**uni.lu** | ☐ FACULTY OF SCIENCE,
TECHNOLOGY
AND MEDICINE

# University of Luxembourg

## Thesis for the Bachelor of Mathematics

# High Dimensional Regression Models

by Joris LIMONIER

*Supervisor:* Mark PODOLSKIJ
*Submitted:* June 4, 2021

ii

# Abstract

Abstract goes here.

# Contents

# Chapter 1

# Introduction

## 1.1 Notes for chapter 1

# Chapter 2

# Classical theory of Linear Regression

To be added

- how to get $\hat{b}$ on page 101.

- where the $\chi^2$ distribution comes from in page 101

## 2.1 Linear models

We consider the setting of having a sample of $n$ observations

$$(\mathbf{X}_1, \mathbf{Y}_1), \ldots, (\mathbf{X}_n, \mathbf{Y}_n)$$

where $X_i \in \mathcal{X} \subseteq \mathbb{R}^p$, $i = 1, \ldots, n$ and $Y_i \in \mathcal{Y} \subseteq \mathbb{R}$, $i = 1, \ldots, n$. In other words, each of the observations contains $p$ covariates. In the real world this could mean having $n$ patients, $p$ observations per patient and trying to predict an outcome such as having a certain type of cancer.

*A bit out of context*

**Definition 2.1** (The linear model). *The relationship between an observation $\mathbf{X}_i \in \mathcal{X}$ and its outcome $\mathbf{Y}_i \in \mathcal{Y}$ can be established by a linear model, that is*

$$i = 1, \ldots, n \qquad \mathbf{Y}_i = \sum_{j=1}^{p} \beta_j \mathbf{X}_i^{(j)} + \varepsilon_i \qquad (2.1)$$

*Any assumptions on $\varepsilon_i$ ?*

Instead of seeing each observation individually we can deal with all of them together by expressing the linear model in matrix notation

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \qquad (2.2)$$

3

**Definition 2.2.** *(a)* **X** *is called the* **design matrix***. It has dimension $n \times p$.*

**X** *consists of stacking the vectors relative to each observation inside of a matrix*

$$X = \begin{bmatrix} - & X_1^T & - \\ & \vdots & \\ - & X_n^T & - \end{bmatrix}$$

*(b)* $\boldsymbol{\beta}$ *is called the* **parameter vector***. It has dimension $p \times 1$.*

*(c)* $\boldsymbol{\varepsilon}$ *is called the* **error vector***. It has dimension $n \times 1$.*

*(d)* **Y** *is called the* **response vector***. It has dimension $n \times 1$.*

## 2.2   The least squares method

We define the objective function $S(\boldsymbol{\beta})$ as follows

$$S(\boldsymbol{\beta}) = \sum_{i=1}^{n} \varepsilon_i^2 = \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} = (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \tag{2.3}$$

*I would skip it.*

which may be rewritten as

$$\begin{aligned} S(\boldsymbol{\beta}) &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \mathbf{Y}^T\mathbf{Y} - \mathbf{Y}^T\mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^T\mathbf{X}^T\mathbf{Y} + \boldsymbol{\beta}^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta} \\ &= \mathbf{Y}^T\mathbf{Y} - 2\boldsymbol{\beta}^T\mathbf{X}^T\mathbf{Y} + \boldsymbol{\beta}^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta} \end{aligned}$$

The least squares method aims at finding the vector $\hat{\boldsymbol{\beta}}$ minimizing $S$, that is

$$\hat{\boldsymbol{\beta}} := \arg\min_{\boldsymbol{\beta}} S(\boldsymbol{\beta})$$

We find $\hat{\boldsymbol{\beta}}$ by differentiating $S$ with respect to $\boldsymbol{\beta}$ and setting the result to 0.

$$\frac{\partial}{\partial \boldsymbol{\beta}} S(\hat{\boldsymbol{\beta}}) = 0$$

$$\implies \frac{\partial}{\partial \hat{\boldsymbol{\beta}}} \left( \mathbf{Y}^T\mathbf{Y} - 2\hat{\boldsymbol{\beta}}^T\mathbf{X}^T\mathbf{Y} + \hat{\boldsymbol{\beta}}^T\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} \right) = 0$$

$$\implies -2\mathbf{X}^T\mathbf{Y} + 2\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = 0$$

$$\implies \mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{Y} \tag{2.4}$$

where equation (2.4) is called the least squares normal equations.

