

Instructor : Wenlong Mou.

Course website : mouwenlong.github.io/teaching/sta3000f25

Querans: homework submission / lecture recording.

Grading: 3 homeworks 20% each.

Take-home final 40%.

Poll: 1. ✓ 75 min lecture 10 min break
75 min lecture.

(end 10 min earlier).

2. 80 min 10 min (10 min 10, 50).

-
- Basic prob tools and review.
 - Decision theory. (optimality, classical theory).
 - Modern tools: empirical process theory (basics).
(application to asymptotics).

Basics of nonparametric estimation.
(including minimax lower bounds).

Probability review.

Radon-Nikodym derivative.

Idea: generalization of pdf.

μ : measure of interest

λ : base measure. (e.g. Lebesgue mrr).

If $\mu < \lambda$ (i.e., $\forall A$ measurable, $\lambda(A) = 0 \Rightarrow \mu(A) = 0$)

then \exists a density function.

$$p(x) = \frac{d\mu}{d\lambda}(x).$$

For any measurable set A .

$$\mu(A) = \int_A \frac{d\mu}{d\lambda}(x) \cdot d\lambda(x).$$

Basic inequalities.

• Markov Inequality.

$$X \geq 0 \quad E[X] < +\infty$$

then $P(X \geq a) \leq \frac{E(X)}{a}$.

Proof: $P(X \geq a) = E[1_{X \geq a}] \leq E\left[\frac{X}{a} 1_{X \geq a}\right] \leq \frac{E(X)}{a}$.

Decaying only at $\frac{1}{a}$ rate.

Higher moment leads to better tail

Chebychev inequality.

Assuming $E[X^2] < +\infty$.

$$P(|X - E(X)| \geq a) \leq \frac{\text{var}(X)}{a^2}.$$

Proof: Let $Y = (X - E(X))^2$.
Apply Markov inequality to Y .

Extension to higher moments.

If $E[X^p] < +\infty$.

$$P(|X - E(X)| \geq a) \leq \frac{E[(X - E(X))^p]}{a^p}.$$

Larger p leads to faster-decaying tail behavior.

What if all moments exist for X ?

Moment generating functions.

Def. $m_X(\lambda) := \mathbb{E}[e^{\lambda X}]$ (for $\lambda \in \mathbb{R}$).

(Assuming existence). $= \sum_{n=0}^{+\infty} \frac{\lambda^n}{n!} \mathbb{E}[X^n].$

e.g. $X \sim N(\mu, \sigma^2)$ $m_X(\lambda) = \exp\left(\lambda\mu + \frac{\lambda^2\sigma^2}{2}\right)$.

e.g. $X \sim \text{Ber}(p)$. $m_X(\lambda) = 1-p+pe^\lambda$.

Fact (direct consequence of Markov ineq).

$\forall \lambda > 0$.

$P(X \geq \mathbb{E}[X] + t) = P(e^{\lambda X} \geq e^{\lambda \mathbb{E}[X] + \lambda t})$

(By Markov). $\leq \frac{m_X(\lambda)}{\exp(\lambda \mathbb{E}[X] + \lambda t)}$

(λ is up to our choice. optimize λ).

e.g. for $X \sim N(\mu, \sigma^2)$.

$$P(X \geq \mu + t) \leq \frac{\exp(\mu\lambda + \frac{\lambda^2\sigma^2}{2})}{\exp(\mu\lambda + \lambda t)} = \exp\left(-\lambda t + \frac{\lambda^2\sigma^2}{2}\right).$$

(Take $\lambda = \frac{t}{\sigma^2}$)

$$P(X \geq \mu + t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

(Exponentially decaying tail).

Def. X "sub Gaussian" random variable if.

$\forall \lambda$, $m_X(\lambda)$ exists; and

$$E[\exp(\lambda(X - E[X]))] \leq \exp\left(\frac{\lambda^2\sigma^2}{2}\right).$$

We call σ "the sub Gaussian parameter".

e.g. Rademacher r.v. is sub-Gaussian.

$$X = \begin{cases} 1 & \text{w.p. } \frac{1}{2} \\ -1 & \text{w.p. } \frac{1}{2}. \end{cases}$$

Proof. $m_X(\lambda) = \frac{1}{2}(e^\lambda + e^{-\lambda})$.

$$= 1 + \frac{\lambda^2}{2!} + \frac{\lambda^4}{4!} + \frac{\lambda^6}{6!} + \dots$$

On the other hand,

$$e^{\frac{1}{2}\lambda^2} = 1 + \frac{\lambda^2}{2} + \frac{1}{2!} \cdot \left(\frac{\lambda^2}{2}\right)^2 + \frac{1}{3!} \cdot \left(\frac{\lambda^2}{2}\right)^3 + \dots$$

Each term:

$$\frac{1}{(2k)!} \quad \text{v.s.} \quad \frac{1}{k!} \cdot \frac{1}{2^k},$$

$$\text{So } m_X(\lambda) \leq \exp\left(\frac{\lambda^2}{2}\right) \quad (\forall \lambda).$$

X is $\{-$ -subGaussian.

