

# STAT5010. Adv. Statistical Inference

---

Term 1, 2022/23

---

Department of Statistics

CUHK

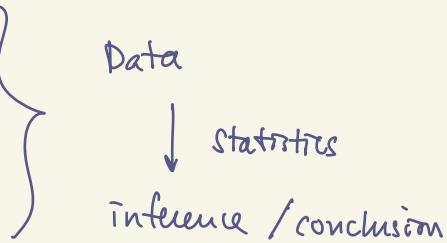
---



# Lecture 1

Different disciplines:

- genetic, medical, earth sciences  
Cox PH ...
- Econ, finance
- Networks, images, recordings



Questions concerned:

- \* Capture the uncertainty.
- \* Methodology → point estimator.  $\hat{\beta} = (X^T X)^{-1} X^T Q$  e.g. random.  $se(\hat{\beta}) = 0.00001$
- \* optimality / optimal inference.  $H_0: \beta = \beta_0 \text{ vs } H_1: \beta = \beta_1$   $\beta_0 = 0$

→ finite-sample optimality (obs:  $X_1, \dots, X_n$ ) ←

→ Asymptotic properties ( $n \rightarrow \infty$ ) ✓

## 1. Decision Theory (Wald, 1939) $E(X) = \int x dF(x)$

$(X_1, X_2, \dots, X_n) \stackrel{iid}{\sim} h(F)$  (distribution) (density:  $f$ )

A statistical model is a family of distributions  $\mathbb{P}$  indexed by a parameter  $\theta$ . We denote

$$\mathbb{P} = \{P_\theta : \theta \in \Omega\}$$

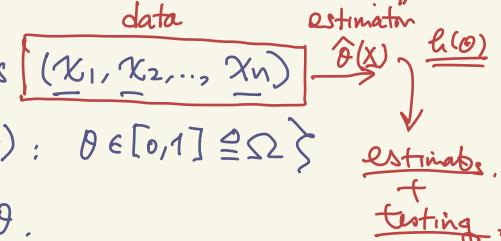
$\Omega \subset \mathbb{R}^k$ : parameter space.  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

Example.

obs: a sequence of coin flips

one can write  $\mathbb{P} = \{\text{Bernoulli}(\theta) : \theta \in [0,1] \cong \Omega\}$

$$P_\theta(X_i=1) = \theta.$$



## A Decision Procedure

Define  $\delta$  (estimator) as a map from the sample space  $X$  to the decision space  $\mathbb{D}$ .

Example. Take  $\mathbb{P} = \{\text{Bernoulli}(\theta)\}$  as shown earlier,

(a) Estimating  $\theta$ :

the decision space is  $\mathbb{D} = [0, 1]$  and one possible decision procedure can be  $\delta(\tilde{x}) = \frac{1}{n} \sum_{i=1}^n x_i \triangleq \bar{X}_n$ .

(b) Hypothesis testing:

Accept / Reject the null hypothesis that  $\theta = 1/2$ .

Correspondingly, the decision space is  $\mathbb{D} = \{\text{Accept}, \text{Reject}\}$ . One possible procedure is

$\delta(\tilde{x}) = \begin{cases} \text{Reject } H_0 & \text{if } \bar{x}_n > 0.5 \text{ and} \\ & \text{Accept it otherwise.} \end{cases}$

(c) A loss function is mapping  $L: \Omega \times \mathbb{D} \rightarrow \underline{\mathbb{R}}^+ \cup \{0\}$

$L(\theta, d)$  represents the penalty for making the decision  $d$  when  $\theta$  is in fact the true parameter for the distribution generating the data.

e.g. Squared-error loss function:  $a) L(\theta, d) = (\theta - d)^2$ .  
 $b) L(\theta, d) = |\theta - d|$ ,  $c) L(\theta, d) = w(\theta - d) \frac{d}{d\theta} L(\theta)$ .

- Frequentist:  $\theta$  is fixed yet unknown.

- Bayesian:  $\Theta \sim \text{distribution.}$

Risk function:  $R(\theta, \delta) = \underbrace{E_\theta}_{\substack{\text{Overall} \\ \text{weighted avg} \\ \text{of } L(\theta, \delta)}} \left( L(\theta, \delta(x)) \right)$

Admissibility:

$\delta$  is inadmissible if there exists  $\delta'$  such that

$R(\theta, \delta') \leq R(\theta, \delta) \forall \theta \in \Omega$  and  $R(\theta', \delta') < R(\theta', \delta)$  for some  $\theta' \in \Omega$ .

- 1) Parametric models: Bernoulli( $\theta$ ),  $N(\mu, \sigma^2)$ ,  $t(k)$  ...  
 - Structure:
- 
- 
- 2) Nonparametric models:  
 $(x_1, \dots, x_n)$   $\hat{F}(x)$   
 $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$ .  
 empirical cdf.
- 3) Semiparametric models:

$$\lambda(t) = \int_0^t \lambda(s) ds$$

Cox proportional hazard model:  
 $\lambda(t|z(t)) = \lambda_0(t) \exp(\beta^T z(t))$ , parametric  
 where  $\lambda(t|z) = \frac{f(t|z)}{1 - F(t|z)}$ ,  $T \sim F_f(t)$   
 $P(T \geq t | z)$ .

## Large-Sample Theory

Convergence in probability:

$$\bar{\gamma}_n(x) = \frac{1}{n} \sum_{i=1}^n \gamma_i(x)$$

Definition: A sequence of random variables  $Y_n$  converges in probability to a random variable  $Y$  as  $n \rightarrow \infty$ , written as  $Y_n \xrightarrow{P} Y$   
 if  $\forall \varepsilon > 0$ ,  $P(|Y_n - Y| \geq \varepsilon) \xrightarrow{n \rightarrow \infty} 0$ .

$$\rightarrow (\hat{\theta}_n \xrightarrow{P} \theta_0)$$

Chebychev's inequality ...

Example: Weak Law of Large Numbers: (WLLN)

$$X_1, \dots, X_n \stackrel{iid}{\sim} (\mu, \sigma^2), \sigma^2 < \infty, \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i,$$

$$E(\bar{X}_n - \mu)^2 = \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n} \rightarrow 0 \text{ as } n \rightarrow \infty,$$

in which case  $\bar{X}_n \xrightarrow{P} \mu$  as  $n \rightarrow \infty$ . ||

[Proposition] If  $f(\cdot)$  is continuous at  $c$ , and if  $Y_n \xrightarrow{P} c$ ,  
 then  $f(Y_n) \xrightarrow{P} f(c)$ .

- Proof. ...

# Lecture 3

19 Sept 2022.

## Large-Sample Theory

- Proof. Since  $f$  is continuous at  $c$ , given any  $\varepsilon > 0$ , there exists  $\delta_\varepsilon > 0$  such that  $|f(y) - f(c)| < \varepsilon$  whenever  $|y - c| < \delta_\varepsilon$ . Thus

$$P(|Y_n - c| < \delta_\varepsilon) \leq P(|f(Y_n) - f(c)| < \varepsilon)$$



Con. in prob.

which implies

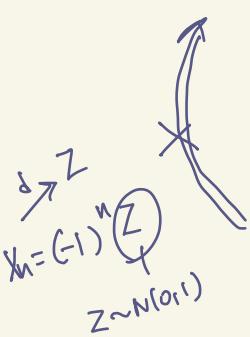
$$P(|f(Y_n) - f(c)| \geq \varepsilon) \leq P(|Y_n - c| > \delta_\varepsilon) \rightarrow 0$$

as  $Y_n \xrightarrow{P} c$  as  $n \rightarrow \infty$ .

□

Stein's method...

Convergence in distribution. (CLT)



Definition. A sequence of r.v.'s  $Y_n$  ( $n \geq 1$ ) with cdfs  $H_n$  converges in distribution (in law) to a random variable  $Y$  with cdf  $H$  if

$$H_n(y) \rightarrow H(y)$$

as  $n \rightarrow \infty$  whenever  $H$  is continuous at  $y$ . For notation, we write  $Y_n \Rightarrow Y$  or  $Y_n \xrightarrow{D/L} Y$  ...

(Example). Suppose  $Y_n = \frac{1}{n}$ , a degenerate r.v. and that  $Y$  is always 0.

Then  $H_n(y) = P(Y_n \leq y) = I\left(\frac{1}{n} \leq y\right)$ .

If  $y > 0$ , then  $H_n(y) = I\left(\frac{1}{n} \leq y\right) \rightarrow 1$  as  $n \rightarrow \infty$  for eventually  $1/n$  will be less than  $y$ . If  $y \leq 0$ , then  $H_n(y) = I(1/n \leq y) = 0$  for all  $n$  and so  $H_n(y) \rightarrow 0$  as  $n \rightarrow \infty$ .

Note that  $H(y) = P(Y \leq y) = I(0 \leq y)$ . Comparisons with the limits just obtained show that  $H_n(y) \rightarrow H(y)$  if  $y \neq 0$ . But  $H_n(0) = 0 \rightarrow 0 \neq 1 = H(0)$ . So, in this example,  $Y_n \xrightarrow{d} Y$  but the cumulative distribution function do not converge to  $H(y)$  as  $y = 0$ .

10 ways of looking at a r.v.

$\sigma$ -algebra /  $\sigma$ -field

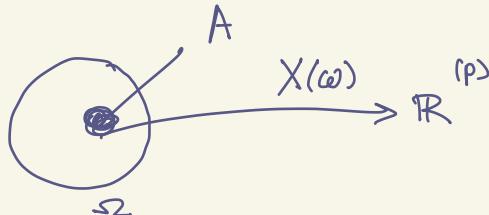
1) Prob. Space  $(\Omega, \mathcal{F}, \mathbb{P})$ .  
sample space ( $\omega$ )

prob. measure, (Measure Theory)

→ Billingsley (1995)

Prob. and Measure.

2) Random variable.



3) Distribution ...

... ↴ Mgf.

8) Characteristic functions.  $\longleftrightarrow$  dist. fn.

$$\phi_x(t) = E(e^{itX}) = \int e^{itx} dF(x),$$

where  $i = \sqrt{-1}$  ( $i^2 = -1$ ) and  $e^{itx} = \cos(tx) + i \sin(tx)$ .

(Example). The characteristic function of  $Z \sim N(0,1)$ :

$$\boxed{\phi_z(t) = \exp\left\{-\frac{1}{2}t^2\right\}}.$$

$$\begin{aligned} |E(e^{itX})| &\leq E|e^{itX}| = E|\cos(tx) + i \sin(tx)| \\ &= E\left(\sqrt{\cos^2(tx) + \sin^2(tx)}\right) = 1. \end{aligned}$$

Fourier transform:

Inversion formula:  $f_x(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-itx} \phi_x(t) dt.$

$$F_x(x) - F_x(y) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{e^{-itx} - e^{-ity}}{it} \phi_x(t) dt.$$

for the cont. points of  $F_x(x)$  and  $F_x(y)$ .

[THEOREM] Uniqueness.

Let  $X$  and  $Y$  be random  $\mathbb{R}^k$ -vectors. Then

(1) If  $\phi_X(t) = \phi_Y(t) \quad \forall t \in \mathbb{R}^k$ , then  $F_X = F_Y$ .  $\checkmark$

(2) If  $M_X(t) = M_Y(t) < \infty$  for all  $t$  in the neighbourhood of 0,  
then  $F_X = F_Y$   $\xrightarrow{\text{E}(e^{tY}) \text{ mgf.}}$  Non-unique moments...  
(C & B, ch.2.X).

Proof.

(1) For any  $\underline{a} = (a_1, \dots, a_k)^T \in \mathbb{R}^k$

and  $\underline{b} = (b_1, \dots, b_k)^T \in \mathbb{R}^k$ ,

and  $[a, b] = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_k, b_k]$

satisfying  $P_X(\text{the boundary of } [a, b]) = 0$ ,

$$P([a, b]) = \lim_{c \rightarrow \infty} \int_{-c}^c \dots \int_{-c}^c \frac{\phi_X(t_1, \dots, t_k)}{(-1)^{k/2} (2\pi)^k} \prod_{j=1}^k \frac{e^{-it_j a_j} - e^{-it_j b_j}}{it_j} dt_j,$$

(inversion formula for the k-dim version; Billingsley (95)).

(2) First, consider the case where  $k=1$ .

$$\text{Since } e^{s|x|} \leq e^{sx} + e^{-sx},$$

we conclude that  $|X|$  has an mgf that is finite in the neighbourhood  $(-c, c)$  for some  $c > 0$ , and that  $|X|$  has finite moments of all orders.

$$\text{Observe that } (\text{Lagrange error bound: } |R_n| \leq \frac{f^{(n+1)}(z) |x-a|^{n+1}}{(n+1)!})$$

$$\left| e^{itx} \left\{ e^{iax} - \sum_{j=0}^n \frac{(iax)^j}{j!} \right\} \right| \leq \frac{|ax|^{n+1}}{(n+1)!},$$

We can write

$$\left| \phi_X(t+a) - \sum_{j=0}^n \frac{a^j}{j!} E((ix)^j e^{iax}) \right| \leq \frac{|a|^{n+1} E|x|^{n+1}}{(n+1)!}$$

Facts:

- 1) If  $M_X(t)$  is finite in a neighbourhood of 0, then  $E(X_1^{r_1} X_2^{r_2} \dots X_k^{r_k})$  is finite for any integers  $r_1, r_2, \dots, r_k \geq 0$  and  $M_X(t)$  has the power series expansion

$$M_X(t) = \sum_{(r_1, r_2, \dots, r_k) \in \mathbb{Z}_+^k} \frac{E(X_1^{r_1} \dots X_k^{r_k}) t_1^{r_1} \dots t_k^{r_k}}{r_1! r_2! \dots r_k!}$$

for any  $k$ -dim random vector  $X$ .

2)  $\frac{\partial^r \phi_X(t)}{\partial t_1^{r_1} \dots \partial t_k^{r_k}} = (-1)^{r/2} E(X_1^{r_1} \dots X_k^{r_k} e^{it^T X})$

with  $\left. \frac{\partial^r \phi_X(t)}{\partial t_1^{r_1} \dots \partial t_k^{r_k}} \right|_{t=0} = (-1)^{r/2} E(X_1^{r_1} \dots X_k^{r_k})$

We can write, for any  $t \in \mathbb{R}$ ,

$$\phi_X(t+a) = \sum_{j=0}^{\infty} \underbrace{\frac{\phi_X^{(j)}(t)}{j!} a^j}_{\text{jth der. of } \phi \text{ wrt } t.} \quad |a| < 0 \quad (*)$$

$$\begin{aligned} (\text{b/c: } E(e^{i(t+a)X})) &= E\left(e^{itX} \sum_{j=0}^{\infty} \frac{a^j}{j!} (iX)^j\right) \\ &= \sum_{j=0}^{\infty} \frac{a^j}{j!} \underbrace{E(e^{itX} \cdot (iX)^j)}_{\text{by (2)} = \phi_X^{(j)}(t)} \end{aligned}$$

The result/expansion also holds for  $Y$ . Under the assumption that  $M_X = M_Y < \infty$  in a neighbourhood of 0,  $X$  and  $Y$  have the same moments of all orders. By (F2),  $\phi_X^{(j)}(0) = \phi_Y^{(j)}(0)$  for all  $j = 1, 2, \dots$ , which and (\*) with  $t=0$  imply that  $\phi_X$  and  $\phi_Y$  are the same on the interval  $(-c, c)$ .

