

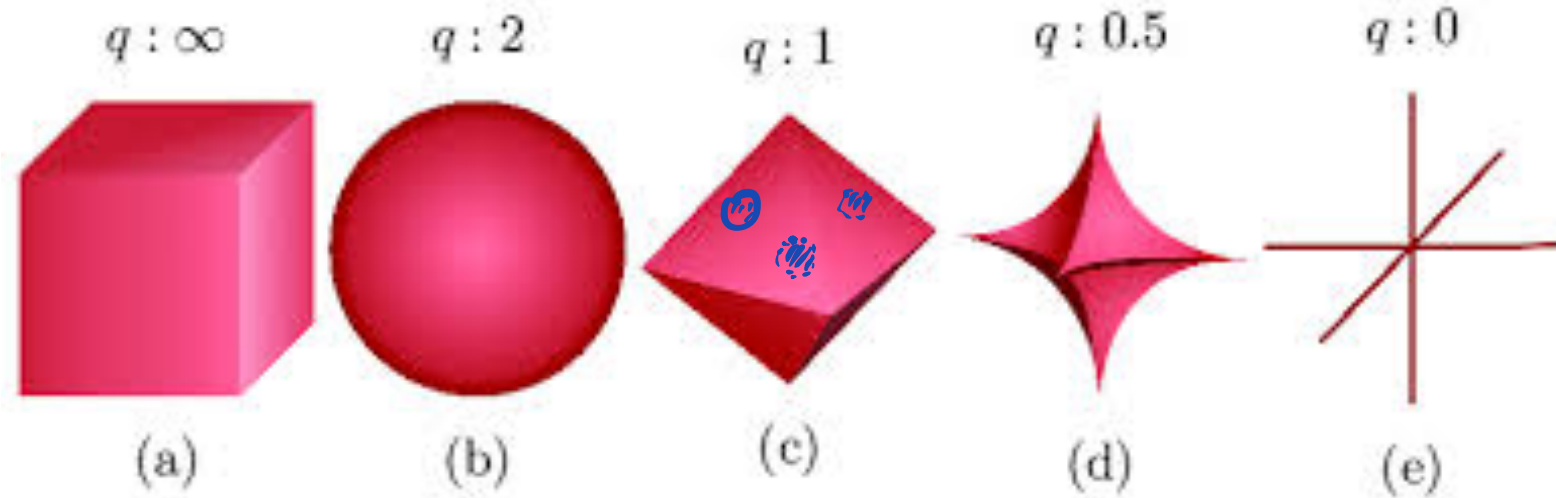
# Lecture 11 Localized Complexity

IEMS 402 Statistical Learning

Northwestern

# Empirical Method of Maurey

# L1 Ball



# Volume Based Bound

Last lecture, we discussed the problem of getting a covering number  $N$  for  $L_1$  balls using  $L_2$  balls.

$$N(\epsilon, B_1^d, \|\cdot\|_2) \tag{1}$$

Using a volume argument, we were able to establish the following result.

$$N(\epsilon, B_1^d, \|\cdot\|_2) \leq N(\epsilon, B_1^d, \|\cdot\|_1) \tag{2}$$

$$N(\epsilon, B_1^d, \|\cdot\|_1) \leq \left(1 + \frac{2}{\epsilon}\right)^d \tag{3}$$

$$\Rightarrow \log N \leq d \log \frac{2}{\epsilon}$$

# Empirical Method of Maurey

**Theorem 1.** When  $\epsilon > \frac{1}{\sqrt{d}}$ ,  $N \leq (2d + 1)^{O(1/\epsilon^2)}$

As a result,  $\log N \lesssim \frac{1}{\epsilon^2} \log(d)$ .

*Proof.* Let's cover the following set:

$$B_1^{d,+} = \{x \in \mathcal{R}^d \mid \|x\|_1 \leq 1 \text{ and } x_i \geq 0 \forall i\}$$

The above set means that  $\sum x_i \leq 1 \forall x_i \geq 0$ .