Moving to i.i.d. sums.

Suppose $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} P$.

P is σ -subGaussian.

How about $\frac{1}{n} \sum_{i=1}^n X_i$. (compared to $E(X)$)?

For $x \in \mathbb{R}$

$$\begin{aligned} & \mathbb{E}\left[\exp\left(\frac{\lambda}{n} \sum_{i=1}^n (X_i - \mathbb{E}(X_i))\right)\right] \\ &= \left\{\mathbb{E}\left[\exp\left(\frac{\lambda}{n}(X - \mathbb{E}(X))\right)\right]\right\}^n \\ &\stackrel{\text{(sub-Gaussian)}}{\leq} \left(\exp\left(\frac{\lambda^2 \sigma^2}{2n^2}\right)\right)^n = \exp\left(\frac{\lambda^2 \sigma^2}{2n}\right). \end{aligned}$$

So $\frac{1}{n} \sum_i X_i$ is $\frac{\sigma}{\sqrt{n}}$ -subGaussian.

Corollary. $P\left(\left|\frac{1}{n} \sum_i X_i - \mathbb{E}(X)\right| \geq t\right) \leq 2 \exp\left(-\frac{nt^2}{2\sigma^2}\right)$.

($t > 0$)
e.g. $X_i \stackrel{iid}{\sim}$ Rademacher.

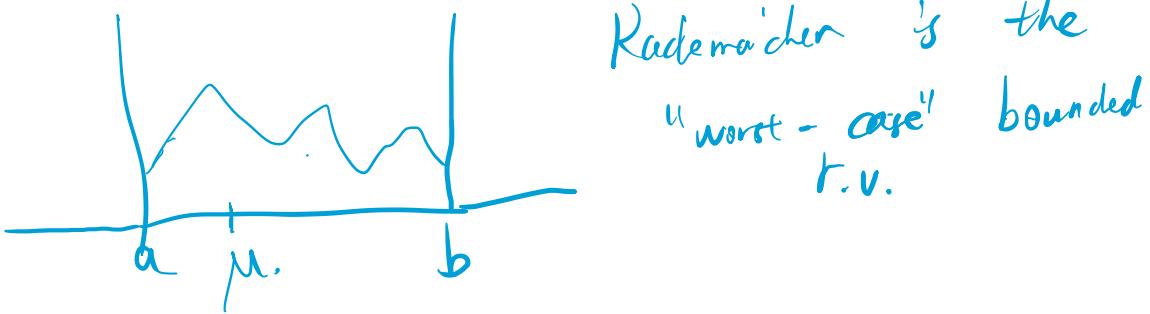
$$P\left(\left|\frac{1}{n} \sum_i X_i\right|\right) \leq 2 \exp\left(-\frac{nt^2}{2}\right).$$

Matching the tail for λ -Gaussian r.v.

This can be extended to general bounded r.v.

X supported on $[a, b]$. ($a < b$).

then X is sub-Gaussian w. parameter $\frac{b-a}{2}$.



Rademacher is the
"worst-case" bounded
r.v.

Corollary (Hoeffding bound).

$X_i \in [\bar{a}, \bar{b}]$ a.s.

X_1, X_2, \dots, X_n i.i.d.

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - E(X)\right| \geq t\right) \leq \exp\left(-\frac{2nt^2}{(\bar{b}-\bar{a})^2}\right).$$

(Chernoff bound).

e.g. $X \sim \text{Ber}(p)$.

$$P(|\bar{X}_n - p| \geq t) \leq \exp(-2nt^2).$$

Remarks: independence is needed,
"identically distributed" can be relaxed.

"union bound".

$$P\left(\bigcup_{i \in I} A_i\right) \leq \sum_{i \in I} P(A_i).$$

Usually, we use A_i to model "bad thing at certain direction".

Tail behavior matters — # of terms.

e.g. Probabilistic method.

Consider Hamming cube $\{-1, 1\}^n$

$$d_{\text{Ham}}(x, y) = \sum_{i=1}^n \mathbb{1}_{x_i \neq y_i}.$$

Question: Can we pack a lot of points in the Hamming cube, s.t. each one is far from each other.

i.e. Does there exist $\{x_1, x_2, \dots, x_M\} \subseteq \{-1, 1\}^n$.

s.t. $d_{\text{Ham}}(x_i, x_j) \geq \frac{1}{4}n$.