If we assume that $\mathbf{X}^T\mathbf{X}$ is invertible, then (2.4) yields that our least squares estimator $\hat{\boldsymbol{\beta}}$ is ~given by~

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} \tag{2.5}$$

We are interested in estimating the quality of our prediction. The residuals can help us do that.

**Definition 2.3** (Residuals). *For a given set of observations* $\mathbf{Y}$*, the **residuals** (or **vector of residuals**) is the difference between the prediction of our model and the observed value, that is*

$$X(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \in \mathbb{R}^n$$

*Delete this sentence*

However, since the residuals take into account the sign of the difference, they may partially cancel out. We would like a measure that indicates how far our predictions are from the measurements. We will use the prediction error for for that purpose.

**Definition 2.4** (Prediction error). *For a given set of observations* $\mathbf{Y}$*, the **prediction error** is the squared* $\ell^2$*-norm of the difference between the prediction of our model and the observed value. In other words, it is the squared residuals, that is*

$$\|X(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\|_2^2 \in \mathbb{R}$$

- What about asymptotic properties of $\hat{\beta}$ or $X\hat{\beta}$ ?
- What about optimality ?

# Chapter 3

# Theory for LASSO in high dimensions

## 3.1 Assuming the truth is linear

*[handwritten: Here you need a long motivation]*

We work with an underdetermined system : there are more variables than equations, or in our context, there are more parameters than observations $(p > n)$.

We define $\hat{\boldsymbol{\beta}}$ as follows

$$\hat{\boldsymbol{\beta}} := \arg\min_{\boldsymbol{\beta}} \left\{ \frac{\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2}{n} + \lambda\|\boldsymbol{\beta}\|_1 \right\} \tag{3.1}$$

*[handwritten: Why exactly this definition ?]*

**Lemma 3.1** (Basic Inequality).

$$\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq 2\frac{\varepsilon^T \mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)}{n} + \lambda\|\boldsymbol{\beta}^0\|_1$$

*Proof.* By definition of $\hat{\boldsymbol{\beta}}$, we have that

$$\forall \boldsymbol{\beta} \quad \frac{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2}{n} + \lambda\|\boldsymbol{\beta}\|_1$$

In particular for $\boldsymbol{\beta} = \boldsymbol{\beta}^0$ we have

*[handwritten: mention that this is the true parameter of the model]*

$$\frac{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^0\|_2^2}{n} + \lambda\|\boldsymbol{\beta}^0\|_1$$

7

{ We now replace $\mathbf{Y}$ using equation (2.2). :

$$\frac{\|\mathbf{Y} - \mathbf{X}\hat{\beta}\|_2^2}{n} + \lambda\|\hat{\beta}\|_1 \leq \frac{\|\mathbf{Y} - \mathbf{X}\beta^0\|_2^2}{n} + \lambda\|\beta^0\|_1$$

$$\Longrightarrow \frac{\|(\mathbf{X}\beta^0 + \varepsilon) - \mathbf{X}\hat{\beta}\|_2^2}{n} + \lambda\|\hat{\beta}\|_1 \leq \frac{\|(\mathbf{X}\beta^0 + \varepsilon) - \mathbf{X}\beta^0\|_2^2}{n} + \lambda\|\beta^0\|_1$$

$$\Longrightarrow \left( \frac{\|\mathbf{X}(\beta^0 - \hat{\beta}) + \varepsilon\|_2^2}{n} + \lambda\|\hat{\beta}\|_1 \leq \frac{\|\mathbf{X}(\beta^0 - \beta^0) + \varepsilon\|_2^2}{n} + \lambda\|\beta^0\|_1 \right) \quad \text{omit this line}$$

$$\Longrightarrow \frac{\langle \mathbf{X}(\beta^0 - \hat{\beta}) + \varepsilon, \mathbf{X}(\beta^0 - \hat{\beta}) + \varepsilon \rangle}{n} + \lambda\|\hat{\beta}\|_1 \leq \frac{\|\varepsilon\|_2^2}{n} + \lambda\|\beta^0\|_1$$

$$\Longrightarrow \frac{\|\mathbf{X}(\beta^0 - \hat{\beta})\|_2^2 + \|\varepsilon\|_2^2 + 2\langle \mathbf{X}(\beta^0 - \hat{\beta}), \varepsilon \rangle}{n} + \lambda\|\hat{\beta}\|_1 \leq \frac{\|\varepsilon\|_2^2}{n} + \lambda\|\beta^0\|_1$$

