

Learning with Sparsity

Zhenwei Lin

日期: 2020 年 8 月 14 日

摘 要

本文主要为几种稀疏模型的介绍:

关键词: Lasso; Ridge

1 Linear Regression

dataset is $\{(x_i, y_i)\}_{i=1}^N$, and the model is

$$\begin{aligned} y &= x^T \beta + \epsilon \\ E(\epsilon) &= 0 \\ \text{Var}(\epsilon) &< \infty \end{aligned} \tag{1}$$

2 Computational perspective

$$\min_{\beta} \frac{1}{2} \sum_{i=1}^N (y_i - x_i^T \beta)^2 = \frac{1}{2} \|Y - X\beta\|^2$$

where

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \in R^N, x_{n \times p} = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} \tag{2}$$

为什么这里不把 β_0 单独给出, 因为该数据是被标准化过后的数据, 可以看到标准化后的数据其 β_0 为 0.

$$\nabla L(\beta) = x^T (Y - X\beta) = 0 \Rightarrow (x^T x) \beta = x^T Y$$

1. if $X^T X$ is invertible, we have $\hat{\beta} = (X^T X)^{-1} X^T Y$
2. if $n < p$ means the size of sample is too small, we have $\text{rank}(X^T X) = \text{rank}(X X^T) \leq n$
3. collinearity

To solve these problems, using **Tikhonov Regulation** $\hat{\beta}_{\lambda} = (X^T X + \lambda I)^{-1} X^T Y$ where $X^T X$ is a semidefinite matrix, and I is a positive definitive matrix.

Tikhonov Regulation in model means Ridge Regression.

$$\begin{aligned} \min_{\beta} \quad & \frac{1}{2} \|Y - X\beta\|^2 + \lambda \|\beta\|^2 \\ \Rightarrow \nabla L(\beta) = & -X^T(Y - X\beta) + 2\lambda\beta = 0 \\ \Rightarrow \hat{\beta}_{\lambda} = & (X^T X + \lambda I)^{-1} X^T Y \end{aligned}$$

3 Bayesian Perspective

下面从贝叶斯的角度来对其进行一定的解释。

$$\begin{aligned} y_i &= x_i^T \beta + \epsilon_i \quad \epsilon_i \sim N(0, 1) \\ \therefore f(y_i|x_i) &= \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{(y_i - x_i^T \beta)^2}{2} \right\} \\ \Rightarrow \max_{\beta} \log \prod_{i=1}^N f(y_i|x_i) &\iff \min_{\beta} \frac{1}{2} \|Y - X\beta\|^2 \end{aligned}$$

by $P(A|B) = \frac{P(A,B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$

prior(先验信息)(对 β 的理解)

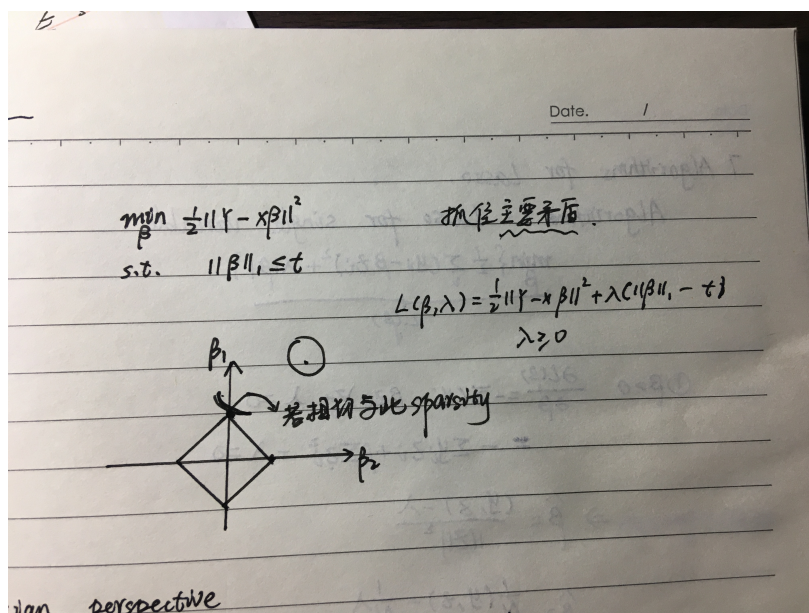
$$\begin{aligned} f(\beta) &= \frac{\lambda^{\frac{p}{2}}}{\sqrt{2\pi}} \exp \{ -\lambda \|\beta\|^2 \} \\ \therefore f(y_i|x_i) &= \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{(y_i - x_i^T \beta)^2}{2} \right\} \\ \therefore p(\beta|y, x) &= \frac{P(Y|\beta, x)P(\beta|x)P(x)}{P(Y|x)P(x)} = \frac{P(Y|\beta, x)P(\beta|x)}{P(Y|x)} \propto P(Y|\beta, x)P(\beta|x) = P(Y|\beta, x)P(\beta) \\ &\iff \max_{\beta} P(Y|x, \beta)P(\beta) \iff \min_{\beta} \frac{1}{2} \|Y - X\beta\|^2 + \tilde{\lambda} \|\beta\|^2 \quad \text{where } \tilde{\lambda} = \frac{\lambda}{2} \end{aligned}$$

