

Lecture 1: January 18

Lecturer: Alessandro Rinaldo

Scribes: Jaehyeok Shin

Note: *LaTeX template courtesy of UC Berkeley EECS dept.*

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

1.1 Minimax lower bounds

For statistical problems such as estimation, testing, confidence set, model selection, etc, we start with a given procedure, and establish

- Consistency
- Rates of convergence (as a function of n and other parameters)

The last one upper bound the minimax risk which we will define later. The *minimax theory* is about a quantifying how hard a problem is by producing rates that lower bound the convergence rates of any procedure.

Example Normal means $Y = \theta^* + \epsilon \in \mathbb{R}^d$, $\epsilon \sim N(0, \sigma^2 I_d)$

Observation : $Y_1, \dots, Y_n \stackrel{i.i.d.}{\sim} N(\theta^*, \sigma^2 I_d)$

Let $\hat{\theta}(Y_1, \dots, Y_n) \in \mathbb{R}^d$. We are interested in the expected squared loss $\mathbb{E} [\|\hat{\theta} - \theta^*\|^2]$

For the penalized regression : $\hat{\theta} \in \arg \min_{\theta \in \Theta} \frac{1}{2n} \sum_{i=1}^n \|Y_i - \theta\|^2$, we get

$$\mathbb{E} [\|\hat{\theta} - \theta^*\|^2] \lesssim \begin{cases} \frac{\sigma^2 d}{n} & \Theta = \mathbb{R}^d \\ \frac{\sigma^2 \log d}{n} & \Theta = B_1 := \{\theta \in \mathbb{R}^d, \|\theta\|_1 \leq 1\} \\ \frac{\sigma^2 k}{n} \log \frac{ed}{n} & \Theta = B_0(k) := \{\theta \in \mathbb{R}^d, \|\theta\|_0 \leq k\} \end{cases}$$

Example Covariance estimation $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} P$ (e.g., $P = N(0, \Sigma_{d \times d})$)

$$\Sigma_{d \times d} \in C(\alpha, M_0, M_1) = \{\Sigma_{d \times d} = (\sigma_{ij})_{i,j=1,\dots,d} \text{ Symmetric \& P.S.D} \\ \text{such that } \max_j \sum_{|i-j|>k} |\sigma_{ij}| \leq M_0 k^{-\alpha}, \forall k \lambda_{\max}(\Sigma) \leq M_1\}$$

Bickel, Levina (2008) estimated Σ with $\hat{\Sigma}_{BL} = \hat{\Sigma}_{BL}(X_1, \dots, X_n)$, and got

$$\sup_{\Sigma \in C(\alpha, M_0, M_1)} \mathbb{E} \|\Sigma - \hat{\Sigma}_{BL}\|_{op} \lesssim \left(\frac{\log d}{n} \right)^{\frac{\alpha}{\alpha+1}}$$

Cai, Zhang, and Zhou (2010) obtained a different estimator $\hat{\Sigma}_{CZZ}$ with

$$\sup_{\Sigma \in C(\alpha, M_0, M_1)} \mathbb{E} \|\Sigma - \hat{\Sigma}_{CZZ}\|_{op} \lesssim \min \left\{ \left(\frac{1}{n} \right)^{\frac{2\alpha}{2\alpha+1}} + \frac{\log d}{n}, \frac{p}{n} \right\}$$

Better rates, in fact, they are optimal.

$$\inf_{\hat{\Sigma}} \sup_{\Sigma \in C(\alpha, M_0, M_1)} \mathbb{E} \|\Sigma - \hat{\Sigma}\|_{op} \text{ is the same order as } \hat{\Sigma}_{CZZ}$$

1.1.1 Set-up

Let \mathcal{P} collection of probability distributions on $(\mathcal{X}, \mathcal{A})$

Let $\theta : \mathcal{P} \rightarrow \Theta$ functional, and $\theta(P)$ parameter.

Example

- $\theta(P) = \mathbb{E}_P[X]$, $X \sim P$
- If P has a density f , $\theta(P) = f(x_0)$, or $\theta(P) = \int \left(f'(x) \right)^2 dx$

In particular, we may have that $\theta(P) = \theta(Q)$ for $P, Q \in \mathcal{P}$

Simple case

θ parametrize \mathcal{P} , in which case $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$, and $\theta(P) \neq \theta(Q)$ iff $P \neq Q$.

