# STA3000F Advanced theory of stats.

## Wenlong Mou

Grading: 3 hw × 20% + 40% take-home final.

10:10 am – 1 pm.

---

$\int$ Decision theory (classical theory).

Asymptotic theory (n large, or $n \to +\infty$)

(a non-asymptotic approach).

nonparametric estimation.

---

Probability Recap.

$(\mathcal{X}, \mathcal{F}, \mu)$ $\qquad$ $\mu: \mathcal{F} \longrightarrow \mathbb{R}$.

$\uparrow$ $\quad$ $\uparrow$ $\nwarrow$ prob msr.

state $\sigma$-field

---

Def (Radon–Nikodym derivative).

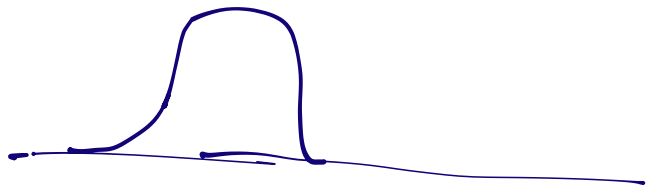$\mu \ll \lambda$. ($\mu$ abs. cts w.r.t. $\lambda$)

$\forall A \in \mathcal{F}$ $\qquad$ $\lambda(A) = 0 \implies \mu(A) = 0$.

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

$$\forall A \in \mathcal{F}. \quad \mu(A) = \int_A p(x)\, d\lambda(x).$$

$$\lambda : \text{Leb}$$

$$\mu : \text{density.}$$

---

Markov ineq.

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

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

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

Chebyshev ineq

Assume $\mathbb{E}[X^2] < +\infty$.

$$Y = (X - \mathbb{E}[X])^2 \geq 0$$

$$\mathbb{P}\left(|X - \mathbb{E}[X]| \geq a\right) = \mathbb{P}(Y \geq a^2) \leq \frac{\text{var}(X)}{a^2}$$

Extension: $\mathbb{E}\left[|X|^p\right] < +\infty$  $(p \geq 2)$

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

$$\mathbb{E}\left[e^{\lambda X}\right] =: m_X(\lambda).$$

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

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

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

by $\lambda > 0$.

$$\mathbb{P}\left(X \geq \mathbb{E}[X] + a\right) \leq \exp\left(-(a + \mathbb{E}X) \cdot \lambda\right) \cdot m_X(\lambda).$$

e.g. $X \sim \mathcal{N}(\mu, \sigma^2)$

$a > 0$  $\underline{\mathbb{P}(X \geq \mu + a)} \leq \exp\left(-(a + \mu) \cdot \lambda\right) \cdot \exp\left(\lambda\mu + \frac{\lambda^2\sigma^2}{2}\right)$

$$= \exp\left(-a\lambda + \frac{\lambda^2\sigma^2}{2}\right)$$

$$\left(\lambda = \frac{a}{\sigma^2}\right)$$

$$= \exp\left(-\frac{a^2}{2\sigma^2}\right).$$

$$\mathbb{P}\left(|\bar{X}_n - \mu| \geq a\right) \leq 2 \cdot \exp\left(-\frac{a^2}{2\sigma^2}\right).$$

eg. $X_1, X_2, \ldots, X_n \overset{iid}{\sim}$ Rade $\begin{cases} +1 & \text{w.p. } \frac{1}{2} \\ -1 & \text{w.p. } \frac{1}{2}. \end{cases}$

$$\mathbb{P}\left(\left|\frac{1}{n}\sum_{i=1}^{n} X_i\right| \geq a\right).$$

$$= 2\mathbb{P}\left(\frac{1}{n}\sum_{i=1}^{n} X_i \geq a\right) \qquad \left(Y = \sum_{i=1}^{n} X_i\right)$$

$$\leq 2\exp(-\lambda a n) \cdot m_Y(\lambda).$$

---

$$m_Y(\lambda) = \mathbb{E}\left[\exp\left(\lambda(X_1 + X_2 + \cdots X_n)\right)\right].$$

$$= \mathbb{E}\left[\exp(\lambda X_1)\right]^n$$

$$= \left(\frac{1}{2}e^{\lambda} + \frac{1}{2}e^{-\lambda}\right)^n.$$

$$\frac{1}{2}(e^{\lambda} + e^{-\lambda}) = 1 + \frac{\lambda^2}{2!} + \frac{\lambda^4}{4!} + \frac{\lambda^6}{6!} + \cdots$$

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

$$\frac{\lambda^{2n}}{(2n)!} \leq \frac{\lambda^{2n}}{n! \cdot 2^n} .$$

So $\quad \frac{1}{2}\left(e^{\lambda} + e^{-\lambda}\right) \leq e^{\frac{1}{2}\lambda^2}$

