

2.2 The Lasso estimator

The Lasso (Tibshirani 1996) performs

$$(\hat{\beta}_\lambda^L, \hat{\beta}_\lambda^L) = \underset{(\mu, \beta) \in \mathbb{R} \times \mathbb{R}^p}{\operatorname{argmin}} \left\{ \frac{1}{2n} \|Y - \mu \mathbf{1} - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

Here $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$

After centering and scaling the cols of X (to have ℓ_2 -norm \sqrt{n}) and centering Y , we can remove μ from the objective

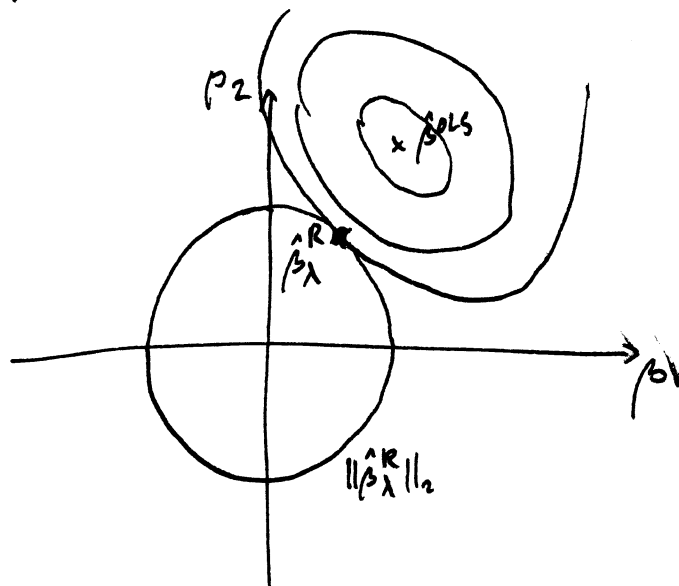
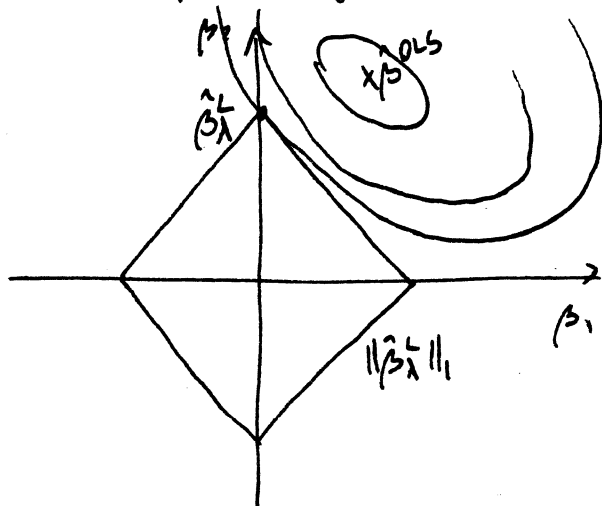
$$Q_\lambda(\beta) = \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

Any minimizer of Q_λ , $\hat{\beta}_\lambda^L$ must also minimize

$$\|Y - X\beta\|_2^2 \text{ subj to } \|\beta\|_1 \leq \|\hat{\beta}_\lambda^L\|_1$$

Indeed if β has $\|\beta\|_1 \leq \|\hat{\beta}_\lambda^L\|_1$ and $\|Y - X\beta\|_2^2 < \|Y - X\hat{\beta}_\lambda^L\|_2^2$ then $Q_\lambda(\beta) < Q_\lambda(\hat{\beta}_\lambda^L)$, contradicting minimality of $\hat{\beta}_\lambda^L$. Similarly, $\hat{\beta}_\lambda^R$ must minimize

$$\|Y - X\beta\|_2^2 \text{ subj to } \|\beta\|_2 \leq \|\hat{\beta}_\lambda^R\|_2$$



The contours of $\beta \mapsto \|Y - X\beta\|_2^2$ are ellipsoids centered at $\hat{\beta}^{OLS}$ and the Lasso and ridge solutions occur when these intersect the ℓ_1 -ball $\{\beta: \|\beta\|_1 \leq \|\hat{\beta}_\lambda^L\|_1\}$ and the ℓ_2 -ball $\{\beta: \|\beta\|_2 \leq \|\hat{\beta}_\lambda^R\|_2\}$ respectively. The fact that the ℓ_1 -ball has 'corners' means some components of $\hat{\beta}_\lambda^L$ are exactly zero.

