

Group Study and Seminar Series (Summer 20)

Minimax Lower Bounds

Presenter: Hao Liang

The Chinese University of Hong Kong, Shenzhen, China

July 2, 2020

Mainly based on:

Wainwright, M. J. (2019). High-dimensional statistics: A non-asymptotic viewpoint (Vol. 48). Cambridge University Press.
John, Duchi (2019). Lecture Notes for Statistics 311/Electrical Engineering 377.

Outline

Introduction

Why study lower bounds?

Preliminaries

Minimax lower bounds

Le Cam's method

Fano's method

Outline

Introduction

Why study lower bounds?

Preliminaries

Minimax lower bounds

Le Cam's method

Fano's method

Example 0: Gaussian location family

- ▶ $\{\mathcal{N}(\theta, \sigma^2), \theta \in \mathbb{R}\}$: a normal distribution family with fixed variance σ^2
- ▶ Data: a collection $Z = (Y_1, \dots, Y_n)$ with Y_i i.i.d. $\sim \mathcal{N}(\theta, \sigma^2)$
- ▶ Method: estimate unknown θ^* via an estimator $\hat{\theta}(Z)$
- ▶ Performance measure: **risk** $R(\hat{\theta}, \theta^*)$
- ▶ How does $\tilde{\theta}_n := \frac{1}{n} \sum_{i=1}^n Y_i$ perform?
- ▶ **Upper bound** provides **worst-case** performance guarantee

$$\sup_{\theta \in \mathbb{R}} R(\tilde{\theta}_n, \theta) \leq \frac{\sigma^2}{n}$$

Example 0: Gaussian location family

- ▶ But how to answer the following questions?
 - Can this analysis be improved? Or does $\tilde{\theta}_n$ actually satisfy better bounds?
 - Can any estimator improve upon the bound?
- ▶ Both questions ask about some form of **optimality**(switch orders?)
 - Optimality of an estimator
 - Optimality of a bound
- ▶ A **positive** answer consists in
 - Finding a **better proof** for $\tilde{\theta}_n$
 - Finding a **better estimator**, together with a proof that it performs better

Example 0: Gaussian location family

- ▶ Lower bound may provide negative answer to both questions

$$\inf_{\hat{\theta}} \sup_{\theta \in \mathbb{R}} \mathbb{E}_{\theta} \left[(\hat{\theta} - \theta)^2 \right] \geq \Theta\left(\frac{\sigma^2}{n}\right)$$

- ▶ Any estimator suffers risk at least $\Theta\left(\frac{\sigma^2}{n}\right)$ in the worst case
- ▶ Recall that $\tilde{\theta}_n$ suffers risk at most $\Theta\left(\frac{\sigma^2}{n}\right)$ in the worst case
- ▶ Both the upper bound $\Theta\left(\frac{\sigma^2}{n}\right)$ and the estimator $\tilde{\theta}_n$ are not improvable!

Outline

Introduction

Why study lower bounds?

Preliminaries

Minimax lower bounds

Le Cam's method

Fano's method

Statistical decision theory

- ▶ $\mathcal{P} = \{\mathbb{P}_\theta | \theta \in \Omega\}$: A parametric family with parameter θ
- ▶ Data/samples: Y_i i.i.d. $\sim \mathbb{P}_\theta$ or $Y^n = (Y_1, \dots, Y_n) \sim \mathbb{P}_\theta^n$
- ▶ Decision rule
 - (Point) **Estimation**: estimate θ^* via an estimator $\hat{\theta}(Y^n)$, $\hat{\theta} : \mathcal{X}^n \rightarrow \Omega$
 - (Hypothesis) **Test**: nature **randomly** choose index $J = j$, decide $j \in \{1, 2, \dots, M\}$ via an test function $\psi(Y^n)$, where $Y^n \sim \mathbb{P}_{\theta_j}^n$
- ▶ Loss function $\rho(\hat{\theta}, \theta^*)$
 - Absolute loss $\rho(\hat{\theta}, \theta^*) = |\hat{\theta} - \theta^*|$
 - Squared loss $\rho(\hat{\theta}, \theta^*) = (\hat{\theta} - \theta^*)^2$
- ▶ Risk $R(\hat{\theta}, \theta^*) = \mathbb{E}_{\mathbb{P}} \left[\rho(\hat{\theta}(Y^n), \theta^*) \right]$