---

$$m_Y(\lambda) \leq \exp\left(\frac{n\lambda^2}{2}\right).$$

$$\mathbb{P}\left(\left|\frac{1}{n}\sum_{i=1}^{n} X_i\right| \geq a\right)$$

$$\leq 2 \cdot \exp\left(-\lambda a n\right) \cdot \exp\left(\frac{n\lambda^2}{2}\right)$$

$$(\lambda = a)$$

$$= 2 \cdot \exp\left(-na^2/2\right).$$

---

11 Union bound ?.

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

e.g — $X_1, X_2, \cdots , X_n \sim \mathcal{N}(0,1) \quad$ marginally.

$$Y = \max_{1 \leq i \leq N} X_i$$

$$\mathbb{P}(Y \geq a) = \mathbb{P}(\exists i, X_i \geq a)$$

$$\leq \sum_{i=1}^{r} \mathbb{P}(X_i \geq a)$$

$$\leq n \cdot \exp\left(-\frac{a^2}{2}\right), \text{ (want } \leq \delta).$$

$$a = \sqrt{2 \log(n/\delta)}.$$

$$\mathbb{P}\left(Y \geq \sqrt{2 \log(n/\delta)}\right) \leq \delta.$$

---

Stats model $\mathcal{P} = \{\mathbb{P}_\theta : \theta \in \Theta\}$

Observe $X \sim \mathbb{P}_{\theta^*}$

— Decision rule $\delta$ maps from $X$ to $a \in A$.

— Loss function

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

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

$$\mathbb{E}_\theta[\cdots] : \text{Expectation under } \mathbb{P}_\theta.$$

e.g. $g: \Theta \to \mathbb{R}$. $\delta: X \to \mathbb{R}$.

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

$R(\theta, \delta) = \mathbb{E}_\theta \left( g(\theta) - \delta(X) \right)^2$.

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

$\Theta_0 \subseteq \Theta$.

$$L(\theta, a) := \begin{cases} 1, & \theta \in \Theta_0, a = 1. \\ 1, & \theta \notin \Theta_0, a = 0. \\ 0, & \text{otherwise.} \end{cases}$$

$$R(\theta, \delta) = \begin{cases} \text{Type-I err}, & \theta \in \Theta_0 \\ \text{Type-II err}, & \theta \notin \Theta_0. \end{cases}$$

e.g. Statistical learning.

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

$A$: set of functions

from $Z$ to $\{0, 1\}$.

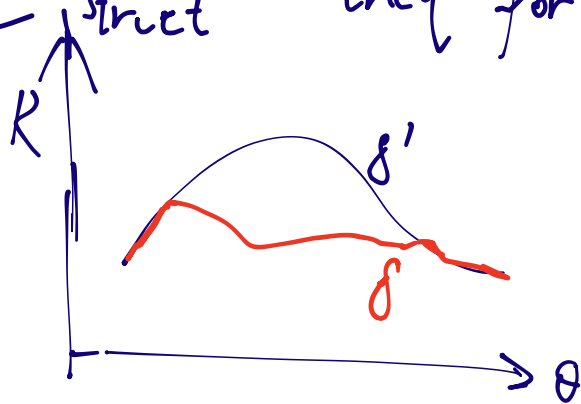$L(\theta, a) = \mathbb{P}_\theta \left( a(z) \neq Y \right)$

$R(\theta, \delta) = \mathbb{E}_\theta \left[ L(\theta, \delta(X)) \right]$

Criteria.
—— Admissibility.

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

with strict ineq for some $\theta$



eg. $X \sim N(\theta, 1)$ $\delta(X) = 0$.

—— Bayes risk. $\pi$ : prior on $\Theta$

$$r_\pi(\delta) = \int_\Theta R(\theta, \delta) \, \pi(d\theta)$$

$$\delta_{Bayes, \pi} = \arg\min \{ r_\pi(\delta) \}.$$

Suppose, $P_\theta(x) = \dfrac{d\mathbb{P}_\theta}{d\nu}(x)$.

$$r_\pi(\delta) = \int_\Theta \int_X L(\theta, \delta(x)) \cdot P_\theta(x) \, \nu(dx) \cdot \pi(d\theta)$$

$$= \int_X \left( \underbrace{\int_\Theta L(\theta, \delta(x)) \cdot p_\theta(x) \, \pi(d\theta)}_{} \right) \lambda(dx)$$

Pointwise minimization.

$$\delta_{Bayes, \pi}(x) = \underset{a \in A}{\arg\min} \left| \frac{\int_\Theta L(\theta, a) \, p_\theta(x) \, \pi(d\theta)}{\int_\Theta p_\theta(x) \, \pi(d\theta)} \right.$$

$$\pi(\cdot | x) := \frac{\pi(\cdot) \, p_\cdot(x)}{\int_\Theta p_\theta(x) \cdot \pi(d\theta)}.$$

$$\delta_{Bayes, \pi}(x) = \underset{a \in A}{\arg\min} \; \mathbb{E}_\pi \left[ L(\theta, a) \right].$$