$$\Longrightarrow \frac{\|\mathbf{X}(\hat{\beta} - \beta^0)\|_2^2}{n} + \lambda\|\hat{\beta}\|_1 \leq \frac{2\langle \mathbf{X}(\hat{\beta} - \beta^0), \varepsilon \rangle}{n} + \lambda\|\beta^0\|_1$$

$$\Longrightarrow \frac{\|\mathbf{X}(\hat{\beta} - \beta^0)\|_2^2}{n} + \lambda\|\hat{\beta}\|_1 \leq \frac{2\varepsilon^T \mathbf{X}(\hat{\beta} - \beta^0)}{n} + \lambda\|\beta^0\|_1$$

This completes the proof.

$\square$

Some explanations are needed

Let

$$\mathcal{T} := \left\{ \max_{1 \leq j \leq p} 2\frac{|\varepsilon^T \mathbf{X}^{(j)}|}{n} \leq \lambda_0 \right\}$$

**Lemma 3.2** (Lemma 6.2.). *For all $t > 0$ and*   ← assumption missing ?

$$\lambda_0 := 2\sigma\sqrt{\frac{t^2 + 2\log p}{n}}$$

*we have*

$$\mathbb{P}(\mathcal{T}) \geq 1 - 2\exp\left[-t^2/2\right]$$

*Proof.* We define

$$V_j := \frac{\varepsilon^T \mathbf{X}^{(j)}}{\sqrt{n\sigma^2}}$$

Then we have

$$\mathbb{P}(\mathscr{T}) = \mathbb{P}\left(\max_{1 \le j \le p} 2\frac{\left|\varepsilon^T \mathbf{X}^{(j)}\right|}{n} \le 2\sigma\sqrt{\frac{t^2 + 2\log p}{n}}\right)$$

$$= \mathbb{P}\left(\max_{1 \le j \le p} \left|\frac{\varepsilon^T \mathbf{X}^{(j)}}{\sqrt{n\sigma^2}}\right| \le \sqrt{t^2 + 2\log p}\right)$$

$$= \mathbb{P}\left(\max_{1 \le j \le p} |V_j| \le \sqrt{t^2 + 2\log p}\right)$$

$$= 1 - \mathbb{P}\left(\max_{1 \le j \le p} |V_j| > \sqrt{t^2 + 2\log p}\right)$$

$$= 1 - \mathbb{P}\left(\bigcup_{j=1}^{p} |V_j| > \sqrt{t^2 + 2\log p}\right)$$

$$\ge 1 - \sum_{j=1}^{p} \mathbb{P}\left(|V_j| > \sqrt{t^2 + 2\log p}\right)$$

$$\ge 1 - p\,\mathbb{P}\left(|V_j| > \sqrt{t^2 + 2\log p}\right) \tag{3.2}$$

Now, let us define $\zeta := \sqrt{t^2 + 2\log p}$.
Since $V_j$ is $\mathscr{N}(0,1)$-distributed and $\zeta > 0$.

$$\mathbb{P}(V_j > \zeta) = \frac{1}{\sqrt{2\pi}}\int_{\zeta}^{\infty} e^{-y^2/2}dy$$

$$< \frac{1}{\sqrt{2\pi}}\int_{\zeta}^{\infty} \frac{y}{\zeta}\, e^{-y^2/2}dy$$

$$= \frac{1}{\zeta\sqrt{2\pi}}\int_{\zeta}^{\infty} y\, e^{-y^2/2}dy$$

$$= \frac{1}{\zeta\sqrt{2\pi}}e^{-\zeta^2/2}$$

We note that $p \ge 2 \implies \zeta\sqrt{2\pi} \ge 1$ therefore

$$\mathbb{P}(V_j > \zeta) < e^{-\zeta^2/2}$$

Moreover by symmetry of the $\mathscr{N}(0,1)$ distribution,

$$\mathbb{P}(|V_j| > \zeta) \Leftarrow \mathbb{P}(V_j > \zeta) + \mathbb{P}(-V_j < -\zeta) \Big) \text{ omit this line}$$
$$= 2\mathbb{P}(V_j > \zeta)$$
$$< 2e^{-\zeta^2/2}$$

