

## Lecture 13: November 8, 2022

Lecturer: Madhur Tulsiani

## 1 Sparse mean estimation

We will conclude our discussion of minimax rates, with this final example of estimating the mean, when we are given the additional condition that the mean is a *sparse* vector. Consider the set of normal distributions, where the mean has only *one* non-zero coordinate.

$$\Pi = \left\{ N(\mu, I_d) \mid \mu \in \mathbb{R}^d, \|\mu\|_0 \leq 1 \right\}.$$

Let  $\theta(P) = \mathbb{E}_{x \sim P}[x]$  be the mean, and let  $\ell(\hat{\theta}, \theta) = \|\hat{\theta} - \theta\|_2^2$  as before. From the previous examples, it seems like the empirical mean estimator is always the best one, and the role of information theory is primarily for proving lower bounds. However, it can also serve as a guide for the right bound to aim for. For this problem, it will be much easier to prove a lower bound. We will then show an estimator which matches this bound.

### 1.1 Lower bound

Let  $\mathcal{V} = \{e_1, \dots, e_d\}$  be the set of standard basis vectors in  $\mathbb{R}^d$ . Consider the set of distributions  $P_v = N(\sqrt{2\delta} \cdot v, I_d)$  for all  $v \in \mathcal{V}$ . Note that the means  $\mu_v = \sqrt{2\delta} \cdot v$  satisfy  $\|\mu_{v_1} - \mu_{v_2}\| = 2\delta$  for all  $v_1 \neq v_2$ . Using the bound from the previous lecture, we get

$$\begin{aligned} \mathcal{M}_n(\Pi, \ell) &\geq \delta^2 \cdot \left( 1 - \frac{n \cdot \mathbb{E}_{v_1, v_2 \in \mathcal{V}} [D(P_{v_1} \| P_{v_2})] + 1}{\log |\mathcal{V}|} \right) \\ &\geq \delta^2 \cdot \left( 1 - \frac{n \cdot (4\delta^2 / (2 \ln 2)) + 1}{\log d} \right) \\ &\geq c \cdot \frac{\log d}{n}, \end{aligned}$$

for an appropriate constant  $c > 0$ , using a choice of  $\delta^2 = c' \cdot \frac{\log d}{n}$ . We will now show that this lower bound is actually tight.

## 1.2 Upper bound

The optimal estimator for the above problem actually extends the definition of the mean as the minimizer of the total square distance (from the sample points). Recall the following.

**Exercise 1.1.** Let  $x_1, \dots, x_n \in \mathbb{R}^d$ . Then the empirical mean  $\eta = \frac{1}{n} \cdot \sum_{i=1}^n x_i$  satisfies

$$\sum_{i=1}^n \|x_i - \eta\|_2^2 = \inf_{v \in \mathbb{R}^d} \left\{ \sum_{i=1}^n \|x_i - v\|_2^2 \right\}.$$

Given a sequence of samples  $\bar{x} = (x_1, \dots, x_n)$ , let the  $\eta$  denote the empirical mean

$$\eta := \frac{1}{n} \cdot \sum_{i=1}^n x_i.$$

As we saw above, the empirical mean is the minimizer of the least square distance. However, it is not sparse. We take our estimator  $\hat{\mu}$  to only consist of the largest entry (in absolute value) of  $\eta$ , and set all other entries to zero i.e.,

$$\hat{\mu}_j := \begin{cases} \eta_j & \text{if } j = \operatorname{argmax}_{k \in [d]} |\eta_k| \\ 0 & \text{otherwise} \end{cases}.$$

Note that the above definition does not make sense if the the coordinate maximizing  $|\eta_k|$  is not unique. In such a case, we arbitrarily pick one of the maximizing coordinates. Check that this definition is a constrained version of the above definition for empirical mean. While the empirical mean  $\eta$  is the minimizer over all of  $\mathbb{R}^d$ , of the average squared distance from the sample points, the estimator above is the minimizer over all sparse vectors.

