

程序代写代做 CS编程辅导

CS861: Theoretical Foundations of Machine Learning

Lecture 13 - 10/04/2023

University of Wisconsin–Madison, Fall 2023

Lecture 13: Hoeffding's lemma, Nonparametric regression

Lecturer: Kirthy

Scribed by: Elliot Pickens, Yuya Shimizu



Disclaimer: These notes are subject to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the instructor.

In this lecture, we provide further methods to derive a lower bound for the minimax risk. First, we will continue our previous discussion on constructing alternatives via **tight packings**. Then, we will introduce the **Varshamov–Gilbert lemma**, which is another method to construct well-separated alternatives. We also briefly mention other methods for lower bounds. Finally, we will discuss **nonparametric regression**.

1 Method 1: Constructing alternatives via tight packings (continued)

In the previous lecture, we have learned ε -packing numbers. We will see ε -covering numbers that behave in the equivalent order to packing numbers as $\varepsilon \downarrow 0$.

Definition 1. (Covering number, metric entropy)

- An ε -**covering** of set \mathcal{X} with respect to a metric ρ is a set $\{x_1, \dots, x_N\}$ such that for all $x \in \mathcal{X}$, there exists some $x_i \in \{x_1, \dots, x_N\}$ s.t. $\rho(x, x_i) \leq \varepsilon$.
- The ε -**covering number** $N(\varepsilon, \mathcal{X}, \rho)$ is the size of the smallest covering.
- The **metric entropy** is $\log N(\varepsilon, \mathcal{X}, \rho)$.

We have the following lemma that relates covering numbers and packing numbers.

Lemma 1. A covering number $N(\cdot, \mathcal{X}, \rho)$ and a packing number $M(\cdot, \mathcal{X}, \rho)$ satisfy

$$M(2\varepsilon, \mathcal{X}, \rho) \leq N(\varepsilon, \mathcal{X}, \rho) \leq M(\varepsilon, \mathcal{X}, \rho).$$

Remark This lemma is useful since we can apply prior work on bounding the metric entropy $\log(N(\varepsilon, \mathcal{X}, \rho))$.

2 Method 2: Varshamov–Gilbert Lemma

To get a lower bound for the minimax risk, it is often convenient to consider alternatives indexed with a hypercube:

$$\{P_\omega; \omega = (\omega_1, \dots, \omega_d) \in \{0, 1\}^d\}.$$

Example 1. (Normal mean estimation in \mathbb{R}^d) Consider a hypercube:

$$\{N(\delta\omega, \sigma^2 I); \omega \in \{0, 1\}^d\}.$$

The diagram shows a QR code with a 'Tutor CS' logo. A graph G is constructed over the QR code. Vertices $v_1, v_2, v_3, v_4, v_6, v_7$ are marked. A dashed line segment is labeled $\sqrt{\delta}$ and a bracket on the right is labeled δ .

WeChat: estuorcs

$$\min_{\omega \neq \omega'} \rho(\theta(P_\omega), \theta(P_{\omega'})) = \min_{\omega \neq \omega'} \|\delta\omega - \delta\omega'\|_2$$

$$\max_{\omega, \omega'} \text{KL}(P_\omega, P_{\omega'}) = \frac{\max_{\omega, \omega'} \|\delta\omega - \delta\omega'\|_2}{2\sigma^2}$$

Email: tutorcs@163.com (ω and ν differ on all coordinates)

Q: 749389476

This example motivates us to introduce Varshamov-Gilbert Lemma. The Varshamov-Gilbert lemma states that we can find a *large* subset of $\{0, 1\}^a$ such that the minimum distance between any two points in the subset is also *large*. Before stating the lemma, we define the Hamming distance.

Definition 2. (Hamming distance) The hamming distance between two binary vectors ω, ω' is $H(\omega, \omega') = \sum_{i=1}^d \mathbf{1}\{\omega_i \neq \omega'_i\}$ for $\omega, \omega' \in \mathbb{R}^d$. It counts the number of coordinate where ω_i and ω'_i differ.

Lemma 2. (Varshamov-Gilbert) Let $m \geq 8$. Then there exists $\Omega_m \subseteq \{0, 1\}^m$ such that the followings are true: (i) $|\Omega_m| \geq 2^{m/8}$. (ii) $\forall \omega, \omega' \in \Omega_m, H(\omega, \omega') \geq m/8$. We will call Ω_m the Varshamov-Gilbert pruned hypercube of $\{0, 1\}^m$.

We will revisit the normal mean estimation example to illustrate an application of the Varshamov-Gilbert lemma.

