36-709: Advanced Statistical Theory I

Spring 2020

Lecture 12: February 27

Lecturer: Siva Balakrishnan

Scribe: Lantian Xu

### 12.1Recap

1. Basis pursuit:  $\hat{\theta} = \arg \min \|\theta\|_1$ , want to find  $y = X\theta^*$ ,  $X \in \mathbb{R}^{n \times d}$ ,  $\theta^*$  s-sparse (Thm 7.8) MJW book, etc)

2. Restricted Nullspace Property(RN): {null(X)  $\cap$  { $\|\Delta_{S^c}\|_1 \leq \|\Delta_S\|_1$ } = {0}

3. Pairwise Incoherence(PI):  $\|\frac{X^TX}{n} - I\|_{\infty} \leq \delta$ 

- 4. Restricted Isometry Property(RIP): For any  $|S| \leq s$ ,  $\|\frac{X_S^T X_S}{n} I_S\|_{op} \leq \delta$

5. If PI holds with  $\delta \leq \frac{1}{2s} \Rightarrow$  RN holds 6. If RIP holds with  $\delta \leq \frac{1}{3} \Rightarrow$  RN holds 7. Notice that for a matrix  $X \in \mathbb{R}^{n \times d}$  with N(0, 1) entries, PI holds with high probability if  $n\gtrsim s^2\log d,$  RIP holds with high probability if  $n\gtrsim s\log(\frac{ed}{s})$ 

## Today 12.2

#### 12.2.1Intuition for RIP

In general we want to solve

$$\widehat{\theta} = \underset{\text{subject to } y = X\theta}{\arg\min \|\theta\|_0}$$
(12.1)

but it is hard to compute. We replace the troublesome  $L_0$  norm by  $L_1$  norm to get  $\hat{\theta}$ . For  $\|\hat{\theta}\|_0 \leq \|\theta^*\|_0 \leq s$ ,  $\|\Delta\|_0 \leq 2s$ . In order to argue that  $\theta^*$  is unique, it suffices to show that  $||X\Delta||_2^2 > 0$  for any (at most) 2s-sparse vector  $\Delta$ .

Let us now see that RIP in fact implies this condition. Observe that,

$$\left\|\frac{X\Delta}{\sqrt{n}}\right\|_{2}^{2} = \frac{\Delta^{T}X^{T}X\Delta}{n} = \Delta^{T}(\frac{X^{T}X}{n} - I)\Delta + \|\Delta\|_{2}^{2} \ge -\delta\|\Delta\|_{2}^{2} + \|\Delta\|_{2}^{2} > 0$$

holds when  $\|\frac{X_S^T X_S}{n} - I_S\|_{op} < 1, \forall |S| \le 2s$ . This shows that an RIP condition is sufficient for the success of the  $L_0$  minimization program.

# 12.2.2 LASSO

Suppose now we observe a vector  $y \in \mathbb{R}^n$  and a matrix  $X \in \mathbb{R}^{n \times d}$  that are linked via the standard linear model  $y = X\theta^* + \epsilon$ , where  $\theta^*$  is s-sparse and  $\epsilon$  is a vector of noise variables whose entries have distribution  $N(0, \sigma^2)$ . Our goal is to find a parameter estimation error bound  $\|\Delta\|_2 := \|\widehat{\theta} - \theta^*\|_2$  between the LASSO solution  $\widehat{\theta}$  and the unknown regression vector  $\theta^*$ .

We have different constrained forms of the LASSO, for example

$$\widehat{\theta} = \underset{\|\beta\|_1 \le t}{\operatorname{arg\,min}} \frac{1}{2n} \|y - X\beta\|_2^2 \tag{12.2}$$

or

$$\widehat{\theta} = \operatorname*{arg\,min}_{\beta} \frac{1}{2n} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \tag{12.3}$$

These convex programs are equivalent due to Lagrangian duality theory.

In low dimension setting, we can obtain a bound without s-sparse condition:

$$\lambda_{\min}(\frac{X^T X}{n}) \|\Delta\|_2^2 \le \frac{\|X\Delta\|_2^2}{2n} \le \left|\frac{(X^T \epsilon)^T \Delta}{n}\right| \le \frac{\|X^T \epsilon\|_2 \|\Delta\|_2}{n}$$

if  $\lambda_{\min}(\frac{X^T X}{n}) > 0$ . While in high dimension setting the eigenvalue condition not always holds (for example, when d > n the rank of  $\frac{X^T X}{n}$  would be at most n). For this reason, we have to come up with some useful conditions in high dimension.

**Definition 12.1 (Restricted Eigenvalue condition)** We say a matrix X satisfies the restricted eigenvalue (RE) condition over S with parameters  $(\kappa, \alpha)$  if for any  $\|\Delta_{S^c}\|_1 \leq \alpha \|\Delta_S\|_1$ ,

$$\left\|\frac{X\Delta}{\sqrt{n}}\right\|_2^2 \ge \kappa \|\Delta\|_2^2$$

In class we discussed  $\alpha = 3$  case (as 7.3.2 (A2), MJW book).