Thus by definition of $\zeta$

$$\mathbb{P}(|V_j| > \zeta) < 2e^{-\zeta^2/2}$$
$$= 2\exp\left[\frac{-\sqrt{t^2 + 2\log p}^2}{2}\right]$$
$$= 2\exp\left[\frac{-t^2}{2} - \log p\right]$$
$$= 2\exp\left[\frac{-t^2}{2}\right]\exp\left[\log\frac{1}{p}\right]$$
$$= \frac{2}{p}\exp\left[\frac{-t^2}{2}\right]$$

omit this

Inserting this result into (3.2) we obtain

$$\mathbb{P}(\mathscr{T}) \geq 1 - p\,\mathbb{P}\left(|V_j| > \sqrt{t^2 + 2\log p}\right)$$
$$\geq 1 - p\frac{2}{p}\exp\left[\frac{-t^2}{2}\right]$$
$$= 1 - 2\exp\left[\frac{-t^2}{2}\right]$$

$\square$

**Corollary 3.3** (Consistency of the LASSO). *Assume $\sigma^2 = 1$ for all $j$. We define the regularization parameter as*

$$\lambda = 4\hat{\sigma}^2\sqrt{\frac{t^2 + 2\log p}{n}}$$

where $\hat{\sigma}$ is some estimator of $\boldsymbol{\sigma}$.

Then with probability at least $1 - \alpha$, where $\alpha := 2\exp(-t^2/2) + \mathbb{P}(\hat{\sigma} \leq \boldsymbol{\sigma})$ we have

$$2\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} \leq 3\lambda\|\boldsymbol{\beta}^0\|_1$$

*proof of it ?*

**Lemma 3.4** (Lemma 6.3.)**.** *We have on $\mathscr{T}$, with $\lambda \geq 2\lambda_0$,*

$$2\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}_{S_0^c}\|_1 \leq 3\lambda\|\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}^0_{S_0}\|_1 \qquad (3.3)$$

*Proof.* We start with the Basic Inequality

$$\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq 2\frac{\varepsilon^T\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)}{n} + \lambda\|\boldsymbol{\beta}^0\|_1$$

Now since we are on $\mathscr{T}$ and since $2\lambda_0 \leq \lambda$

*why does it hold ?*

$$\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \lambda_0\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0\|_1 + \lambda\|\boldsymbol{\beta}^0\|_1$$

$$2\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + 2\lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \lambda\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0\|_1 + 2\lambda\|\boldsymbol{\beta}^0\|_1$$

Let $\boldsymbol{\beta}_{j,S} := \boldsymbol{\beta}_j 1\{j \in S\}$. We use the triangle inequality on the left hand side

$$\begin{aligned}
\|\hat{\boldsymbol{\beta}}\|_1 &= \|\hat{\boldsymbol{\beta}}_{S_0}\|_1 + \|\hat{\boldsymbol{\beta}}_{S_0^c}\|_1 \\
&= \|\boldsymbol{\beta}^0_{S_0} - \boldsymbol{\beta}^0_{S_0} + \hat{\boldsymbol{\beta}}_{S_0}\|_1 + \|\hat{\boldsymbol{\beta}}_{S_0^c}\|_1 \\
&\geq \|\boldsymbol{\beta}^0_{S_0}\|_1 - \|\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}^0_{S_0}\|_1 + \|\hat{\boldsymbol{\beta}}_{S_0^c}\|_1
\end{aligned}$$

whereas on the right hand side

$$\begin{aligned}
\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0\|_1 &= \|(\hat{\boldsymbol{\beta}}_{S_0} + \hat{\boldsymbol{\beta}}_{S_0^c}) - (\boldsymbol{\beta}^0_{S_0} + \underbrace{\boldsymbol{\beta}^0_{S_0^c}}_{=0})\|_1 \\
&= \|\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}^0_{S_0}\|_1 + \|\hat{\boldsymbol{\beta}}_{S_0^c}\|_1
\end{aligned}$$