Information theory

- ▶ Entropy $H(X) := \int_{\mathcal{X}} p_X(u) \log \frac{1}{p_X(u)} du$
- ▶ Relative entropy/KL divergence $D(\mathbb{P}_X \parallel \mathbb{P}_Y) := \int_{\mathcal{X}} p_X(u) \log \frac{p_X(u)}{p_Y(u)} du$
 - $D(\mathcal{N}(\theta_1, \sigma_1^2), \mathcal{N}(\theta_2, \sigma_2^2)) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$
- ▶ Mutual information $I(X; Y) := D(\mathbb{P}_{X,Y} \parallel \mathbb{P}_X \mathbb{P}_Y)$
 - KL divergence form: $I(X; Y) = \mathbb{E}_X D(\mathbb{P}_{Y|X} \parallel \mathbb{P}_Y) = \sum_x \mathbb{P}(x) D(\mathbb{P}_{Y|X=x} \parallel \mathbb{P}_Y)$
- ▶ **Fano's inequality** provides a lower bound on the error in a M -ary testing problem

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

Outline

Introduction

Why study lower bounds?

Preliminaries

Minimax lower bounds

Le Cam's method

Fano's method

Basic framework

- ▶ Given a class of distributions \mathcal{P} and $\theta : \mathcal{P} \rightarrow \Omega$ is a **functional** mapping distributions to a parameter $\theta(\mathbb{P})$
- ▶ For **parametric** classes, $\theta(\mathbb{P})$ uniquely determines the underlying distribution \mathbb{P} , write $\mathcal{P} = \{\mathbb{P}_\theta | \theta \in \Omega\}$ (e.g. Gaussian location family)
- ▶ The viewpoint of estimating functionals here is more general than a parametric family (e.g. estimating the mode of the density $\theta(\mathbb{P}) = \arg \max_{x \in [0,1]} f(x)$)

Minimax risk

- ▶ Given a sample $X \sim \mathbb{P}_{\theta^*}$, θ^* fixed but unknown
- ▶ The goal of an estimator $\hat{\theta}$ is to estimate θ^* based on X , write also $\hat{\theta} \equiv \hat{\theta}(X)$
- ▶ Let $\rho : \Omega \times \Omega \rightarrow [0, \infty)$ be a **semi-metric**, consider r.v. $\rho(\hat{\theta}, \theta^*)$
- ▶ Taking expectations over X yields the **deterministic** quantity $R(\hat{\theta}, \theta^*) := \mathbb{E}_{\mathbb{P}} \left[\rho(\hat{\theta}, \theta^*) \right]$
- ▶ Typically referred to as the **risk function** associated with $\hat{\theta}$

Minimax risk

- ▶ Goal: $\min_{\hat{\theta}} R(\hat{\theta}, \theta^*), \forall \theta^*$?
- ▶ **Multi-objective** optimization problem
- ▶ Two ways to deal with this issue: **Bayesian** approach and **minimax** approach
 - Bayesian approach: taking average over parameters

$$\inf_{\hat{\theta}} \mathbb{E}_{\theta^* \sim \pi} [R(\hat{\theta}, \theta^*)]$$

- Minimax approach: adversarial perspective

$$\inf_{\hat{\theta}} \sup_{\theta^*} R(\hat{\theta}, \theta^*)$$

Minimax risk

- ▶ More generally

$$\inf_{\hat{\theta}} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}}[\rho(\hat{\theta}, \theta(\mathbb{P}))]$$

- ▶ The ρ -minimax risk

$$\mathfrak{M}(\theta(\mathcal{P}); \rho) := \inf_{\hat{\theta}} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}}[\rho(\hat{\theta}, \theta(\mathbb{P}))] \quad (2)$$

- ▶ Introduce a non-decreasing function $\Phi : [0, \infty) \rightarrow [0, \infty)$,

$$\mathfrak{M}(\theta(\mathcal{P}); \Phi \circ \rho) := \inf_{\hat{\theta}} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}}[\Phi(\rho(\hat{\theta}, \theta(\mathbb{P})))] \quad (3)$$