**Exercise 1.2.** Check that for  $\hat{\mu}$  defined as above

$$\sum_{i=1}^n \|x_i - \hat{\mu}\|_2^2 = \inf_{\|v\|_0 \leq 1} \left\{ \sum_{i=1}^n \|x_i - v\|_2^2 \right\}.$$

While we will use the above estimator, the operation of picking the largest coordinate does not combine well with analytic expressions such as expectation etc. For this reason, we will use the empirical mean  $\eta$  as an intermediate object in the analysis. We need the following basic properties

**Proposition 1.3.** Let  $\bar{x} \sim (N(\mu, I_d))^n$  be a sequence of  $n$  independent samples, and let  $\eta = \frac{1}{n} \cdot \sum_{i=1}^n x_i$  be the empirical mean. Then  $\eta - \mu$  is distributed according to the Gaussian distribution  $N(0, \frac{1}{n} \cdot I_d)$ .

**Proof:** Since different coordinates are independent in each of  $x_1, \dots, x_n$ , they are also independent in  $\delta - \mu$ . For any single coordinate  $j \in [d]$ , we have

$$(\eta - \mu)_j = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)_j.$$

By definition of  $(x_1, \dots, x_n)$ , each term  $(x_i - \mu)_j$  is independently distributed according to  $N(0, 1)$ . Since a linear combination of independent Gaussians is still a Gaussian, and variances add for independent variables, we get

$$\text{Var} [(\eta - \mu)_j] = \frac{1}{n^2} \sum_{i=1}^n \text{Var} [(x_i - \mu)_j] = \frac{1}{n^2} \cdot n = \frac{1}{n}.$$

Combined with  $\mathbb{E} [x_i - \mu] = 0$ , this completes the proof. ■

**Corollary 1.4.** Let  $\bar{x} = (x_1, \dots, x_n) \sim (N((\mu, I_d)))^n$  as above. Then,

$$\mathbb{P} [\exists j \in [d] \quad |\mu_j - \eta_j| \geq t] \leq 2d \cdot \exp(-nt^2/2).$$

**Proof:** Using the standard Gaussian tail bound, we know that for  $y \sim N(0, \sigma^2)$ , we have

$$\mathbb{P} [|y| \geq t] \leq 2 \cdot \exp(-t^2/(2\sigma^2)).$$

Using [Proposition 1.3](#) for each coordinate  $\eta_j - \mu_j$ , and taking a union bound over all  $j \in [d]$  gives the desired bound. ■

Recall the our goal is to bound the expected loss  $\mathbb{E}_{\bar{x} \sim (N(\mu, I_d))^n} [\|\mu - \hat{\mu}(\bar{x})\|_2^2]$ . Using the above, we can first prove a tail bound: the probability that the loss is too large, is small.

**Claim 1.5.** For the estimator  $\hat{\mu}$  as above

$$\mathbb{P} [\|\mu - \hat{\mu}\|_2 \geq t] \leq 2d \cdot \exp(-nt^2/18).$$

**Proof:** We will prove that

$$\|\mu - \hat{\mu}\|_2 \geq t \quad \Rightarrow \quad \exists j \in [d] \quad |\eta_j - \mu_j| \geq t/3.$$

Using this, together with [Corollary 1.4](#) will prove the claim. Recall that both  $\mu$  and  $\hat{\mu}$  have at most one non-zero coordinate. If  $\mu = 0$  and  $\hat{\mu}_j \neq 0$ , then we must have  $|\hat{\mu}_j| = |\eta_j - \mu_j| \geq t$ . The case when  $\hat{\mu} = 0$  can be handled similarly.

If  $\mu \neq 0$ , then let unique the non-zero coordinate be 1 (without loss of generality) i.e.,  $|\mu_1| > 0$ . If  $\hat{\mu}_1 \neq 0$ , then we again have

$$|\mu_1 - \eta_1| = |\mu_1 - \hat{\mu}_1| = \|\mu - \hat{\mu}\|_2 \geq t,$$

and we are done. So let's assume  $\hat{\mu}_1 = 0$  and  $\hat{\mu}_j \neq 0$  for some  $j > 1$ . Since we must have  $\hat{\mu}_j = \eta_j$  in this case, we have

$$|\mu_1| + |\eta_j| = |(\mu - \hat{\mu})_1| + |(\mu - \hat{\mu})_j| \geq \|\mu - \hat{\mu}\|_2 \geq t.$$