**Theorem 12.2** ( $\ell_2$  error between  $\hat{\theta}(12.2)$  and  $\theta^*$ ) If  $\theta^*$  is s-sparse, X satisfies RE condition with parameter ( $\kappa$ ,3) and we pick  $t = \|\theta^*\|_1$  in (12.2), then

$$\|\widehat{\theta} - \theta^*\|_2 \le 4\|\frac{X^T \epsilon}{n}\|_{\infty} \frac{\sqrt{s}}{\kappa} \tag{12.4}$$

For  $\|\frac{X^T \epsilon}{n}\|_{\infty} \lesssim \sqrt{\frac{\log d}{n}}, \|\widehat{\theta} - \theta^*\|_2 \lesssim \sqrt{\frac{s \log d}{n}}.$ 

**Proof:** First given  $t = \|\theta^*\|_1$  the target vector  $\theta^*$  is feasible; from RE condition we have

$$\kappa \|\Delta\|_2^2 \le \left\|\frac{X\Delta}{\sqrt{n}}\right\|_2^2$$

From  $\frac{1}{2n} \|y - X\hat{\theta}\|_2^2 \le \frac{1}{2n} \|y - X\theta^*\|_2^2$  and Holder inequality one can derive

$$\frac{\|X\Delta\|_2^2}{n} \le \left|\frac{2(X^T\epsilon)^T\Delta}{n}\right| \le 2\|\frac{X^T\epsilon}{n}\|_{\infty}\|\Delta\|_1$$

On the other hand, following the proof of the analysis of basis pursuit we obtain that  $\|\Delta_{S^c}\|_1 \leq \|\Delta_S\|_1$ , under our constraint on  $\theta$ . This in turn gives  $\|\Delta\|_1 \leq 2\|\Delta_S\|_1 \leq 2\sqrt{s}\|\Delta\|_2$ ; Putting together these pieces yields  $\|\widehat{\theta} - \theta^*\|_2 \leq 4\|\frac{X^T\epsilon}{n}\|_{\infty}\frac{\sqrt{s}}{\kappa}$ .

**Theorem 12.3** ( $\ell_2$  error between  $\hat{\theta}(12.3)$  and  $\theta^*$ ) If  $\theta^*$  is s-sparse, X satisfies RE condition with parameter ( $\kappa$ ,3) and we have a regularization parameter lower bounded as  $\lambda \geq 2 \|\frac{X^T \epsilon}{n}\|_{\infty}$ , then

$$\|\widehat{\theta} - \theta^*\|_2 \le \frac{3\lambda\sqrt{s}}{\kappa} \tag{12.5}$$

**Proof:** Condition (12.3) gives

$$\frac{1}{2n} \|y - X\widehat{\theta}\|_{2}^{2} + \lambda \|\widehat{\theta}\|_{1} \le \frac{1}{2n} \|y - X\theta^{*}\|_{2}^{2} + \lambda \|\theta^{*}\|_{1}$$

Rearranging yields

$$\frac{1}{2n} \|X\Delta\|_2^2 \le \left|\frac{2(X^T\epsilon)^T\Delta}{n}\right| + \lambda\{\|\theta^*\|_1 - \|\widehat{\theta}\|_1\}$$

Rewriting  $\theta^*$  under s-sparse condition and applying triangle inequality, we have

$$\|\theta^*\|_1 - \|\widehat{\theta}\|_1 = \|\theta^*_S\|_1 - \|\theta^*_S + \Delta_S\|_1 - \|\Delta_{S^c}\|_1 \le \|\Delta_S\|_1 - \|\Delta_{S^c}\|_1$$

Plugging in above inequality yields

$$\frac{1}{2n} \|X\Delta\|_2^2 \le \frac{\lambda}{2} \{3\|\Delta_S\|_1 - \|\Delta_{S^c}\|_1\}$$
(12.6)

since  $\lambda \geq 2 \|\frac{X^T \epsilon}{n}\|_{\infty}$ . Here we are in a situation to apply the RE condition:  $\frac{1}{n} \|X\Delta\|_2^2 \geq \kappa \|\Delta\|_2^2$ . Thus from the inequality in Theorem 12.2,

$$\kappa \|\Delta\|_2^2 \le \frac{3\lambda}{2} \|\Delta_S\|_1 \le \frac{3\lambda\sqrt{s}}{2} \|\Delta\|_2 \tag{12.7}$$

as desired.

12-3

# 12.2.3 In-sample Prediction Error

We want to come up with some methods which do not need strong assumptions to predict well. In general, we want to get the bound

$$\frac{\|X\Delta\|_2^2}{n} \le C\lambda \|\theta^*\|_1$$

without assumptions on X or  $\theta^*$  (remains for later discussion).

On the other hand, a simple inspection of the above proof for the  $\ell_2$  error shows that under RE and if  $\theta^*$  is s-sparse we obtain the faster rates for the in-sample prediction error of:

$$\frac{\|X\Delta\|_2^2}{n} \le \frac{Cs\log p}{n},$$

when  $\lambda$  in the Lagrangian LASSO is chosen as  $\lambda \geq 2 \|\frac{X^T \epsilon}{n}\|_{\infty}$ . In more detail, from (12.6) we obtain the bound:

$$\frac{1}{n} \|X\Delta\|_2^2 \le 3\lambda \|\Delta_S\|_1 \le 3\lambda \sqrt{s} \|\Delta\|_2.$$

Combining this with the RE condition, to obtain that,

$$\frac{1}{n} \| X \Delta \|_2^2 \leq \frac{3\lambda \sqrt{s} \| X \Delta \|_2}{\sqrt{n\kappa}}$$

Re-arranging this we obtain that,

$$\frac{\|X\Delta\|_2^2}{n} \le \frac{Cs\lambda^2}{\kappa},$$

as desired.