We can think about a probability distribution over  $\{e_1, \dots, e_d, 0\}$ :

$$z = \sum_{i=1}^d x_i e_i + (1 - \|x\|_1) \cdot 0$$

$\Rightarrow$  sum of  $\begin{pmatrix} 1 \\ \vdots \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ \vdots \\ 1 \end{pmatrix}$   
 $\dots \begin{pmatrix} 0 \\ \vdots \\ 1 \end{pmatrix}, 0$

number of (random) basis that you need  $\leq d^{1/\epsilon^2}$   
 $\rightarrow$  Li norm of sample + data is  $\binom{t}{d}$   
 Randomly sample + basis.  $\Rightarrow$  error:  $s = O(\frac{1}{\sqrt{\epsilon}})$

# Empirical Method of Maurey

This implies the following probabilities.

$$\mathbb{P}[z = e_j] = x_j \forall j \in [d]$$

$$\mathbb{P}[z = 0] = 1 - \|x\|_1$$

With these, we can get a mean of the probability distribution.

$$\mathbb{E}[z] = \sum \mathbb{P}[z = e_j] \cdot e_j + \mathbb{P}[z = 0] \cdot 0 = \sum x_j \cdot e_j = x$$

We will draw  $t$  samples  $z_1, \dots, z_t$  from the distribution where each  $z$  is some  $e_i$ . After drawing the samples, we can take the average of the samples:

$$\bar{z} = \frac{1}{t} \sum_{i=1}^t z_i$$

We want to show that  $\mathbb{E}[\|\bar{z} - x\|_2^2] \leq \epsilon^2$ . If we can do this, then if we take all possible  $\bar{z}$ , we get an  $\epsilon$ -cover of the space using those  $\bar{z}$  since then all  $x$  we can choose will be within  $\epsilon$  of some point in the cover by what we argue above.

# Empirical Method of Maurey vs Volume

- Volume :  $N = (\frac{1}{\epsilon})^d \Rightarrow \log N = d \cdot \log(1/\epsilon)$

This is better when  $\epsilon$  is small

- Empirical method of Maurey :  $N = d^{(\frac{1}{\epsilon^2})}$


$\hookrightarrow$  the optimal dependency over  $\epsilon$

$\Rightarrow \log N = \frac{1}{\epsilon^2} \log(d)$

The dimension dependency  $\log(d)$  is better.

Lasso

The dependency on  $\epsilon$  is bad



# Localized Complexity



# Example: Mean Estimation

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n x_i = \operatorname{argmin}_{\theta} \underbrace{\frac{1}{n} \sum_{i=1}^n (\theta - x_i)^2}_{\hat{L}(\theta)}$$

$x_i \sim \text{Unif}([0, 1])$

$$L(\theta) = \mathbb{E} (\theta - x)^2$$

$\parallel \theta - \theta^* \parallel^2 + \text{Var}(x)$

①  $L(\hat{\theta}) - L(\theta) \leq \text{Rad} = o\left(\frac{\text{Complexity}}{\sqrt{n}}\right)$

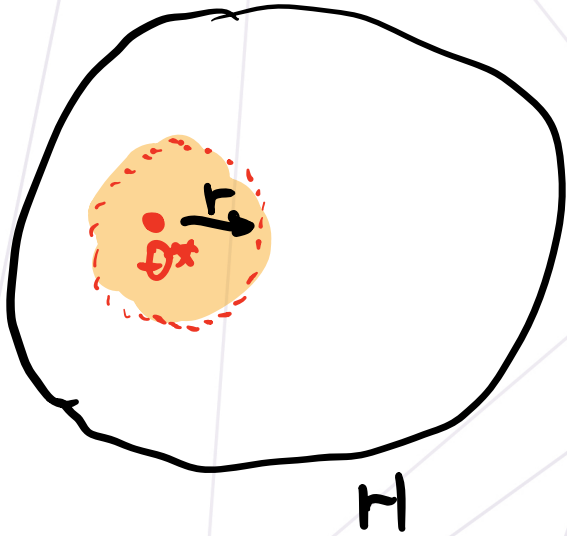
② Right Bound

$$\|\hat{\theta} - \theta^*\| = o\left(\frac{1}{\sqrt{n}}\right)$$



$$\Rightarrow L(\hat{\theta}) - L(\theta^*) = \|\hat{\theta} - \theta^*\|^2 = o\left(\frac{1}{n}\right) \Leftarrow \text{Fast Rate}$$

