

Lecture 20: Minimax Lower Bounds - Global Fano's Method

Lecturer: Yudong Chen

Scribe: Xumei Xi

## 1 Recap: Local Fano's Method

Recall the two key theorems from previous lectures.

**Theorem 1** (Estimation to testing). *If  $\{\theta_1, \theta_2, \dots, \theta_M\}$  is a  $2\delta$ -packing of parameter space  $\Theta$  in  $\rho(\cdot, \cdot)$ , then*

$$\inf_{\hat{\theta}} \sup_{\theta \in \Theta} \mathbb{E}_{\mathbb{P}_{\theta}} \left[ \rho(\hat{\theta}(X), \theta) \right] \geq \delta \inf_{\psi} \frac{1}{M} \sum_{j=1}^M \mathbb{P}[\psi(X) \neq j | \mathbb{P}_{\theta_j}]. \quad (1)$$

**Theorem 2** (Fano's inequality). *For any testing procedure  $\psi$ , we have*

$$\mathbb{P}[\psi(X) \neq J] \geq 1 - \frac{I(X; J) + \log 2}{\log M}. \quad (2)$$

In the local Fano's method, we upper bound the mutual information by the maximum of pairwise KL-divergence,

$$I(X; J) \leq \max_{i,j} D(\mathbb{P}_{\theta_i} \| \mathbb{P}_{\theta_j}) \leq g(\delta).$$

This upper bound, based on convexity of KL, is relatively crude and not tight under some settings, especially for non-parametric problems. In particular, with this upper bound we only makes use of a *local* packing of  $\Theta$ . In order to capture the full capacity of the entire parameter space, we will develop the so-called global Fano's method.

## 2 Global Fano's method

The global Fano's method is based on the following better upper bound on the mutual information.

### 2.1 Yang-Barron's upper bound

**Lemma 1** (Yang-Barron's upper bound). *Let  $N_{\text{KL}}(\varepsilon, \Theta)$  denote the  $\varepsilon$ -covering number of  $\Theta$  in the pseudo-distance  $\rho_{\text{KL}}(\theta, \theta') := \sqrt{D(\mathbb{P}_{\theta} \| \mathbb{P}_{\theta'})}$ . Then we have*

$$I(X; J) \leq \varepsilon^2 + \log N_{\text{KL}}(\varepsilon, \Theta), \quad \forall \varepsilon > 0. \quad (3)$$

**Proof** Recall the notation  $Q_X := \frac{1}{M} \sum_{j=1}^M \mathbb{P}_{\theta_j}$  for the marginal distribution of  $X$ . Then for any distribution of  $X$ , denoted by  $Q'$ , we have

$$\begin{aligned} I(X; J) &= \frac{1}{M} \sum_{j=1}^M D(\mathbb{P}_{\theta_j} \| Q_X) && \text{definition of KL} \\ &\leq \frac{1}{M} \sum_{j=1}^M D(\mathbb{P}_{\theta_j} \| Q'), && \text{the mixture distribution } Q_X \text{ minimizes the average KL} \\ &\leq \max_{j=1, \dots, M} D(\mathbb{P}_{\theta_j} \| Q'). && (4) \end{aligned}$$

Note that the second step above is an analogue of the fact that

$$\arg \min_{\theta'} \left\{ \frac{1}{M} \sum_{j=1}^M \|\theta_j - \theta'\|_2^2 \right\} = \frac{1}{M} \sum_{j=1}^M \theta_j.$$

Let  $\{\beta_1, \beta_2, \dots, \beta_N\}$  be a minimal  $\varepsilon$ -covering of  $\Theta$  w.r.t.  $\rho_{\text{KL}}$ , where  $N = N_{\text{KL}}(\varepsilon, \Theta)$ . Note that we are free to choose any  $Q'$ . Here we set  $Q' = \frac{1}{N} \sum_{\ell=1}^N \mathbb{P}_{\beta_\ell}$ .

Fix an arbitrary index  $j \in \{1, 2, \dots, M\}$ . By definition of  $\varepsilon$ -covering, there exists some  $\beta_i$  such that  $\rho_{\text{KL}}(\theta_j, \beta_i) \leq \varepsilon$ . We therefore have

$$\begin{aligned} D(\mathbb{P}_{\theta_j} \| Q') &= \mathbb{E}_{\mathbb{P}_{\theta_j}} \left[ \log \frac{d\mathbb{P}_{\theta_j}}{\frac{1}{N} \sum_{l=1}^N d\mathbb{P}_{\beta_l}} \right], & d\mathbb{P}_{\theta_j} \text{ denotes the density of } \mathbb{P}_{\theta_j} \\ &\leq \mathbb{E}_{\mathbb{P}_{\theta_j}} \left[ \log \frac{d\mathbb{P}_{\theta_j}}{\frac{1}{N} d\mathbb{P}_{\beta_i}} \right] & \text{sum of } N \text{ terms is lower bounded by any one term} \\ &= D(\mathbb{P}_{\theta_j} \| \mathbb{P}_{\beta_i}) + \log N \\ &\leq \varepsilon^2 + \log N. \end{aligned}$$

Combining this upper bound (valid for all  $j = 1, 2, \dots, M$ ) with equation (4) gives the desired inequality  $I(X; J) \leq \varepsilon^2 + \log N_{\text{KL}}(\varepsilon, \Theta)$ .  $\square$

**Remark** Note that we have two sets of points here:

- $\{\theta_1, \theta_2, \dots, \theta_M\} \subset \Theta$  is a  $\delta$ -packing in  $\rho$ , with cardinality  $M = M(\delta, \Theta, \rho)$ .
- $\{\beta_1, \beta_2, \dots, \beta_N\} \subset \Theta$  is an  $\varepsilon$ -packing in  $\rho_{\text{KL}}$ , with cardinality  $N = N(\varepsilon, \Theta, \rho_{\text{KL}}) = N_{\text{KL}}(\varepsilon, \Theta)$ .

