

Let us now apply the results in the previous week to some concrete applications.

### Example 1 (LASSO with exactly sparse models)

Consider the standard linear model

$$y = X\theta^* + w, \quad y, w \in \mathbb{R}^n; \quad X \in \mathbb{R}^{n \times d}; \quad \theta^* \in \mathbb{R}^d.$$

We assume that the ground-truth  $\theta^*$  is exactly sparse with sparsity  $s$ . The LASSO estimator takes the form

$$\hat{\theta} \in \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \underbrace{\frac{1}{2n} \|y - X\theta\|_2^2}_{\mathcal{L}(\theta)} + \lambda \underbrace{\|\theta\|_1}_{R(\theta)} \right\}$$

By taking  $S$  to be the support of  $\theta^*$  (in particular,  $|S| = s$ ),  $\mathcal{C}$  takes the form

$$\mathcal{C} = \left\{ \Delta \in \mathbb{R}^d : \|\Delta_S\|_1 \leq 3\|\Delta_S\|_1 \right\}.$$

Moreover, to show that  $\mathcal{C}$  satisfies the RSC property, it suffices to prove that

$$\frac{\|X\Delta\|_2^2}{n} \geq k \|\Delta\|_2^2 \quad \forall \Delta \in \mathcal{C} \quad (*)$$

for some  $k > 0$ . Before we do that, let us use (\*) to obtain a bound on the statistical error of the LASSO estimate:

Proposition 1: Suppose that

(i)  $X$  satisfies (\*);

(ii)  $\frac{\|X_j\|_2}{\sqrt{n}} \leq 1 \quad \forall j$ , where  $X_j$  is the  $j^{\text{th}}$  column of  $X$ ;

(iii)  $w$  is Sub-Gaussian with parameter  $\sigma > 0$ ; i.e.,  $w$  has zero mean and for any fixed  $v \in \mathbb{R}^n$  with  $\|v\|_2 = 1$ ,

$$\Pr \left[ |v^T w| \geq t \right] \leq 2 \exp \left( -\frac{t^2}{2\sigma^2} \right) \quad \forall t > 0;$$

(iv)  $\theta^*$  is  $s$ -sparse.

Then, by taking  $\lambda = 4\sigma\sqrt{\frac{\log d}{n}}$  in the LASSO estimator,

1/2

we have

$$\|\hat{\theta} - \theta^*\|_2^2 \leq 36 \frac{\sigma^2}{K^2} \cdot \frac{s \log d}{n}$$

with probability at least  $1 - \frac{2}{d}$ .

Proof: Consider the model subspace

$$m \equiv m(S) = \bar{m}(S) = \{\theta \in \mathbb{R}^d : \theta_j = 0 \ \forall j \notin S\},$$

where  $S = \text{supp}(\theta^*)$ . Then, we have  $\theta_{m^c}^* = 0$ . By our general statistical error bound result (Theorem 4 in Week 2's notes), we need to bound  $\Psi^2(\bar{m})$  and  $\lambda$ . For the former, we have

$$\Psi^2(\bar{m}) = \sup_{u \in \bar{m} \setminus \{0\}} \frac{\|u\|_1^2}{\|u\|_2^2} = s.$$

For the latter, recall that we need  $\lambda \geq 2R^*(\nabla \mathcal{L}(\theta^*))$ . Since

$$\nabla \mathcal{L}(\theta^*) = \frac{1}{n} X^T (X\theta^* - y) = -\frac{1}{n} X^T w \quad \text{and} \quad R^*(v) = \|v\|_\infty,$$

we need to bound  $\|X^T w\|_\infty$ . Towards that end, we utilize our statistical assumptions on the model. Specifically, we have

$$\begin{aligned} \Pr\left[ \left| \frac{X_j^T w}{n} \right| \geq t \right] &= \Pr\left[ \left| \left( \frac{X_j}{n} \right)^T w \right| \geq \sqrt{n} t \right] \\ &\leq 2 \exp\left(-\frac{t^2 n}{2\sigma^2}\right) \quad \text{for } j=1, \dots, d. \end{aligned}$$