From estimation to testing

- ▶ Developing methods for lower bounding the minimax risk
- ▶ **Reduction** to the problem of obtaining lower bounds for certain testing problems
- ▶ Start with constructing such testing problems as follows:
- ▶ Suppose that $\{\theta^1, \dots, \theta^M\} \subseteq \theta(\mathcal{P})$ is a **2δ -separated set**, i.e., $\rho(\theta^j, \theta^k) \geq 2\delta$ for all $j \neq k$
- ▶ For each θ^j , choose some representative distribution \mathbb{P}_{θ^j} for which $\theta(\mathbb{P}_{\theta^j}) = \theta^j$

From estimation to testing

- ▶ Generate a random variable Z by the following procedure:
 - Sample a random integer J from the **uniform** distribution over the index set $[M] := \{1, \dots, M\}$
 - Given $J = j$, sample $Z \sim \mathbb{P}_{\theta^j}$
- ▶ let \mathbb{Q} denote the **joint** distribution of the pair (Z, J) , then the **marginal** distribution over Z is $\bar{\mathbb{Q}} := \frac{1}{M} \sum_{j=1}^M \mathbb{P}_{\theta^j}$
- ▶ Consider the M -ary hypothesis testing problem of determining J based on a sample Z
- ▶ A **testing function** for this problem is a mapping $\psi : \mathcal{Z} \rightarrow [M]$

From estimation to testing

- ▶ The probability of error of ψ is $\mathbb{Q}[\psi(Z) \neq J]$, can be used to obtain lower bound

Proposition 1.

(From estimation to testing) For any non-decreasing function Φ and choice of 2δ -separated set, the minimax risk is lower bounded as

$$\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho) \geq \Phi(\delta) \inf_{\psi} \mathbb{Q}[\psi(Z) \neq J], \quad (4)$$

where the infimum ranges over all test functions.

Remarks

- ▶ The r.h.s. of the bound involves two terms, and both of them depends δ
 - The function Φ is decreasing in δ
 - As δ increases, M decreases
 - The underlying testing problem becomes easier, $\mathbb{Q}[\psi(Z) \neq J]$ decreases
- ▶ Choose a sufficiently small δ^* to ensure that this testing error is at least 0.5,

$$\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho) \geq \frac{1}{2} \Phi(\delta^*) \quad (5)$$

Proof

- ▶ For any $\mathbb{P} \in \mathcal{P}$ with parameter $\theta = \theta(\mathbb{P})$,

$$\mathbb{E}_{\mathbb{P}}[\Phi(\rho(\hat{\theta}, \theta))] \stackrel{(i)}{\geq} \Phi(\delta)\mathbb{P}[\Phi(\rho(\hat{\theta}, \theta)) \geq \Phi(\delta)] \stackrel{(ii)}{\geq} \Phi(\delta)\mathbb{P}[\rho(\hat{\theta}, \theta) \geq \delta]. \quad (6)$$

- ▶ It suffices to lower bound $\sup_{\mathbb{P} \in \mathcal{P}} \mathbb{P}[\rho(\hat{\theta}, \theta(\mathbb{P})) \geq \delta]$
- ▶ Recall that \mathbb{Q} denotes the joint distribution over the pair (Z, J) ,

$$\sup_{\mathbb{P} \in \mathcal{P}} \mathbb{P}[\rho(\hat{\theta}, \theta(\mathbb{P})) \geq \delta] \geq \frac{1}{M} \sum_{j=1}^M \mathbb{P}_{\theta^j} [\rho(\hat{\theta}, \theta^j) \geq \delta] = \mathbb{Q}[\rho(\hat{\theta}, \theta^J) \geq \delta]. \quad (7)$$

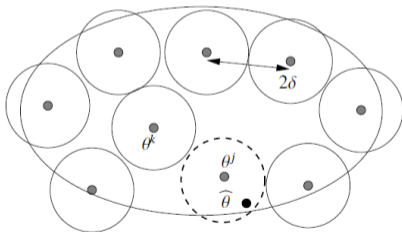
- ▶ Reduced to lower bounding $\mathbb{Q}[\rho(\hat{\theta}, \theta^J) \geq \delta]$

Proof

- ▶ Construct a test based on the estimator $\hat{\theta}$ via

$$\psi(Z) := \arg \min_{\ell \in [M]} \rho(\theta^\ell, \hat{\theta})$$

- ▶ Suppose that the true parameter is θ^j , then the event $\{\rho(\theta^j, \hat{\theta}) < \delta\}$ ensures that the test is correct