因为 x 与 β 无关, 所以 $P(\beta|x) = P(\beta)$ 成立。

4 Lasso(least absolute shrinkage and selection operator)

$$\begin{aligned} \min_{\beta} \quad & \frac{1}{2} \|Y - X\beta\|^2 \\ \text{s.t.} \quad & \|\beta\|_1 \leq t \end{aligned}$$

and the intuition could be shown in the following graph.



In Bayesian Perspective: the prior is $P(B) = \text{constant} \times \exp \{-\lambda \|\beta\|_1\}$ (called laplace distribution)

$$\begin{aligned}
 p(B) &\propto p(Y|x, \beta)p(\beta) \\
 \max_{\beta} \log \prod_{i=1}^N \exp \left\{ -\frac{(y_i - x_i^T \beta)^2}{2} \right\} \exp \{-\lambda \|\beta\|_1\} \\
 &\max_{\beta} -\frac{1}{2} \|Y - x\beta\|^2 - \lambda \|\beta\|_1 \\
 &\min_{\beta} \frac{1}{2} \|Y - x\beta\|^2 + \lambda \|\beta\|_1 \quad (\lambda \geq 0)
 \end{aligned}$$

5 Algorithms for lasso

5.1 Algorithm 1(lasso for single variable)

$$\min_{\beta} \underbrace{\left\{ \frac{1}{2} \sum_i (y_i - \beta z_i)^2 + \lambda \|\beta\|_1 \right\}}_{L(\beta)} \quad (3)$$

1. when $\beta > 0$,

$$\begin{aligned}
 \frac{\partial L(\beta)}{\partial \beta} &= - \sum_i (y_i - \beta z_i) z_i + \lambda = 0 \\
 &= - \sum_i y_i z_i + \beta \sum_i z_i^2 + \lambda = 0 \\
 \Rightarrow \hat{\beta} &= \frac{(y, z) - \lambda}{\|z\|^2} \Rightarrow \hat{\beta} = \frac{\frac{1}{N}(y, z) - \frac{1}{N}\lambda}{\frac{1}{N} \|z\|^2} \\
 &\because \frac{1}{N} \|z\|^2 = 1 \\
 \therefore \hat{\beta} &= \frac{1}{N}(y, z) - \frac{1}{N}\lambda = \frac{1}{N}(y, z) - \tilde{\lambda}
 \end{aligned} \quad (4)$$

2. when $\beta < 0$, similarly we have

$$\hat{\beta} = \frac{1}{N}(y, z) + \tilde{\lambda}$$

3. when $\beta = 0$

$$\begin{aligned} 0 &\in \partial L(0) \\ 0 &\in -\frac{1}{N} \sum (y_i - \beta z_i) z_i + \lambda \partial |\beta| \\ \stackrel{\text{by } \beta=0}{\iff} 0 &\in -\frac{1}{N} (y, z) + \lambda \partial |\beta| \\ \Rightarrow \frac{1}{N} (y, z) &= \lambda \partial |0| = [-\lambda, \lambda] \end{aligned}$$

summary: we have

$$\hat{\beta} = \begin{cases} \frac{1}{N}(y, z) - \tilde{\lambda} & \text{if } \frac{1}{N}(y, z) > \tilde{\lambda} \\ 0 & \text{if } |\frac{1}{N}(y, z)| \leq \tilde{\lambda} \\ \frac{1}{N}(y, z) + \tilde{\lambda} & \text{if } \frac{1}{N}(y, z) < -\tilde{\lambda} \end{cases}$$