Hence, by the union bound,

$$\Pr\left[ \left\| \frac{X^T w}{n} \right\|_\infty \geq t \right] \leq 2d \cdot \exp\left(-\frac{t^2 n}{2\sigma^2}\right).$$

Upon setting  $t^2 = \frac{4\sigma^2 \log d}{n}$ , the RHS equals  $\frac{2}{d}$ . It follows that

with probability at least  $1 - \frac{2}{d}$ , we have  $\lambda \geq \frac{2}{n} \|X^T w\|_\infty$ .  $\square$

Now, let us return to (\*). Key to proving it is the following result:

Theorem 1: Suppose that the rows of  $X$  are iid  $\mathcal{N}(0, \Sigma)$  with  $\Sigma \succ 0$ .

Then, there exist constants  $C_1, C_2 > 0$  such that

$$\frac{\|X\Delta\|_2}{\sqrt{n}} \geq \frac{1}{4} \|\Sigma^{1/2}\Delta\|_2 - 9\rho(\Sigma) \sqrt{\frac{\log d}{n}} \|\Delta\|_2 \quad \forall \Delta \in \mathbb{R}^d \quad (**)$$

with probability at least  $1 - C_1 \exp(-C_2 n)$ , where

$$\rho^2(\Sigma) = \max_{1 \leq j \leq d} \Sigma_{jj}.$$

### Implications and Remarks

1) Note that  $\forall \Delta \in \mathbb{C}$ ,

$$\|\Delta\|_1 = \|\Delta_s\|_1 + \|\Delta_{s^c}\|_1 \leq 4\|\Delta_s\|_1 \leq 4\sqrt{s}\|\Delta_s\|_2 \leq 4\sqrt{s}\|\Delta\|_2.$$

Thus, if we take  $\Sigma = I$  for simplicity, then Theorem 1 yields

$$\frac{\|X\Delta\|_2}{\sqrt{n}} \geq \frac{1}{4} \|\Delta\|_2 - 36 \sqrt{\frac{s \log d}{n}} \|\Delta\|_2$$

In particular, as long as

$$k = \frac{1}{4} - 36 \sqrt{\frac{s \log d}{n}} > 0 \iff n > 144^2 s \log d,$$

we have (\*) holds.

2) Theorem 1 essentially says that the nullspace of  $X$  cannot have vectors that are too sparse.

3) The subtlety in proving (\*\*) is that it has to hold for all  $\Delta \in \mathbb{R}^d$ , not just a single  $\Delta$ . This requires tools from random matrix theory and empirical process theory.

### Example 2 (LASSO with weakly sparse models)

Let us now consider the setup of Example 1, except that  $\theta^*$  is not necessarily exactly sparse but can be well approximated by a sparse vector. To formalize this, one way is to stipulate that

$$\theta^* \in B_g(R_g) \triangleq \left\{ \theta \in \mathbb{R}^d : \|\theta\|_g^2 \leq R_g \right\}, \quad g \in [0, 1].$$

When  $g=0$ ,  $\theta^*$  is exactly sparse with sparsity  $R_0$ . For  $g \in (0, 1]$ , the ball  $B_g(R_g)$  imposes a certain decay rate on the ordered absolute values of  $\theta^*$ .

In this weakly sparse case,  $\mathcal{C}$  takes the form

$$\mathcal{C} = \left\{ \Delta \in \mathbb{R}^d : \|\Delta_{S^c}\|_1 \leq 3 \|\Delta_S\|_1 + 4 \|\theta_{S^c}^*\|_1 \right\}, \quad S \subseteq \{1, \dots, d\}$$

Note that  $\mathcal{C}$  contains a ball centered at the origin. Hence, it is impossible to have a  $\kappa > 0$  such that

$$\frac{\|X\Delta\|_2^2}{2n} \geq \kappa \|\Delta\|_2^2 \quad \forall \Delta \in \mathcal{C}.$$

This is why in the definition of RSC, we need the extra tolerance term  $\tau(\cdot)$ .