If $\Theta \subset \mathbb{R}^d$, \mathcal{P} is a parametric family.

Examples

1. $\mathcal{P} = \{N(\theta, I) : \theta \in \mathbb{R}^d\}$
2. $\Theta = \{\text{Set of smooth functions on } [0, 1]^d\}$,
 \mathcal{P} consists of probabilities P for $(Y, X) \in \mathbb{R} \times \mathbb{R}^d$ such that

$$Y = f(X) + \epsilon, \quad f \in \Theta, \quad \epsilon \sim (0, \sigma^2) \perp\!\!\!\perp X$$

\implies Non-parametric function estimation problem.

3. \mathcal{P} of $(Y, X) \in \mathbb{R} \times \mathbb{R}^d$, $Y = X^T \theta + \epsilon$, $\theta \in \Theta \subset \mathbb{R}^d$

Let $\tilde{X} = (X_1, \dots, X_n) \stackrel{i.i.d.}{\sim} P \in \mathcal{P}$. We will estimate $\theta(P)$ using \tilde{X} .

You may allow \mathcal{P} to change with n , i.e., for each n , we will have \mathcal{P}_n, Θ_n depending on n

Let $d : \Theta \times \Theta \rightarrow [0, \infty)$ be a metric

1. $d \geq 0$
2. $d(x, y) = d(y, x)$
3. $d(x, z) \leq d(x, y) + d(y, z)$

4. $d(x, y) = 0 \Leftrightarrow x = y$

More generally, we will consider $w\left(d(\theta, \theta')\right)$ where

$$w : [0, \infty) \rightarrow [0, \infty) \text{ non-decreasing, } w(0) = 0, \quad w \not\equiv 0$$

Example $d(\theta, \theta') = \|\theta - \theta'\|$

- $w(x) = x^2$: square error
- $w(x) = \mathbf{1}\{x > c\}$, $c > 0$

Assuming $\tilde{X} = (X_1, \dots, X_n) \stackrel{i.i.d.}{\sim} P \in \mathcal{P}$. We are concerned with $\theta(P) \in \Theta$.

For a given procedure $\hat{\theta} : \mathcal{X}^n \rightarrow \Theta$, its (point-wise) risk at P is

$$\mathbb{E}_{\tilde{X} \sim P^n} \left[w\left(d(\hat{\theta}(\tilde{X}), \theta(P))\right) \right]$$

Example

- $\mathcal{P} = \{N(\theta, I), \theta \in \mathbb{R}^d\}$, $d(x, y) = \|x - y\|_2$, $w(x) = x^2$
Risk : $\mathbb{E}\|\hat{\theta} - \theta\|_2^2$
- $\mathcal{P} = \{P_0, P_1\}$, $\theta(\tilde{X}) = 1 \text{ or } 0$
Risk : $l_i P_i \left(\hat{X} \neq i \right)$, $i = 0, 1$, $l_i > 0$

To measure how well $\hat{\theta}$ does, let's take sup over all $P \in \mathcal{P}_n$

$$r_n(\hat{\theta}, \mathcal{P}_n) = \sup_{P \in \mathcal{P}_n} \mathbb{E}_P \left[w\left(d(\hat{\theta}, \theta(P))\right) \right]$$

Upper bound calculations entail finding a constant $C = C(\mathcal{P}_n)$ and a sequence $\psi_n \rightarrow 0$ as $n \rightarrow \infty$ such that

$$r_n(\hat{\theta}, \mathcal{P}_n) \leq C\psi_n$$

The minimax risk is

$$R_n(\mathcal{P}_n) = \inf_{\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)} \sup_{P \in \mathcal{P}_n} \mathbb{E}_P \left[w\left(d(\hat{\theta}, \theta(P))\right) \right]$$

If we can show that

$$R_n(\mathcal{P}_n) \geq C' \psi'_n, \quad C' = C'(\mathcal{P}_n), \quad \psi'_n \rightarrow 0$$

Then, $C' \psi'_n$ is a lower bound on minimax risk, $\forall n$

If $\hat{\theta}$ is such that $\frac{\psi_n}{\psi'_n} = \Theta(1)$, then $\hat{\theta}$ is minimax rate optimal.