Proof

- ▶ Conditioned on $J = j$, $\{\rho(\theta^j, \hat{\theta}) < \delta\} \subseteq \{\psi(Z) = j\}$, implying

$$\mathbb{P}_{\theta^j}[\rho(\hat{\theta}, \theta^j) \geq \delta] \geq \mathbb{P}_{\theta^j}[\psi(Z) \neq j] \quad (8)$$

- ▶ Taking averages over index j ,

$$\mathbb{Q}[\rho(\hat{\theta}, \theta^J) \geq \delta] = \frac{1}{M} \sum_{j=1}^M \mathbb{P}_{\theta^j}[\rho(\hat{\theta}, \theta^j) \geq \delta] \geq \mathbb{Q}[\psi(Z) \neq J] \quad (9)$$

- ▶ Combined with the earlier argument, $\sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}}[\Phi(\rho(\hat{\theta}, \theta))] \geq \Phi(\delta) \mathbb{Q}[\psi(Z) \neq J]$
- ▶ Take the infimum over all estimators $\hat{\theta}$ on the l.h.s., and the the infimum over the **induced set of tests** on the r.h.s.
- ▶ Finally notice that the full infimum over all tests can only be smaller, from which the claim follows

Some divergence measures

▶ Three important measures

- **Total variation (TV) distance** $\|\mathbb{P} - \mathbb{Q}\|_{\text{TV}} := \sup_{A \subseteq \mathcal{X}} |\mathbb{P}(A) - \mathbb{Q}(A)| = \frac{1}{2} \int_{\mathcal{X}} |p(x) - q(x)| dx$
- **KL divergence** $D(\mathbb{P} \parallel \mathbb{Q}) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} dx$
- **Squared Hellinger distance** $H^2(\mathbb{P} \parallel \mathbb{Q}) := \int_{\mathcal{X}} (\sqrt{p(x)} - \sqrt{q(x)})^2 dx = 2 - 2 \int_{\mathcal{X}} \sqrt{p(x)q(x)} dx$

▶ The second and third distance can be used to upper bound TV distance

- **Pinsker's inequality** $\|\mathbb{P} - \mathbb{Q}\|_{\text{TV}} \leq \sqrt{\frac{1}{2} D(\mathbb{Q} \parallel \mathbb{P})}$
- **Le Cam's inequality** $\|\mathbb{P} - \mathbb{Q}\|_{\text{TV}} \leq H(\mathbb{P} \parallel \mathbb{Q}) \sqrt{1 - \frac{H^2(\mathbb{P} \parallel \mathbb{Q})}{4}}$

▶ These inequalities are useful when dealing with **product distributions**

Some divergence measures

- ▶ Let $\mathbb{P}^{1:n} = \bigotimes_{i=1}^n \mathbb{P}_i$ be the product distribution of $(\mathbb{P}_1, \dots, \mathbb{P}_n)$ defined on \mathcal{X}^n
- ▶ What's the expression of $\text{Div}(\mathbb{Q}^{1:n} \parallel \mathbb{P}^{1:n})$ in terms of $\text{Div}(\mathbb{Q}_i \parallel \mathbb{P}_i)$?
- ▶ The TV distance behaves badly: difficult to decouple
- ▶ The KL divergence exhibits a very attractive **decoupling** property,

$$D(\mathbb{P}^{1:n} \parallel \mathbb{Q}^{1:n}) = \sum_{i=1}^n D(\mathbb{P}_i \parallel \mathbb{Q}_i) \quad (10)$$

- ▶ The squared Hellinger distance does not decouple in a simple way, but

$$\frac{1}{2} H^2(\mathbb{P}^{1:n} \parallel \mathbb{Q}^{1:n}) = 1 - \prod_{i=1}^n \left(1 - \frac{1}{2} H^2(\mathbb{P}_i \parallel \mathbb{Q}_i) \right) \quad (11)$$

Some divergence measures

- ▶ In the **i.i.d.** case where $\mathbb{P}_i = \mathbb{P}_1$ and $\mathbb{Q}_i = \mathbb{Q}_1$ for all i

$$D(\mathbb{P}^{1:n} \parallel \mathbb{Q}^{1:n}) = nD(\mathbb{P}_1 \parallel \mathbb{Q}_1) \quad (12)$$

$$\frac{1}{2}H^2(\mathbb{P}^{1:n} \parallel \mathbb{Q}^{1:n}) = 1 - \left(1 - \frac{1}{2}H^2(\mathbb{P}_1 \parallel \mathbb{Q}_1)\right)^n \leq \frac{1}{2}nH^2(\mathbb{P}_1 \parallel \mathbb{Q}_1) \quad (13)$$