Similar to Proposition 1, we have the following statistical error bound for the weakly sparse case:

Proposition 2: Suppose that

- (i) There exist constants  $\kappa_1, \kappa_2 > 0$  such that

$$\frac{\|X\Delta\|_2^2}{n} \geq \kappa_1 \|\Delta\|_2^2 - \kappa_2 \frac{\log d}{n} \|\Delta\|_1^2 \quad \forall \Delta \in \mathbb{R}^d; \quad \text{--- (***)}$$

- (ii)  $\frac{\|X_j\|_2}{\sqrt{n}} \leq 1 \quad \forall j$ , where  $X_j$  is the  $j^{\text{th}}$  column of  $X$ ;

- (iii)  $w$  is sub-Gaussian with parameter  $\sigma > 0$ ;

- (iv)  $\theta^* \in B_g(R_g)$ , where  $\sqrt{R_g} \left( \frac{\log d}{n} \right)^{\frac{1}{2} - \frac{g}{4}} \leq 1$ .

Then, by taking  $\lambda = 4\sigma\sqrt{\frac{\log d}{n}}$  in the LASSO estimator, we have

$$\|\hat{\Theta} - \Theta^*\|_2^2 \leq \left[ 17 + 128 \frac{k_2}{k_1} \left( \frac{4\sigma}{k_1} \right)^{-8} \right] \cdot \left( \frac{16\sigma^2}{k_1^2} \right)^{1-\frac{8}{2}} R_8 \cdot \left( \frac{\log d}{n} \right)^{1-\frac{8}{2}}$$

with probability at least  $1 - \frac{2}{d}$ .

Proof: For a threshold  $\eta > 0$ , define

$$S_\eta = \{j : |\theta_j^*| > \eta\}$$

and consider the model subspace

$$\mathcal{M} \equiv \mathcal{M}(S_\eta) = \bar{\mathcal{M}}(S_\eta) = \{\theta \in \mathbb{R}^d : \theta_j = 0 \ \forall j \notin S_\eta\}$$

Claim 2: Under the assumptions of Proposition 2 and the assumption

that  $n \geq 64 \frac{k_2}{k_1} |S_\eta| \log d$ , the following RSC property of

$\theta \mapsto \mathcal{L}(\theta) = \frac{1}{2n} \|y - X\theta\|_2^2$  holds:

$$\mathcal{L}(\theta^* + \Delta) \geq \mathcal{L}(\theta^*) + \nabla \mathcal{L}(\theta^*)^\top \Delta + \frac{k_1}{2} \|\Delta\|_2^2 - \underbrace{32k_2 \frac{\log d}{n} \|\theta_{S_\eta^c}^*\|_1^2}_{\tau^2(\theta^*)}$$

$\forall \Delta \in \mathcal{E}$ .

Assuming the claim, by invoking our general statistical error bound result (Theorem 1 in Week 2's notes), we get

$$\|\hat{\Theta} - \Theta^*\|_2^2 \leq \frac{9\lambda^2}{k_1^2} |S_\eta| + \frac{4}{k_1} \left( 32k_2 \frac{\log d}{n} \|\theta_{S_\eta^c}^*\|_1^2 + 2\lambda \|\theta_{S_\eta^c}^*\|_1 \right), \quad (\Delta)$$

where we use the facts that  $\Psi^2(\bar{\mathcal{M}}) = |S_\eta|$  and  $\lambda \geq \frac{2}{n} \|X^\top w\|_\infty$  with probability at least  $1 - \frac{2}{d}$  by our calculation in the proof of Proposition 1. To obtain the final result, it remains to bound  $|S_\eta|$  and  $\|\theta_{S_\eta^c}^*\|_1$ .

