

Penalized/constrained loss minimization: The linear regression case

Lasso ("least absolute shrinkage & selection operator"), Tibshirani (1996, JRSSB)

Linear model: $y = X\beta + \varepsilon$
 $n \times 1$ $n \times p$ $p \times 1$ $p \times 1$ $n \times 1$
 response \rightarrow predictor/covariate/design
 $E(\varepsilon|X) = 0$
 $\text{Var}(\varepsilon|X) = \sigma^2$

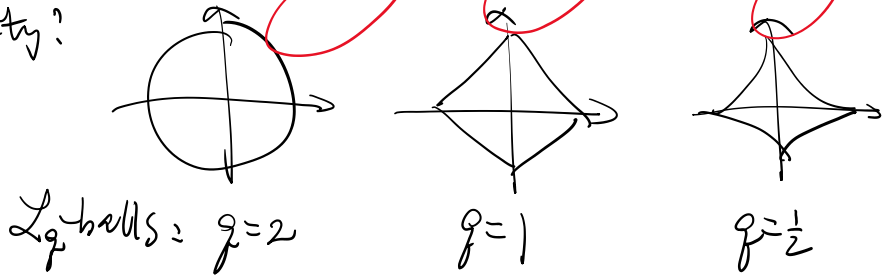
Lasso (L_1 regularization):

minimize $\frac{1}{2n} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$ — penalized form

or

minimize $\|y - X\beta\|_2^2$
 subject to $\|\beta\|_1 \leq t$ — constrained form

Why L_1 induces sparsity?

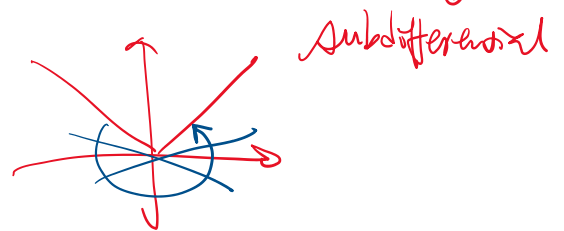


Prop. The Lasso estimator $\hat{\beta}^\lambda$ satisfies

$$\left\| \frac{1}{n} X^T (y - X\hat{\beta}^\lambda) \right\|_\infty = \lambda.$$

Pf. The optimality condition for the Lasso problem is $0 \in \frac{1}{n} X^T (y - X\hat{\beta}^\lambda) + \lambda \partial \|\hat{\beta}^\lambda\|_1$

where $\partial \|\cdot\|_1 = [-1, 1]$.



Related: Dantzig selector $\hat{\beta}^{DS}$ (Candes & Tao, 2007, AoS)

minimize $\|\beta\|_1$

subject to $\left\| \frac{1}{n} X^T (y - X\beta) \right\|_\infty \leq \lambda$

Goal. Derive a nonasymptotic bound on the estimate/prediction error of $\hat{\beta}^\lambda$

& $\hat{\beta}^{DS}$ \forall finite p, n

0, ..., n, p, lambda, sigma

& β^*

\forall finite p, n

High dimensional setting: $p \gg n$, $\lim \| \hat{\beta} \|_0 = \#\{j: \beta_j^* \neq 0\} \triangleq s \ll n$.

ambient dimension

intrinsic dimension

Identifiability issue. When $p > n$, the Gram matrix $\Psi_n = \frac{1}{n} X^T X$ is singular, & hence

β^* is not identifiable

linearly dependent

Remedy. Assume Ψ_n is nondegenerate in certain "sparse" directions.

Restricted eigenvalue condition: $RE(s, c_0)$

$$R(s, c_0) \triangleq \min_{\substack{1 \leq j \leq p \\ |S| \leq s}} \min_{\substack{\delta \neq 0 \\ \|\delta_j\| \leq c_0 \|\delta\|}} \frac{\|X \delta\|_2}{\sqrt{n} \|\delta_S\|_2} > 0$$

