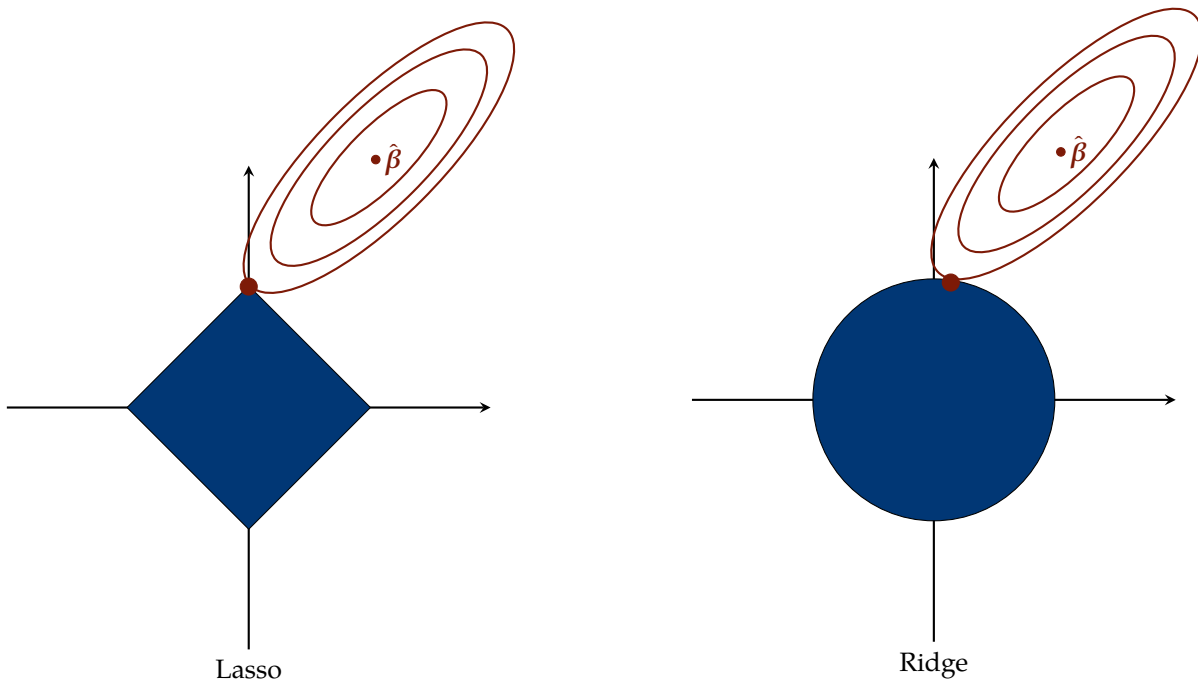
$$\min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \sum_{j=1}^p p_{\lambda}(|\beta_j|) \right\}$$

where  $p_{\lambda}(\cdot)$  is a penalty function indexed by the regularization parameter  $\lambda \geq 0$ . [Antoniadis and Fan \(2001\)](#) showed that the PLS estimator  $\hat{\beta}$  has the following properties:

- **sparsity**: if  $\min_{t \geq 0} \{t + p'_{\lambda}(t)\} > 0$
- **approximate unbiasedness**: if  $p'_{\lambda}(t) = 0$  for  $t$  large enough
- **continuity**: iff  $\arg \min_{t \geq 0} \{t + p'_{\lambda}(t)\} = 0$

In general

- the **sigularity** of penalty function at the origin,  $p'_{\lambda}(0_+) > 0$  is needed for generating **sparsity** in variable selection
- the **concavity** is needed to reduce the bias



## References

- Anestis Antoniadis and Jianqing Fan. Regularization of wavelet approximations. *Journal of the American Statistical Association*, 96(455):939–967, 2001.
- Emmanuel Candes and Terence Tao. The dantzig selector: Statistical estimation when  $p$  is much larger than  $n$ . *The Annals of Statistics*, 35(6):2313–2351, 2007.
- Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Least angle regression. *The Annals of Statistics*, 32(2):407 – 499, 2004. doi: 10.1214/0090536040000000067. URL <https://doi.org/10.1214/0090536040000000067>.
- Gareth M James, Peter Radchenko, and Jinchi Lv. Dasso: connections between the dantzig selector and lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(1):127–142, 2009.
- Nicolai Meinshausen and Peter Bühlmann. Variable selection and high-dimensional graphs with the lasso. *Ann Stat*, 34:1436–1462, 2006.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.