- ▶ Combined with the previous inequalities,
 - Pinsker's inequality in the i.i.d. case

$$\|\mathbb{P}^{1:n} - \mathbb{Q}^{1:n}\|_{\text{TV}} \leq \sqrt{\frac{n}{2}D(\mathbb{P}_1 \parallel \mathbb{Q}_1)} \quad (14)$$

- Le Cam's inequality in the i.i.d. case

$$\|\mathbb{P}^{1:n} - \mathbb{Q}^{1:n}\|_{\text{TV}} \leq H(\mathbb{P}^{1:n} \parallel \mathbb{Q}^{1:n}) \sqrt{1 - \frac{H^2(\mathbb{P}^{1:n} \parallel \mathbb{Q}^{1:n})}{4}} \leq H(\mathbb{P}^{1:n} \parallel \mathbb{Q}^{1:n}) \leq \sqrt{n}H(\mathbb{P}_1 \parallel \mathbb{Q}_1)$$

Outline

Introduction

Why study lower bounds?

Preliminaries

Minimax lower bounds

Le Cam's method

Fano's method

Binary test and Le Cam's method

- ▶ Recall the reduction from estimation to testing

$$\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho) \geq \Phi(\delta) \inf_{\psi} \mathbb{Q}[\psi(Z) \neq J]$$

- ▶ Le Cam's method is to consider the simplest type of testing problem - **binary hypothesis test**, which involves only two distributions
- ▶ In a binary testing problem with **equally weighted** hypotheses, $Z \sim \bar{Q} := \frac{1}{2}\mathbb{P}_0 + \frac{1}{2}\mathbb{P}_1$
- ▶ For a given decision rule $\psi : \mathcal{Z} \rightarrow \{0, 1\}$, the associated probability of error is

$$\mathbb{Q}[\psi(Z) \neq J] = \frac{1}{2}\mathbb{P}_0[\psi(Z) \neq 0] + \frac{1}{2}\mathbb{P}_1[\psi(Z) \neq 1] \quad (16)$$

Bayes error and TV distance

- ▶ Take the infimum over all decision rules yields **Bayes error**
- ▶ Recall the definition of TV distance $\|\mathbb{P} - \mathbb{Q}\|_{\text{TV}} := \sup_{A \subseteq X} |\mathbb{P}(A) - \mathbb{Q}(A)|$
- ▶ There is a **one-to-one correspondence** between ψ and partitions (A, A^c) of the space \mathcal{X}
 $A = \{x \in X \mid \psi(x) = 1\}$
- ▶ The Bayes risk can be expressed in terms of $\|\mathbb{P}_1 - \mathbb{P}_0\|_{\text{TV}}$

$$\begin{aligned} \inf_{\psi} \mathbb{Q}[\psi(Z) \neq J] &= \frac{1}{2} \inf_{\psi} (\mathbb{P}_0[\psi(Z) \neq 0] + \mathbb{P}_1[\psi(Z) \neq 1]) \\ &= \frac{1}{2} \inf_{A \subseteq \mathcal{X}} (\mathbb{P}_0[A] + \mathbb{P}_1[A^c]) \\ &= \frac{1}{2} \{1 - \|\mathbb{P}_1 - \mathbb{P}_0\|_{\text{TV}}\} \end{aligned}$$

Le Cam's method

Proposition 2.

(Le Cam's bound) For any pair of distributions $\mathbb{P}_0, \mathbb{P}_1 \in \mathcal{P}$ s.t. $\rho(\theta(\mathbb{P}_0), \theta(\mathbb{P}_1)) \geq 2\delta$,

$$\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho) \geq \frac{\Phi(\delta)}{2} \{1 - \|\mathbb{P}_1 - \mathbb{P}_0\|_{\text{TV}}\} \quad (17)$$

► Two extremes

- Worst case: $\mathbb{P}_1 = \mathbb{P}_0$, hypotheses are completely indistinguishable
- Best case: $\|\mathbb{P}_1 - \mathbb{P}_0\|_{\text{TV}} = 1$, \mathbb{P}_0 and \mathbb{P}_1 have disjoint supports