Lemma (Basic inequality) Let $S = \text{supp}(\beta^*)$, $\lambda = C \sigma \sqrt{\frac{\log p}{n}}$ w/ $C > 2\sqrt{2}$. Then

it holds w/ probability $\geq 1 - p^{-c_1}$ that

$$\frac{1}{n} \|X(\hat{\beta}^Z - \beta^*)\|_2^2 + \lambda \|\hat{\beta}^Z - \beta^*\|_1 \leq 4\lambda \|\hat{\beta}_S^Z - \beta_S^*\|_1 \leq 4\lambda \sqrt{s} \|\hat{\beta}_S^Z - \beta_S^*\|_2.$$

prediction error

estimate error

can be estimated well on sparse set

pf. The optimality of $\hat{\beta}^Z$ implies

$$\frac{1}{2n} \|y - X \hat{\beta}^Z\|_2^2 + \lambda \|\hat{\beta}^Z\|_1 \leq \frac{1}{2n} \|y - X \beta^*\|_2^2 + \lambda \|\beta^*\|_1.$$

Substituting $y = X \beta^* + \varepsilon$ gives

$$\frac{1}{2n} \|\varepsilon - X(\hat{\beta}^Z - \beta^*)\|_2^2 + \lambda \|\hat{\beta}^Z\|_1 \leq \frac{1}{2n} \|\varepsilon\|_2^2 + \lambda \|\beta^*\|_1,$$

$$\cancel{\frac{1}{2n} \|\varepsilon\|_2^2} + \lambda \|\hat{\beta}^Z - \beta^*\|_1^2 - 2\varepsilon^T X(\hat{\beta}^Z - \beta^*)$$

$$\text{or } \frac{1}{2n} \|X(\hat{\beta}^Z - \beta^*)\|_2^2 \leq \frac{1}{n} \varepsilon^T X(\hat{\beta}^Z - \beta^*) + \lambda \|\beta^*\|_1 - \lambda \|\hat{\beta}^Z\|_1.$$

By Gaussian concentration inequality,

$$\mathbb{P}\left(\| \frac{1}{n} X^T \varepsilon \|_\infty \geq \frac{\lambda}{2}\right) \leq \sum_{j=1}^p \mathbb{P}\left(\left| \frac{1}{n} X_j^T \varepsilon \right| \geq \frac{\lambda}{2}\right) \leq p e^{-\frac{n \lambda^2}{8 \sigma^2}}.$$

$$\mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \epsilon_j^2 \leq \frac{\sigma^2}{2}\right) \leq \sum_{j=1}^n \mathbb{P}\left(\frac{1}{n} \epsilon_j^2 \leq \frac{\sigma^2}{2}\right) \leq p e^{-n \lambda^2 / (2 \sigma^2)}$$

union bound *standardize ϵ_j to have $\|\epsilon_j\|_2 = \sqrt{n}$*

Take $\lambda = C \sqrt{\frac{\log p}{n}}$. The probability becomes

$$p e^{-\frac{C^2}{2} \log p} = p^{1 - C^2/8} \rightarrow 0 \text{ since } C > 2\sqrt{2}.$$

Thus, w/ high probability the event $A = \left\{ \frac{1}{n} \sum_{j=1}^n \epsilon_j^2 \leq \frac{\sigma^2}{2} \right\}$ holds. On event A ,

$$\frac{1}{n} \|\lambda(\hat{\beta}^Z - \beta^*)\|_2^2 \leq 2\|\hat{\beta}^Z - \beta^*\|_1 + 2\lambda(\|\beta^*\|_1 - \|\hat{\beta}^Z\|_1).$$

Adding $2\|\hat{\beta}^Z - \beta^*\|_1$ to both sides yields

$$\begin{aligned} \frac{1}{n} \|\lambda(\hat{\beta}^Z - \beta^*)\|_2^2 + 2\|\hat{\beta}^Z - \beta^*\|_1 &\leq 2\lambda(\|\beta^*\|_1 + \|\beta^*\|_1 - \|\hat{\beta}^Z\|_1) \\ &= 2\lambda(\|\beta^*\|_1 + \|\beta^*\|_1 - \|\hat{\beta}^Z\|_1) \leq 4\lambda \|\beta^*\|_1. \end{aligned}$$

triangle inequality *= 0 on S^c*

Consequently, on event A , $\|\delta_S\|_1 \leq 3\|\delta_S\|_1$, where $\delta = \hat{\beta}^Z - \beta^*$.