从计算结果来说，其内积即相关性越小，则其数值更偏向于 0。

so $\hat{\beta} = \text{Soft}_{\lambda}(\frac{1}{N}(y, z))$

5.2 Algorithm 2(Multivariable of orthogonal Design)

$$\begin{aligned} X_{n \times p}^T X &= I \\ \min_{\beta} \quad & \frac{1}{2N} \|Y - X\beta\|^2 + \lambda \|\beta\|_1 \\ \|Y - X\beta\|^2 &= \|Y\|^2 - 2(Y, X\beta) + \beta^T X^T X \beta \\ &= \|\beta\|^2 - 2(X^T Y, \beta) + \|Y\|^2 \\ &= \|\beta - X^T Y\|^2 + \|Y\|^2 - \|X^T Y\|^2 \\ \iff \min_{\beta} \quad & \left\{ \frac{1}{2N} \|\beta - X^T Y\|^2 + \lambda \|\beta\|_1 \right\} \end{aligned}$$

can be decomposition.

5.3 Algorithm 3 (Cyclic Coordinate Descent)

$$\beta \in \mathbf{R}^p = (\beta_1, \dots, \beta_p)^T$$

and β_2, \dots, β_p is known.

$$\begin{aligned} \min_{\beta} \quad & \left\{ \frac{1}{2N} \sum_{i=1}^N (y_i - \sum_{j=2}^p \beta_j x_{ij} - \beta_1 x_{i1})^2 + \lambda \sum_{i=2}^p |\beta_i| + \lambda |\beta_1| \right\} \\ \min_{\beta} \quad & \left\{ \frac{1}{2N} \sum_{i=1}^N (r_{i1} - \beta_1 x_{i1})^2 + \lambda |\beta_1| \right\} \\ \Rightarrow \hat{\beta}_1 &= \text{Soft}_{\lambda}(\frac{1}{N}(r_i, x_{i1})) \end{aligned} \tag{4}$$

按照一个一个坐标轴分开进行优化。

5.4 Algorithm 4 Proximal gradient descent

$$\min \left\{ \frac{1}{2N} \|Y - X\beta\|^2 + \lambda \|\beta\|_1 \right\}$$

and we know $\frac{1}{2N} \|Y - X\beta\|^2$ is a β -smooth function.

so

$$\begin{aligned} & \|\nabla f(\beta_1) - \nabla f(\beta_2)\| \\ &= \left\| \frac{1}{N} X^T (Y - X\beta_1) - \frac{1}{N} X^T (Y - X\beta_2) \right\| \\ &= \frac{1}{N} \|X^T X(\beta_1 - \beta_2)\| \leq \frac{\lambda_{\max}(X^T X)}{N} \|\beta_1 - \beta_2\| \\ \Rightarrow b &= \frac{\lambda_{\max}(X^T X)}{N} \end{aligned}$$

$$\therefore \beta_{t+1} = \text{prox}_{g/b}(\beta_t - \frac{1}{b} \nabla f(\beta_t))$$

recall in definition, we have

$$\text{prox}_{g/b}(z) = \underset{x}{\operatorname{argmin}} \left\{ \frac{b}{2} \|x - z\|^2 + g(z) \right\}$$

$$\min_x \left\{ \frac{1}{2} \|x - z\|^2 + \frac{\lambda}{b} \|x\|_1 \right\} \Rightarrow \text{prox}_{g/b} = \text{Soft}_{\lambda/b}(z)$$

$$\begin{aligned} \therefore \hat{\beta}_{t+1} &= \text{Soft}_{\lambda/b}(\beta_t + \frac{N}{\lambda_{\max}(X^T X)} \frac{1}{N} X^T (Y - X\beta_t)) \\ &= \text{Soft}_{\lambda/b}(\beta_t - \frac{X^T X}{\lambda_{\max}(X^T X)} \beta_t + \frac{X^T Y}{\lambda_{\max}(X^T X)}) \\ &= \text{Soft}_{\lambda/b} \left\{ (I - \frac{X^T X}{\lambda_{\max}(X^T X)}) \beta_t + \frac{X^T Y}{\lambda_{\max}(X^T X)} \right\} \end{aligned}$$