Considering  $t=c-\varepsilon$ , and  $-c+\varepsilon$  for an arbitrarily small  $\varepsilon > 0$  in (\*) shows that  $\phi_X$  and  $\phi_Y$  also agree on  $(-2c+\varepsilon, 2c-\varepsilon)$  and hence  $(-2c, 2c)$ . Likewise, the same argument  $\phi_X$  and  $\phi_Y$  are the same on  $(-3c, 3c)$  and so on. Hence,  $\phi_X(t) = \phi_Y(t)$  for all  $t$  and by (i)  $F_X = F_Y$ .

For the case  $k \geq 2$ , if  $F_x \neq F_y$ , then by part (i), there exists  $t \in \mathbb{R}^k$  such that  $\phi_x(t) \neq \phi_y(t)$ . Then  $\phi_{t^T x}(1) \neq \phi_{t^T y}(1)$ , which implies that  $F_{t^T x} = F_{t^T y}$ .

But  $M_x = M_y < \infty$  in a neighbourhood of  $0 \in \mathbb{R}^k$  implies that  $M_{t^T x} = M_{t^T y} < \infty$  in a neighbourhood of  $0 \in \mathbb{R}$  and by the proved result for  $k=1$ ,  $M_{t^T x} = M_{t^T y}$ . The contradiction shows that  $F_x = F_y$ .  $\square$

a) Conditional Probability  $\xrightarrow{\text{Ch. 6 of Keener}}$

$$P(B|A) = \frac{P(A \cap B)}{P(A)} \quad \text{when } P(A) \neq 0,$$

where by convention  $P(B|A) = 0$  when  $P(A) = 0$ .

Suppose  $E|f(x, Y)| < \infty$ ,

$$E(f(x, Y)) = E\left(E(f(x, Y) | X)\right) \quad \checkmark \quad \begin{matrix} \text{lower expectation.} \\ \text{c.f.} \end{matrix} \quad (***) \quad E(Y|X) = \beta_0^T X$$

with  $E(f(x, Y) | X) \equiv h(X)$ , and

$$\rightarrow h(x) = E(f(x, Y) | X=x) = \int f(x, y) dQ_x(y),$$

where  $Q_x(y)$  denotes the conditional distribution of  $Y$  given  $X=x$ .

(\*\*\*) becomes

$$E(f(x, Y)) = \underbrace{\int H(x) dP_x(x)}_{\substack{\text{dist. of } X}} = \iint f(x, y) dQ_x(y) dP_x(x).$$

### Definition



The function  $Q$  is a conditional distribution of  $Y$  given  $X$ , written as  $Y | X=x \sim Q_x$  if

(i)  $Q_x(\cdot)$  is a prob. measure for all  $x$ ,

(ii)  $Q_x(B)$  is a measurable function of  $x$  for any Borel set  $B$ ,

(iii) for any Borel sets  $A$  and  $B$ ,

$$P(X \in A, Y \in B) = \int_A Q_x(B) dP_x(x).$$

## 10) Tail behaviour:

For a scalar random variable  $X$  with pdf  $f$ , we say  $X$  has

i) an exponential tail if

$$\lim_{a \rightarrow \infty} \frac{-\log(1-F(a))}{ca^r} = 1 \quad \text{for some } c > 0, r > 0.$$

ii) an algebraic tail if

$$\lim_{a \rightarrow \infty} \frac{-\log(1-F(a))}{m \log a} = 1 \quad \text{for some } m > 0,$$

Examples:

Exponential  $F(a) = 1 - e^{-\lambda a} \Rightarrow c = \lambda, r = 1$

Gaussian  $F(a) = \dots \Rightarrow c = 2, r = 2$ .

Student t-distr. heavy-tail dist. for small d.f.,  $m=1$ .

## → Sufficiency (data reduction/compression)

Def. A statistic  $T: X \mapsto T$  is a function of data

e.g.  $T(X) = \frac{1}{n} \sum_{i=1}^n X_i \triangleq \bar{X}_n$ ,

Def. A statistic is sufficient for a model  $P = \{P_\theta : \theta \in \Omega\}$  if for all  $t$ , the conditional distribution  $[X | T(x)=t]$  does not depend on  $\theta$ .

(Example) Let  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$ . Is  $\sum_{i=1}^n X_i$  sufficient?

$$\begin{aligned} P_\theta \left( X=x \mid T(x)=t \right) &= \frac{P_\theta \left( X=x, T(X)=t \right)}{P_\theta(T(X)=t)} \\ &= \frac{\mathbb{I}(t=\sum_{i=1}^n X_i) \theta^t (1-\theta)^{n-t}}{\binom{n}{t} \theta^t (1-\theta)^{n-t}} \\ &= \frac{\mathbb{I}(\sum_{i=1}^n X_i = t)}{\binom{n}{t}} = \frac{1}{\binom{n}{t}} \end{aligned}$$

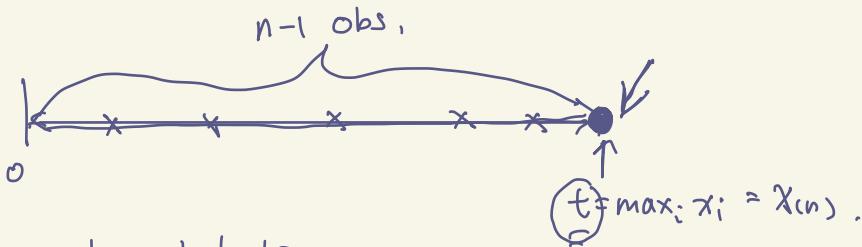
(Example). Let  $X_1, \dots, X_n \stackrel{iid}{\sim} U(0, \theta)$ . Show that  $T(X) = \max_{1 \leq i \leq n} X_i \stackrel{\text{def}}{=} X_{(n)}$  is sufficient.

$$F_T(t) = P(T \leq t) = P\left(\max_{1 \leq i \leq n} X_i \leq t\right) = P(X_1 \leq t, \dots, X_n \leq t)$$

$$= \left\{ P(X_i \leq t) \right\}^n = \left(\frac{t}{\theta}\right)^n$$

$$f_T(t) = \frac{n t^{n-1}}{\theta^n} \quad \text{and} \quad P(X|T) \perp\!\!\!\perp \theta.$$

Intuition:



(Example) The order statistics

$T(\underline{x}) = (X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)})$  from a random sample of  $X_i \stackrel{iid}{\sim} f$  are sufficient. Given  $T(\underline{x})$ , the possible values of  $X$  are  $n!$  permutations of  $T$ :

$$P(X_1 = X_{(1)}, \dots, X_n = X_{(n)}) = \frac{1}{n!}$$

[THEOREM]. If  $X \sim P_\theta \in \mathcal{P}$  and  $T$  is sufficient for  $P$ , then for any decision procedure  $\delta$ , there is a (possibly randomized) decision procedure of equal risk that depends on  $X$  through  $T(X)$  only.

Remark: Suppose we are given an independent source of randomness, say  $U \perp\!\!\!\perp \theta$ , we can generate a new dataset  $X'$  from the cond. dist.  $P(X|T(X))$  and define a randomised procedure:

$$\delta^*(x, u) \stackrel{\text{def}}{=} \delta(f(T(x)), u) = \delta(x') \stackrel{d}{=} \delta(x)$$

(Example)  $X$  and  $Y \stackrel{\text{iid}}{\sim} f_\theta$ , where  $f_\theta(x) = \begin{cases} \theta e^{-\theta x} & , x \geq 0 \\ 0 & , \text{otherwise.} \end{cases}$

Let  $U \sim \text{uniform}(0, 1) \perp\!\!\!\perp (X, Y)$ . Define  $T = X + Y$ .

$$\rightarrow \tilde{X} = UT \text{ and } \tilde{Y} = (1-U)T. \quad \boxed{\tilde{X} + \tilde{Y} = X + Y}$$

$$\begin{aligned} P(T \leq t | Y=y) &= P(X+Y \leq t | Y=y) \\ &= E(I(X+Y \leq t) | Y=y) \\ &= \int I(x \leq t-y) dF_x(x) \\ &= F_x(t-y) \end{aligned}$$

We can write  $P(T \leq t | Y) = F_x(t-Y)$ , a r.v. of  $Y$ .

$$\begin{aligned} P_+(t) &= P(T \leq t) = E_Y(F_x(t-Y)) \\ &= \int_0^t \underbrace{1 - e^{-\theta(t-y)}}_{\theta e^{-\theta y}} dy, \\ &= 1 - e^{-\theta t} + \theta e^{-\theta t}. \end{aligned}$$

$$\text{which gives } f_T(t) = \frac{\partial F_T(t)}{\partial t} = t\theta^2 e^{-\theta t}, t \geq 0,$$

Note<sup>that</sup>  $(U, T)$  has the joint density

$$p_\theta(t, u) = \begin{cases} t\theta^2 e^{-\theta t} & , t \geq 0, u \in (0, 1) \\ 0 & , \text{otherwise} \end{cases}$$

$$P_\theta \left( \begin{pmatrix} \tilde{X} \\ \tilde{Y} \end{pmatrix} \in B \right) = \iint_B \begin{pmatrix} \tilde{X} \\ \tilde{Y} \end{pmatrix} \begin{pmatrix} tu \\ (1-u)t \end{pmatrix} p_\theta(t, u) du dt$$

Hence,  $\tilde{X}$  and  $\tilde{Y}$  have the joint density of  $t$

$$\frac{P_\theta(x+y, \frac{x}{x+y})}{x+y} = \begin{cases} \theta^2 e^{-\theta \underline{x+y}} & , x \geq 0, y \geq 0 \\ 0 & , \text{otherwise.} \end{cases}$$

which is the same joint density as  $(X, Y)$ . ||

[THEOREM]. Neyman-Fisher Factorisation Criterion.

Suppose each  $P_\theta \in \mathcal{P}$  has density  $p(x; \theta)$  with respect to a common  $\sigma$ -finite measure  $\mu$ , i.e.  $dP_\theta / d\mu = p(x; \theta)$ . Then,  $T(X)$  is sufficient if and only if

e.g.  $(\mu, \sigma^2)$  for normal,

$$p(x; \theta) = g_\theta(T(x)) h(x)$$

for some functions  $g_\theta$  and  $h$  where  $h(x)$  is free of  $\theta$ .

Proof. (Discrete: Refer to Casella & Berger, 2001) ...

For the continuous case,  $T$  and  $X$  will not have a joint density with respect to any product measure. To begin, suppose  $P_\theta$ ,  $\theta \in \Omega$  has density

$$\rightarrow p_\theta(x) = g_\theta(T(x)) \underbrace{h(x)}_{\text{support of } x} \quad \leftarrow \text{assumed}$$

with respect to  $\mu$ . Modifying  $h$ , we can assume wlog that  $\mu$  is a prob. measure equivalent to the family  $\mathcal{P} = \{P_\theta : \theta \in \Omega\}$ .

$$\mu(N) = 0 \text{ if and only if } \int_N p_\theta(x) dx = 0 \quad \forall \theta \in \Omega.$$

Let  $E^*$  and  $P^*$  denote the expectation and probability when  $X \sim \mu$ ; let  $G^*$  and  $g(\theta)$  denote the marginal distributions for  $T(x)$  when  $X \sim \mu$  and  $X \sim P_\theta$ ; and let  $\varrho$  be the conditional distribution for  $X$  given  $T$  when  $X \sim \mu$ .

To find the densities for  $T$ :

$$\begin{aligned} (E_\theta(f(T))) &= \int f(T(x)) g_\theta(T(x)) h(x) d\mu(x) \\ &= E^*(f(T) g_\theta(T) h(x)) \\ &\stackrel{\text{tower exp.}}{=} \iint f(t) g_\theta(t) h(x) d\varrho_t(x) dG^*(t) && \text{cond. dist. of } x \text{ given } T. \\ &\triangleq \int \boxed{f(t)} \underbrace{g_\theta(t) w(t)}_{\substack{\text{fixed } t \\ \text{density}}} dG^*(t), \end{aligned}$$

$$\text{where } w(t) = \frac{\int h(x) d\varrho_t(x)}{\int h(x) d\mu(x)}$$

↑  $\int g_\theta(x) dx$   
↑  $\text{lea. mea.}$   
↑  $\text{density}$

If  $f(\cdot)$  is an indicator function, this shows that  $\hat{Q}_\theta$  has density  $g_\theta(t) w(t)$  with respect to  $G^*$

Next, define  $\tilde{Q}_t$  to have the density  $\frac{h}{w(t)}$  with respect to  $Q_t$  so that

$$\tilde{Q}_t(B) = \int_B \frac{h(x)}{w(t)} dQ_t(x)$$

Then,

$$\begin{aligned} E_\theta(f(x, T)) &= E^*(f(x, T) g_\theta(T) h(x)) \\ &= \iint f(x, t) g_\theta(t) h(x) d\tilde{Q}_t(x) dG^*(t) \\ &= \iint f(x, t) \left[ \frac{h(x)}{w(t)} dQ_t(x) \right] [g_\theta(t) w(t) dG^*(t)] \\ &= \iint f(x, t) d\tilde{Q}_t(x) dG_\theta(t) \end{aligned}$$

By the definition, this shows that  $\hat{Q}$  is a conditional distribution of  $X$  given  $T$  under  $P_\theta$ . Because  $\hat{Q}$  does not depend on  $\theta$ ,  $T$  is sufficient.

To assume that  $T$  is suff... (next lecture).

## Lecture 4

26 Sept.

### NFFC proof (cont'd)

- Proof. Mixture distributions: Given a marginal distribution  $G^*$  and a conditional distribution  $Q$ , we can define a mixture distribution

$$\hat{P} \text{ by } \rightarrow \hat{P}(B) = \int \underbrace{Q_t(B)}_{\text{cond. "given"} t} dG^*(t) = \iint 1_B(x) dQ_t(x) dG^*(t).$$

Then for any integrable  $f$ ,

$$\iint f d\hat{P} = \iint f(x) \underbrace{dQ_t(x)}_{d\hat{P}} dG^*(t).$$

Suppose now that  $T$  is sufficient (with  $Q_T$  the conditional distribution for  $X$  given  $T$ ). Let  $g_\theta$  be the  $G^*$  density of  $T$  when  $X \sim P_\theta$ .

Ch. 1 of  
Keener

Existence: If  $G^*(N) = 0$ ,  $\mu(N_0) = 0$ , where  $N_0 = \overline{T}(N)$  null sets

$$\begin{aligned} \text{and so } G_\theta(N) &= P_\theta(T \in N) \\ &= P_\theta(X \in N_0) = \int_{N_0} g_\theta(y) d\mu = 0 \end{aligned}$$

$$\begin{aligned} \text{Then } P_\theta(X \in B) &= E_\theta(P_\theta(X \in B | T)) \\ &= E_\theta(Q_T(B)) \\ &= \int Q_t(B) g_\theta(t) dG^*(t) \\ &= \iint 1_B(x) dQ_t(x) g_\theta(t) dG^*(t) \\ &= \iint 1_B(x) g_\theta(T(x)) dQ_t(x) dG^*(t) \\ &= \int_B g_\theta(T(x)) d\hat{P}(x), \quad P(X \in B) = \int_B f(x) dx. \end{aligned}$$

which shows that  $P_\theta$  has the density  $\underline{g}_\theta(T(\cdot))$  with respect to  $\hat{P}$ ,

The mixture distribution  $(\hat{P})$  is absolutely continuous with respect to  $\mu$ .