2.2.1 Prediction error with no design assumptions

Suppose cols of X are centered and scaled to have l_2 -norm \sqrt{n} . Assume normal linear model (after centering)

$$Y = X\beta^0 + \varepsilon - \bar{\varepsilon} \mathbf{1} \quad \varepsilon \sim N_n(0, \sigma^2 I)$$

Theorem 9 Let $\hat{\beta}$ be the Lasso estimate when $\lambda = A\sigma\sqrt{\log p/n}$

With probability of at least $1 - 2p^{-(A^2/2-1)}$

$$\frac{1}{n} \|X\beta^0 - X\hat{\beta}\|_2^2 \leq 4A\sigma\sqrt{\frac{\log p}{n}} \|\beta^0\|_1$$

Proof From defⁿ of $\hat{\beta}$

$$\frac{1}{2n} \underbrace{\|Y - X\hat{\beta}\|_2^2}_{X(\beta^0 - \hat{\beta}) + (\varepsilon - \bar{\varepsilon}\mathbf{1})} + \lambda \|\hat{\beta}\|_1 \leq \frac{1}{2n} \underbrace{\|Y - X\beta^0\|_2^2}_{\varepsilon - \bar{\varepsilon}\mathbf{1}} + \lambda \|\beta^0\|_1$$

Rearranging

$$\frac{1}{2n} \|X(\beta^0 - \hat{\beta})\|_2^2 \leq \frac{1}{n} (\varepsilon - \bar{\varepsilon}\mathbf{1})^T X(\hat{\beta} - \beta^0) + \lambda \|\beta^0\|_1 - \lambda \|\hat{\beta}\|_1$$

$\rightarrow 0$ as cols X centered

$$|\varepsilon^T X(\hat{\beta} - \beta^0)| = |(X^T \varepsilon)^T (\hat{\beta} - \beta^0)| \stackrel{\text{Holder}}{\leq} \|X^T \varepsilon\|_\infty \|\hat{\beta} - \beta^0\|_1$$

[If $v \in \mathbb{R}^p$, $\|v\|_\infty = \max_{j=1, \dots, p} |v_j|$,
Hölder's inequality
if $u, v \in \mathbb{R}^p$ $|u^T v| \leq \|u\|_1 \|v\|_\infty$]

Let $\Omega_\lambda = \{\|X^T \varepsilon\|_\infty / n \leq \lambda\}$. We will show later

(Lemma 13) that $P(\Omega_\lambda) \geq 1 - 2p^{-(A^2/2-1)}$. Working on Ω_λ

$$\frac{1}{2n} \|X(\beta^0 - \hat{\beta})\|_2^2 \leq \lambda \|\hat{\beta} - \beta^0\|_1 + \lambda \|\beta^0\|_1 - \lambda \|\hat{\beta}\|_1 \leq 2\lambda \|\beta^0\|_1$$

by Δ inequality. \square

2.7.2 Basic concentration inequalities

Consider the event Ω

$$P(\|X^T \varepsilon\|_\infty / n > \lambda) = P\left(\bigcup_{j=1}^p \{ |X_j^T \varepsilon| / n > \lambda \}\right)$$

union bound

$$\leq \sum_{j=1}^p P(|X_j^T \varepsilon| / n > \lambda)$$

Now $X_j^T \varepsilon / n \sim N(0, \sigma^2/n)$.

Markov's inequality: Given a non-negative r.v. W , $t \mathbb{1}_{\{W \geq t\}} \leq W$

$$P(W \geq t) \leq \frac{1}{t} EW \quad (t > 0)$$

Given any strictly increasing function $\varphi: \mathbb{R} \rightarrow [0, \infty)$ and any r.v. W

$$P(W \geq t) = P(\varphi(W) \geq \varphi(t)) \leq \frac{E\varphi(W)}{\varphi(t)} \quad (\varphi(t) > 0)$$

Apply this with $\varphi(t) = e^{\alpha t}$ yields the Chernoff bd

$$P(W \geq t) \leq \inf_{\alpha} e^{-\alpha t} \underbrace{E e^{\alpha W}}_{\text{mgf } W \text{ (moment generating fn)}}$$

When $W \sim N(0, \sigma^2)$, $E e^{\alpha W} = e^{\alpha^2 \sigma^2 / 2}$, so

$$\begin{aligned} P(W \geq t) &\leq \inf_{\alpha} \exp\left(\underbrace{\frac{\alpha^2 \sigma^2}{2} - \alpha t}_{\frac{\sigma^2}{2} \left\{ \left(\alpha - \frac{t}{\sigma^2} \right)^2 - \frac{t^2}{\sigma^2} \right\}}\right) \\ &= e^{-t^2 / (2\sigma^2)} \end{aligned}$$

Defn 3 A r.v. W with $\mu = EW$ is sub-Gaussian if $\exists \sigma > 0$ s.t.

$$E e^{\alpha(W-\mu)} \leq e^{\alpha^2 \sigma^2 / 2}$$

We say that W is sub-G with parameter σ .

Prop 10 If W is sub-G with parameter σ , then $P(W - EW \geq t) \leq e^{-t^2 / (2\sigma^2)}$

Bounded r.v.-s are sub-G.

Lemma 11 (Hoeffding's Lemma) If W has mean zero and takes values in $[a, b]$ then

W is a sub-G with param $\frac{b-a}{2}$.