5.5 Algorithm 5 : AGD \rightarrow FISTA

$$\beta_0 = 0, \alpha_1 = \beta_0, t = 0$$

$$\beta_t = \text{Soft}_{\lambda/b} \left\{ (I - \frac{X^T X}{\lambda_{\max}(X^T X)}) \alpha_t + \frac{X^T Y}{\lambda_{\max}(X^T X)} \right\}$$

$$a_{t+1} = \frac{1 + \sqrt{1 + 4a_t^2}}{2}$$

$$\alpha_{t+1} = \beta_t + \frac{a_t - 1}{a_{t+1}} (\beta_t - \beta_{t+1})$$

and the convergence speed is : $O(\frac{1}{T^2})$

5.6 Algorithm6:ADMM

$$\begin{aligned}
& \min_{\beta} \left\{ \frac{1}{2N} \|Y - X\beta\|^2 + \lambda \|\beta\|_1 \right\} \\
& \iff \min_{\beta} \left\{ \frac{1}{2N} \|Y - X\beta\|^2 \right\} \\
& \quad s.t. \quad \beta - \alpha = 0 \\
L(\beta, \alpha, v) &= \frac{1}{2N} \|Y - X\beta\|^2 + \lambda \|\alpha\|_1 + v^T(\beta - \alpha) + \frac{\rho}{2} \|\beta - \alpha\|^2 \\
\frac{\partial L}{\partial \beta} &= -\frac{1}{N} X^T(Y - X\beta) + v + \rho(\beta - \alpha) \\
&= \left(\frac{1}{N} X^T X + \rho I\right)\beta + v - \frac{1}{N} X^T Y - \rho\alpha = 0 \\
\Rightarrow \hat{\beta} &= \left(\frac{1}{N} X^T X + \rho I\right)^{-1} \left(\frac{1}{N} X^T Y - v + \rho\alpha\right)
\end{aligned}$$

for α :

$$\begin{aligned}
& \min_{\alpha} \left\{ \frac{\rho}{2} \|\alpha - \beta\|^2 - v^T(\alpha - \beta) + \lambda \|\alpha\|_1 \right\} \\
& \iff \min_{\alpha} \left\{ \frac{\rho}{2} \left\| \alpha - \beta - \frac{v}{\rho} \right\|^2 + \lambda \|\alpha\|_1 \right\} \\
& \iff \min_{\alpha} \left\{ \frac{1}{2} \left\| \alpha - \beta - \frac{v}{\rho} \right\|^2 + \frac{\lambda}{\rho} \|\alpha\|_1 \right\} \\
& \Rightarrow \hat{\alpha} = \text{Soft}_{\lambda/\rho}(\beta + \frac{v}{\rho}) \\
\therefore \beta_{t+1} &= \left(\frac{1}{N} X^T X + \rho I\right)^{-1} \left(\frac{1}{N} X^T Y - v + \rho\alpha_t\right) \\
\alpha_{t+1} &= \text{Soft}_{\lambda/\rho}(\beta_{t+1} - \frac{v_t}{\rho}) \\
v_{t+1} &= v_t + \rho(\beta_{t+1} - \alpha_{t+1})
\end{aligned}$$

5.7 Alogrithm 7: Dual Method

$$\begin{aligned}
& \min_{\beta} \left\{ \frac{1}{2N} \|Y - X\beta\|^2 + \lambda \|\beta\|_1 \right\} \\
& \min_{\beta} \left\{ \sup_{\alpha} \left\{ (\alpha, Y - X\beta) - \frac{1}{2} \|\alpha\|^2 \right\} + \lambda \sup_{\lambda} \{(\gamma, \beta) - \delta_{B_{\infty}}(\gamma)\} \right\}
\end{aligned}$$

6 Analysis of LASSO

当我们使用数值方法求解出了 $\hat{\beta} = \text{argmin} \left\{ \frac{1}{2N} \|Y - X\beta\|^2 + \lambda \|\beta\|_1 \right\}$
 但是我们考虑真实的模型为 $Y = X^T \beta^* + \epsilon$, where β^* is real coefficient.
 这时我们对这个问题有三个方法可以来阐述其误差大小