Also, since  $\eta_j$  must be the largest coordinate in absolute value, we have

$$|\eta_j| \geq |\eta_1| \geq |\mu_1| - |\mu_1 - \eta_1|.$$

Adding the above inequalities gives

$$|\mu_1 - \eta_1| + 2 \cdot |\eta_j| = |\mu_1 - \eta_1| + 2 \cdot |\mu_j - \eta_j| \geq t.$$

Hence, either  $|\mu_1 - \eta_1| \geq t/3$  or  $|\mu_j - \eta_j| \geq t/3$ , which is what we wanted to prove. ■

We can now finish the computation of the expected loss, using the above tail bound. Using  $s = t^2$  in the above bound, we can write it as

$$\mathbb{P} \left[ \|\mu - \hat{\mu}\|_2^2 \geq s \right] \leq 2d \cdot \exp(-ns/18).$$

This yields the following bound.

**Claim 1.6.** *For the estimator  $\hat{\mu}$  as above*

$$\mathbb{E}_{\bar{\mathbf{x}} \sim (N(\mu, I_d))^n} \left[ \|\mu - \hat{\mu}(\bar{\mathbf{x}})\|_2^2 \right] = O\left(\frac{\log d}{n}\right).$$

**Proof:** We use the fact that for a non-negative random variable  $Z$ ,  $\mathbb{E}[Z] = \int_s^\infty \mathbb{P}[Z \geq s] ds$ . Using this, we get

$$\begin{aligned} \mathbb{E}_{\bar{\mathbf{x}} \sim (N(\mu, I_d))^n} \left[ \|\mu - \hat{\mu}(\bar{\mathbf{x}})\|_2^2 \right] &= \int_0^\infty \mathbb{P} \left[ \|\mu - \hat{\mu}\|_2^2 \geq s \right] ds \\ &= \int_0^u \mathbb{P} \left[ \|\mu - \hat{\mu}\|_2^2 \geq s \right] ds + \int_u^\infty \mathbb{P} \left[ \|\mu - \hat{\mu}\|_2^2 \geq s \right] ds \\ &\leq \int_0^u 1 ds + \int_u^\infty 2d \cdot \exp(-ns/18) ds \\ &= u + \frac{36d}{n} \cdot \exp(-nu/18). \end{aligned}$$

Choosing  $u = c \cdot \frac{\log d}{n}$  for an appropriate constant  $c$ , then finishes the proof. ■

## 2 I-Projections and applications

We will now talk more about finding a distribution in a set  $\Pi$  that minimizes  $D(P\|Q)$  for a fixed distribution  $Q$ . We encountered this when discussing Sanov's theorem and hypothesis testing, and will now discuss its properties in some detail. When  $Q$  is the uniform distribution on  $\mathcal{X}$ . Then we also have,

$$D(P\|Q) = \log |\mathcal{X}| - H(P)$$

Hence, in this case  $P^*$  is a distribution that maximizes entropy. In general, when the given information does not uniquely determine a distribution, we choose  $P^*$  that maximizes entropy. This can be thought of as picking  $P^*$  in the set of distributions  $\Pi$ , subject to the least amount of additional assumptions. This is sometimes called the *Maximum Entropy Principle*. In this lecture, we will characterize the distributions obtained by minimizing KL-divergence (or maximizing entropy).

For closed convex set  $\Pi$ , such a  $P$  is called the I-projection of  $Q$  onto  $\Pi$ .

**Definition 2.1.** Let  $\Pi$  be a closed convex set of distributions over  $\mathcal{X}$ . In addition, assume that  $\text{Supp}(Q) = \mathcal{X}$ . Then

$$\text{Proj}_{\Pi}(Q) := \arg \min_{P \in \Pi} D(P\|Q) = P^*$$

Note that the assumption  $\text{Supp}(Q) = \mathcal{X}$  above is without loss of generality since  $D(P\|Q) = \infty$  for any  $P$  such that  $\text{Supp}(P) \not\subseteq \text{Supp}(Q)$ . Use the (strict) convexity of KL-divergence to check the following.

**Exercise 2.2.** For a closed, convex set  $\Pi$ , the projection  $P^* = \text{Proj}_{\Pi}(Q)$  exists and is unique.

It is immediate from definition that if  $P \in \Pi$ , then  $D(P\|Q) \geq D(P^*\|Q)$ . In fact,  $P^*$  tells us more. It also tells us how "far"  $P$  is away from  $Q$  in KL-divergence measure.

