

Lecture 32 — 19th November, 2021

Sebastien Roch, UW-Madison

Scribe: Gokcan Tatli

1 Overview

In the last lecture, we lower bounded the minimax risk \mathfrak{M} using Proposition 15.12 from W [1]. Then, we used the relation between mutual information and Kullback Leibler divergence (15.34 in W [1]), to find the following result:

Recall: If $(\theta^1), \dots, \theta^M$ are 2δ separated under P , then

$$\mathfrak{M} \geq \phi(\delta) \left(1 - \frac{\frac{1}{\mu^2} \sum_{i \neq j} KL(P_{\theta^i} \| P_{\theta^j}) + 1}{\log_2 M} \right)$$

In this lecture, we continue to linear regression example and prove the lower bound for minimax MSE for that example. After that, we cover some basics of information theory.

2 Example: Linear Regression (continued)

- Suppose $\mathbb{Y} \sim N(\mathbb{X}\theta^*, \sigma^2 I_n)$ when (\mathbb{Y}, \mathbb{X}) are observed and θ^* is unknown.
- We want a lower bound on minimax MSE $(\mathbb{X}\theta^*)$, where

$$\text{MSE}(\mathbb{X}\theta^*) = \frac{1}{n} \|\mathbb{X}\theta^* - \mathbb{X}\hat{\theta}\|_2^2.$$

That is why we need following claims.

- **Claim 1.** If $\gamma^1, \dots, \gamma^M \in \text{range}(\mathbb{X})$ are such that

$$\|\gamma^i - \gamma^j\| > 2\delta\sqrt{n} \quad \forall i \neq j,$$

then $\exists \theta^1, \dots, \theta^M$ such that

$$\rho(\theta^i, \theta^j) \frac{1}{\sqrt{n}} \|\mathbb{X}\theta^i - \mathbb{X}\theta^j\|_2 > 2\delta$$

- **Claim 2.** If $\gamma^1, \dots, \gamma^M \in \text{range}(\mathbb{X})$ are such that

$$\|\gamma^i\| \leq 4\delta\sqrt{n} \quad \forall i,$$

then

$$KL(P_{\theta^i} \| P_{\theta^j}) \leq \frac{32n\delta^2}{\sigma^2}$$

Proof of Claim 2: We know that $P_{\theta^j} = N(\mathbb{X}\theta^j, \sigma^2 I_n)$. That is why,

$$KL(P_{\theta^i} \| P_{\theta^j}) = \frac{1}{2\sigma^2} \|\mathbb{X}\theta^i - \mathbb{X}\theta^j\|^2 \quad (\text{by explicit calculation}) \quad (1)$$

$$\leq \frac{1}{2\sigma^2} \left[\|\gamma^i\|^2 + \|\gamma^j\|^2 \right] \quad (\text{by triangle inequality}) \quad (2)$$

Then, we use the norm bound, i.e., $4\delta\sqrt{n}$, in Claim 2 to get the desired KL-divergence bound $\frac{32n\delta^2}{\sigma^2}$.

Recall: The largest cardinality of an ϵ -separated set in a subset K of a metric space (T, d) is called the packing number of K , i.e., $\mathcal{P}(K, \epsilon)$, which is a function of ϵ .

We proved that

$$\frac{\text{vol}(K)}{\text{vol}(\epsilon B_2^r)} \leq N(K, \epsilon) \leq \mathcal{P}(K, \epsilon),$$

where $N(K, \epsilon)$ is the covering number and B_2^r is the unit norm-ball.

In our case (linear regression example), $K \subseteq \mathbb{R}^r$, $\text{rank}(\mathbb{X})=r$, $\epsilon = 2\delta\sqrt{n}$, $K = "4\delta\sqrt{n}B_2^r"$ (inside the $\text{range}(\mathbb{X})$). Therefore,

$$M \geq \left(\frac{4\delta\sqrt{n}}{2\delta\sqrt{n}} \right)^r = 2^r \Rightarrow \log_2 M \geq r$$

So, we can write

$$\mathfrak{M} \geq \delta^2 \left(1 - \frac{32n\delta^2}{\sigma^2} + 1 \right) \quad (3)$$

$$\approx \frac{\delta^2}{2} \quad (4)$$

where (3) comes from combining above bound for r and the bound we got from Claim 2, and plugging them in the lower bound of minimax risk given in Section 1. Then, we can write the lower bound in (4) by taking $\delta^2 = \frac{\sigma^2}{64} \frac{r}{n}$ when r is sufficiently large.

3 Quick Tour of Information Theory

- Note that in the interest of keeping things simple, we will do everything in discrete spaces
- Mostly from Chapter 2 of Cover-Thomas [2]

Definition 1. If $X \in \mathcal{X}$ is a discrete random variable with pmf $p(x)$, the entropy of X is

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log p(x),$$

with the convention that $0 \log 0 = 0$ and $0 \log \frac{0}{0} = 0$.

Claim 2. If \mathcal{X} is finite, then

$$H(X) \leq \log |\mathcal{X}|$$

with equality if and only if X is uniform over \mathcal{X} .

Proof: Suppose U is uniform over \mathcal{X} , then

$$H(U) = - \sum_{u \in \mathcal{X}} p(u) \log p(u) \tag{5}$$

$$= - \sum_{u \in \mathcal{X}} \frac{1}{|\mathcal{X}|} \log \frac{1}{|\mathcal{X}|} \tag{6}$$

$$= \log |\mathcal{X}| \tag{7}$$

On the other hand,

$$0 \leq KL(X||U) \tag{8}$$

$$= \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{1/|\mathcal{X}|} \tag{9}$$

$$= \log |\mathcal{X}| - \left(- \sum_{x \in \mathcal{X}} p(x) \log p(x) \right) \tag{10}$$

$$= \log |\mathcal{X}| - H(X) \tag{11}$$

which completes the proof.

Definition 3. For a pair of random variables (X, Y) , the conditional entropy of $X|Y$ is

$$H(X|Y) = - \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log p(x, y),$$

where $p(x, y)$ can be also written as $p(y)p(x|y)$. This can be considered as expectation of conditional distribution $x|y$.

References

- [1] Wainwright, Martin J., *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*, Cambridge Series in Statistical and Probabilistic Mathematics, 2019.
- [2] Cover, Thomas M., *Elements of information theory*, John Wiley & Sons, 1999.