- (a) l_2 error $L_2(\hat{\beta}, \beta^*) = \|\hat{\beta} - \beta^*\|_2$
- (b) Prediction error : $L_p(\hat{\beta}, \beta^*) = \frac{1}{N} \|X\beta^* - X\hat{\beta}\|^2$
- (c) Variable Selection Error : $\{j | \text{sign}(\hat{\beta}) \neq \text{sign}(\beta^*)\}$

下面考虑第一种从 l_2 -error 的角度来看问题。

denote $\hat{v} = \hat{\beta} - \beta^*$

$$\Rightarrow L_2(\hat{\beta}, \beta^*) = \|\hat{v}\|$$

$$\Rightarrow L_p(\hat{\beta}, \beta^*) = \frac{1}{N} \|x\hat{v}\|_2^2$$

recall: α -strong convex function

$$f(y) - f(x) - (\nabla f(x), y - x) \geq \frac{\alpha}{2} \|y - x\|^2$$

这里假设 $f(\beta)$ 是一个 α -strong convex function

再对 $f(\hat{\beta})$ 在 β^* 处进行 Taylor 展开, 则有

$$f(\hat{\beta}) - f(\beta^*) - (\nabla f(\beta^*), \hat{\beta} - \beta^*) = \frac{1}{2N} \hat{v}^T x^T x \hat{v} \geq \frac{\gamma}{2} \|\hat{v}\|^2 \quad (5)$$

由于无法认为 $X^T X$ 满秩, 故在某些区域上满足该条件进行思考, 鉴于此, 我们有 restricted strong convex 成立。

定义 6.1. Restricted Strong convex

$$f(y) - f(x) - (\nabla f(x), y - x) \geq \frac{\gamma}{2} \|y - x\|^2 \quad \forall x, y, y - x \in C$$

定义 6.2. Restricted Eigenvalue Condition(REC)

$$\frac{\hat{v}^T x^T x \hat{v}}{N} \geq \gamma \|\hat{v}\|^2 \quad \hat{v} \in C$$

现在我们先假设 REC 条件成立, 考虑以下两个 lasso 形式

1.

$$\begin{aligned} \min \quad & \frac{1}{2} \|Y - X\beta\|^2 \\ \text{s.t.} \quad & \|\beta\|_1 \leq R \end{aligned} \quad (6)$$

2. Lagrangian Lasso

$$\min_{\beta} \left\{ \frac{1}{2N} \|Y - X\beta\|^2 + \lambda_N \|\beta\|_1 \right\} \quad (7)$$

那么 C 是哪个集合呢? 我们先有如下符号定义:

$$\|\beta^*\|_0 = k$$

说明其有 k 个位置不为 0

$$S \subseteq \{1, 2, \dots, p\}$$

$$v_s = \{v | v_j = 0, j \notin s\}$$

$$v_{s^c} = \{v | v_j = 0, j \in s\}$$

例 6.1.

$$v = (1, 2, 3)$$

$$s = \{2\}$$

$$\Rightarrow v_s = \{(0, 2, 0)\}$$

$$\Rightarrow v_{s^c} = \{(1, 0, 3)\}$$

(8)

因此我们特定的取

$$s = \{i | \beta_i^* \neq 0\} \Rightarrow \|\beta^*\|_1 = \|\beta_s^*\|_1$$

引理 6.1. Assume

$$\hat{\beta} \in \underset{\|\beta\|_1 \leq R}{\operatorname{argmin}} \left\{ \frac{1}{2N} \|Y - X\beta\|^2 \right\}$$

$$\hat{v} = \hat{\beta} - \beta^*$$

then $\|\hat{v}_{s^c}\|_1 \leq \|\hat{v}_s\|_1$ where $s = \{j | \beta_j^* \neq 0\}$

证明. 假设其在边界上达到最优

$$\begin{aligned} R = \|\beta^*\|_1 &\geq \|\hat{\beta}\|_1 = \|\hat{v} + \beta^*\|_1 = \|\hat{v}_s + \hat{v}_{s^c} + \beta^*\|_1 \\ &= \underbrace{\|\hat{v}_{s^c}\|_1}_{\text{取出 } (\beta - \beta^*) \text{ 中 } \beta^* = 0 \text{ 的部分}} + \underbrace{\|\beta^* + \hat{v}_s\|_1}_{\text{取出 } (\beta - \beta^*) \text{ 中 } \beta^* \neq 0 \text{ 的部分}} \stackrel{\text{triangle inequality}}{\geq} \|\beta^*\|_1 - \|\hat{v}_s\|_1 + \|\hat{v}_{s^c}\|_1 \\ &\Rightarrow \|\beta^*\|_1 \geq \|\beta^*\|_1 - \|\hat{v}_s\|_1 + \|\hat{v}_{s^c}\|_1 \end{aligned}$$