Example 2. (Normal mean estimation in \mathbb{R}^d) Let an i.i.d. data $S = \{x_1, \dots, x_n\} \sim P^n$ where $P \in \mathcal{P}$ and $\mathcal{P} = \{N(\mu, \Sigma); \mu \in \mathbb{R}^d, \Sigma \preceq \sigma^2 I\}$. Consider $\Phi \circ \rho(\theta_1, \theta_2) = \|\theta_1 - \theta_2\|_2^2$. Let Ω_d be the Varshamov-Gilbert pruned hypercube of $\{0, 1\}^d$. Define

$$\mathcal{P}' = \left\{ N \left(\sqrt{\frac{8}{d}} \delta \omega, \sigma^2 I \right); \omega \in \Omega_d \right\}.$$

程序代写代做 CS编程辅导

For these alternatives, we have the following bound:

$$\begin{aligned} P_{\omega \neq \omega'}(P_{\omega'}) &= \min_{\omega \neq \omega'} \sqrt{\sum_{j=1}^d \left(\sqrt{\frac{8}{d}} \delta \omega_j - \sqrt{\frac{8}{d}} \delta \omega'_j \right)^2} \\ &= \sqrt{\frac{8}{d}} \delta \min_{\omega \neq \omega'} \sqrt{H(\omega, \omega')} \\ &\geq \sqrt{\frac{8}{d}} \delta \sqrt{\frac{d}{8}} = \delta, \end{aligned}$$

where the inequality follows from the property (ii) of the Varshamov-Gilbert pruned hypercube. Since the maximum ℓ_2 -distance over a hypercube is the length of a diagonal, we also have

$$\max_{P_{\omega}, P_{\omega'} \in \mathcal{P}'} \text{KL}(P_{\omega}, P_{\omega'}) = \frac{\left(\sqrt{d} \times \left(\sqrt{\frac{8}{d}} \delta \right) \right)^2}{2\sigma^2} = \frac{4\delta^2}{\sigma^2}.$$

Choose $\delta = \sigma \sqrt{\frac{d \log(2)}{128n}}$. Then,

Assignment Project Exam Help

$$\begin{aligned} \max_{P_{\omega}, P_{\omega'} \in \mathcal{P}'} \text{KL}(P_{\omega}, P_{\omega'}) &= \frac{4\delta^2}{\sigma^2} = \frac{d \log(2)}{32n} \\ &= \frac{\log(2^{d/8})}{4n} \\ &\leq \frac{\log(|\mathcal{P}'|)}{4n}, \end{aligned}$$

where the inequality follows from the property (i) of the Varshamov-Gilbert pruned hypercube: $|\mathcal{P}'| = |\Omega_d| \geq 2^{d/8}$. Therefore, by the local Fano's method (here we also require $d \geq 32$),

QQ: 749389476

$$R_n^* \geq \frac{1}{2} \Phi\left(\frac{\delta}{2}\right) = \frac{\log(2)}{1024} \cdot \frac{d\sigma^2}{n}.$$

https://tutorcs.com

3 Other methods for lower bounds

Theorem 3. *Informal theorem (Ch 2, Tsybakov)*

Let $S = \{(x_1, y_1), \dots, (x_n, y_n)\} \sim P \in \mathcal{P}$, $\{P_0, \dots, P_N\} \subseteq \mathcal{P}$, and $\delta = \max_{j \neq k} \Phi \circ \rho(\theta(P_j), \theta(P_k))$. Then if

$$\frac{1}{N} \sum_{j=1}^N \text{KL}(P_j, P_0) \leq C_1 \log(N)$$

we can say that

$$R_n^* \geq C_2 \Phi\left(\frac{\delta}{2}\right).$$

Roughly, this informal theorem says that if the average KL distance between $P_j \forall j$ and some "central" distribution P_0 is small enough we get a the lower bound of on the minimax risk seen above.

A related method, *Assouad's method*, can be found in chapter 9 of John Duchi's "Lecture Notes on Information Theory." This method applies when there is additional structure in the problem.