# Idea: Localized Complexity



$$\phi(r) = \text{Rad} \left( \{ \theta \in H, \|\theta - \hat{\theta}\| \leq r \} \right)$$

$$= \mathbb{E} \sup_{\{ \theta \in H, \|\theta - \hat{\theta}\| \leq r \}} \sum \sigma_i \ell_{\theta}(x_i)$$

expect  $r$  is the right result:

$$\|\hat{\theta} - \theta^*\| \leq r$$

$\Rightarrow$  using Radmecher

previous:

$$\sup_{\theta \in H} | \mathbb{E} - \hat{\mathbb{E}} f |$$

Now :

$$\sup_{\theta \text{ lies in the neighbor of } \theta^*} | \mathbb{E} - \hat{\mathbb{E}} f |$$

$\phi(r) = r \Rightarrow r$  is the final error ??

# Localize Leads to Fast Rate

localize  $\|\theta - \theta^*\| \leq \delta$

Assume Loss  $\leq \|\theta - \theta^*\|^\alpha$

- localized Rademacher complexity linear respect to  $\delta$

$$\frac{\delta}{\sqrt{n}} = \delta^\alpha \Rightarrow \delta = (\sqrt{n})^{\frac{1}{1-\alpha}}$$

$$\Rightarrow \text{final loss} = \delta^\alpha = \left(\frac{1}{\sqrt{n}}\right)^{\frac{\alpha}{1-\alpha}}$$

# Non-parametric Least Square

To estimate the unknown regression function  $f^*$ , we consider the empirical risk minimizer (ERM), which is given by

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2. \quad (2)$$

# Method 1

**Proof of Theorem 1:** Since  $\hat{f}$  is optimal to the ERM problem (2) and  $f^* \in \mathcal{F}$  is feasible, we have

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \leq \frac{1}{n} \sum_{i=1}^n (y_i - f^*(x_i))^2. \quad (3)$$

Also recall that

$$y_i = f^*(x_i) + \sigma w_i, \quad 1 \leq i \leq n.$$

We plug this expression into  $y_i$ 's in equation (3), open the squares and rearrange terms. Doing so gives the "basic inequality"

$$\frac{1}{2} \|\hat{f} - f^*\|_n^2 \leq \frac{\sigma}{n} \sum_{i=1}^n w_i (\hat{f}(x_i) - f^*(x_i)) \leq \text{Red Complexity}. \quad (4)$$

*loss is quadratic growth* ← **Loss**      **Complexity** → *is a linear growth*

Introducing the shorthand  $\Delta := \hat{f} - f^* \in \mathcal{F}^*$ , we rewrite the above basic inequality compactly as

$$\|\Delta\|_n^2 \leq C \|\Delta\| \quad \Leftrightarrow \quad \frac{1}{2} \|\Delta\|_n^2 \leq \frac{\sigma}{n} \sum_{i=1}^n w_i \Delta(x_i).$$

*⇒ ||Δ|| ≤ C*

**square term**      **linear term**      **||Δ|| can be bounded by the linear growth<sup>(5)</sup> speed of the localized complexity**

# We need star shape

localized  
complexity

$$G_n(\delta; \mathcal{F}^*) := \text{Rad}(\{\| \theta - \theta^* \| \leq \delta, \theta \in \mathcal{H}\})$$

**Lemma 1.** If  $\mathcal{F}^*$  is star-shaped, then the function  $\delta \mapsto \frac{G_n(\delta; \mathcal{F}^*)}{\delta}$  is non-increasing on  $(0, \infty)$ . Hence  $\delta^*$  exists and is finite.

**Proof** For any  $0 < \delta < t$ , we want to show that  $\frac{G_n(t; \mathcal{F}^*)}{t} \leq \frac{G_n(\delta; \mathcal{F}^*)}{\delta}$ .