Example 1: Gaussian location family

- ▶ $\{\mathcal{N}(\theta, \sigma^2), \theta \in \mathbb{R}\}$: a normal distribution family with **fixed** variance σ^2
- ▶ Goal: Estimate θ
 - Metric: either $|\hat{\theta} - \theta|$ or $(\hat{\theta} - \theta)^2$
 - Data: a collection $Z = (Y_1, \dots, Y_n) \sim \mathcal{N}(\theta, \sigma^2)^{1:n} = P_\theta^n$
- ▶ Apply the two-point Le Cam bound with the distributions \mathbb{P}_0^n and \mathbb{P}_θ^n
 - Set $\theta = 2\delta$ s.t. the two means are 2δ -separated
 - Bound $\|\mathbb{P}_\theta^n - \mathbb{P}_0^n\|_{\text{TV}}$

$$\|\mathbb{P}_\theta^n - \mathbb{P}_0^n\|_{\text{TV}}^2 \leq \frac{n}{2} D(\mathbb{P}_\theta \| \mathbb{P}_0) = \frac{n}{2} \frac{\theta^2}{2\sigma^2} \leq \frac{1}{4} \left\{ e^{n\theta^2/\sigma^2} - 1 \right\} = \frac{1}{4} \left\{ e^{4n\delta^2/\sigma^2} - 1 \right\}$$

Example 1: Gaussian location family

- ▶ Apply the two-point Le Cam bound with the distributions \mathbb{P}_0^n and \mathbb{P}_θ^n
 - Setting $\delta = \frac{1}{2} \frac{\sigma}{\sqrt{n}}$ thus yields

$$\inf_{\hat{\theta}} \sup_{\theta \in \mathbb{R}} \mathbb{E}_\theta [|\hat{\theta} - \theta|] \geq \frac{\delta}{2} \left\{ 1 - \frac{1}{2} \sqrt{e-1} \right\} \geq \frac{\delta}{6} = \frac{1}{12} \frac{\sigma}{\sqrt{n}}$$

$$\inf_{\hat{\theta}} \sup_{\theta \in \mathbb{R}} \mathbb{E}_\theta [(\hat{\theta} - \theta)^2] \geq \frac{\delta^2}{2} \left\{ 1 - \frac{1}{2} \sqrt{e-1} \right\} \geq \frac{\delta^2}{6} = \frac{1}{24} \frac{\sigma^2}{n}$$

- ▶ Although the pre-factors $1/12$ and $1/24$ are not optimal, the scalings σ/\sqrt{n} and σ^2/n are sharp/**order optimal**
- ▶ Matching upper bound: the sample mean $\tilde{\theta}_n := \frac{1}{n} \sum_{i=1}^n Y_i$ satisfies the bounds

$$\sup_{\theta \in \mathbb{R}} \mathbb{E}_\theta [|\tilde{\theta}_n - \theta|] = \sqrt{\frac{2}{\pi}} \frac{\sigma}{\sqrt{n}} \quad \text{and} \quad \sup_{\theta \in \mathbb{R}} \mathbb{E}_\theta [(\tilde{\theta}_n - \theta)^2] = \frac{\sigma^2}{n}$$

Example 2: Uniform location family

- ▶ Mean-squared error decaying as n^{-1} is typical for parametric problems, but faster rates is possible for some other problems
- ▶ $\{\mathbb{U}_\theta, \theta \in \mathbb{R}\}$: \mathbb{U}_θ is uniform over the interval $[\theta, \theta + 1]$
- ▶ Impossible to use Pinsker's inequality to control the TV norm!
- ▶ Consider $H^2(\mathbb{U}_\theta \parallel \mathbb{U}_{\theta'})$, recall that

$$H^2(\mathbb{P} \parallel \mathbb{Q}) = 2 - 2 \int_{\mathcal{X}} \sqrt{p(x)q(x)} dx$$

- It suffices to consider the case $\theta' > \theta$
- If $\theta' > \theta + 1$, then $H^2(\mathbb{U}_\theta \parallel \mathbb{U}_{\theta'}) = 2$
- Otherwise, $H^2(\mathbb{U}_\theta \parallel \mathbb{U}_{\theta'}) = 2 - 2 \int_{\theta'}^{\theta+1} dt = 2|\theta' - \theta|$

