

Theorem 11.1 b

Derives error bounds on Minimum of the objective function:

$$J(\beta) = \frac{1}{2N} \|\vec{y} - X\vec{\beta}\|^2 + \lambda_N \|\beta\|_1$$

λ_N : Regularisation Parameter

$\beta \in \mathbb{R}^p$: unknown sparse signal.

$$\vec{y} = X\vec{\beta} + \vec{w} \quad \text{Measurement Vector}$$

$$\vec{w} \sim \mathcal{N}(0, \sigma^2)$$

Side-stuff

if $N > p$ high N Error

if $N < p$

assume more structure
(here sparseness)

$$X \in \mathbb{R}^{N \times p}$$

↓ covariates

$X \in \mathbb{R}^{N \times p}$: Sensing Matrix with Normalised Columns

minimizer of $J(\vec{\beta})$: known by LASSO in statistics.

a) Define Restricted Eigenvalue Condition

Least Squares Objective Function $f_N(\beta) = \frac{1}{2N} \|y - X\beta\|^2$

$$\frac{1}{N} \frac{v^T X^T X v}{\|v\|_2^2} \geq \gamma \quad \forall \text{ non-zero } v \in C \quad (\text{Constrain Set})$$

↓
some const

b) Starting from 11.20 on page 309 explain why $G(\hat{v}) \leq G(v)$

$$G(v) := \frac{1}{2N} \|y - X(\beta^* + v)\|_2^2 + \lambda_N \|\beta^* + v\|_1$$

β^* original

$$G(v) = \frac{\|y - X\beta^*\|_2^2}{2N} + \lambda \|\beta^*\|_1$$

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

$$G(\hat{v}) = \frac{1}{2N} \|y - X\hat{\beta}\|_2^2 + \lambda_N \|\hat{\beta}\|_1$$

Since $\hat{\beta}$ is the minimiser of LASSO

we have that

$$G(\hat{v}) \leq G(v)$$

$$G(\hat{v}) \leq G(v)$$

c) Do the algebra to obtain 11.21

↓

$$\frac{\|X\hat{v}\|_2^2}{2N} \leq \frac{w^T X \hat{v}}{N} + \lambda_N \{ \|\beta^*\|_1 - \|\beta^* + \hat{v}\|_1 \}$$

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

$$G(\hat{v}) \leq G(v)$$

$$\therefore \frac{1}{2N} \|\vec{y} - X\beta^* - X\hat{v}\|_2^2 + \lambda_N \|\beta^* + \hat{v}\|_1 \leq \frac{1}{2N} \|\vec{y} - X\beta^*\|_2^2 + \lambda_N \|\beta^*\|_1$$

$$\frac{1}{2N} \|\vec{y} - X\beta^* - X\hat{v}\|_2^2 \leq \frac{1}{2N} \|\vec{y} - X\beta^*\|_2^2 + \lambda_N \{ \|\beta^*\|_1 - \|\beta^* + \hat{v}\|_1 \}$$

→ Last Eqn on pg 308 of PDF states $y = X\beta^* + w$

$$\therefore \frac{1}{2N} \|w - X\hat{v}\|_2^2 \leq \frac{1}{2N} \|w\|_2^2 + \lambda_N \{ \|\beta^*\|_1 - \|\beta^* + \hat{v}\|_1 \}$$

$$\therefore \frac{\|X\hat{v}\|_2^2}{2N} \leq \frac{(w^T)(X\hat{v})}{N} + \lambda_N \{ \|\beta^*\|_1 - \|\beta^* + \hat{v}\|_1 \}$$

c) Holder's Inequality as taught in EE734

$$\|X\|_p = E[\|X\|^p]^{\frac{1}{p}}$$

$$\text{For } 1 \leq p \leq \infty \text{ and } \frac{1}{p} + \frac{1}{q} = 1$$

if Random Variables X and Y satisfy $0 < \|X\|_p, \|Y\|_q < \infty$

$$\text{then } |E[XY]| \leq \|X\|_p \|Y\|_q$$

$$(w^T X)(\hat{v})$$

where \hat{v}_i : i^{th} entry.

Let $\{B_1, B_2, \dots, B_p\}$ partition Ω

$$\text{Let } Y = \sum \frac{\hat{v}_i I(w \in B_i)}{N} \text{ and } Z = \sum \frac{(w^T X) I(w \in B_i)}{N}$$

$$\therefore E[YZ] = \frac{(w^T X)(\hat{v})}{N}$$

$$\|Z\|_\infty = \frac{\|w^T X\|_\infty}{1}$$

$$\|Y\|_1 = \frac{\|\hat{v}\|_1}{N}$$

$$\therefore E[YZ] \leq \|Y\|_1 \|Z\|_\infty$$

hence 11.22

I was earlier going to state "Trivially follows from Holder's" but then decided to elaborate ;)

c) Derive 11.23

$$\frac{\|X^T \hat{v}\|_2^2}{2N} \leq \frac{\lambda_N}{2} \left\{ \|\hat{v}_S\|_1 + \|\hat{v}_{S^c}\|_1 \right\} + \lambda_N \left\{ \|\hat{v}_S\|_1 - \|\hat{v}_{S^c}\|_1 \right\} \leq \frac{3}{2} \sqrt{k} \lambda_N \|\hat{v}\|_2$$

$$\text{Theorem 11.1 b)} \quad \text{assumes} \quad \lambda_N \geq \frac{2 \|X^T w\|_\infty}{N} \quad (\text{A})$$