Keener (2010) Thm 1.10.

If a finite measure  $P$  is abs. cont. w.r.t. a  $\sigma$ -finite measure  $\mu$ , then there exists a non-negative measurable function  $f$  such that

$$P(A) = \int_A d\underline{P} = \int_A \left( \frac{dP}{d\mu} \right) d\mu \stackrel{\text{def}}{=} \int_A f d\mu$$

$$P(A) = \int_A f d\mu = \int f 1_A d\mu.$$

The function  $f$  is called the Radon-Nikodym derivative of  $P$  with respect to  $\mu$ , or the density, denoted as  $f = dP/d\mu$ .

Suppose  $\mu(N) = 0$ , then  $P_\theta(N) = \int_{Q_t(N)} dG_\theta(t) = 0$ , which implies that  $G_\theta(\tilde{N}) = 0$ , where  $\tilde{N} = \{t : Q_t(N) > 0\}$ . Because  $\mu$  is equivalent to  $P$  and  $G_\theta(\tilde{N}) = P_\theta(T \in \tilde{N}) = 0 \forall \theta \in \Omega$ ,  $P^*(T \in \tilde{N}) = G^*(\tilde{N}) = 0$ . Thus,  $Q_t(N) = 0$  (a.e.  $G^*$ ) and so  $\hat{P}(N) = \int_{Q_t(N)} dG^*(t) = 0$ . Taking  $\hat{\mu} = \frac{d\hat{P}}{dP^*}$ ,  $P_\theta$  has density  $\underline{g_\theta(T(x)) \hat{\mu}(x)}$  with respect to  $P^*$ .

(Example) let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $\theta = (\mu, \sigma^2)$ .

$\mu \quad G_\theta \quad \square$

The joint distribution of  $(X_1, \dots, X_n)$  is

$$\begin{aligned} p(\underline{x}; \theta) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(x_i - \mu)^2}{2\sigma^2} \right\} \\ &= \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp \left\{ +\frac{1}{2\sigma^2} \left( -\sum_{i=1}^n x_i^2 + 2\mu \sum_{i=1}^n x_i - n\mu^2 \right) \right\} \end{aligned}$$

$\overbrace{\qquad\qquad\qquad}^{\text{"S}_{\theta}(\underline{T})\text{"}}$

$\prod_{i=1}^n I(-\infty < x_i < \infty)$

$\hat{\mu}(x)$

$$\text{We have } T(\underline{x}) = \left( \sum_{i=1}^n x_i^2, \sum_{i=1}^n x_i \right).$$

By NFFC, we claim that  $T$  is sufficient for  $\theta$  (or  $P_\theta$ ).  $\square$

### Exponential families

The model  $\mathcal{P} = \{P_\theta : \theta \in \Omega\}$  forms an  $n$ -dimensional exponential family if each  $P_\theta$  has the density of the form:

$$p(x; \theta) = \exp \left\{ \sum_{i=1}^n \eta_i(\theta) T_i(x) - B(\theta) \right\} h(x)$$

suff. stat. by NFFC  
base measure.  
e.g.  $I(x \in R)$ .

"natural parameters" "standardiser"

$$B(\theta) = \log \left( \int \exp \left\{ \sum_{i=1}^n \eta_i(\theta) T_i(x) \right\} h(x) d\mu(x) \right) \in \mathbb{R},$$

Examples include: Normal, Binomial, Poisson...



(Example). Exponential distribution  $\mathcal{P} = \{\text{Exp}(\theta) : \theta > 0\}$

$$\begin{aligned} \text{The densities take the form } p(x; \theta) &= \theta e^{-\theta x} I(x \geq 0) \\ &= \exp(-\theta x + \log \theta) I(x \geq 0) \end{aligned}$$

We have:  $\eta_1(\theta) = -\theta$ ;  $T_1(x) = x$ ,  $B(\theta) = -\log \theta$ ,  $\theta(x) = I(x \geq 0)$   
 $\tilde{\eta}_1(\theta) = \theta$ ;  $\tilde{T}_1(x) = -x$ ,  $B(\theta) = -\log \theta$ ,  $\theta(x) = I(x \geq 0)$

(Example) Beta distribution  $P = \{Beta(\alpha, \beta) : \alpha, \beta > 0\}; \theta = (\alpha, \beta)$

$$p(x; \theta) = x^{\alpha-1}(1-x)^{\beta-1} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} I(0 \leq x \leq 1)$$

$$= \exp \left\{ \underbrace{(\alpha-1) \log x}_{\eta_1(\theta)} + \underbrace{(\beta-1) \log(1-x)}_{\eta_2(\theta)} + \log \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \right\} I(0 \leq x \leq 1)$$

$\sim B(\theta)$

$$= \exp \left\{ \alpha \log x + \beta \log(1-x) + \log \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \right\} \frac{I(0 \leq x \leq 1)}{x(1-x)}$$

### Definition

An exponential family is in canonical form when density has the form:

$$p(x; \eta) = \exp \left\{ \sum_{i=1}^s \eta_i T_i(x) - A(\eta) \right\} \theta(x).$$

### Definition

The set of all valid natural parameters  $\Omega$  is called the natural parameter space for each  $\eta \in \Omega$ , there exists a normalizing constant  $A(\eta)$  such that  $\int p(x; \eta) dx = 1$ .

Equivalently,

$$\Omega = \left\{ \eta : 0 < \int \exp \left\{ \sum_{i=1}^s \eta_i T_i(x) \right\} \theta(x) d\mu(x) < \infty \right\}.$$

For any canonical exponential family,  $P = \{f_\eta : \eta \in H\}$ , we have  $H \subseteq \Omega$ ,  $\Omega$  is convex (see TSH, §2.7.1).

### Definition

If  $P = \{P_\theta : \theta \in \Sigma\}$ , then  $\theta$  is unidentifiable if for two parameters  $\theta_1 \neq \theta_2$ ,  $f_{\theta_1} = f_{\theta_2}$ .

### Remark

Two cases when the superficial dimension of an  $n$ -exponential family  $P$  can be reduced:

1)  $X \sim \text{Exp}(\eta_1, \eta_2)$  with density

$$p(x; \eta_1, \eta_2) = \exp\{-\eta_1 x - \eta_2 x + \log(\eta_1 + \eta_2)\} I(x \geq 0),$$

where  $T_1(X) = T_2(X) = X$  (are linearly dependent). We can actually combine  $(\eta_1, \eta_2)$  into  $\eta_1 + \eta_2$  and rewrite (\*) as

$$p(x; \eta_1, \eta_2) = \exp\left\{-\underbrace{(\eta_1 + \eta_2)}_{\eta} x + \log\underbrace{(\eta_1 + \eta_2)}_{\eta}\right\} I(x \geq 0)$$

The  $T_i(x)$ 's satisfy an affine equality (linearly dep.) constraint for all  $x \in X$  which will result in unidentifiability, i.e.

$$p(x; \eta_1 + c, \eta_2 - c) = p(x; \eta_1, \eta_2) \text{ for all } c < \eta_2$$

2. The  $\eta_i$ 's satisfy an affine equality constraint for all  $y \in H$

$p(x; \eta) \propto \exp(\eta_1 x + \eta_2 x^2)$  for all  $(\eta_1, \eta_2)$  satisfying  $\eta_1 + \eta_2 = 1$ . Then we can write

$$p(x; \eta) \propto \exp(\eta_1 x + \eta_2 x^2) = \exp\left(\eta_1 \ln x + x^2\right).$$

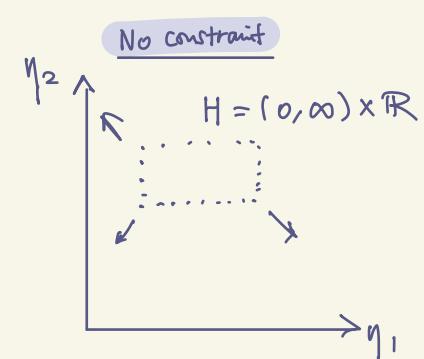
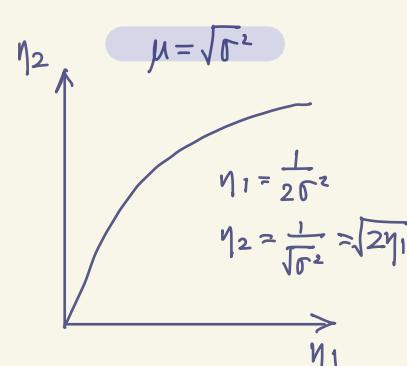
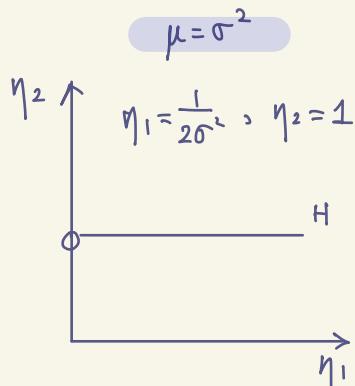
Definition A canonical exponential family  $P = \{P_\eta : \eta \in H\}$  is minimal if

(a) No affine  $T_i$ 's:  $\sum_{i=1}^s \lambda_i T_i(x) = \lambda_0 \quad \forall x \in X \Rightarrow \lambda_i = 0 \quad \forall i \in \{0, \dots, s\}$

(b) No affine  $\eta_i$ 's:  $\sum_{i=1}^s \lambda_i \eta_i = \lambda_0 \quad \forall \eta \in H \Rightarrow \lambda_i = 0 \quad \forall i \in \{0, \dots, s\}$

Definition Suppose  $P = \{P_\eta : \eta \in H\}$  is an  $s$ -dimensional minimal exponential family. If  $H$  contains an open rectangle ( $s$ -dim), then  $P$  is called full-rank, otherwise,  $P$  is curved.

(Example) Consider  $N(\mu, \sigma^2)$ , where  $\eta_1 = \frac{1}{2\sigma^2}$ ,  $\eta_2 = \frac{\mu}{\sigma^2}$ ,  $T_1(x) = -x^2$ ,  $T_2(x) = x$



Curved exp. family (Ch.5 of Keener, 2010)

Remarks: Properties of Exponential Families.

1) If  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p(x; \theta) = \exp \left\{ \sum_{i=1}^s \eta_i(\theta) T_i(x) - A(\theta) \right\} h(x)$

$$\text{then, by NFFC, } T = \left( \sum_{j=1}^n T_1(X_j), \dots, \sum_{j=1}^n T_s(X_j) \right)$$

' is sufficient. The exponential family data is highly compressible.

2) If  $f$  is integrable, and  $\eta \in \Omega$ , then

$$G(f, \eta) = \int f(x) \exp \left\{ \sum_{i=1}^s \eta_i T_i(x) \right\} h(x) d\mu(x)$$

is infinitely differentiable with respect to  $\eta$  and the derivatives can be obtained by differentiating it under the integral sign.

3) Moments of  $T_i$ 's...

Take  $f(x) \equiv 1$ , then

$$G(1, \eta) = \int \exp \left\{ \sum_{i=1}^s \eta_i T_i(x) \right\} h(x) d\mu(x) = \exp \{ A(\eta) \},$$

$$\frac{\partial G(1, \eta)}{\partial \eta_i} = \int T_i(x) \exp \left\{ \sum_{i=1}^s \eta_i T_i(x) \right\} h(x) d\mu(x) = \frac{\partial A(\eta)}{\partial \eta_i} = \exp \{ A(\eta) \}$$

$$\begin{aligned} \frac{\partial A(\eta)}{\partial \eta_i} &= \int T_i(x) \exp \left\{ \sum_{i=1}^s \eta_i T_i(x) - A(\eta) \right\} h(x) d\mu(x) \\ &= E_\eta(T_i(x)) \end{aligned}$$

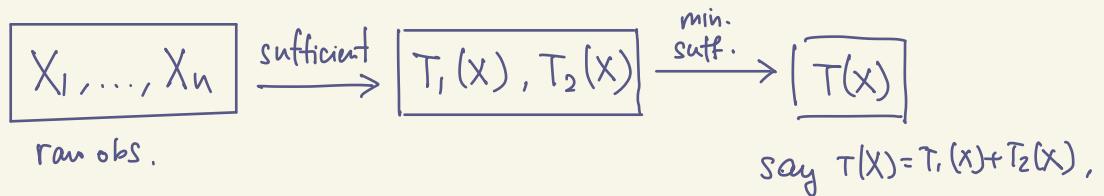
$$\begin{aligned} \frac{\partial^2 A(\eta)}{\partial \eta_i \partial \eta_j} &= \int T_i(x) \left( T_j(x) - \frac{\partial A(\eta)}{\partial \eta_j} \right) \exp \left\{ \sum_{l=1}^s \eta_l T_l(x) - A(\eta) \right\} \\ &\quad h(x) d\mu(x) \end{aligned}$$

$$= E_\eta(T_i(x) T_j(x)) - E_\eta(T_i(x)) E_\eta(T_j(x))$$

$$= \text{cov}_\eta(T_i(x), T_j(x)),$$

### Minimal Sufficiency

Definition A sufficient statistic  $T$  is minimal if for every sufficient statistic  $T'$  and for every  $x, y \in \mathcal{X}$ ,  $T(x) = T(y)$  when  $T'(x) = T'(y)$ . In other words,  $T$  is a function of  $T'$ , i.e. there exists a function  $f$  such that  $T(x) = f(T'(x))$  for any  $x \in \mathcal{X}$ .



[THEOREM]. Let  $\{p(x; \theta) : \theta \in \Omega\}$  be a family of densities with respect to some measure  $\mu$  (Lebesgue measure for continuous distribution  $dx$ ; counting measure for discrete distribution  $(x; -x_{i-1})$ ). Suppose that there exists a statistic  $T$  such that, for every  $x, y \in X$ ,

$$\rightarrow \boxed{p(x; \theta) = C_{x,y} p(y; \theta)} \Leftrightarrow \boxed{T(x) = T(y)}$$

for every  $\theta$  and some  $C_{x,y} \in \mathbb{R}$ . Then  $T$  is a minimal sufficient statistic.

• Proof. We first prove that  $T$  is sufficient and then  $T$  is minimal.

$\Rightarrow$  sufficiency for T:

Start with  $T(X) = \{t : t = T(x) \text{ for some } x \in X\} = \text{range of } T.$

For each  $t \in T(X)$ , we consider the preimage  $A_t = \{X : T(X) = t\}$ .

Then, for any  $y \in X$ , we have  $y \in A_{T(y)}$  and  $X_{T(y)} \in A_{T(y)}$ .  
 By the definition of  $A_T$ , we can see that  $T(y) = T(X_{T(y)})$ .

Because of the assumption / condition imposed in the statement ,

$$\text{ply}(\theta) = \underbrace{C_{y, X_{T(y)}}}_{\hat{p}} p(X_{T(y)}; \theta) \triangleq \hat{p}(y) g_\theta(T(y)),$$

which implies that  $T$  is sufficient as a result of the NFFC.