**Theorem 2.3.** Let  $P^* = \text{Proj}_{\Pi}(Q)$ . Then, for all  $P \in \Pi$ ,

$$\begin{aligned} \text{Supp}(P) &\subseteq \text{Supp}(P^*) \\ D(P\|Q) &\geq D(P\|P^*) + D(P^*\|Q) \end{aligned}$$

**Proof:** Define  $P_t = tP + (1-t)P^*$ , where  $t \in [0, 1]$ . By minimality of  $P^*$ , it is clear that  $D(P_t\|Q) - D(P^*\|Q) \geq 0$ . By the mean value theorem, we also have that

$$0 \leq \frac{1}{t} \cdot (D(P_t\|Q) - D(P^*\|Q)) \leq \left. \frac{d}{dt} D(P_t\|Q) \right|_{t=t' \in [0, t]}$$

Since  $t' \rightarrow 0$  as  $t \rightarrow 0$ , we get

$$\lim_{t \downarrow 0} \frac{d}{dt} D(P_t\|Q) \geq 0.$$

We now compute  $\frac{d}{dt}D(P_t||Q)$ .

$$\frac{d}{dt}D(P_t||Q) = \sum_{x \in \mathcal{X}} \frac{d}{dt} p_t(x) \log \frac{p_t(x)}{q(x)} + \sum_{x \in \mathcal{X}} p_t(x) \frac{d}{dt} (\log p_t(x) - \log q(x))$$

Note that

$$\begin{aligned} \frac{d}{dt} p_t(x) &= p(x) - p^*(x) \\ \frac{d}{dt} \log p_t(x) &= \frac{1}{\ln 2} \frac{1}{p_t(x)} (p(x) - p^*(x)) \end{aligned}$$

Using these facts, we have

$$\begin{aligned} \frac{d}{dt}D(P_t||Q) &= \sum_{x \in \mathcal{X}} (p(x) - p^*(x)) \log \frac{p_t(x)}{q(x)} + \sum_{x \in \mathcal{X}} \frac{1}{\ln 2} (p(x) - p^*(x)) \\ &= \sum_{x \in \mathcal{X}} (p(x) - p^*(x)) \log \frac{p_t(x)}{q(x)} \end{aligned}$$

Here, note that if  $(\exists x)$  such that  $p(x) > 0$  and  $p^*(x) = 0$ , then  $\lim_{t \downarrow 0} \frac{d}{dt}D(P_t||Q) \rightarrow -\infty$ , which contradicts the fact that  $\frac{d}{dt}D(P_t||Q) \geq 0$ . Hence, if  $p(x) > 0$ , then  $p^*(x) > 0$  and therefore,  $\text{Supp}(P) \subseteq \text{Supp}(P^*)$ . This proves the first part of the theorem. Now we evaluate  $\frac{d}{dt}D(P_t||Q)$  at  $t = 0$ .

$$\begin{aligned} \frac{d}{dt}D(P_t||Q)|_{t=0} &= \sum_{x \in \mathcal{X}} p(x) \log \frac{p^*(x)}{q(x)} - p^*(x) \log \frac{p^*(x)}{q(x)} \\ &= \sum_{x \in \mathcal{X}} p(x) \log \frac{p^*(x)}{q(x)} \frac{p(x)}{p(x)} - D(P^*||Q) \\ &= \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} - \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{p^*(x)} - D(P^*||Q) \\ &= D(P||Q) - D(P||P^*) - D(P^*||Q) \geq 0 \end{aligned}$$

Hence,  $D(P||Q) \geq D(P||P^*) + D(P^*||Q)$ . ■

Consider the following example, which shows that the inequality can in fact be strict.

**Exercise 2.4.** Let  $\mathcal{X} = \{0, 1\}$  and  $\Pi = \{P : p(1) \leq 1/2\}$ . Let  $Q$  be defined as

$$Q = \begin{cases} 1 & \text{with prob. } 3/4 \\ 0 & \text{with prob. } 1/4 \end{cases}$$

1. Show that

$$P^* = \begin{cases} 1 & \text{with prob. } 1/2 \\ 0 & \text{with prob. } 1/2 \end{cases}$$

2. Show that  $D(P||Q) > D(P||P^*) + D(P^*||Q)$  for the above example.