Example 2: Uniform location family

- ▶ Apply the Le Cam bound with the distributions \mathbb{U}_θ^n and $\mathbb{U}_{\theta'}^n$
 - Take a pair θ, θ' s.t. $|\theta' - \theta| = 2\delta := \frac{1}{4n}$
 - $\|\mathbb{U}_\theta^n - \mathbb{U}_{\theta'}^n\|_{\text{TV}}^2 \leq H^2(\mathbb{U}_\theta^n \|\mathbb{U}_{\theta'}^n) \leq nH^2(\mathbb{U}_\theta \|\mathbb{U}_{\theta'}) = n2|\theta' - \theta| = \frac{1}{2}$
 - $\inf_{\hat{\theta}} \sup_{\theta \in \mathbb{R}} \mathbb{E}_\theta \left[(\hat{\theta} - \theta)^2 \right] \geq \frac{\delta^2}{2} \left\{ 1 - \sqrt{\frac{1}{2}} \right\} = \frac{(1 - \frac{1}{\sqrt{2}})}{128} \frac{1}{n^2}$
- ▶ Contrasted with the n^{-1} rate, this lower bound has faster n^{-2} rate!
- ▶ Matching upper bound: the estimator $\tilde{\theta} = \min \{Y_1, \dots, Y_n\}$ satisfies the bound $\sup_{\theta \in \mathbb{R}} \mathbb{E} \left[(\tilde{\theta} - \theta)^2 \right] \leq \frac{2}{n^2}$
- ▶ Estimating the location parameter of uniform location family is easier.

Outline

Introduction

Why study lower bounds?

Preliminaries

Minimax lower bounds

Le Cam's method

Fano's method

Fano's method

- ▶ Le Cam's method reduces the estimation problem to binary test, how about M -ary hypothesis testing problem?
- ▶ In information theory, **Fano's inequality** lower bounds the error probability in such problems

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

- ▶ Combined with the reduction in Proposition 1

Proposition 3.

(Fano's bound) Let $\{\theta^1, \dots, \theta^M\}$ be a 2δ -separated set in the ρ semi-metric on $\Theta(\mathcal{P})$,

$$\mathfrak{M}(\theta(\mathcal{P}); \Phi \circ \rho) \geq \Phi(\delta) \left\{ 1 - \frac{I(Z; J) + \log 2}{\log M} \right\} \quad (18)$$

Remarks on Fano's method

- ▶ Consider the behavior of the different terms of r.h.s. as $\delta \rightarrow 0^+$
 - The 2δ -separation criterion becomes milder, $M \equiv M(2\delta)$ increases
 - $J \in [M(2\delta)]$ can take on a larger number of potential values, $I(Z; J)$ decreases
 - Decreasing δ sufficiently may ensure that

$$\frac{I(Z; J) + \log 2}{\log M} \leq \frac{1}{2} \quad (19)$$

- $\mathfrak{N}(\theta(\mathcal{P}); \Phi \circ \rho) \geq \frac{1}{2}\Phi(\delta)$
- ▶ Two technical and possibly challenging steps
 - Specify 2δ -separated sets with large cardinality $M(2\delta)$, metric entropy theory
 - Compute or upper bound $I(Z; J)$, non-trivial
- ▶ Using convexity of KL divergence and the mixture representation

$$I(Z; J) = \frac{1}{M} \sum_{j=1}^M D(\mathbb{P}_{\theta^j} \|\overline{\mathbb{Q}}) \leq \frac{1}{M^2} \sum_{j,k=1}^M D(\mathbb{P}_{\theta^j} \|\mathbb{P}_{\theta^k}) \quad (20)$$

Example 3: Gaussian location model via Fano's method

- ▶ Consider the 2δ -separated set $\{\theta^1, \theta^2, \theta^3\} = \{0, 2\delta, -2\delta\}$
- ▶ $D(\mathbb{P}_{\theta^j}^{1:n} \parallel \mathbb{P}_{\theta^k}^{1:n}) = \frac{n}{2\sigma^2} (\theta^j - \theta^k)^2 \leq \frac{8n\delta^2}{\sigma^2}$ for all $j, k = 1, 2, 3$
- ▶ $I(Z; J) \leq \frac{8n\delta^2}{\sigma^2}$
- ▶ Choosing $\delta^2 = \frac{\sigma^2}{80n}$ ensures that $\frac{8n\delta^2/\sigma^2 + \log 2}{\log 3} = \frac{0.1 + \log 2}{\log 3} < 0.75$
- ▶ The Fano's bound with $\Phi(t) = t^2$ implies

$$\sup_{\theta \in \mathbb{R}} \mathbb{E}_{\theta} \left[(\hat{\theta} - \theta)^2 \right] \geq \frac{\delta^2}{4} = \frac{1}{320} \frac{\sigma^2}{n}$$