程序代写代做 CS编程辅导

4 Nonparametric regression

The model: Let \mathcal{F} be the set of Lipschitz functions in $[0, 1]$. That is



$$B] ; |f(x_1) - f(x_2)| \leq L|x_1 - x_2|$$

where L is the Lipschitz constant

We observe some $S = \{(x_i, y_i)\}_{i=1}^n$ drawn i.i.d. from $P_{xy} = \mathcal{P}$ where

$$\mathcal{P} = \{P_{xy} ; 0 < \alpha_0 \leq p(x) \leq \alpha_1 < \infty$$

the regression function $f(x) \triangleq \mathbb{E}[Y|X = x] \in \mathcal{F}$

WeChat: cstutorcs

We wish to estimate the regression function via the following loss

Assignment Project Exam Help

where f is the regression function and $g : [0, 1] \rightarrow \mathbb{R}$. The risk of \hat{f} is

Email: tutors@163.com

$$R(P_{xy}, \hat{f}) = \mathbb{E}_S [\ell(P_{xy}, \hat{f})] \\ = \mathbb{E}_S \left[\int (f(x) - \hat{f}(x))^2 p(x) dx \right]$$

QQ: 749389476

and the minimax risk of \hat{f} is

https://tutorcs.com

We want to show that the minimax risk is $\Theta(n^{-\frac{2}{3}})$. We will do this in two steps:

1. Establish a lower bound with Fano's method
2. Get an upper bound using Nadaraya-Watson Estimation

4.1 Lower Bound

To begin it should be noted that we have a problem where the loss does not can be written as $\ell = \Phi \circ \rho$. To circumvent this, denote $\mathcal{P}'' = \{P_{xy} \in \mathcal{P} ; p(x) = 1\}$.¹ Then

$$R_n^* = \inf_{\hat{f}} \sup_{P_{xy} \in \mathcal{P}''} \mathbb{E}_S \left[\underbrace{\int (f(x) - \hat{f}(x))^2 p(x) dx}_{\Phi \circ \rho} \right]$$

where $\Phi \circ \rho(f_1, f_2) = \|f_1 - f_2\|_2^2$.

Now we proceed with proving the lower bound via the following three steps:

¹We choose $p(x) = 1$ here for simplicity, but any uniform pdf will work.

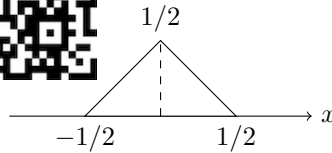
程序代写代做 CS编程辅导

1. Constructing the alternatives

2. Lower bounding μ

3. Upper bounding μ

4.1.1 Constructing



WeChat: cstutorcs

Figure 2: $U(x)$

Let

Assignment Project Exam Help

$$\Psi(x) = \begin{cases} x + \frac{1}{2} & \text{if } x \in [-1/2, 0] \\ -x + \frac{1}{2} & \text{if } x \in [0, 1/2] \end{cases}$$

Email: tutores@163.com

where Ψ is 1-Lipschitz and $\int \Psi^2(x)dx = \frac{1}{12} < \infty$.

Now let $h > 0$ (we will choose it more precisely later) and let $m = \frac{1}{h}$. Define \mathcal{F}'

QQ: 749389476

$$\mathcal{F}' = \left\{ f_w; f_w(\cdot) = \sum_{j=1}^m w_j \phi_j(\cdot), w \in \Omega_m \right\}$$

https://tutores.com

where Ω_m is the Vashamov-Gilbert pruned subset of $\{0, 1\}^d$. And ϕ_j as

$$\phi_j(x) = Lh\Psi\left(\frac{(x - (j - 1/2)h)}{h}\right).$$

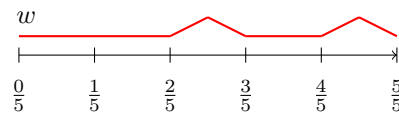
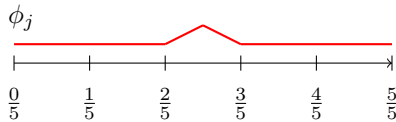


Figure 3: Depiction of ϕ_j and w when $w = \{0, 0, 1, 0, 1\}$.

Now we need to check that $\mathcal{F}' \subseteq \mathcal{F}$ (to show that f_w is L -Lipschitz). To do this, it is sufficient to check within one of the "bumps." By the chain rule,

$$|\phi'_j| = |Lh\Psi'\left(\frac{(x - (j - 1/2)h)}{h}\right)| \frac{1}{h} = L.$$

程序代写代做 CS编程辅导

Finally we define our set of alternatives \mathcal{P}' as follows.

$$\mathcal{P}' = \left\{ \mathcal{P} ; \mathcal{P} \triangleq \mathbb{E}[Y|X = x] \in \mathcal{F}', \quad Y|X = x \sim \mathcal{N}(f(x), \sigma^2) \right\}.$$

We see that $\mathcal{P}' \subseteq \mathcal{P}'' \subseteq$



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>