□

Definition 6.1.

$$C(S, 1) = \{v | \|v_{s^c}\|_1 \leq 1 \|v_s\|_1\} \quad (9)$$

引理 6.2.

$$p \geq q \geq 1, \quad x \in R^d$$

$$\|x\|_p \leq \|x\|_q \leq d^{\frac{1}{q} - \frac{1}{p}} \|x\|_p \quad (10)$$

等价范数，范数可以被同样上下限控制住。

定理 6.3. Consider Constrained LASSO. Assume REC over $C(S, 1)$ then

- (a) $\|\hat{\beta} - \beta^*\| \leq \frac{4\sqrt{k}}{\gamma} \left\| \frac{x^T \epsilon}{N} \right\|_\infty$
- (b) if we further assume that $\epsilon \stackrel{i.i.d}{\sim} N(0, \sigma^2)$ then for any $0 < \delta < 1$ we have

$$\|\hat{\beta} - \beta^*\|_2 \leq C(\delta) \frac{4\sigma}{\gamma} \sqrt{\frac{k \log P}{N}}$$

依概率 $1 - \delta$ 成立，其中 P 为维数 P .

(a)proof: Basci inequality

$$\hat{v} = \hat{\beta} - \beta^*$$

$$\frac{1}{2N} \|Y - X\beta\|^2 \leq \frac{1}{2N} \|Y - X\beta^*\|^2 \quad Y = X\beta^* + \epsilon$$

basic inequality 成立的原因在于 $\hat{\beta}$ 是其最小化后的结果，而 β^* 却只是其中符合模型的最优解

$$\begin{aligned}
& \Rightarrow \|Y - X\beta^* + X\beta^* - X\hat{\beta}\|^2 \leq \|\epsilon\|^2 \\
& \|\epsilon + X(\beta^* - \hat{\beta})\|^2 \leq \|\epsilon\|^2 \\
& \Rightarrow \|X\hat{v}\|^2 \leq 2(X^T \epsilon, \hat{v}) \\
& \Rightarrow \frac{1}{N} \|X\hat{v}\|^2 \stackrel{\text{Cauchy-Schwarz inequality 对任意对偶范数成立}}{\leq} 2 \left\| \frac{X^T \epsilon}{N} \right\|_{\infty} \|\hat{v}\|_1 \\
& \stackrel{\text{by REC}}{\rightarrow} \frac{1}{N} \|X\hat{v}\|_2^2 \geq \gamma \|\hat{v}\|_2^2 \\
& \|\hat{v}\|_1 = \|\hat{v}_s + \hat{v}_{s^c}\|_1 \leq 2 \|\hat{v}_s\|_1 \leq 2\sqrt{k} \|\hat{v}_s\|_2 \leq 2\sqrt{k} \|\hat{v}\|_2 \\
& \Rightarrow \gamma \|\hat{v}\|_2^2 \leq 2 \left\| \frac{X^T \epsilon}{N} \right\|_{\infty} \|\hat{v}\|_1 \leq 4\sqrt{k} \left\| \frac{X^T \epsilon}{N} \right\|_{\infty} \|\hat{v}\|_2 \\
& \Rightarrow \|\hat{v}\|_2 \leq \frac{4\sqrt{k}}{\gamma} \left\| \frac{X^T \epsilon}{N} \right\|_{\infty}
\end{aligned}$$