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

$$\delta_{Bayes, \pi} = \int_\Theta g(\theta) \, \pi(d\theta | x).$$

Conjugate prior.

e.g. $X_1, X_2, \cdots X_n \overset{i.i.d}{\sim} N(\theta, 1)$.

$$\pi = N(0, \tau^2)$$

$$\pi(\theta) \cdot p_\theta(X_1) \cdot p_\theta(\dot{X}_2) \cdots \cdots p_\theta(X_n)$$

$$\propto \exp\left(-\frac{\theta^2}{2\tau^2} - \frac{(X_1-\theta)^2}{2} - \frac{(X_2-\theta)^2}{2} - \cdots - \frac{(X_n-\theta)^2}{2}\right)$$

$$\propto \exp\left(-\left(\frac{1}{2\tau^2} + \frac{n}{2}\right)\theta^2 + \left(\sum_{i=1}^{n} X_i\right)\theta\right)$$

$$\pi(\theta \mid x_1^n) = \mathcal{N}\left(\frac{\tau^2 n \bar{X}_n}{\tau^2 n + 1}, \frac{\tau^2}{\tau^2 n + 1}\right)$$

Another e.g.

Beta-Bernoulli.

$$\theta \sim Beta(\alpha, \beta)$$

$$\pi(\theta) \propto \theta^{\alpha-1} \cdot (1-\theta)^{\beta-1} \quad \text{for} \quad \theta \in [0,1].$$

$$X_1, X_2, \cdots, X_n \overset{iid}{\sim} Ber(\theta)$$

$$\pi(\theta \mid x_1^n) = Beta\left(\alpha + \sum_{i=1}^{n} X_i, \ \beta + \left(n - \sum_{i=1}^{n} X_i\right)\right).$$

___

— Minimax rule

$$\text{minimize} \quad \max_{\theta \in \Theta} R(\theta, \delta).$$

i.e. $\quad \inf_{\delta} \sup_{\theta \in \Theta} \quad R(\theta, \delta).$

---

## Relation between criteria.

— Unique Bayes rules are admissible.

Proof: If $\delta'$ satisfying $R(\theta, \delta') \leq R(\theta, \delta) \; (\forall \theta \in \Theta)$

then $\quad r_\pi(\delta') \leq r_\pi(\delta) \qquad \delta'$ is also Bayes.
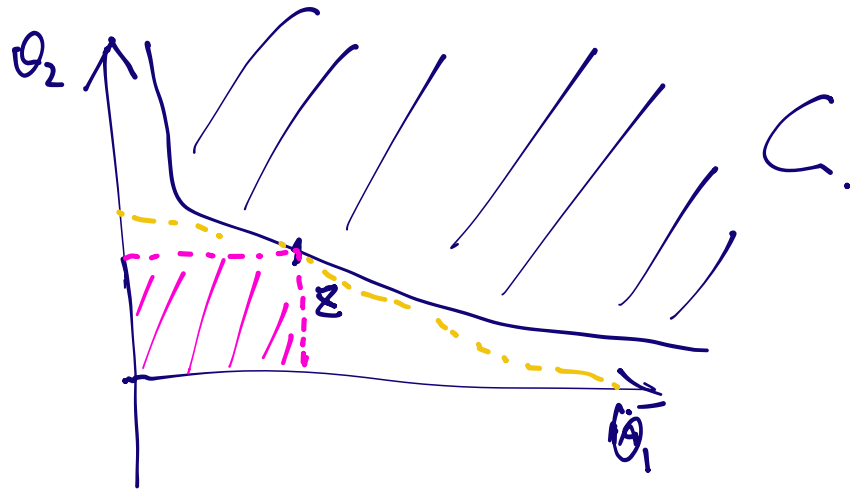
— Suppose $\Theta$ is finite.

All admissible rules are Bayes.

Proof: $\Theta = \{\theta_1, \theta_2, \ldots, \theta_K\}.$

$C := \{(R(\theta_j, \delta))_{j=1}^{K} : \delta \text{ is a rule}\}$

$\delta$ adm. risk vector $z \in \mathbb{R}^K.$

$\theta_2$

$C$

$\theta_1$

$C$ is convex

Separating hyperplane thm.

$\exists \; \lambda \in \mathbb{R}^k \quad \text{s.t.} \quad \pi^T \lambda f \;\begin{cases} \geq a \\ \leq a \end{cases} \qquad \begin{aligned} & x \in C \\ & 0 \leq x \leq B \end{aligned}$