## 2.2 Procedure for using Yang-Barron's upper bound

We can employ the following two steps when choosing the parameters  $\varepsilon$  and  $\delta$ .

1. Choose  $\varepsilon > 0$  such that

$$\varepsilon^2 \geq \log N_{\text{KL}}(\varepsilon, \Theta). \quad (5)$$

2. Choose the largest  $\delta > 0$  such that

$$\log M(\delta, \Theta, \rho) \geq 4\varepsilon^2 + 2 \log 2. \quad (6)$$

With the above choice, we have the following lower bound the minimax error:

$$\begin{aligned} \inf_{\hat{\theta}} \sup_{\theta \in \Theta} \mathbb{E} \left[ \rho(\hat{\theta}, \theta) \right] &\geq \delta \left( 1 - \frac{I(X; J) + \log 2}{\log M} \right), & \text{by combining (1) and (2)} \\ &\geq \delta \left( 1 - \frac{\varepsilon^2 + \log N_{\text{KL}}(\varepsilon, \Theta) + \log 2}{\log M} \right), & \text{by Lemma 1} \\ &\geq \delta \left( 1 - \frac{2\varepsilon^2 + \log 2}{4\varepsilon^2 + 2 \log 2} \right), & \text{by what we just did in (5) and (6)} \\ &= \frac{1}{2} \delta. & (7) \end{aligned}$$

### 3 Application: Lipschitz regression

For application, we consider a non-parametric regression problem over Lipschitz functions. We observe

$$y_i = f(x_i) + e_i, \quad i = 1, 2, \dots, n,$$

where  $x_i$ 's are fixed,  $e_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$  for all  $i$ , and  $f$  is an unknown function from the function class

$$\mathcal{F} := \{f : [0, 1] \rightarrow \mathbb{R} \mid f(0) = 0, f \text{ is 1-Lipchitz}\}.$$

Here  $\mathcal{F}$  acts as the parameter space  $\Theta$  in the non-parametric setting.

We have proved in Homework 1 that

$$\log N(\delta, \mathcal{F}, \|\cdot\|_\infty) \asymp \log M(\delta, \mathcal{F}, \|\cdot\|_\infty) \asymp \frac{1}{\delta}.$$

So we can find a desired  $\delta$ -packing of  $\mathcal{F}$  in  $\|\cdot\|_\infty$ .

Next we need an  $\varepsilon$ -covering of  $\mathcal{F}$  in  $\rho_{\text{KL}}$ . Observe that the distribution of the data  $(y_i)_{i=1}^n$  is

$$\begin{aligned} \mathbb{P}_f &= \mathcal{N}(f(x_1), \sigma^2) \times \dots \times \mathcal{N}(f(x_n), \sigma^2) \\ &= \mathcal{N}(f(x_1^n), \sigma^2 I_n), \end{aligned}$$

which represents an  $n$ -dimensional Gaussian distribution with mean vector  $f(x_1^n) = (f(x_1), \dots, f(x_n))$ . Then we can calculate the pairwise KL-divergence:

$$\begin{aligned} D(\mathbb{P}_f \parallel \mathbb{P}_g) &= \frac{1}{2\sigma^2} \|f(x_1^n) - g(x_1^n)\|_2^2 \\ &\leq \frac{n}{2\sigma^2} \|f - g\|_\infty^2. \end{aligned}$$

Hence the metric entropy of  $\mathcal{F}$  in KL-divergence is proportional to that in  $\|\cdot\|_\infty$ :

$$\begin{aligned} \log N_{\text{KL}}(\varepsilon, \mathcal{F}) &\asymp \log N\left(\sqrt{\frac{2\sigma^2}{n}}\varepsilon, \mathcal{F}, \|\cdot\|_\infty\right) \\ &\asymp \frac{\sqrt{n}}{\varepsilon\sigma}. \end{aligned}$$

Now we are ready to set the parameters  $\varepsilon$  and  $\delta$  according to our two-step procedure:

1. Choose  $\varepsilon = \left(\frac{\sqrt{n}}{\sigma}\right)^{\frac{1}{3}}$  so that  $\varepsilon^2 \geq \frac{\sqrt{n}}{\sigma\varepsilon} \gtrsim \log N_{\text{KL}}(\varepsilon, \mathcal{F})$ .
2. Choose  $\delta \asymp \frac{1}{\varepsilon^2}$  so that  $\log M(\delta, \Theta, \|\cdot\|_\infty) \gtrsim \frac{1}{\delta} \geq 4\varepsilon^2 + 2\log 2$ .

Thus we satisfy the requirements in (5) and (6), and the minimax lower bound (7) holds. In particular, we have

$$\inf_{\hat{f}} \sup_{f \in \mathcal{F}} \mathbb{E} \left[ \|\hat{f} - f\|_\infty \right] \geq \frac{1}{2}\delta \asymp \left(\frac{\sigma^2}{n}\right)^{\frac{1}{3}}.$$

Note that this tight lower bound cannot be achieved using local Fano's method instead.

We may compare this lower bound with the upper bound we derived in Lectures 14-15:

$$\|\hat{f} - f\|_n := \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f(x_i))^2} \lesssim \left(\frac{1}{n}\right)^{\frac{1}{3}}.$$

We can do some extra work to match the norms. In this case, the upper and lower bound match.