Given  $h \in \mathcal{F}^*$  with  $\|h\|_n \leq t$ , define the rescaled function  $\tilde{h} = \frac{\delta}{t}h$ . We have  $\tilde{h} \in \mathcal{F}^*$  by definition with  $\|\tilde{h}\|_n \leq \delta$ . It is easy to see that

$$\frac{1}{n} \left( \frac{\delta}{t} \sum_{i=1}^n w_i h(x_i) \right) = \frac{1}{n} \sum_{i=1}^n w_i \tilde{h}(x_i).$$

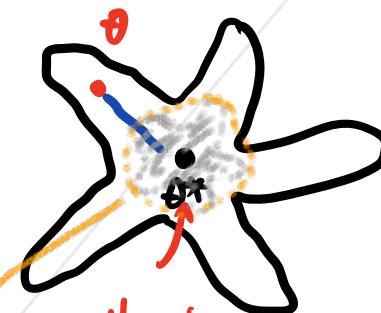
Taking the supreme and expectation on both side over  $h$ , we obtain that

$$\frac{\delta}{t} \mathbb{E} \left[ \sup_{h \in \mathcal{F}^*: \|h\|_n \leq t} \frac{1}{n} \sum_{i=1}^n w_i h(x_i) \right] \leq \mathbb{E} \left[ \sup_{\tilde{h} \in \mathcal{F}^*: \|\tilde{h}\|_n \leq \delta} \frac{1}{n} \sum_{i=1}^n w_i \tilde{h}(x_i) \right].$$

This is equivalent to desired inequality

$$\frac{G_n(t; \mathcal{F}^*)}{t} \leq \frac{G_n(\delta; \mathcal{F}^*)}{\delta}$$

$\delta_1 \leq \delta_2 \Rightarrow \frac{G_n(\delta_1)}{\delta_1} \geq \frac{G_n(\delta_2)}{\delta_2}$



the true  
solution

orange:  
localized complexity  
at  $\delta_1$

gray:  
localized complexity  
at  $\delta_2$

the line connects  
 $\theta$  and  $\theta^*$  should  
lie in the hypothesis  
space!

localized complexity

# Final Error

$G_n(\delta; \mathcal{F}^*) = \frac{\delta^2}{2\sigma} \rightarrow$  the bias.

linear growth speed at  $\delta^*$

$$\delta^* := \min_{\delta > 0} \left\{ \delta : \frac{G_n(\delta; \mathcal{F}^*)}{\delta} \leq \frac{\delta}{2\sigma} \right\}$$

$$\mathbb{E}_\sigma \sup_{\|g\|_n \leq u} \frac{\sigma}{n} \sum \sigma_i g(x_i) \leq u \delta^*$$

Then, if  $u \geq \delta^*$ , then  $\mathbb{E}_\sigma \sup_{\|g\|_n \leq u} \frac{\sigma}{n} \sum \sigma_i g(x_i) \leq u \delta^*$ .

$$G_n(u; \mathcal{F}^*) = u \cdot \frac{G_n(u; \mathcal{F}^*)}{u} \leq u \cdot \frac{G_n(\delta^*; \mathcal{F}^*)}{\delta^*} = u \cdot \delta^*$$

because  $u \geq \delta^*$

Then, proof of  $\|\hat{f} - f^*\| \leq o(\delta^*)$

- ① if  $\|\hat{f} - f^*\| \leq \delta^*$  ✓
- ② if  $\|\hat{f} - f^*\| \geq \delta^*$ . wlog  $\Delta$ :  $\|\Delta\|^2 \leq G_n(\|\Delta\|) \leq \|\Delta\| \cdot \delta^* \Rightarrow \|\Delta\| \leq \delta^*$ .

Not Required

# Method 2: Peeling

**Lemma 1 (Peeling Technique)** *If there is a function  $\phi : [0, \infty) \rightarrow [0, \infty)$  and  $r^* > 0$  s.t.  $\forall r > \hat{r}^*$ , we have*

- $\phi(4r) \leq 2\phi(r)$
- $R_n(G_r) \leq \phi(r)$

Then we have for all  $r > \hat{r}^*$  we have

$$\mathbb{E}_{\sigma_i, z_i} \left[ \frac{\frac{1}{n} \sum_{i=1}^n \sigma_i g(z_i)}{\mathbb{E}g + r} \right] \leq \frac{4\phi(r)}{r}$$

$\Rightarrow$  find  $r^*$  such that  $\phi(r^*) = r^*$   
- normalized Empirical Process

empirical process  
 $\sup_{g \in \mathcal{F}} |\mathbb{E} - \hat{\mathbb{E}}| g$