Bounds based on local packings

- ▶ Construct a 2δ -separated set contained within Ω s.t. for some $c > 0$

$$\sqrt{D(\mathbb{P}_{\theta^j} \|\mathbb{P}_{\theta^k})} \leq c\sqrt{n}\delta \quad \text{for all } j \neq k \quad (21)$$

- ▶ The bound (20) then implies that $I(Z; J) \leq c^2 n \delta^2$, and hence the bound (19) will hold if

$$\log M(2\delta) \geq 2 \{c^2 n \delta^2 + \log 2\} \quad (22)$$

- ▶ The minimax risk is lower bounded as $\mathfrak{M}(\theta(\mathcal{P}), \Phi \circ \rho) \geq \frac{1}{2}\Phi(\delta)$

Example 4: Minimax risks for linear regression

- ▶ Standard linear regression model: $y = \mathbf{X}\theta^* + w$
- ▶ $\mathbf{X} \in \mathbb{R}^{n \times d}$: fixed design matrix, $w \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$: observation noise
- ▶ Metric: prediction norm $\rho_{\mathbf{X}}(\hat{\theta}, \theta^*) := \frac{\|\mathbf{X}(\hat{\theta} - \theta^*)\|_2}{\sqrt{n}}$, θ^* : vary over \mathbb{R}^d
- ▶ Consider the set $\{\gamma \in \text{range}(\mathbf{X}) \mid \|\gamma\|_2 \leq 4\sqrt{n}\delta\}$
- ▶ let $\{\gamma^1, \dots, \gamma^M\}$ be a $2\sqrt{n}\delta$ -packing in the ℓ_2 -norm, $r = \text{rank}(\mathbf{X})$
- ▶ Lemma 5.7 in HDS book implies such a packing with $\log M \geq r \log 2$ elements

Example 4: Minimax risks for linear regression

- ▶ A collection of vectors of the form $\gamma^j = \mathbf{X}\theta^j$ for some $\theta^j \in \mathbb{R}^d$ s.t.

$$\frac{\|\mathbf{X}\theta^j\|_2}{\sqrt{n}} \leq 4\delta, \quad \text{for each } j \in [M] \quad (23)$$

$$2\delta \leq \frac{\|\mathbf{X}(\theta^j - \theta^k)\|_2}{\sqrt{n}} \leq 8\delta, \quad \text{for each } j \neq k \in [M] \times [M] \quad (24)$$

- ▶ Let \mathbb{P}_{θ^j} denote the distribution of y when $\theta^* = \theta^j$, then $\mathbb{P}_{\theta^j} = \mathcal{N}(\mathbf{X}\theta^j, \sigma^2\mathbf{I}_n)$
- ▶ $D(\mathbb{P}_{\theta^j} \|\mathbb{P}_{\theta^k}) = \frac{1}{2\sigma^2} \|\mathbf{X}(\theta^j - \theta^k)\|_2^2 \leq \frac{32n\delta^2}{\sigma^2}$
- ▶ Condition (21) holds with $c = \frac{\sqrt{32}}{\sigma}$

Example 4: Minimax risks for linear regression

- ▶ Need to lower bound $\log M \geq r \log 2 \geq 2(c^2 n \delta^2 + \log 2)$
- ▶ Choose $\delta^2 = \frac{\sigma^2}{64} \frac{r}{n}$, then condition (22) holds for sufficiently large r since
 - $D(\mathbb{P}_{\theta^j} \parallel \mathbb{P}_{\theta^k}) \leq \frac{r}{2}$
 - $2(c^2 n \delta^2 + \log 2) = 2(\frac{r}{2} + \log 2) = r + \log 2$
- ▶ Set $\Phi(t) = t^2$

$$\inf_{\hat{\theta}} \sup_{\theta \in \mathbb{R}^d} \mathbb{E} \left[\frac{1}{n} \|\mathbf{X}(\hat{\theta} - \theta)\|_2^2 \right] \geq \frac{1}{128} \frac{\text{rank}(\mathbf{X}) \sigma^2}{n}$$

- ▶ This lower bound is sharp up to constant pre-factors