Injecting these two results, we get that

$$2\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + 2\lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \lambda\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0\|_1 + 2\lambda\|\boldsymbol{\beta}^0\|_1$$

$$\implies 2\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + 2\lambda\left(\|\boldsymbol{\beta}_{S_0}^0\|_1 - \|\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0\|_1 + \|\hat{\boldsymbol{\beta}}_{S_0^c}\|_1\right)$$

$$\leq \lambda\left(\|\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0\|_1 + \|\hat{\boldsymbol{\beta}}_{S_0^c}\|_1\right) + 2\lambda\|\boldsymbol{\beta}^0\|_1$$

$$\implies 2\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + 2\lambda\|\underbrace{\boldsymbol{\beta}_{S_0}^0}_{=\boldsymbol{\beta}^0}\|_1 + \lambda\|\hat{\boldsymbol{\beta}}_{S_0^c}\|_1 \leq 3\lambda\|\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0\|_1 + 2\lambda\|\boldsymbol{\beta}^0\|_1$$

$$\implies 2\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}_{S_0^c}\|_1 \leq 3\lambda\|\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0\|_1$$

□

*Why do we need this condition?*

**Definition 3.5** (Compatibility condition). *We say that the compatibility condition is met for the set $S_0$, if for some $\phi_0 > 0$, and for all $\boldsymbol{\beta}$ satisfying $\|\boldsymbol{\beta}_{S_0^c}\|_1 \leq 3\|\boldsymbol{\beta}_{S_0}\|_1$, it holds that*

*$\overset{\wedge}{\geq}$ ?*

$$\|\boldsymbol{\beta}_{S_0}\|_1^2 \leq (\boldsymbol{\beta}^T \hat{\partial} \boldsymbol{\beta}) \frac{s_0}{\phi_0^2} \tag{3.4}$$

**Theorem 3.6** (Theorem 6.1.). *Suppose the compatibility condition holds for $S_0$. Then on $\mathscr{T}$, we have for $\lambda \geq 2\lambda_0$,*

$$\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0\|_1 \leq 4\lambda^2\frac{s_0}{\phi_0^2}$$

*what is lemma 3.3 ?*

*Proof.* Using lemma 3.3 we have that

$$2\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0\|_1$$

$$= 2\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}_{S_0} + \hat{\boldsymbol{\beta}}_{S_0^c} - \boldsymbol{\beta}_{S_0}^0 - \underbrace{\boldsymbol{\beta}_{S_0^c}^0}_{=0}\|_1$$

$$= 2\frac{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0)\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0\|_1 + \lambda\|\hat{\boldsymbol{\beta}}_{S_0^c}\|_1 \quad \text{(by lemma 3.3)}$$

$$\leq 4\lambda\|\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0\|_1$$

*? why*

$$\leq 4\lambda\sqrt{\left(\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0\right)^T \hat{\boldsymbol{\sigma}} \left(\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0\right) s_0/\phi_0^2}$$

*$\hat{\Sigma}$ ?*

$$\leq \sqrt{\left(\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0\right)^T \mathbf{X}^T\mathbf{X} \left(\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0\right)} \frac{4\lambda\sqrt{s_0}}{\phi_0\sqrt{n}}$$

$$\leq \|\mathbf{X}(\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0)\|_2 \frac{4\lambda\sqrt{s_0}}{\phi_0\sqrt{n}}$$

$$\leq \|\mathbf{X}(\hat{\boldsymbol{\beta}}_{S_0} - \boldsymbol{\beta}_{S_0}^0)\|_2^2 + \frac{4\lambda^2 s_0}{\phi_0^2 n}$$

Where the last inequality follows from $4uv \leq u^2 + 4v^2$. $\qquad\square$

## 3.2 Linear approximation of the truth

*Introduction is needed here.*

Now $\mathbf{Y} := \mathbf{f}^0 + \boldsymbol{\varepsilon}$, therefore $\mathbb{E}[\mathbf{Y}] := \mathbf{f}^0$.

**Lemma 3.7** (New version of the Basic Inequality). $\forall \boldsymbol{\beta}^* \in \mathbb{R}^p$ *we have*

$$\frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{2\boldsymbol{\varepsilon}^T\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)}{n} + \lambda\|\boldsymbol{\beta}^*\|_1 + \frac{\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n} \tag{3.5}$$

*Proof.* By definition of $\hat{\boldsymbol{\beta}}$, we have that

$$\forall \boldsymbol{\beta} \quad \frac{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2}{n} + \lambda\|\boldsymbol{\beta}\|_1$$

In particular for $\boldsymbol{\beta} = \boldsymbol{\beta}^*$ we have

$$\forall \boldsymbol{\beta}^* \quad \frac{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^*\|_2^2}{n} + \lambda\|\boldsymbol{\beta}^*\|_1$$