$\Leftrightarrow$ ) Minimality of  $T$ :

Consider another sufficient statistic, say  $T'$ . By NFFC,

$$p(x; \theta) = \tilde{q}_\theta(T'(x)) \tilde{h}(x)$$

Take any  $x$  and  $y$  such that  $T'(x) = T'(y)$ , then

$$p(x; \theta) = \tilde{g}_{\theta}(T'(x)) \hat{h}(x)$$

$$= \underbrace{\tilde{g}_\theta(T'(y))}_{\text{wavy line}} \widehat{h}(y) \cdot \frac{\tilde{h}(x)}{\tilde{h}(y)} \stackrel{\text{def}}{=} \text{ply}(\theta)[x,y]$$

- Hence,  $T(x) = T(y)$  by the assumption of the minimal sufficient statistic theorem. So  $T'(x) = T'(y)$  implies  $T(x) = T(y)$  for any sufficient statistic  $T'$  and any  $x$  and  $y$ . Hence,  $T$  is a minimal sufficient statistic.  $\square$

(Example) Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$  with both  $\mu$  and  $\sigma^2$  unknown. Let  $\tilde{x}$  and  $\tilde{y}$  be any two sample points, and let  $(\bar{x}_n, S_x^2)$  and  $(\bar{y}_n, S_y^2)$  be the sample means and variances corresponding to the  $\tilde{x}$  and  $\tilde{y}$  samples, respectively.

$$\begin{aligned} \frac{f(\tilde{x}; \mu, \sigma^2)}{f(\tilde{y}; \mu, \sigma^2)} &= \frac{(2\pi\sigma^2)^{-n/2} \exp\left\{-\frac{n(\bar{x}_n - \mu)^2 + (n-1)S_x^2}{2\sigma^2}\right\}}{(2\pi\sigma^2)^{-n/2} \exp\left\{-\frac{n(\bar{y}_n - \mu)^2 + (n-1)S_y^2}{2\sigma^2}\right\}} \\ &= \exp\left\{-\frac{1}{2\sigma^2} \left\{ -n(\bar{x}_n^2 - \bar{y}_n^2) + 2n\mu(\bar{x}_n - \bar{y}_n) - (n-1)(S_x^2 - S_y^2) \right\} \right\} \end{aligned}$$

The ratio will be constant as a function of  $\theta = (\mu, \sigma^2)$  iff  $\bar{x}_n = \bar{y}_n$  and  $S_x^2 = S_y^2$ . Thus, by the theorem " $(\bar{x}_n, S_x^2)$  is a minimum sufficient statistic for  $\theta$ ", where  $S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ .

(Example) (curved exponential family)

Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\sigma, \sigma^2)$ , where  $\sigma > 0$ . Then  $\theta = \sigma$

$$\frac{p(\tilde{x}; \theta)}{p(\tilde{y}; \theta)} = \exp\left\{-\frac{1}{2\sigma^2} \left( \sum_{i=1}^n x_i^2 - \sum_{i=1}^n y_i^2 \right) + \frac{1}{\sigma} \left( \sum_{i=1}^n x_i - \sum_{i=1}^n y_i \right) \right\}.$$

This ratio will be constant as a function of  $\theta = \sigma$  iff  $\sum_{i=1}^n x_i = \sum_{i=1}^n y_i$  and  $\sum_{i=1}^n x_i^2 = \sum_{i=1}^n y_i^2$ . Thus, by the above theorem,  $T(X) = (\sum_{i=1}^n x_i, \sum_{i=1}^n x_i^2)$  is minimal sufficient.

Remark: If  $p(x; \theta) = C_{x,y} p(y; \theta)$  and  $x$  and  $y$  must be supported by the sample  $\theta$ . (Support of  $X$ :  $\{x \in \mathcal{X} : p(x; \theta) > 0\}$ ). Otherwise, the constant "C<sub>x,y</sub> will be  $\theta$ -dependent.

(Example) Let  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Uniform}(0, \theta)$  and  $T(X) = \max_{1 \leq i \leq n} X_i = X_{(n)}$

In that case, for  $\underline{x} = (x_1, \dots, x_n)$  such that  $x_i > 0$ , for  $i=1, \dots, n$ ,

$$p(\underline{x}; \theta) = \prod_{i=1}^n \frac{1}{\theta} I(x_i \leq \theta) = \frac{1}{\theta^n} I\left(\max_{1 \leq i \leq n} x_i \leq \theta\right)$$

$\underbrace{\max_{1 \leq i \leq n} x_i}_{X_{(n)}}$

If  $T(x) = T(y)$ , then  $p(\underline{x}; \theta) = 1 \times p(\underline{y}; \theta)$ . The ratio between the two distributions does not depend on  $\theta$ , so  $T$  is sufficient.

Conversely, if  $x, y > 0$ , i.e.  $x_i, y_i > 0$  for  $i=1, \dots, n$ , are supported by the same  $\theta$ 's, then

$$\{\theta \text{ supporting } x\} = (T(x), \infty) = (T(y), \infty) = \{\theta \text{ supporting } y\}$$

Therefore,  $T(x) = T(y)$  and  $\underline{x}$  is a min. suff. statistic.

II.

## Mid term test

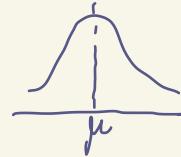
- \* 5-6 questions (shorter than long asg problems)
- \* Not open-book / -notes exam. ; have 1 double-sided, A4 size info sheet.
- \* 2 hr exam..

final

## Lecture 5

3 OCTOBER 2022,

Ancillarity and Completeness.



(Example) Consider  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{CauchyLoc}(\theta)$ , whose distribution is given by

$$p(x; \theta) = \frac{1}{\pi} \cdot \frac{1}{1+(x-\theta)^2} = f(x-\theta)$$

(see TPE.1.5: we can see that  $(X_{(1)}, X_{(2)}, \dots, X_{(n)})$  is minimal sufficient)

(This is also true for the double exponential location model:  $p(x; \theta) \propto \underline{\exp(|x-\theta|)}$ ,  $\downarrow \rightarrow$  Kou (2002) MS. jump diffusion.)

Question: What determines the level of compressibility?

Definition A statistic  $A$  is ancillary for  $X \sim P_\theta \in \mathcal{P}$  if the distribution of  $A(X)$  does not depend on  $\theta$ .

(Example) Consider again  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{CauchyLoc}(\theta)$ .

Define  $A(X) = \underline{X_{(n)} - X_{(1)}}^{\underline{(Z_{(n)} + \theta)} \underline{(Z_{(1)} + \theta)}}$ , we can see that  $A(X)$  is ancillary w.r.t.  $\theta$ .

Let  $X_i = Z_i + \underline{\theta}$ , for  $Z_i \stackrel{iid}{\sim} \text{CauchyLoc}(0)$

So  $X_{(i)} = Z_{(i)} + \theta$ ,  $A(X) = A(Z)$

Hence  $A(X)$  does not depend on  $\theta$ . ||

! The statistic that we make use of should include as little ancillary information as possible.

Definition A statistic  $A$  is first-order ancillary for  $X \sim P_\theta \in \mathcal{P}$  if  $E_\theta(A(X))$  does not depend on  $\theta$ .

Definition A statistic  $T$  is complete for  $X \sim P_\theta \in \mathcal{P}$  if no non-constant function of  $T$  is first-order ancillary. In other words, if  $E_\theta(f(T(X))) = 0$  for all  $\theta$ , then  $f(T(X)) = 0$  with prob. 1 for all  $\theta$ .

Completeness formalises our ideal concept of optimal data reduction. Minimal sufficiency is our achievable notion of optimal data reduction.

Properties:

- 1) If  $T$  is complete sufficient, then  $T$  is minimal sufficient. This is known as Bahadur's theorem.
- 2) Complete sufficient statistics yield optimal unbiased estimator. (More results later)...

$$\uparrow E_\theta(\cdot)$$

(Example)

Let  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$ ,  $\theta \in (0, 1)$ . Then we know that  $T(X) = \sum_{i=1}^n X_i$  is sufficient. Suppose  $E_\theta(f(T(X))) = 0$  for all  $\theta \in (0, 1)$ .

Observe:

$$E_\theta(f(T(X))) = \sum_{j=0}^n f(j) \binom{n}{j} \theta^j (1-\theta)^{n-j} = 0 \quad \forall \theta \in (0, 1) \quad (*)$$

Dividing through by  $\theta^n$  and let  $\beta = \frac{\theta}{1-\theta}$ , we have

$$\sum_{j=0}^n f(j) \binom{n}{j} \beta^j = 0 \quad \leftarrow, \quad \forall \beta > 0.$$

If  $f$  are non-zero, then the LHS is a polynomial of degree at most  $n$ . However, an  $n$ th degree polynomial can have at most  $n$  roots. Hence, it is impossible for the LHS to equal 0 for EVERY  $\beta > 0$  unless  $f = 0$ . We can conclude that  $T$  is complete.

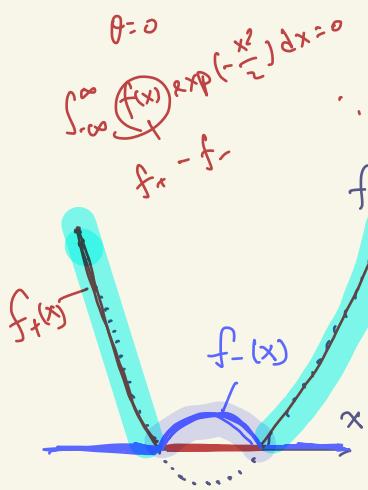
(Example) Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, \sigma^2)$ , with  $\theta$  unknown and a known  $\sigma^2 > 0$ . Is  $\bar{X}_n = \frac{\sum_{i=1}^n X_i}{n}$  complete for this model?

min. suff.

To simplify the notation, we consider a special case where  $n = 1$  and  $\sigma = 1$ , so that  $T(X) = X \sim N(\theta, 1)$ .

Suppose

$$E_\theta(f(x)) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(x) \exp\left(-\frac{(x-\theta)^2}{2}\right) dx = 0 \quad \forall \theta \in \mathbb{R}. \quad (\#*)$$



Multiplying both sides by  $\sqrt{2\pi} e^{\theta^2/2}$ , we can rewrite  $(\#*)$  as

$$\int_{-\infty}^{\infty} \underline{f(x)} \exp\left(-\frac{x^2}{2}\right) \exp(\theta x) dx = 0 \quad \forall \theta \in \mathbb{R}. \quad (\#\#)$$

We decompose  $f$  into its positive and negative parts as

$$f(x) = \underline{[f_+(x) - f_-(x)]}$$

where  $f_+(x) = \max(f(x), 0)$  and  $f_-(x) = \max(-f(x), 0)$ .

Observe that  $f_+(x) \geq 0$  and  $f_-(x) \geq 0$  for all  $x \in \mathbb{R}$  and  $f_+(x) = f_-(x)$  if and only if  $f_+(x) = f_-(x) = 0$ .

If  $f(x) \geq 0$  a.e. or  $f(x) \leq 0$  a.e. then  $(\#\#)$  implies that  $f(x) = 0$  a.e. because setting  $\theta = 0$  gives us an integral of a non-neg (or a non-pos) function being zero.

$$m_x(t) = E(e^{tx}) = \int e^{tx} \cdot f(x) dx$$

We may write

$$M_f(t) = M_{f_-}(t)$$

$$\begin{aligned} M_f(t) &= \int_{-\infty}^{\infty} f(x) e^{tx} dx \\ &= \int_{-\infty}^{\infty} f_+(x) e^{-\frac{x^2}{2}} e^{\theta x} dx \\ &= \int_{-\infty}^{\infty} f_+(x) \exp\left(-\frac{x^2}{2}\right) dx \end{aligned}$$

$$\frac{\int_{-\infty}^{\infty} f_-(x) e^{-\frac{x^2}{2}} e^{\theta x} dx}{\int_{-\infty}^{\infty} f_-(x) \exp\left(-\frac{x^2}{2}\right) dx} \quad (\#)$$

Notice that in  $(\#)$ ,

$$\frac{f_-(x) \exp\left(-\frac{x^2}{2}\right)}{\int_{-\infty}^{\infty} f_-(x) \exp\left(-\frac{x^2}{2}\right) dx} \quad \text{defines a legitimate}$$

probability density. Notice also that  $\text{exp}(-x^2/2)$  is the mgf of this density. Similarly, the RHS is the mgf of the density  $f_-(x) \exp(-x^2/2) / \int_{-\infty}^{\infty} f_-(x) e^{-x^2/2} dx$ .

It implies that  $f_+(x) = f_-(x)$  a.e or, in other words,  $f_+(x) = f_-(x) = 0$  a.e., i.e.  $f(x) = 0$  a.s. Hence T is sufficient and complete. ||

(Example) / Exercise.

If  $X \sim p(x; \theta) \propto h(x) \exp(\theta x)$ , then the statistic  $T(X) = X$  is complete.

Key steps:

- 1). Suppose  $\int f(x) h(x) \exp(\theta x) dx = 0 \forall \theta$ .
- 2) Decompose  $f(x) = f_+(x) - f_-(x)$  with  $f_+, f_- \geq 0$
- 3)  $f_+$  and  $f_-$  can be viewed as legitimate densities
- 4) The decomposition shows that the mgf's of  $f_+$  and  $f_-$  coincide and hence  $f(T) \equiv 0$ .

[THEOREM].  $(T_1, \dots, T_s)$  is complete for any  $s$ -dimensional full rank exponential family. (TSR)

[THEOREM] (Basu's Theorem). If T is complete and sufficient for  $P = \{P_\theta : \theta \in \Sigma\}$  and V is ancillary, then  $T(X) \perp\!\!\!\perp V(X)$ .

• Proof. Define  $g_A(t) = P_\theta(V \in A | T=t)$

$$\text{or } g_A(T) = P_\theta(V \in A | T)$$

$$\text{and } p_A = P_\theta(V \in A).$$

$\text{S} \subseteq \perp\!\!\!\perp \theta$

$\frac{1}{X}$

By sufficiency and ancillarity, neither  $p_A$  nor  $g_A(t)$  depends on  $\theta$ . By smoothing / tower expectation, we can write

$$p_A = P_\theta(V \in A) = E_\theta(P_\theta(V \in A | T)) = E_\theta(g_A(T)).$$

By completeness,

$$E_\theta(g_A(T) - p_A) \stackrel{\Delta}{=} E_\theta(f(\underline{T})) = 0$$

Hence  $f(T) = g_A(T) - p_A = 0$ , i.e.  $\underline{g_A(T)} = p_A$  a.s. P

Again, by smoothing / tower expectation,

$$\begin{aligned} P_\theta(T \in B, V \in A) &= E_\theta(1_B(T) 1_A(V)) \quad 1_B(T) = I(T \in B). \\ &= E_\theta(E_\theta(1_B(T) 1_A(V) | \underline{T})) \\ &= E_\theta(1_B(T) \underline{E_\theta(1_A(V) | T)}) \\ &= E_\theta(1_B(T) \underline{g_A(T)}) \\ &= E_\theta(1_B(T) \cdot p_A) \\ &= p_A \cdot E_\theta(1_B(T)) \\ &= P_\theta(V \in A) P_\theta(T \in B). \end{aligned}$$

Hence, T and V are independent as A and B are arbitrary Borel sets.

(Example) Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , where both  $\mu$  and  $\sigma^2$  are unknown. Then  $\bar{X}_n \perp\!\!\!\perp \sum_{i=1}^n (X_i - \bar{X}_n)^2$  and  $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$ .