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

$$E[f(X)] \geq f(E[X])$$

$$f(x) = x^2$$

$$\frac{\|\hat{v}\|_2^2}{k} \geq \left(\frac{\|\hat{v}_S\|_1}{k} \right)^2$$

$$\begin{aligned} \|\hat{v}\|_2 &\geq \|\hat{v}_S\|_2 \\ &\geq \frac{\|\hat{v}_S\|_1}{\sqrt{k}} \end{aligned} \quad (\text{B})$$

hence 11.23 follows trivially from 11.22, A and B

(f)

Lemma 11-1

Suppose that $\lambda_N \geq 2 \left\| \frac{X^T w}{N} \right\|_{\infty} > 0$ then error $\hat{\gamma} := \hat{\beta} - \beta^*$ associated with

any LASSO solution $\hat{\beta}$ belongs to the cone set $C(S; 3)$

and 11-23 is derived completes proof for error bound of 11-14b

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{3}{8} \sqrt{\frac{k}{N}} \sqrt{N} \lambda_N$$

$$\lambda_N \geq 2 \left\| \frac{X^T w}{N} \right\|_{\infty}$$

Using Lemma 11-1,

Using $\gamma = \hat{\gamma}$ in (11-10) we get

$$\frac{1}{N} \frac{\|X^T \hat{\gamma}\|_2^2}{\|\hat{\gamma}\|_2^2} \geq \gamma \stackrel{(11-23)}{\Rightarrow} \gamma \frac{\|\hat{\gamma}\|_2^2}{2} \leq \frac{3}{2} \lambda_N \sqrt{k} \|\hat{\gamma}\|_2$$

$$\Rightarrow \|\hat{\gamma}\|_2 \leq \frac{3 \lambda_N \sqrt{k}}{8}$$

g) In Lemma 11-1 $\lambda_N \geq 2 \frac{\|X^T w\|_{\infty}}{N}$

is used to show that the error $\hat{\gamma} := \hat{\beta} - \beta^*$ associated with Lasso soln $\hat{\beta}$ belongs to Cone set $C(S; 3)$

(ii)

In RHS of 11-22

$$\frac{\|X^T w\|_{\infty}}{N} \|\hat{\gamma}\|_1 + \lambda_N \{ \|\hat{\gamma}_S\|_1 - \|\hat{\gamma}_{S^c}\|_1 \}$$

In upper bounding $\lambda_N \geq 2 \frac{\|X^T w\|_{\infty}}{N}$ is used.

h) Importance of cone constraint

Consider the un-numbered equation written between 11.19 and 11.20

$$\|\hat{\gamma}\|_1 = \|\hat{\gamma}_S\|_1 + \|\hat{\gamma}_{S^c}\|_1 \leq 2\|\hat{\gamma}_S\|_1 \leq 2\sqrt{k}\|\hat{\gamma}\|_2$$

↓

Cone constraint is used here $[\hat{\gamma} \in C(S, 1)]$
 And this bound is used in final upper bounding of 11.23.

S : support set for original sparse β^*

Cone constraint also appears in eq. 11.12 $\|\hat{\gamma}_{S^c}\| \leq 3\|\hat{\gamma}_S\|_1$

Cone constraint can be thought of as proxy for strong convexity for vectors in some constraint set

i) The particular theorem for special case of Gaussian Noise vector gives bound

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{C\sigma}{\sqrt{S}} \sqrt{\frac{IK \log p}{N}} \quad \begin{matrix} \text{with probability} \\ w.p \geq 1 - e^{-\frac{1}{2}(I-2)\log p} \end{matrix} \quad (11.15)$$

So we get tighter bounds here (albeit with high probability) as compared to theorem 3.

Advantage of Theorem 3

$$\text{Bounds for theorem 3} \quad \|\theta^* - \theta\|_2 \leq \frac{C_0}{\sqrt{S}} \|\theta - \theta_S\|_1 + C\sigma$$

Theorem 3 does not require hard sparsity (even if non-zero entries are more than some 's' entries we would get nice bounds as long as those entries are small (relative term)), whereas theorem in book requires hard sparsity.

j) Common Thread between bounds on 'Dantzig Selector' and LASSO

$$\text{LASSO} \quad \|\hat{x} - x\|_2 \leq \frac{C_0 \|\sigma_k(x)\|_1}{\sqrt{k}} + C_2 \epsilon$$

A: RIP δ_{2k} with $\delta_{2k} < \sqrt{2}-1$

$$y = Ax + e \quad \|e\|_2 \leq \epsilon$$

$$B(y) = \{z : \|Az - y\|_2 \leq \epsilon\}$$

here k : sparsity.

'Dantzig Selector'

A RIP δ_{2k} ($< \sqrt{2}-1$)

$$y = Ax + e$$

$$\|A^T e\|_{\infty} \leq \lambda$$

$$B(y) = \{z : \|A^T(Az - y)\|_{\infty} \leq \lambda\}$$

A different recovery algorithm based on Dantzig selector in the case when $\|A^T e\|_{\infty}$ is small.

Both bounds are minimax estimators

Both bounds have roughly the same flavor of derivation and nearly same looking results. If the case of Gaussian Noise is considered and if m and n are fixed, and we consider the effect of varying k . then the bound based on Dantzig selector (adaptive to k) improves whereas bound based on LASSO does not.