We since $\mathbf{Y} = \mathbf{f}^0 + \boldsymbol{\varepsilon}$

$$\frac{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^*\|_2^2}{n} + \lambda\|\boldsymbol{\beta}^0\|_1 \qquad \beta^*$$

$$\implies \left( \frac{\|(\mathbf{f}^0 + \boldsymbol{\varepsilon}) - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{\|(\mathbf{f}^0 + \boldsymbol{\varepsilon}) - \mathbf{X}\boldsymbol{\beta}^*\|_2^2}{n} + \lambda\|\boldsymbol{\beta}^*\|_1 \right) \quad omit$$

$$\implies \frac{\|(\mathbf{f}^0 - \mathbf{X}\hat{\boldsymbol{\beta}}) + \boldsymbol{\varepsilon}\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{\|(\mathbf{f}^0 - \mathbf{X}\boldsymbol{\beta}^*) + \boldsymbol{\varepsilon}\|_2^2}{n} + \lambda\|\boldsymbol{\beta}^*\|_1$$

$$\implies \frac{\langle (\mathbf{f}^0 - \mathbf{X}\hat{\boldsymbol{\beta}}) + \boldsymbol{\varepsilon}, (\mathbf{f}^0 - \mathbf{X}\hat{\boldsymbol{\beta}}) + \boldsymbol{\varepsilon} \rangle}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1$$

$$\leq \frac{\langle (\mathbf{f}^0 - \mathbf{X}\boldsymbol{\beta}^*) + \boldsymbol{\varepsilon}, (\mathbf{f}^0 - \mathbf{X}\boldsymbol{\beta}^*) + \boldsymbol{\varepsilon} \rangle}{n} + \lambda\|\boldsymbol{\beta}^*\|_1$$

$$\implies \frac{\|\mathbf{f}^0 - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 + 2\langle \mathbf{f}^0 - \mathbf{X}\hat{\boldsymbol{\beta}}, \boldsymbol{\varepsilon} \rangle}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1$$

$$\leq \frac{\|\mathbf{f}^0 - \mathbf{X}\boldsymbol{\beta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 + 2\langle \mathbf{f}^0 - \mathbf{X}\boldsymbol{\beta}^*, \boldsymbol{\varepsilon} \rangle}{n} + \lambda\|\boldsymbol{\beta}^*\|_1$$

$$\implies \left( \frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{2\langle \mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*), \boldsymbol{\varepsilon} \rangle}{n} + \lambda\|\boldsymbol{\beta}^*\|_1 + \frac{\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n} \right) \quad omit$$

$$\implies \frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{2\boldsymbol{\varepsilon}^T\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)}{n} + \lambda\|\boldsymbol{\beta}^*\|_1 + \frac{\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n}$$

$\square$

**Lemma 3.8** (New version of Lemma 6.3.). *We have on $\mathscr{T}$, with $\lambda \geq 4\lambda_0$,*

$$\frac{4\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + 3\lambda\|\hat{\boldsymbol{\beta}}_{S_*^c}\|_1 \leq 5\lambda\|\hat{\boldsymbol{\beta}}_{S_*} - \boldsymbol{\beta}_{S_*}^*\|_1 + \frac{4\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n} \qquad (3.6)$$

*where $S_* := \{j : \boldsymbol{\beta}_j^* \neq 0\}$.*

*Proof.* We start with the Basic Inequality

$$\frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{2\boldsymbol{\varepsilon}^T\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)}{n} + \lambda\|\boldsymbol{\beta}^*\|_1 + \frac{\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n}$$

Now since we are on $\mathscr{T}$ and since $4\lambda_0 \leq \lambda$

$$\frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \frac{2\boldsymbol{\varepsilon}^T\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)}{n} + \lambda\|\boldsymbol{\beta}^*\|_1 + \frac{\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n}$$

$$\implies \left( \frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + \lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \lambda_0\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_1 + \lambda\|\boldsymbol{\beta}^*\|_1 + \frac{\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n} \right) \quad omit$$

$$\implies 4\frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + 4\lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \lambda\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_1 + 4\lambda\|\boldsymbol{\beta}^*\|_1 + 4\frac{\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n}$$