Fact: We can make M exponentially large

Proof. Let $X_1, X_2, \dots, X_n \sim \text{Unif}(\{-1, 1\}^n)$.

$$\Pr\left(\underbrace{d_{\text{Ham}}(x_i, x_j)}_{\text{sum of } n \text{ iid Bernoulli r.v.'s.}} \leq \frac{n}{4}\right) \leq \exp\left(-2n \cdot \left(\frac{1}{4}\right)^2\right)$$

$$\leq \exp\left(-\frac{n}{8}\right).$$

$$\begin{aligned} & \mathbb{P}\left(\exists i, j \in [M], d_{Ham}(x_i, x_j) \leq \frac{n}{4}\right) \\ & \leq M^2 \cdot \mathbb{P}\left(d_{Ham}(x_i, x_j) \leq \frac{n}{4}\right) \\ & \leq M^2 \cdot \exp\left(-\frac{n}{8}\right). \end{aligned}$$

Let $M = \left\lfloor \frac{1}{2} \exp\left(\frac{n}{16}\right) \right\rfloor$.

$$\mathbb{P}(\text{bad event}) \leq \frac{1}{4}.$$

Classical decision theory.

Key question: compare stats methods and prove optimality.

A statistical model

$$\mathcal{P} = \{P_\theta : \theta \in \Theta\}$$

Observe $X \sim P_{\theta^*}$ for some $\theta^* \in \Theta$.

Goal: to make decisions (about θ^*)
based on data.

— A decision rule: mapping from data
to action $a \in A$.

— loss function:

$$L(\theta, a) \in \mathbb{R}.$$

— Risk $R(\theta; \delta) = \mathbb{E}_{\theta} [L(\theta; \delta(X))]$.

(\mathbb{E}_{θ} , P_{θ} denote expectation, prob under θ).

e.g. estimation. $g: \mathbb{H} \rightarrow \mathbb{R}^k$ $A = \mathbb{R}^k$.
 $\hat{g} = g: \text{data} \rightarrow \mathbb{R}^k$.

$$L(\theta, a) = \|g(\theta) - a\|_2^2.$$

$$R(\theta, \hat{g}) = \mathbb{E}_{\theta} [\|g(\theta) - \hat{g}(X)\|_2^2].$$

e.g. Testing. $A = \{0, 1\}$.

$$\mathbb{H}_0 \subseteq \mathbb{H}.$$

$$L(\theta; a) = \begin{cases} 1 & (\theta \in \mathbb{H}_0, a=1) \text{ or } (\theta \notin \mathbb{H}_0, a=0) \\ 0 & \text{otherwise.} \end{cases}$$

$$R(\theta; \delta) = \begin{cases} P(\text{type-I error}) & \theta \in \mathbb{H}_0 \\ P(\text{type-II error}) & \theta \notin \mathbb{H}_0. \end{cases}$$

e.g. Learning a binary classifier.

A : the set of functions from
input domain $\rightarrow \{0, 1\}$.

$$X = \{f(z_i, Y_i)\}_{i=1}^n$$

$$L(\theta, f) = P_\theta(f(z) \neq Y).$$

$$(f \in A) \quad R(\theta; \delta) = \mathbb{E}_\theta [L(\theta; f(x))].$$

(usually evaluated on a test set)

Comparing decision rules

Idea: minimize $R(\theta; \delta)$.

Multi-objective: there could be f_1, f_2

each performing better on a subset of medals.

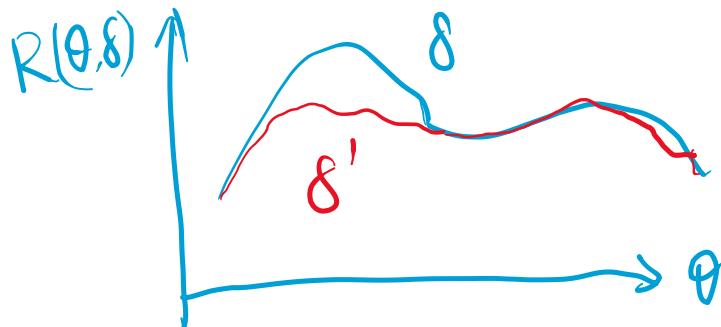
— Admissibility.

If $\exists f'$ s.t.

$$R(\theta; f') \leq R(\theta; \delta) \quad (\forall \theta \in \Theta)$$

with strict inequality for some θ .

then we call δ inadmissible.
 A decision rule δ is ^Vadmissible if not called inadmissible.



Problem: stupid decision rules can be admissible.

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$ $\theta \in \mathbb{R}$.

Goal: estimate θ .

$$\delta(X) = \theta.$$

$R(\theta, \delta) = 0$ if $\theta = 0$.
 (But the risk is very large when $\theta \neq 0$).

Some "good" estimators are inadmissible.

"Stein's phenomenon":

$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\theta, I_d)$ $\theta \in \mathbb{R}^d$.

Goal: estimate θ under MSE less.

$$\delta(x_1 \dots x_n) = \frac{1}{n} \sum_1^n x_i$$

Fact: δ is not admissible when ($d \geq 3$).