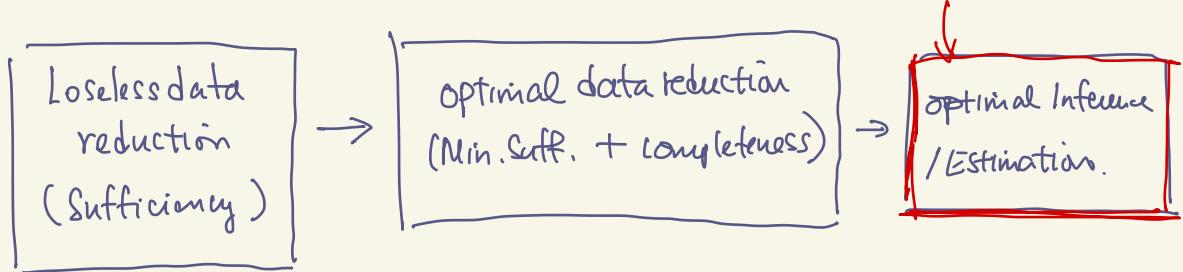
For any  $\sigma > 0$ , we consider the submodel  $P_\sigma = \{N(\mu, \sigma^2) : \mu \in \mathbb{R}\}$

In each submodel,  $\bar{X}_n$  is complete and sufficient, and  $n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  is ancillary. (Let  $X_i = Z_i + \mu$ , then  $X_i - \bar{X}_n = Z_i - \bar{Z}_n$ , which does not depend on  $\mu$ .)

By Basu's theorem,  $\bar{X}_n \perp\!\!\!\perp \sum_{i=1}^n (X_i - \bar{X}_n)^2$  under  $N(\mu, \sigma^2)$  for any  $\mu$ . Since  $\sigma$  is arbitrary, we can conclude that

$\bar{X}_n \perp\!\!\!\perp n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  holds for the full model  $P = \{N(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0\}$ .

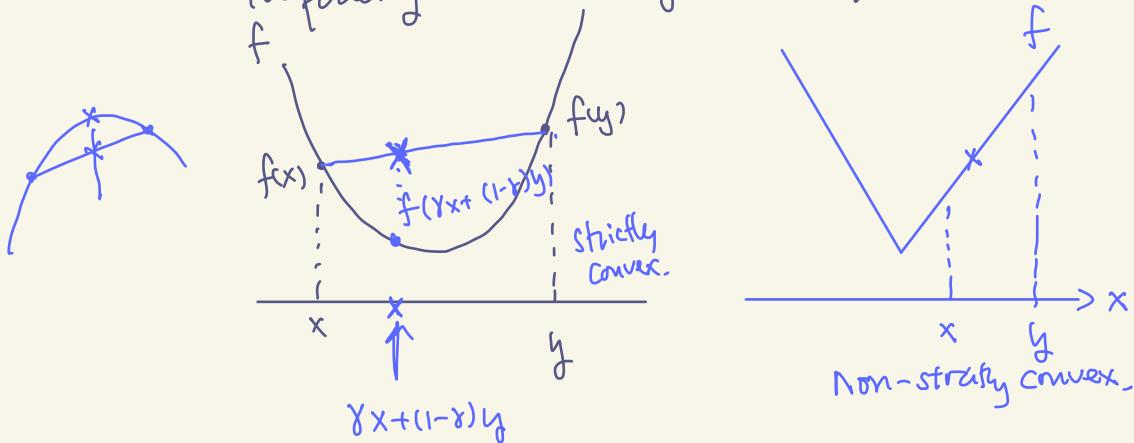
- From data reduction to optimal inference.



Definition A function  $f \cdot$  is a convex function if  $x \neq y$  and  $\gamma \in (0, 1)$

$$f(\gamma x + (1-\gamma)y) \leq \gamma f(x) + (1-\gamma)f(y)$$

The function  $f$  is said to be strictly convex if the above inequality holds strictly (i.e. ' $<$ ' )



(Examples) a) For any  $\theta \in \Omega$ , the function  $f(d) = (d-\theta)^2$  is strictly convex on  $\mathbb{R}$ .

b) For any  $\theta \in \Omega$ , the function  $\bar{f}(d) = |d-\theta|$  is convex but not strictly convex.

[THEOREM] (Jensen's inequality). For  $\stackrel{\text{a}}{\text{convex}} f$ , and  $E(X)$  exists, then

$$f(E(X)) \leq E(f(X)).$$

- If  $f$  is strictly convex, then the above inequality holds strictly unless  $X = E(X)$  with probability 1.

[THEOREM] (Rao-Blackwell Theorem).

Suppose  $T$  is sufficient for  $P = \{P_\theta : \theta \in \Omega\}$ , that  $\delta(X)$  is an estimator for  $g(\theta)$  for which  $E(\delta(X))$  exists, and that  $R(\theta, \delta) = E_\theta(L(\theta, \delta(X))) < \infty$ . If, in particular,  $L(\theta, \cdot)$  is convex, then

$$\rightarrow R(\theta, \eta) \leq R(\theta, \underline{\delta})$$

for  $\boxed{\eta = E(\delta(X) | T(X))}$ .

If  $L(\theta, \cdot)$  is strictly convex, then  $R(\theta, \eta) < R(\theta, \underline{\delta})$  for any  $\theta$  unless  $\eta(T(X)) = \underline{\delta}(X)$  with prob. 1.

(Example). If  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$ , where  $\theta \in (0, 1)$ ,

Consider  $L(\theta, d) = (\theta - d)^2$ . Suppose we have a naive estimator  $\delta(X) = X_1$ . We know that  $T(X) = \bar{X}_n$  is sufficient. So, we can apply the Rao-Blackwell Theorem to improve our estimator  $\delta$ :

$$\begin{aligned}\eta(T(X)) &= E(\delta(X) | T(X)) \\ &\stackrel{x_1}{=} \underbrace{E(\delta(X) | \bar{X}_n)}_{\frac{1}{n} \sum_{i=1}^n} \\ &= E(\bar{X}_n | \bar{X}_n) = \bar{X}_n.\end{aligned}$$

Recall that

$$R(\theta, \eta) = \frac{\theta(1-\theta)}{n} < \underline{\theta(1-\theta)} = R(\theta, \underline{\delta}).$$

Next lecture:

$\boxed{\text{UMVUE}}$

11

## Lecture 6

10/ October / 2022.

### Unbiased estimation

Definition An estimator is said to be unbiased if  $E_\theta(\delta(X)) = g(\theta)$  for all  $\theta$ .

? uniformly best estimator (challenging)

Bu. we can find an unbiased estimator that gives minimum risk uniformly, i.e.  $R(\theta, \delta) \leq R(\theta, \delta')$  for all  $\theta \in \Omega$  and any other unbiased estimator.

Such an estimator is called a uniformly minimum risk unbiased estimator (UMRUE).

In particular, if we adopt  $L(\theta, d) = (\theta - d)^2$ , an UMRUE will become a UMVUE, because

$$E_\theta((g(\theta) - \delta(X))^2) = \underbrace{\{E_\theta(\delta(X)) - g(\theta)\}^2}_{\text{Bias}^2} + \underbrace{E_\theta(\{\delta(X) - E_\theta(\delta(X))\}^2)}_{\text{Variance of } \delta(X)}$$

↑ Variance  
MSE.  
↓ Bias

and if  $\delta(X)$  is unbiased, then the MSE will be reduced to

$$E_\theta((g(\theta) - \delta(X))^2) = E_\theta(\{\delta(X) - E_\theta(\delta(X))\}^2)$$

Definition If an unbiased estimator exists, then  $g(\cdot)$  is called U-estimable.

(Example) Suppose  $X \sim \text{uniform}(0, \theta)$ . Then  $\delta(X)$  is unbiased if

$$\int_0^\theta \frac{\delta(x)}{\theta} dx = g(\theta), \quad \forall \theta > 0$$

$$\text{or if } \int_0^\theta \delta(x) dx = \theta g(\theta), \quad \forall \theta > 0. \quad (*)$$

So,  $g$  cannot be U-estimable unless  $\theta g(\theta) \rightarrow 0$  as  $\theta \downarrow 0$ . If  $g'$  exists, then differentiating  $(*)$ , and also by the fundamental theorem of calculus, we have

$$\delta(x) = \frac{\partial}{\partial x} \{x g(x)\} = g(x) + x g'(x).$$

For example, say  $g(\theta) = 0$ , then  $\delta(X) = X + \underline{X} \cdot 1 = 2X$ .

(Example) Suppose  $X \sim \text{Binomial}(n, \theta)$ . If  $g(\theta) = \sin \theta$ , then  $\delta(X)$  will be unbiased if

$$\sum_{k=0}^n \delta(k) \binom{n}{k} \theta^k (1-\theta)^{n-k} = \sin \theta \quad \forall \theta \in (0,1), \quad (**)$$

The LHS  $(**)$  is a polynomial in  $\theta$  with degree at most  $n$ . The sin function cannot be written as a polynomial of degree  $n$ . Therefore,  $\sin \theta$  is not U-estimable. ||

Definition An unbiased estimator  $\delta$  is uniformly minimum variance unbiased (UMVUE) if  $\text{Var}_\theta(\delta) \leq \text{Var}_\theta(\delta')$   $\forall \theta \in \Omega$  and for any other competing unbiased estimator  $\delta'$ .

[THEOREM]. (Lehmann-Scheffé Theorem). If  $T$  is a complete and sufficient statistic, and  $E_\theta(\tilde{h}(T(X))) = g(\theta)$  [i.e.  $\tilde{h}(T(X))$  is unbiased for  $g(\theta)$ ], then  $\tilde{h}(T(X))$  is

- (a) the only function of  $T(X)$  that is unbiased for  $g(\theta)$ ,
- (b) an UMVUE under any convex loss function,
- (c) the unique UMVUE (hence UMVUE), up to a P-null set, under any strictly convex loss function.

. Proof. (a) Suppose  $E_\theta(\tilde{h}(T(X))) = g(\theta)$ , then

$$E_\theta(\tilde{h}(T(X)) - h(T(X))) = 0 \quad \forall \theta \in \Omega.$$

Thus,  $\tilde{h}(T(X)) = h(T(X))$  almost surely for all  $\theta \in \Omega$  by Completeness.

(b) Consider any unbiased estimator  $\delta(X)$  and let

$$\tilde{h}(T(X)) = E_\theta(\delta(X) | T(X)). \quad \text{Then}$$

$$E_\theta(\tilde{h}(T(X))) = E_\theta(E_\theta(\delta(X) | T(X))) = E_\theta(\delta(X)) = g(\theta)$$

(by "smoothing"). By (a),  $\tilde{h}(T(X)) = h(T(X))$ , then  $R(\theta, \tilde{h}(T(X))) = R(\theta, h(T(X)))$ . By Rao-Blackwell Theorem,

we have  $R(\theta, \tilde{h}(T(X))) \leq R(\theta, \delta)$  for all  $\theta \in \Omega$  if the loss function is convex. It follows that

$$R(\theta, h(T(X))) \leq R(\theta, \delta) \quad \forall \theta \in \Omega.$$

Therefore,  $h(T(X))$  is an UMVUE for any convex loss function.

(c) If the loss function is strictly convex,  $R(\theta, h(T(X))) < R(\theta, \delta)$  unless  $\delta(X) \stackrel{a.s.}{=} h(T(X))$ . Hence,  $h(T(X))$  is the unique UMRUE (resp. UMVUE if the loss function is the squared error loss function).

Possible solution a: Rao-Blackwellization

(Example) (Rao-Blackwellization). Suppose  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$ .

We know that  $T(X) = \sum_{i=1}^n X_i$  is a complete, sufficient statistic. We also know that  $n^{-1} T(X)$  is an unbiased estimator for  $\theta$ , i.e.

$$E_\theta(T(X)/n) = n^{-1} \sum_{i=1}^n E_\theta(X_i) = \theta.$$

Therefore,  $n^{-1} T(X)$  is an UMRUE for  $\theta$  under any convex loss function.

Suppose that, instead, we are interested in estimating  $\theta^2$ . Let's observe

$$\begin{aligned} \delta(X) &= I(X_1 = X_2 = 1) = X_1 X_2 \\ E_\theta(\delta(X)) &= E_\theta(X_1 \cdot X_2) = \{E_\theta(X_1)\}^2 = \theta^2 \\ \rightarrow E_\theta(\delta(X) \mid T(X) = t) &= P_\theta(X_1 = X_2 = 1 \mid T(X) = t) \\ &= \frac{P_\theta(X_1 = X_2 = 1, \sum_{i=3}^n X_i = t-2)}{P_\theta(T(X) = t)} \\ &= \frac{\theta^2 \binom{n-2}{t-2} \theta^{t-2} (1-\theta)^{n-t} I(t \geq 2)}{\binom{n}{t} \theta^t (1-\theta)^{n-t}} \\ &= \frac{t(t-1) I(t \geq 2)}{n(n-1)}. \end{aligned}$$

Hence, we conclude that a UMVUE for  $\theta^2$  is  $\left[ \frac{T(X)I(T(X) \geq 2)}{n(n-1)} \right]$ . //

(Example) Suppose  $X_1, \dots, X_n \stackrel{iid}{\sim} U(0, \theta)$ . In this case  $T(X) = X_{(n)}$  is a complete and sufficient statistic, and  $\delta(X) = 2X$  is an unbiased estimator for  $\theta$ , i.e.  $E_\theta(2X) = 2 \cdot \left(\frac{\theta}{2}\right) = \theta$ .

Given  $X_{(n)}$ ,  $X_1$  is equal to  $X_{(n)}$  with prob.  $1/n$  and follows uniform  $(0, X_{(n)})$  with prob.  $1 - 1/n$ . Hence

$$P_\theta(X_1 = x_1 \mid T(X)) = \frac{I(T(X) = x_1)}{n} + \underbrace{\frac{I(0 < x_1 < T(X))}{T(X)} \left(1 - \frac{1}{n}\right)}_{\text{uniform.}}$$

To find the UMVUE, we calculate

$$\begin{aligned}
 E_\theta(\delta(X) | T(X)) &= 2E_\theta(X_1 | T(X)) \\
 &= 2\left\{\frac{1}{n}T(X) + (1-\frac{1}{n})\int_0^{T(X)} \frac{x_1 dx_1}{T(X)}\right\} \\
 &= 2\left\{\frac{T(X)}{n} + (1-\frac{1}{n})\frac{T(X)}{2}\right\} \\
 &= \left(\frac{n+1}{n}\right)T(X),
 \end{aligned}$$

Which gives the desired result.

Possible sol'n B: solve for  $\delta$  directly

(Example) Let  $X \sim \text{Poisson}(1/\theta)$ .  $X$  is a complete and sufficient statistic for  $\theta$ .  $X$  is also unbiased and therefore UMVUE for  $\theta$ .

Suppose we are interested in estimating  $g(\theta) = e^{-a\theta}$  for  $a \in \mathbb{R}$ .