We use the triangle inequality on the left hand side

$$\|\hat{\boldsymbol{\beta}}\|_1 = \|\hat{\boldsymbol{\beta}}_{S_*}\|_1 + \|\hat{\boldsymbol{\beta}}_{S_*^c}\|_1$$
$$= \|\boldsymbol{\beta}_{S_*}^* - \boldsymbol{\beta}_{S_*}^* + \hat{\boldsymbol{\beta}}_{S_*}\|_1 + \|\hat{\boldsymbol{\beta}}_{S_*^c}\|_1$$
$$\geq \|\boldsymbol{\beta}_{S_*}^*\|_1 - \|\hat{\boldsymbol{\beta}}_{S_*} - \boldsymbol{\beta}_{S_*}^*\|_1 + \|\hat{\boldsymbol{\beta}}_{S_*^c}\|_1$$

whereas on the right hand side

$$\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_1 = \|(\hat{\boldsymbol{\beta}}_{S_*} + \hat{\boldsymbol{\beta}}_{S_*^c}) - (\boldsymbol{\beta}_{S_*}^* + \underbrace{\boldsymbol{\beta}_{S_*^c}^*}_{=0})\|_1$$
$$= \|\hat{\boldsymbol{\beta}}_{S_*} - \boldsymbol{\beta}_{S_*}^*\|_1 + \|\hat{\boldsymbol{\beta}}_{S_*^c}\|_1$$

Injecting these two results, we get that

$$4\frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + 4\lambda\|\hat{\boldsymbol{\beta}}\|_1 \leq \lambda\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_1 + 4\lambda\|\boldsymbol{\beta}^*\|_1 + 4\frac{\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n}$$

$$\implies 4\frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + 4\lambda\left(\|\boldsymbol{\beta}_{S_*}^*\|_1 - \|\hat{\boldsymbol{\beta}}_{S_*} - \boldsymbol{\beta}_{S_*}^*\|_1 + \|\hat{\boldsymbol{\beta}}_{S_*^c}\|_1\right)$$
$$\leq \lambda\left(\|\hat{\boldsymbol{\beta}}_{S_*} - \boldsymbol{\beta}_{S_*}^*\|_1 + \|\hat{\boldsymbol{\beta}}_{S_*^c}\|_1\right) + 4\lambda\|\boldsymbol{\beta}^*\|_1 + 4\frac{\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n}$$

$$\implies 4\frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + 4\lambda\|\underbrace{\boldsymbol{\beta}_{S_*}^*}_{=\boldsymbol{\beta}^*}\|_1 + 3\lambda\|\hat{\boldsymbol{\beta}}_{S_*^c}\|_1$$
$$\leq 5\lambda\|\hat{\boldsymbol{\beta}}_{S_*} - \boldsymbol{\beta}_{S_*}^*\|_1 + 4\lambda\|\boldsymbol{\beta}^*\|_1 + 4\frac{\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n}$$

$$\implies 4\frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{f}^0\|_2^2}{n} + 3\lambda\|\hat{\boldsymbol{\beta}}_{S_*^c}\|_1 \leq 5\lambda\|\hat{\boldsymbol{\beta}}_{S_*} - \boldsymbol{\beta}_{S_*}^*\|_1 + 4\frac{\|\mathbf{X}\boldsymbol{\beta}^* - \mathbf{f}^0\|_2^2}{n}$$

$\square$

**Definition 3.9** (Compatibility condition for general sets). *We say that the compatibility condition holds for the set $S$, if for some constant $\phi(S) > 0$, and for all $\boldsymbol{\beta}$, with $\|\boldsymbol{\beta}_{S^c}\|_1 \leq 3\|\boldsymbol{\beta}_S\|_1$, one has*

$$\|\boldsymbol{\beta}_S\|_1^2 \leq (\boldsymbol{\beta}^T\hat{\boldsymbol{\sigma}}\boldsymbol{\beta})\frac{|S|}{\phi^2(S)}.$$

We define $\mathscr{S}$ as the collection of sets $S$ for which the compatibility condition holds.

**Definition 3.10** (The oracle). *We define the oracle $\beta^*$ as*

$$\beta^* = \arg \min_{\beta:S_\beta \in \mathscr{S}} \left\{ \frac{\|\mathbf{X}\beta - \mathbf{f}^0\|_2^2}{n} + \frac{4\lambda^2 s_\beta}{\phi^2(S_\beta)} \right\}$$

*where $S_\beta := \{j : \beta_j \neq 0\}$, $s_\beta := |S_\beta|$ denotes the cardinality of $S_\beta$ and the factor 4 in the right hand side comes from choosing $\lambda \geq \lambda_0$.*