First, observe that

$$R_8 \geq \|\theta^*\|_8^8 \geq \sum_{j \in S_\eta} |\theta_j^*|^8 \geq \eta^8 |S_\eta|.$$

It follows that

$$|S_\eta| \leq \eta^{-8} R_8 \quad \text{for any } \eta > 0.$$

Next, we have

$$\|\theta_{S_\eta^c}^*\|_1 = \sum_{j \in S_\eta^c} |\theta_j^*| = \sum_{j \in S_\eta^c} |\theta_j^*|^8 |\theta_j^*|^{-7} \leq R_8 \eta^{1-8}.$$

Substituting these into  $(\Delta)$  and setting  $\eta = \frac{\lambda}{k_1}$  yields

$$\begin{aligned} \|\hat{\Theta} - \Theta^*\|_2^2 &\leq \frac{9\lambda^2}{k_1^2} \eta^{-8} R_g + \frac{8\lambda}{k_1} \cdot \eta^{1-8} R_g \\ &\quad + 128 \frac{k_2}{k_1} \frac{\log d}{n} \eta^{2-28} R_g^2 \\ &= 17 \left(\frac{\lambda^2}{k_1^2}\right)^{1-8/2} R_g + 128 \frac{k_2}{k_1} \left(\frac{\lambda^2}{k_1^2}\right)^{1-8/2} R_g \cdot \left[ R_g \frac{\log d}{n} \left(\frac{\lambda}{k_1}\right)^{-8} \right] \end{aligned}$$

By our choice of  $R_g$  and  $\lambda$ , we have

$$R_g \frac{\log d}{n} \left(\frac{\lambda}{k_1}\right)^{-8} = \underbrace{\left(\frac{4\sigma}{k_1}\right)^{-8} R_g \left(\frac{\log d}{n}\right)^{1-8/2}}_{\leq 1} \leq \left(\frac{4\sigma}{k_1}\right)^{-8}$$

It follows that

$$\begin{aligned} \|\hat{\Theta} - \Theta^*\|_2^2 &\leq \left[ 17 + 128 \frac{k_2}{k_1} \left(\frac{4\sigma}{k_1}\right)^{-8} \right] R_g \left(\frac{\lambda^2}{k_1^2}\right)^{1-8/2} \\ &= \left[ 17 + 128 \frac{k_2}{k_1} \left(\frac{4\sigma}{k_1}\right)^{-8} \right] \cdot \left(\frac{16\sigma^2}{k_1^2}\right)^{1-8/2} \cdot R_g \cdot \left(\frac{\log d}{n}\right)^{1-9/2} \quad \square \end{aligned}$$

Proof of Claim 2: For any  $\Delta \in \mathcal{E}$ , we have

$$\|\Delta\|_1 \leq \|\Delta_{S_\eta}\|_1 + \|\Delta_{S_\eta^c}\|_1 \leq 4 \|\Delta_{S_\eta}\|_1 + 4 \|\Theta_{S_\eta^c}^*\|_1$$

Hence, using the inequality  $(a+b)^2 \leq 2(a^2+b^2) \forall a, b \geq 0$ , we get

$$\begin{aligned} \|\Delta\|_1^2 &\leq 32 \left( \|\Delta_{S_\eta}\|_1^2 + \|\Theta_{S_\eta^c}^*\|_1^2 \right) \\ &\leq 32 \left( |S_\eta| \cdot \|\Delta_{S_\eta}\|_2^2 + \|\Theta_{S_\eta^c}^*\|_1^2 \right) \\ &\leq 32 |S_\eta| \|\Delta\|_2^2 + 32 \|\Theta_{S_\eta^c}^*\|_1^2 \end{aligned}$$

This, together with (\*\*\*), yields

$$\frac{\|X\Delta\|_2^2}{n} \geq \left( k_1 - 32 k_2 \frac{\log d}{n} |S_\eta| \right) \|\Delta\|_2^2 - 32 k_2 \frac{\log d}{n} \|\Theta_{S_\eta^c}^*\|_1^2 \quad \forall \Delta \in \mathcal{E}$$

Our choice of  $n$  implies that  $k_1 - 32 k_2 \frac{\log d}{n} |S_\eta| \geq \frac{k_1}{2}$ . This

completes the proof. □