We need to find an estimator  $\delta$  such that  $E_\theta(\delta(X)) = g(\theta)$  for all  $\theta > 0$ . Under this model, we have

$$\begin{aligned}
 E_\theta(\delta(X)) &= \sum_{x=0}^{\infty} \delta(x) \cdot \frac{e^{-\theta} \theta^x}{x!} = e^{-a\theta}, \quad \forall \theta > 0 \\
 \Leftrightarrow \quad \sum_{x=0}^{\infty} \frac{\delta(x) \theta^x}{x!} &= e^{(1-a)\theta} = \sum_{x=0}^{\infty} \frac{(1-a)^x \theta^x}{x!} \\
 \Rightarrow \quad \delta(X) &= (1-a)^X \text{ is the UMVUE for } g(\theta).
 \end{aligned}$$

Note that this estimator is not ideal in the sense that if  $a=2$ , then the estimator  $\delta(X) = (-1)^X$  will change sign according to the "evenness" of  $X$ .

The estimator is hence inadmissible when  $a > 1$  and will be dominated by  $\max(\delta(X), 0)$ .

Possible sol'n C: Guess

(Example) Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ . Consider the case where  $\theta = (\mu, \sigma^2)$  is unknown.

(a) The UMVUE for  $\sigma^2 = (n-1)^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = S^2$ , [Try].

(b) How about the UMVUE for  $\sigma$ ?

(c) What is the UMVUE for  $\mu$ ?

(b) Observe:

$$X_i - \bar{X}_n \sim N\left(0, \frac{n-1}{n} \sigma^2\right) \quad \checkmark$$

$$E(|X_i - \bar{X}_n|) = \sigma \sqrt{\frac{2}{\pi}} \times \sqrt{\frac{n-1}{n}}. \quad (\text{Verify})$$

$$\text{This implies: } \delta' = \frac{\sqrt{\pi n}}{\sqrt{2(n-1)}} |X_i - \bar{X}_n| \quad \checkmark$$

is unbiased for  $\sigma$ . [Rao-Blackwellisation may not work so easily here..]

Instead, we can observe another fact that

$$S_*^2 = \sum_{i=1}^n (X_i - \bar{X}_n)^2 = (n-1) S^2 \sim \sigma^2 \chi_{n-1}^2$$

Hence,

$$E(S_*) = \sigma E(\chi_{n-1}) \Rightarrow \sigma = \frac{E(S_*)}{E(\chi_{n-1})},$$

meaning that  $\frac{S_*}{E(\chi_{n-1})}$  is unbiased for  $\sigma$  and hence UMVUE.

(c) Take the expectation of the UMVUE for  $\mu$  and squaring it,

$$\text{we obtain } E(\bar{X}_n^2) = \mu^2 + \frac{\sigma^2}{n}$$

↑

and so,  $\delta_n(X) = \bar{X}_n^2 - \underbrace{\frac{S_*^2}{n(n-1)}}_{\text{negative...}}$  is UMVUE for  $\mu^2$ .

$$\max(0, \delta_n(X))$$

H

### Fisher Information

Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

↓ unbiased

Define  $S^2 = \sum_{i=1}^n (X_i - \bar{X})^2$ ,  $\frac{S^2}{n-1}$  is the UMVUE for  $\sigma^2$

MLE for  $\sigma^2$  is  $\frac{S^2}{n}$

{ neither MLE nor UMVUE is admissible }

Shrink estimator:  $\frac{S^2}{n+1}$ . TOWER R&E.

↓ [James-Stein estimator...]

Rabinowitz (?)  
2000  
Sinica.  
Semiparametric efficiency

Suppose we have  $\delta_1$  and  $\delta_2$  as UMVUEs for  $g_1(\theta)$  and  $g_2(\theta)$ , respectively.

Is  $\delta_1 + \delta_2$  an UMVUE for  $g_1(\theta) + g_2(\theta)$ ?

[THEOREM] (TPE 2.1.7) (Characterization of UMVUEs).

Let  $\Delta = \{ \delta : E_\theta(\delta^2) < \infty \}$ . Then  $\underline{\delta}_0 \in \Delta \Leftrightarrow$  UMVU for  $g(\theta) = E(\delta_0)$

If and only if  $E(\delta(\theta) U) = 0$  for every  $U \in \mathcal{U} = \{E(U) = 0\}$ .

unbiased estimators of  $\theta$ .

Proof. If  $\delta_0$  is an UMVUE, let's consider  $\delta_\lambda = \delta_0 + \lambda U$  for a fixed  $\lambda \in \mathbb{R}$  and  $U \in \mathcal{U}$ . Since  $\delta_0$  has minimal variance,

$$\text{Var}_\theta(\delta_\lambda) = \text{Var}_\theta(\delta_0) + \lambda^2 \text{Var}_\theta(U) + 2\lambda \text{Cov}_\theta(\delta_0, U) \geq \underline{\text{Var}}_\theta(\delta_0),$$

or equivalently,

$$\lambda^2 \text{Var}_\theta(U) + 2\lambda \text{Cov}(\delta_0, U) \geq 0.$$

Consider the quadratic form  $g(\lambda) = \lambda^2 \text{Var}_\theta(U) + 2\lambda \text{Cov}(\delta_0, U)$ .

The form  $g$  has the roots

$$\lambda = 0 \text{ and } \lambda = \frac{-2\text{Cov}(\delta_0, U)}{\text{Var}_\theta(U)}.$$

If the roots are distinct, the form must be negative at some point, which violates the inequality above. Hence  $-2\text{Cov}_\theta(\delta_0, U)/\text{Var}_\theta(U) = 0$  in which case  $E_\theta(U\delta_0) = \underline{\text{Cov}}(\delta_0, U) = 0$ .

We assume that  $E_\theta(U\delta_0) = 0$  for all  $U \in \mathcal{U}$ , and consider any  $\delta$  unbiased for  $g(\theta)$ .

It follows that  $\delta - \delta_0 \in \mathcal{U}$ . So  $E_\theta(\delta_0(\delta - \delta_0)) = 0$ . This implies that  $E_\theta(\delta_0 \delta) = E_\theta(\delta_0^2)$  and subtracting  $\{E_\theta(\delta_0)\}^2 (= E_\theta(\delta_0) E_\theta(\delta))$  on both sides, we obtain

$$\text{Var}_\theta(\delta_0) = \text{Cov}_\theta(\delta_0, \delta) \stackrel{\text{CSI}}{\leq} \sqrt{\text{Var}_\theta(\delta_0) \text{Var}_\theta(\delta)}$$

Hence  $\underline{\text{Var}}_\theta(\delta_0) \leq \text{Var}_\theta(\delta)$  for any arbitrary unbiased estimator  $\delta$ . Hence  $\delta_0$  is a UMVUE for  $g(\theta)$ .

$$\underline{E_\theta \left\{ (\delta_1 + \delta_2) u \right\}} = \underline{E_\theta (\delta_1 u)} + \underline{E (\delta_2 u)} = 0 \quad \forall u \in \mathcal{U}$$

$$\Rightarrow \delta_1 + \delta_2 \text{ is a UMVUE of } g_1(\theta) + g_2(\theta).$$

### Variance Bounds and Information

$$\underline{\text{Var}_\theta(\delta)} \geq \underline{\text{LB}}.$$

Recall from the Cauchy-Schwarz inequality,

$$\text{Cov}(X, Y) \leq \sqrt{\text{Var}(X)\text{Var}(Y)} \quad \text{or} \quad |\text{Cov}(X, Y)| \leq \sigma_X \sigma_Y.$$

If  $\delta$  is an unbiased estimator for  $g(\theta)$  and  $\psi$  is an arbitrary random variable, then

$$\underline{\text{Var}_\theta(\delta)} \geq \frac{\text{Cov}_\theta(\delta, \psi)}{\text{Var}_\theta(\psi)}.$$

Suitable  $\psi$  so that the  $\text{Cov}_\theta(\delta, \psi)$  is the same for all  $\delta$  that are unbiased for  $g(\theta)$ .  
(above)

Let  $P = \{P_\theta : \theta \in \Omega\}$  be a dominated family with densities  $f_\theta : \theta \in \Omega = \mathbb{R}$ .

To begin,  $E_{\theta+\Delta}(\delta) - E_\theta(\delta)$  gives  $g(\theta+\Delta) - g(\theta)$  for any unbiased  $\delta$ .

Here  $\Delta$  must be chosen such that  $\Delta + \theta \in \Omega$ . Next, we write

$E_{\theta+\Delta}(\delta) - E_\theta(\delta)$  as a covariance under  $f_\theta$ . This step includes the use of "likelihood ratio." We assume here that  $f_{\theta+\Delta}(x) = 0$  whenever  $f_\theta(x) = 0$ . Define

$$L(x) = \begin{cases} \frac{f_{\theta+\Delta}(x)}{f_\theta(x)}, & f_\theta(x) > 0 \\ 0, & f_\theta(x) = 0 \end{cases}$$

Let  $\underline{f_{\theta+\Delta}(x)} = \frac{f_{\theta+\Delta}(x)}{f_\theta(x)} \times f_\theta(x)$  as.  $x$ .

and so for any function  $h$  integrable under  $P_{\theta+\Delta}$ ,

$$E_{\theta+\Delta}(h(x)) = \int h \underline{f_{\theta+\Delta}} d\mu = \int h L f_\theta d\mu = E_\theta(L(x) h(x)).$$

Take  $h = 1$ , then,

$$E_\theta(L(x)) = \int \frac{f_{\theta+\Delta}(x)}{f_\theta(x)} f_\theta(x) dx = \int f_{\theta+\Delta}(x) dx = 1$$

Take  $h = \delta$ , then

$$E_{\theta+\Delta}(\delta) = E_\theta(L(x)\delta)$$

So, if we define  $\psi(x) = \underline{L(x) - 1}$ ,

$$\text{then we can see that } E_\theta(\psi(x)) = E_\theta(\underline{L(x)-1}) = 0$$

ch. of measure /  
importance  
sampling

and

$$E_{\theta+\Delta}(\delta) - E_\theta(\delta) \approx E_\theta(\underline{\delta}) - E_\theta(\delta) = E_\theta(\underline{\psi}\delta) = \underline{\text{Cov}_\theta(\delta, \psi)}.$$

$$\text{As a result, } \text{Cov}_\theta(\delta, \psi) = g(\theta + \Delta) - g(\theta)$$

for any unbiased estimator  $\delta$ . With this particular choice of  $\psi$ , the above inequality can be rewritten as

$$\text{Var}_\theta(\delta) \geq \frac{\text{Cov}_\theta^2(\delta, \psi)}{\text{Var}_\theta(\psi)} = \frac{(g(\theta + \Delta) - g(\theta))^2}{\text{Var}_\theta(\psi)} = \frac{(g(\theta + \Delta) - g(\theta))^2}{E_\theta \left( \left( \frac{f_{\theta+\Delta}(x)}{f_\theta(x)} - 1 \right)^2 \right)}$$

Hammersley - Chapman - Robinson bounds

Under suitable regularity conditions, we can show that

$$\frac{\left( \frac{g(\theta + \Delta) - g(\theta)}{\Delta} \right)^2}{E_\theta \left( \left( \frac{\{f_{\theta+\Delta}(x) - f_\theta(x)\}/\Delta}{f_\theta(x)} \right)^2 \right)} \xrightarrow{\Delta \rightarrow 0} \left[ \frac{(g'(\theta))^2}{E_\theta \left( \left( \frac{\partial f_\theta(x)/\partial \theta}{f_\theta(x)} \right)^2 \right)} \right]$$

(To be continued...)

Fisher information  
 $I(\theta) = E_\theta \left( \left( \frac{\partial \log f_\theta(x)}{\partial \theta} \right)^2 \right)$

Next lecture: { ① Finish the remaining discussion on FI.  
 ② Bayes estimator. and avg risk optimality (TPE Ch.4).

31/Oct

Lecture 1 - Fisher Information (Lecture 6 - FRB)

FINAL EXAM

Final date

Last teaching day : 28/Nov.

↓  
 5/Dec

OR 12/Dec

Pending.

## Lecture 7

17 Oct / 2022.

Agenda: (1) Fisher Information  
(2) Bayes Rule --

- ~~Equivariance (Ch. 3 of FPE)~~ --

### [Cont'd] Fisher Information / C&R Lower Bound

$$\text{Define } I(\theta) = E_{\theta} \left( \left( \frac{\partial \log f_{\theta}(x)}{\partial \theta} \right)^2 \right),$$

where  $\log f_{\theta}(x)$  is the log-likelihood function, as the Fisher information.

$$0 = \frac{\partial}{\partial \theta} (1) = \frac{\partial}{\partial \theta} \int f_{\theta}(x) d\mu(x) = \int \underbrace{\frac{\partial}{\partial \theta} f_{\theta}(x)}_{(dx)} d\mu(x) = \int \frac{\partial \log f_{\theta}(x)}{\partial \theta} f_{\theta}(x) d\mu(x)$$

$$= E \left( \frac{\partial \log f_{\theta}(x)}{\partial \theta} \right)$$

$$\text{and so, } I(\theta) = E_{\theta} \left( \left( \frac{\partial \log f_{\theta}(x)}{\partial \theta} \right)^2 \right) - \underbrace{\left\{ E_{\theta} \left( \frac{\partial \log f_{\theta}(x)}{\partial \theta} \right) \right\}}_0^2$$

$$= \text{Var}_{\theta} \left( \frac{\partial \log f_{\theta}(x)}{\partial \theta} \right).$$

Furthermore, since

$$\int \frac{\partial^2 f_{\theta}(x)}{\partial \theta^2} d\mu(x) = E \left( \frac{\partial^2 f_{\theta}(x)/\partial \theta^2}{f_{\theta}(x)} \right) = 0$$

We can see that

$$\frac{\partial^2 \log f_{\theta}(x)}{\partial \theta^2} = \frac{\partial^2 f_{\theta}(x)/\partial \theta^2}{f_{\theta}(x)} - \left( \frac{\partial \log f_{\theta}(x)}{\partial \theta} \right)^2$$

$$\Rightarrow I(\theta) = -E \left( \frac{\partial^2 \log f_{\theta}(x)}{\partial \theta^2} \right).$$

[THEOREM]. Let  $P = \{P_{\theta} : \theta \in \Omega\}$  be a dominated family with  $\Omega$  an open set in  $\mathbb{R}$  and densities  $f_{\theta}$  differentiable with respect to  $\theta$ . If  $E_{\theta}(\psi) = 0$  and  $E_{\theta}(\delta^2) < \infty$ , then

$$\text{Var}_{\theta}(\delta) \geq \frac{\{g'(\theta)\}^2}{I(\theta)}, \quad \theta \in \Omega.$$

This is called the Cramér-Rao lower bound, or information bound.

(Example) Let  $P$  be a one-parameter exponential family in canonical form with densities  $f_\theta$  given by

$$f_\theta(x) = \exp\{\eta T(x) - A(\eta)\} h(x).$$

Then

$$\frac{\partial \log f_\theta(x)}{\partial \eta} = T(x) - A'(\eta).$$

By the previous results we have

$$I(\eta) = \text{Var}_\eta(T(X) - A'(\eta)) = \text{Var}_\eta(T(X)) = A''(\eta)$$

because

$$\frac{\partial^2 \log f_\theta(x)}{\partial \eta^2} = \frac{\partial^2}{\partial \eta^2} \left( \eta T(x) - A(\eta) + \log h(x) \right) = -A''(\eta).$$

If the family is parameterised instead by  $\mu = A'(\eta) = E_\eta(T(X))$ , then  $A''(\eta) = I(\mu) \{A''(\eta)\}^2$ .