□

(b)proof:

$$\begin{aligned}
\epsilon_i & \sim N(0, \sigma^2) \quad X = (x_1, \dots, x_p) \quad X^T \epsilon = \begin{pmatrix} x_1^T \epsilon \\ x_2^T \epsilon \\ \vdots \\ x_p^T \epsilon \end{pmatrix} \\
\frac{x_1^T \epsilon}{N} & \sim N(0, \frac{\sigma^2 \|x_1\|^2}{N^2}) \stackrel{\frac{\|x_1\|^2}{N}=1}{=} N(0, \frac{\sigma^2}{N}) \\
p(|z| \geq t) & \stackrel{z \sim N(0,1)}{\leq} 2 \exp(-\frac{t^2}{2\sigma^2}) \\
& \Rightarrow p\left(\left| \frac{x_i^T \epsilon}{N} \right| \geq t\right) \leq 2 \exp(-\frac{Nt^2}{2\sigma^2}) \\
p\left(\left\| \frac{X^T \epsilon}{N} \right\|_{\infty} \geq t\right) & \leq \sum_{j=1}^P p\left(\left| \frac{x_j^T \epsilon}{N} \right| \geq t\right) = \underbrace{2P \exp(-\frac{Nt^2}{2\sigma^2})}_{\delta}
\end{aligned}$$

无穷范数表示其中的最大值

$$\begin{aligned}
& \Rightarrow \log\left(\frac{\delta}{2P}\right) = -\frac{Nt^2}{2\sigma^2} \\
t & = \sqrt{\frac{2\sigma^2}{N} \log\left(\frac{P}{2\delta}\right)} = c(\delta) \cdot \sigma \sqrt{\frac{\log(P)}{N}}
\end{aligned}$$

代入 a 的结论即可得证。

□

这个定理想表达的主要有：

1. σ 越大， $\|\hat{\beta} - \beta^*\|$ 越大不好估计
2. k 越大，信息量越大，则不好找最优解
3. want $\frac{\log P}{N} \rightarrow 0$, we have $O(\log P) < O(N)$ 同时要求 k 小点，eg: $N=100$, 则有 $\Rightarrow P = \exp(10)$

Consider Lagrangian LASSO

引理 6.4. If $\lambda_N \geq \frac{2\|X^T \epsilon\|_\infty}{N}$ then

$$\hat{v} \in C(S, 3) = \{v \mid \|v_{s^c}\|_1 \leq 3 \|v_s\|_1\}$$

证明. By Basic inequality, we have

$$\begin{aligned} & \frac{1}{2N} \|Y - X\hat{\beta}\|^2 + \lambda_N \|\hat{\beta}\|_1 \leq \frac{1}{2N} \|Y - X\beta^*\|^2 + \lambda_N \|\beta^*\|_1 \\ \Rightarrow 0 & \leq \frac{1}{N} \|X\hat{v}\|^2 \leq \frac{\|X^T \epsilon\|_\infty}{N} \|\hat{v}\|_1 + \lambda_N (\|\beta^*\|_1 - \|\hat{\beta}\|_1) \\ \Rightarrow 0 & \leq \frac{\lambda_N}{2} \|\hat{v}\|_1 + \lambda_N (\|\beta^*\|_1 - \|\hat{\beta}\|_1) \\ & \|\hat{\beta}\|_1 = \|\hat{v} + \beta^*\|_1 = \|\hat{v}_s + \beta^* + \hat{v}_{s^c}\|_1 = \|\hat{v}_s + \beta^*\|_1 + \|\hat{v}_{s^c}\|_1 \geq \|\beta^*\|_1 - \|\hat{v}_s\|_1 + \|\hat{v}_{s^c}\|_1 \\ \Rightarrow \|\hat{v}_s\|_1 - \|\hat{v}_{s^c}\|_1 & \geq \|\beta^*\|_1 - \|\hat{\beta}\|_1 \\ \Rightarrow \frac{1}{2} (\|\hat{v}_s\|_1 + \|\hat{v}_{s^c}\|_1) + (\|\hat{v}_s\|_1 - \|\hat{v}_{s^c}\|_1) & \geq 0 \\ \Rightarrow \|\hat{v}_{s^c}\|_1 & \leq 2 \|\hat{v}_s\|_1 \end{aligned}$$

□

定理 6.5. let us consider the **Lagrangian LASSO** problem, and assume that x satisfies REC over $C(S, 3)$, then for any $\lambda_N \geq 2 \left\| \frac{X^T \epsilon}{N} \right\|_\infty$, we have

(a) $\|\hat{\beta} - \beta^*\|_1 \leq \frac{3\sqrt{k}}{\gamma} \lambda_N$

(b) If we further assume $\epsilon \stackrel{i.i.d}{\sim} N(0, \sigma^2)$, then we have $\|\hat{\beta} - \beta^*\| \leq \frac{C(s)}{\gamma} \sigma \sqrt{\frac{k \log P}{N}}$