Pf. The basic inequality implies

$$\lambda \|\delta\|_1 \leq 4\lambda \|\delta_S\|_1, \text{ or } \|\delta_S\|_1 \leq 3\|\delta_S\|_1.$$

$\|\delta_S\|_1 + \|\delta_{S^c}\|_1$

sparsity for the RB condition

Pf of Thm. By the basic inequality,

$$\frac{1}{n} \|\lambda \delta\|_2^2 \leq 4\lambda \|\delta_S\|_1.$$

On the other hand, by $\text{RB}(s, 3)$ condition,

$$\frac{\|\lambda \delta\|_2}{\sqrt{n} \|\delta_S\|_1} \geq k > 0, \text{ or } \|\delta_S\|_1 \leq \frac{\|\lambda \delta\|_2}{\sqrt{n} k}.$$

Combining two inequalities yields

$$\left(\frac{1}{n} \|\lambda \delta\|_2^2\right) \leq \left(4\lambda \sqrt{s} \frac{\|\lambda \delta\|_2}{\sqrt{n} k}\right)^2, \text{ or } \frac{1}{n} \|\lambda \delta\|_2^2 \leq \frac{16\lambda^2 s}{k^2} = \frac{16C^2 \log p}{k^2 n}.$$

Moreover,

$$\|\delta\|_1 \leq 4\|\delta_S\|_1 \leq 4\sqrt{s} \|\delta_S\|_2 \leq \frac{4\sqrt{s}}{k} \frac{\|\lambda \delta\|_2}{\sqrt{n}}$$

$\log p = o(n)$

$$\begin{aligned} \|D\|_1 &\leq 4\|D_S\|_1 = 4\sqrt{S} \|W_S\|_2 = \frac{4\sqrt{S}}{K} \sqrt{W} \leftarrow \text{stop} = o(n) \\ &\leq \frac{4\sqrt{S}}{K} \frac{4C}{K} \sigma \sqrt{\frac{\log p}{n}} = \frac{16C}{K^2} \sigma \sqrt{\frac{\log p}{n}} \quad \# \\ &\quad \text{stop} = o(n) \end{aligned}$$

Remark. Prediction is easier than estimation.

Question. What if our only goal is prediction? Is the Lasso estimator still prediction consistent w/o any assumption on the design matrix?

Thm. Assume $\|B^*\|_1 \leq K$. The solution $\hat{\beta}$ to the Lasso problem

$$\underset{\beta}{\text{minimize}} \|y - X\beta\|_2^2 \quad \text{subject to } \|\beta\|_1 \leq K$$

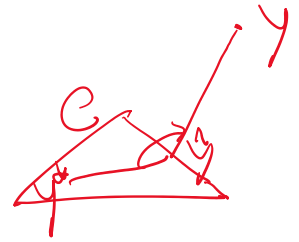
satisfies, w/ high probability, that

$$\frac{1}{n} \|X(\hat{\beta} - \beta^*)\|_2^2 \leq CK \sigma \sqrt{\frac{\log p}{n}}.$$

Prf. By definition, $\hat{y} = X\hat{\beta}$ is the projection of y onto the compact convex set $C = \{X\beta : \|\beta\|_1 \leq K\}$. Also, since $\|B^*\|_1 \leq K$,

we have $y^* = X\beta^* \in C$. Thus,

$$\begin{aligned} 0 &\geq (\hat{y} - y^*)^T (\hat{y} - y) = (\hat{y} - y^*)^T (\hat{y} - y^* - y + y^*) \\ &= \|\hat{y} - y^*\|_2^2 - (\hat{y} - y^*)^T (y - y^*), \end{aligned}$$



$$\text{or } \frac{1}{n} \|\hat{y} - y^*\|_2^2 \leq \frac{1}{n} (\hat{y} - y^*)^T (y - y^*) = \frac{1}{n} \varepsilon^T X(\hat{\beta} - \beta^*)$$

$$\leq \underbrace{\left\| \frac{1}{n} \varepsilon^T X \right\|_\infty}_{\text{w.h.p.}} \underbrace{\|\hat{\beta} - \beta^*\|_1}_{\leq \frac{\lambda}{2}} \leq \frac{\lambda}{2} \cdot 2K = CK \sigma \sqrt{\frac{\log p}{n}} \quad \#$$

$\lambda = C\sigma \sqrt{\frac{\log p}{n}}$

Remark. There is a rate loss in bounding $\|\hat{\beta} - \beta^*\|_1$.