And, so, because  $A''(\eta) = \text{Var}(T)$ , we have  $I(\mu) = \frac{1}{\text{Var}_\eta(T)}$

Observe also that because  $T$  is UMVUE of  $\mu$ , the lower bound

$\text{Var}(\hat{\delta}) \geq \frac{1}{I(\mu)}$  for an unbiased estimator of  $\delta$  if  $\mu$  is sharp.

(Example) Suppose  $E$  is an absolutely continuous random variable with density  $f$ . The family of distributions  $P = \{P_\theta : \theta \in \mathbb{R}\}$  with  $f_\theta$  the distribution of  $\theta + E$  is called a location family.

$$\begin{aligned} \int g(x) dP_\theta(x) &= E_\theta(g(X)) \stackrel{X=\theta+E}{=} E_\theta(g(\theta+E)) \\ &= \int g(\theta+e) f(e) de = \int g(x) f(x-\theta) dx. \end{aligned}$$

So,  $P_\theta$  has density  $f_\theta(x) = f(x-\theta)$ . The corresponding Fisher information for this family is

$$I(\theta) = E_\theta \left( \left( \frac{\partial \log f(X-\theta)}{\partial \theta} \right)^2 \right) = E_\theta \left( \left( -\frac{f'(X-\theta)}{f(X-\theta)} \right)^2 \right)$$

$$= E \left( \left( \frac{f'(\theta)}{f(\theta)} \right)^2 \right) = \int \frac{(f'(x))^2}{f(x)} dx \perp\!\!\!\perp \theta$$

So, for the location families,  $I(\theta)$  is constant with respect to  $\theta$ . ||

If two (or more) independent vectors are observed, then the total Fisher information is the sum of the Fisher information provided by the individual observations:

Suppose  $X$  and  $Y$  are independent,  $X$  has density  $f_\theta$  and  $Y$  has density  $g_\theta$ .

$$\text{The Fisher information from } X \text{ is } I_X(\theta) = \text{Var}_\theta \left( \frac{\partial \log f_\theta(X)}{\partial \theta} \right)$$

$$I_Y(\theta) = \text{Var}_\theta \left( \frac{\partial \log g_\theta(Y)}{\partial \theta} \right)$$

$$I_{(X,Y)}(\theta) = \text{Var}_\theta \left( \frac{\partial \log (f_\theta(X) g_\theta(Y))}{\partial \theta} \right)$$

$$= \text{Var}_\theta \left( \frac{\partial \log f_\theta(X)}{\partial \theta} + \frac{\partial \log g_\theta(Y)}{\partial \theta} \right)$$

$$= \text{Var}_\theta \left( \frac{\partial \log f_\theta(X)}{\partial \theta} \right) + \text{Var}_\theta \left( \frac{\partial \log g_\theta(Y)}{\partial \theta} \right)$$

$$= I_X(\theta) + I_Y(\theta).$$

More generally, for  $X_i \stackrel{iid}{\sim} f_\theta$  ( $i=1, \dots, n$ ), then

$$I_{\tilde{X}}(\theta) = \sum_{i=1}^n I_{X_i}(\theta) = n I_1(\theta).$$

$$\text{Var}_\theta(\delta) \geq \frac{f'g'(\theta)^2}{n I(\theta)}.$$

How about  $\theta \in \Omega \subset \mathbb{R}^k$ ?

Suppose  $\theta$  takes values in  $\mathbb{R}^k$ , then the Fisher information will be a matrix, defined in regular case by

$$\begin{aligned} \{I(\theta)\}_{ij} &= E_\theta \left( \frac{\partial \log f_\theta(X)}{\partial \theta_i} \cdot \frac{\partial \log f_\theta(X)}{\partial \theta_j} \right) \\ &= \text{Cov}_\theta \left( \frac{\partial \log f_\theta(X)}{\partial \theta_i}, \frac{\partial \log f_\theta(X)}{\partial \theta_j} \right) \\ &= -E_\theta \left( \frac{\partial^2 \log f_\theta(X)}{\partial \theta_i \partial \theta_j} \right). \end{aligned}$$

Note that

$$E_{\theta} (\nabla_{\theta} \log f_{\theta}(x)) = 0,$$

$$\begin{aligned} I(\theta) &= E_{\theta} \left( \{\nabla_{\theta} \log f_{\theta}(x)\} \{\nabla_{\theta} \log f_{\theta}(x)\}^T \right) \\ &= \text{Cov}_{\theta} (\nabla_{\theta} \log f_{\theta}(x)) = -E_{\theta} (\nabla_{\theta}^2 \log f_{\theta}(x)). \end{aligned}$$

where  $\nabla_{\theta}$  is the gradient with respect to  $\theta$  and  $\nabla_{\theta}^2$  is the Hessian matrix of second order derivatives, i.e.

$$\nabla_{\theta} = \begin{pmatrix} \frac{\partial}{\partial \theta_1} & \cdots & \frac{\partial}{\partial \theta_K} \end{pmatrix}^T \quad \text{and} \quad \nabla_{\theta}^2 = \begin{pmatrix} \frac{\partial^2 f_{\theta}(x)}{\partial \theta_1^2} & \frac{\partial^2 f_{\theta}(x)}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 f_{\theta}(x)}{\partial \theta_1 \partial \theta_K} \\ \vdots & \ddots & & \\ & & \frac{\partial^2 f_{\theta}(x)}{\partial \theta_K^2} \end{pmatrix}$$

The lower bound for the variance of an unbiased estimator  $\delta$  of  $g(\theta)$ , where  $g: \Omega \rightarrow \mathbb{R}^k$  is

$$\text{Var}_{\theta}(\delta) \geq (\nabla_{\theta} g(\theta))^T I(\theta) (\nabla_{\theta} g(\theta)).$$

## BAYES ESTIMATORS AND AVERAGE RISK OPTIMALITY

Explore an alternative approach to achieve optimality.

Notation: Data:  $X$ ,  $P = \{P_{\theta} \in \Omega\}$ ,  $L(\theta, d)$  loss function,  $R(\theta, \delta) = E_{\theta}(L(\theta, d))$  Risk function.

We need to introduce a measure  $\Lambda$  over  $\Omega$ . This measure  $\Lambda$  can be viewed as an assignment of weights to each parameter value  $\theta \in \Omega$  a priori (prior dist.).

Before any data is observed,  $\theta$  is no longer a constant. It is instead a random variable ( $\mathbb{H}$ ).

Given a measure  $\Lambda$ , our objective is to find an estimator  $\delta_{\Lambda}$  which minimises the average risk:

$$r(\Lambda, \delta) = \int [R(\theta, \delta)] \Lambda(\theta) = E_{\mathbb{H}}(R(\theta, \delta))$$

If  $\Lambda$  is a probability distribution on  $(\mathbb{H})$ , we call  $\Lambda$  the prior distribution.

Correspondingly, the estimator  $\delta_{\Lambda}$ , if exists, is called the Bayes estimator with respect to  $\Lambda$ , and the minimized average risk is called the Bayes risk.

We shall pay attention to  $E(L(\theta, \delta(X)) | X=x)$ , the conditional risk at (almost) every value of  $X$ . Note that the expectation here is taken w.r.t. the conditional distribution of  $\theta$  given  $X$ , i.e.  $\theta | X=x$ .

[THEOREM] Suppose  $\theta \sim \Lambda$  and  $X | \theta=\theta \sim P_\theta$ . If

1. there exists  $\delta_0$ , and estimator of  $g(\theta)$  with finite risk for all  $\theta$ ,
2. there exists a value  $\delta_\Lambda(x)$  that minimises  
 $\rightarrow \bullet E(L(\theta, \delta_\Lambda(x)) | X=x)$  for almost every  $x$ , (\*)  
then  $\delta_\Lambda$  is a Bayes estimator with respect to  $\Lambda$ .

Proof. If (1) and (2) hold, for any other estimator  $\delta'$ , say, and for almost every  $x$ ,

$$\cdot E(L(\theta, \delta_\Lambda(x)) | X=x) \leq E(L(\theta, \delta'(x)) | X=x)$$

Taking expectation over  $X$ , we obtain

$$\rightarrow \underline{E(L(\theta, \delta_\Lambda(x)))} \leq E(L(\theta, \delta'(x))) \quad \forall \delta'$$

Remark: The almost sure statement in (\*) is defined w.r.t. the marginal distribution of  $X$ :  $P(X \in A) = \int P_\theta(X \in A) d\Lambda(\theta)$ . □

(Example) If we consider  $L(\theta, d) = (d-\theta)^2$ , to find the Bayes estimator, we need to minimize  $E(\{g(\theta) - \delta(x)\}^2 | X=x)$

$$= E(\{g(\theta) - E(g(\theta) | X) + E(g(\theta) | X) - \delta(x)\}^2 | X=x)$$

$$= E(\{g(\theta) - E(g(\theta) | X)\}^2 | X=x) \quad \text{+}$$

$$E(\{E(g(\theta) | X) - \delta(x)\}^2 | X=x),$$

or equivalently, minimizing  $E(\{E(g(\theta) | X) - \delta(x)\}^2 | X=x)$ .

The Bayes estimator is then  $\underline{E(g(\theta) | X)}$ .

Posterior mean.

II

Remark: By the Baye's theorem,

$$\text{posterior} = \frac{\text{joint}}{\text{marginal.}} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal}}$$

$$p(\theta | X) = \frac{p(\theta, x)}{\int p(\theta', x) d\theta'} = \frac{\boxed{p(x|\theta)} \times \pi(\theta)}{\int p(\theta' | x) d\theta'}$$

↑ "normalizer"

(Example) Binomial mean estimation.

$E(\theta|X)$   
 $\pi(\theta|x)$   
 Conjugate prior

Suppose  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$  given  $\Theta = \theta$  and that  $\Theta$  has a prior distribution  $\text{Beta}(\alpha, \beta)$ , where  $\alpha$  and  $\beta$  are two fixed prior hyperparameters. The prior density is given by

$$\rightarrow \pi(\theta; \alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} I(0 < \theta < 1) \quad \checkmark$$

The model pmf of  $X = X_1 + \dots + X_n$  is

$$f(x; \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

So the posterior distribution of  $\Theta$  given  $X$  is

$$\begin{aligned} \pi(\theta|x) &\propto \left[ \binom{n}{x} \theta^x (1-\theta)^{n-x} \right] \left[ \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} I(0 < \theta < 1) \right] \\ &\propto \frac{\theta^{(x-\alpha)-1}}{\cdot} \cdot \frac{(1-\theta)^{(n-x+\beta)-1}}{\cdot} I(0 < \theta < 1) / \int_{0}^{1} \theta^{(x-\alpha)-1} (1-\theta)^{(n-x+\beta)-1} d\theta \\ &\sim \text{Beta}(x-\alpha, n-x+\beta). \end{aligned}$$

From textbooks/wiki  
 $E(\text{Beta}(\alpha, \beta)) = \frac{\alpha}{\alpha+\beta}$

We know that the posterior mean of  $\Theta|X=x$  is

$$\begin{aligned} E_{\Theta}(\Theta|X=x) &= \frac{x+\alpha}{(x+\alpha)+(n-x+\beta)} = \frac{x+\alpha}{n+\alpha+\beta} \\ &= \frac{n}{n+\alpha+\beta} \cdot \frac{(x)}{\binom{n}{x}} + \frac{\alpha+\beta}{n+\alpha+\beta} \left( \frac{\alpha}{\alpha+\beta} \right) \\ &\quad \text{"Let the data speak for themselves"} \quad \boxed{\text{Sample mean}} \quad \boxed{\text{"mean of the prior."}} \end{aligned}$$

(Example) Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\Theta, \sigma^2)$  with  $\sigma^2$  known. Furthermore, let  $\Theta \sim N(\mu, b^2)$ , where  $\mu$  and  $b^2$  are two fixed prior hyperparameters. Then the posterior distribution of  $\Theta|X$  is given by

$$\begin{aligned} \pi(\theta|x) &\propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(X_i - \theta)^2}{2\sigma^2}\right) \times \frac{1}{\sqrt{2\pi b^2}} \exp\left(-\frac{(\theta - \mu)^2}{2b^2}\right) \\ &\propto \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \theta)^2 - \frac{1}{2b^2} (\theta - \mu)^2\right\} \\ &= \exp\left\{\frac{1}{\sigma^2} \sum_{i=1}^n X_i \theta - \frac{n\theta^2}{2\sigma^2} - \frac{1}{2b^2} \theta^2 + \frac{\mu\theta}{b^2}\right\} \end{aligned}$$

$$\begin{aligned}
 &= \exp \left\{ -\frac{1}{2} \left( \frac{n}{b} + \frac{1}{b^2} \right) B^2 + \left( \frac{n \bar{x}_n}{\sigma^2} + \frac{\mu}{b^2} \right) B \right\} \\
 &\stackrel{D}{=} \exp \left\{ -\frac{1}{2\sigma^2} (\theta - \tilde{\mu})^2 \right\}, \quad \sim N(\tilde{\mu}, \sigma^2)
 \end{aligned}$$

where  $\hat{\mu} = \frac{n\bar{x}_n/\sigma^2 + \mu/b^2}{n/\sigma^2 + 1/b^2}$  and  $\hat{\sigma}^2 = \frac{1}{n/\sigma^2 + 1/b^2}$ .

Hence, the posterior mean of  $\Theta | X = x$  is

$$F_{\Theta}(\Theta | X=x) = \frac{n\bar{x}_n/\sigma^2 + \mu/b^2}{n/\sigma^2 + 1/b^2} = \frac{\frac{n/\sigma^2}{n/\sigma^2 + 1/b^2} \bar{x}_n + \frac{1/b^2}{n/\sigma^2 + 1/b^2} \mu}{\Phi}$$

Sample mean      prior mean.

Next week:

- 1) How about other loss functions?
  - 2) Relationship b/w  $\hat{X}_n$  and Bayes estimator...
  - 3) Admissibility.
  - 4) Empirical Bayes. -

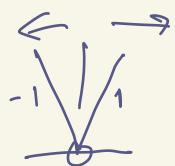
## Lecture 8

24 Oct 2022.

Recall: Average risk optimality:

$$\theta \downarrow \text{fixed constant}$$

$$\Theta \sim \Lambda \text{ prior dist.}$$



Bayes estimator (under the squared error loss function):  $E(g(\Theta) | X)$

Conjugate priors

abs loss function

posterior dist.

$$\text{med}(g(\Theta) | X)$$

$$\frac{\partial}{\partial d} | - | = [\text{Sgn } \underline{\quad}]$$

(Example) Assume that we consider instead  $L(\theta, d) = \omega(\theta)(g(\theta) - d)^2$ , where  $\omega(\theta) \geq 0$ , which can be interpreted as a weight function.

Our goal is to find the corresponding Bayes estimator, which minimises  $\underset{d}{\mathbb{E}} (\omega(\Theta)(g(\Theta) - d)^2 | X=x)$  with respect to  $d$ .

*the quantity*

$\underset{d}{\mathbb{E}}$  can be rewritten as

$$\begin{aligned} & d^2 \mathbb{E}(\omega(\Theta) | X=x) - 2d \mathbb{E}(\omega(\Theta) g(\Theta) | X=x) \\ & + \mathbb{E}(\omega(\Theta) g(\Theta)^2 | X=x) \end{aligned} \quad (*)$$



Taking derivative of (\*) w.r.t.  $d$  gives

$$2d^* \mathbb{E}(\omega(\Theta) | X=x) - 2 \mathbb{E}(\omega(\Theta) g(\Theta) | X=x) = 0$$

$$d^* \left( \stackrel{\cong}{=} \delta_{\Lambda}(x) \right) = \frac{\mathbb{E}(\omega(\Theta) g(\Theta) | X=x)}{\mathbb{E}(\omega(\Theta) | X=x)}.$$

In particular, if  $\omega(\cdot) \equiv 1$ ,  $\delta_{\Lambda}(x) = \mathbb{E}(g(\Theta) | X=x)$

[consistent with our previous result].

Q: Can  $\bar{X}_n (= n^{-1} \sum_{i=1}^n X_i)$  be a Bayes estimator?

No. Except that it degenerates to 0, which is a rare case.



[THEOREM] (TPE 4.2.3) If  $\delta$  is unbiased for  $g(\theta)$  with  $r(\Lambda, \delta) < \infty$  and  $E(g(\theta)^2) < \infty$ , then  $\delta$  is not Bayes under the squared error loss function unless its average risk is zero, i.e.

$$\rightarrow E_{(X, \Theta)} (\{\delta(X) - g(\Theta)\}^2) = 0 .$$

Proof. Let  $\delta$  be an unbiased estimator under the squared error loss function. Then, we know that  $\delta$  is the posterior mean, i.e.

$$\delta(X) = E(g(\Theta) | X) \quad \text{a.s. . } \checkmark$$

Thus we have

$$\begin{aligned} E(\delta(X)g(\Theta)) &= E(E(\delta(X)g(\Theta) | X)) \\ &= E(\delta(X)E(g(\Theta) | X)) = E(\{\delta(X)\}^2). \end{aligned}$$

Also,

$$\begin{aligned} E(\delta(X)g(\Theta)) &= E(E(\delta(X)g(\Theta) | \Theta)) \\ &= E(g(\Theta)E(\delta(X) | \Theta)) = E(\{g(\Theta)\}^2). \end{aligned}$$

Observe that

$$\begin{aligned} E(\{\delta(X) - g(\Theta)\}^2) &= E(\delta^2(X)) - 2E(\delta(X)g(\Theta)) + E(g^2(\Theta)) \\ &= E(\delta^2(X)) - E(\delta(X)g(\Theta)) + E(g^2(\Theta)) \\ &\quad - E(\delta(X)g(\Theta)) \\ &= E(\delta^2(X)) - E(\delta^2(X)) + E(g^2(\Theta)) - E(g^2(\Theta)) \\ &\quad \quad \quad \text{(By ①)} \quad \quad \quad \text{(By ②)} \end{aligned}$$

$$\text{Thus we have } E_{X, \Theta} (\{\delta(X) - g(\Theta)\}^2) = 0 .$$

(Example) Suppose  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\theta, \sigma^2)$  with  $\sigma^2$  known.

$$\hat{\theta} \sim N(\mu, b^2)$$

Is  $\bar{X}_n$  Bayes under the squared error loss function for some choice of the prior?

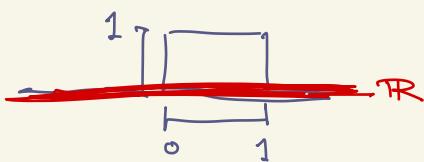
Observe that  $E(\bar{X}_n | \theta) = \theta$ . Hence  $\bar{X}_n$  is unbiased for  $\theta$ .

The corresponding average risk under the squared error loss function is given by

$$E_{(\bar{X}_n, \theta)} \left( \frac{1}{n} \sum_{i=1}^n (\bar{X}_n - \theta)^2 \right) = \frac{\sigma^2}{n} \neq 0$$

$\Rightarrow \bar{X}_n$  is not a Bayes estimator under "any prior distribution".

Remark:



$\bar{X}_n$  does minimize a form of average risk: it minimizes the average risk with respect to the Lebesgue measure, that is, with respect to the density  $\pi(\theta) = 1$  for all  $\theta$ . We call this choice of  $\pi$  an improper prior since the integral  $\int \pi(\theta) d\theta = \infty$  [in which  $\pi(\theta)$  does not define a proper prob. distribution]

We may define a formal posterior for this improper prior:

Formal posterior  $\propto$  likelihood  $\times$  Improper Prior.

(e.g. Normal:

$$\text{Formal posterior} \propto \prod_{i=1}^n \exp \left( -\frac{1}{2\sigma^2} (X_i - \theta)^2 \right) \times \underset{\text{Improper prior}}{\downarrow} 1$$

$$\propto \exp \left( \frac{n\bar{X}_n}{\sigma^2} \theta - \frac{n}{2\sigma^2} \theta^2 \right)$$

$$\sim N(\bar{X}_n, \sigma^2/n)$$

We call

$\bar{X}_n$  as a generalized Bayes estimator)

DBA

non-informative prior.

[THEOREM] A unique Bayes estimator (as for all  $P_\theta$ ) is admissible.

(TPE 5.2.4)

An estimator is admissible if it is not uniformly dominated by some other estimator.  $\delta$  is said to be inadmissible if and only if there exists  $\delta'$  such that

$$R(\theta, \delta') \leq R(\theta, \delta) \quad \forall \theta \in \Omega \text{ and}$$

$$R(\theta, \delta') < R(\theta, \delta) \quad \text{for some } \theta \in \Omega.$$

Proof. Suppose  $\delta_\Lambda$  is Bayes for  $\Lambda$ , and for some  $\delta'$ ,  $R(\theta, \delta') \leq R(\theta, \delta_\Lambda)$  for all  $\theta \in \Omega$ . If we take the expectation with respect to  $\Lambda$ , the inequality above is preserved and can be written as

$$\int_{\theta \in \Omega} R(\theta, \delta') d\Lambda(\theta) \leq \int_{\theta \in \Omega} R(\theta, \delta_\Lambda) d\Lambda(\theta).$$

This implies that  $\delta'$  is also Bayes because  $\delta'$  has less (or equal) risk than  $\delta_\Lambda$  which minimizes the average risk. Hence  $\delta' = \delta_\Lambda$  with prob. 1 for all  $P_\theta$ .

Question: When will a Bayes estimator be unique?

[THEOREM] Let  $Q$  be the marginal distribution of  $X$ , that is

(TPE 4.14)

$$Q(E) = \int P(X \in E | \Lambda) d\Lambda.$$

Then under a strictly convex loss function,  $\delta_\Lambda$  is unique (as for all  $P_\theta$ ) if

(a)  $R(\Lambda, \delta_\Lambda)$  is finite and finiteness for comparison

(b)  $P_\theta \ll Q$  \Omega is defined on the support of \Lambda.  
abs. cont.

Why do we need to consider Bayes estimators?

1). All admissible estimators are limits of Bayes estimators.

[explain/elaborate further later when we discuss sequences of priors]

$$\Lambda_m \xrightarrow{\sim} N(\mu, \frac{\sigma^2}{m})$$

minimax

## z) Prior information

↓ How can we choose a prior?

- a. subjective - prior knowledge / experience
- b. objective - maximally non-informative prior or reference prior (Jeffreys priors etc.).
- c. conjugate prior - for computational convenience
- d. Hierarchical Bayes → TPE §4.5.
- "  
2) Empirical Bayes (\*).

we begin with the following Bayes model :

$$x_i | \theta \sim f(x_i | \theta) \quad i = 1, \dots, n$$

$$\Theta | \gamma \sim \Pi(\theta | \gamma) \quad (\text{prior distribution})$$

One can calculate the marginal distribution of  $\underline{x}$  with density

$$m(\underline{x} | \gamma) = \underbrace{\int \prod_{i=1}^n f(x_i | \theta)}_{\text{likelihood}} \underbrace{\Pi(\theta | \gamma)}_{\text{prior.}} d\theta$$

Based on  $m(\underline{x} | \gamma)$ , we can obtain an estimate, say  $\hat{\gamma}(x)$  of  $\gamma$ .

We can then substitute  $\hat{\gamma}(x)$  for  $\gamma$  in  $\Pi(\theta | \gamma)$  and determine the estimator that minimizes the empirical posterior loss

$$\int L(\theta, \hat{\gamma}(x)) \Pi(\theta | x, \hat{\gamma}(x)) d\theta$$

empirical Bayes estimator  
(minimizer)

(Example). Suppose there are  $K$  different groups of patients, where each group has  $n$  patients. Each group is given a different treatment for the same illness, and in the  $k^{\text{th}}$  group ( $k=1, \dots, K$ ), we count  $X_k$ , the number of successful treatments out of these  $n$  patients.

$$X_k \sim \text{binomial}(n, p_k)$$

$$p_k \sim \text{Beta}(a, b) \quad k=1, 2, \dots, K$$

instead of  $(a_k, b_k)$ .

Share a common prior distribution

The single-prior Bayes estimator  $p_k$  under the squared error loss is

$$\delta^\pi(X_k) = E(p_k | X_k, a, b) = \frac{a + X_k}{a + b + n} \quad \text{[done b+]} \quad (6.7)$$

In the empirical Bayes model, we consider these hyperparameters unknown and we want to estimate them.

To construct an EBE, we calculate

$$m(\tilde{x} | a, b) = \int_0^1 \dots \int_0^1 \prod_{k=1}^K \binom{n}{x_k} p_k^{x_k} (1 - p_k)^{n - x_k}$$

Check Table 6.1 (on Page 265 of TPE §4.6)

$K=10, n=20$

$$\begin{aligned} \text{Prior parameters } a & \quad b \\ \downarrow \text{Bayes Risk} & \\ \delta^\pi & \quad \delta^\pi \\ (6.7) & \quad (6.9) \\ x/n & \\ \hat{a}, \hat{b} & \stackrel{\Delta}{=} \text{argmax}_{\hat{a}, \hat{b}} \end{aligned} = \left[ \prod_{k=1}^K \binom{n}{x_k} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+x_k)\Gamma(n-x_k+b)}{\Gamma(a)\Gamma(b)\Gamma(a+b+n)} \right]$$

product of beta-binomial distribution

The corresponding empirical Bayes estimator is given by:

$$\begin{array}{ccc} 2 & 2 & 0.0833 \quad 0.0850 \quad 0.1000 \\ \vdots & & \hat{\delta}^\pi(X_k) = E(p_k | X_k, \hat{a}, \hat{b}) = \frac{\hat{a} + X_k}{\hat{a} + \hat{b} + n} \end{array} \quad (6.9)$$

(check)

[THEOREM] (TPE 4.6.2) For the situation of Corollary 3.3 with the prior distribution  $\pi(\lambda | \lambda)$ . Suppose  $\hat{\lambda}(x)$  is the MLE of  $m(x | \lambda)$ , then the empirical Bayes estimator is

$$p_\eta(x) = e^{\sum_i \eta_i x_i - A(\eta)} h(x) \quad E(\eta_i | x, \hat{\lambda}) = \frac{\partial}{\partial x_i} \log m(x | \hat{\lambda}(x)) - \frac{\partial}{\partial x_i} \log h(x)$$

(You may verify this result)

(Example) Consider the estimation of

$$X_i | \theta_i \sim N(\theta_i, \sigma^2) \quad i=1, \dots, p, \text{ independent}$$

$$\Theta_i \sim N(\mu, \tau^2) \quad i=1, \dots, P$$

Here  $\mu$  is unknown. We can use the previous theorem to calculate the empirical Bayes estimator,

$$E(\Theta_i | \tilde{x}, \hat{\mu}) = \sigma^2 \left( \frac{\partial}{\partial x_i} \log m(\tilde{x} | \hat{\mu}) - \frac{\partial}{\partial x_i} f_i(x_i) \right),$$

where  $\hat{\mu}$  is the MLE of  $\mu$  from

$$m(\tilde{x} | \mu) = \frac{1}{2\pi (\sigma^2 + \tau^2)^{P/2}} e^{-\frac{\sum_{i=1}^P (x_i - \mu)^2}{2(\sigma^2 + \tau^2)}} \quad \checkmark$$

Here  $\hat{\mu} = \bar{x}$  and

$$\frac{\partial}{\partial x_i} \log m(\tilde{x} | \hat{\mu}) = \frac{\partial}{\partial x_i} \left( -\frac{1}{2(\sigma^2 + \tau^2)} \sum_{i=1}^P (x_i - \bar{x})^2 \right),$$

which gives the elliptical Bayes estimator

$$E(\Theta_i | \tilde{x}, \hat{\mu}) = \frac{\tau^2}{\sigma^2 + \tau^2} x_i + \underbrace{\frac{\sigma^2}{\sigma^2 + \tau^2} \bar{x}}_{\text{corr. to the prior } \pi(\theta | \bar{x})}.$$

$\hookrightarrow$  corr. to the prior  $\pi(\theta | \bar{x})$ .

... (You may also take a look of

Risk calculation for the emp. Bayes estimator | )

$$\begin{aligned} R(\eta, E(\eta | X, \hat{\lambda}(X))) &= R(\eta, -\nabla \log h(X)) \\ &\quad + \sum_{i=1}^P E_\eta \left\{ 2 \frac{\partial^2}{\partial x_i^2} \log m(X | \hat{\lambda}(X)) \right. \\ &\quad \left. + \left( \frac{\partial}{\partial x_i} \log m(X | \hat{\lambda}(X)) \right)^2 \right\} \end{aligned}$$

## Worst-case Optimality

Given  $X \sim P_\theta$ , where  $\theta \in \Omega$  and a loss function  $L(\theta, d)$ , we want to minimise the maximum risk:  $\sup_{\theta \in \Omega} R(\theta, \delta)$ . The minimiser is known as a minimax estimator.

Recall the definition of Bayes risk under an arbitrary prior distribution  $\Lambda$ :

$$\Lambda: r_\Lambda = \inf_{\delta} \underline{r}(\Lambda, \delta) = \inf_{\delta} \int_{\theta \in \Omega} R(\theta, \delta) d\Lambda(\theta)$$

$$E(L(\theta, \delta(x)) | \Theta)$$

Definition A prior distribution  $\Lambda$  is said to be a least favourable prior if  $r_\Lambda \geq r_{\Lambda'}$  for any other prior distribution  $\Lambda'$ .

[THEOREM] Suppose  $\delta_\Lambda$  is Bayes for  $\Lambda$  with  $r_{\Lambda \in \Omega} = \sup_{\theta} R(\theta, \delta_\Lambda)$

Bayes risk is the max risk of  $\delta_\Lambda$

- 1)  $\delta_\Lambda$  is minimax;
- 2)  $\Lambda$  is a least favourable prior and
- 3) If  $\delta_\Lambda$  is unique Bayes estimator for  $\Lambda$ , it is also unique minimax.

Proof:

- i)  $\sup_{\theta \in \Omega} R(\theta, \delta) \geq \int R(\theta, \delta) d\Lambda(\theta)$
- $\geq \int R(\theta, \delta_\Lambda) d\Lambda(\theta)$
- $= \sup_{\theta} R(\theta, \delta_\Lambda)$

$\Rightarrow \delta_\Lambda$  is minimax. 3) ...

2) Let  $\Lambda'$  be any other prior distribution. Then we have

$$\begin{aligned} r_{\Lambda'} &= \inf_{\delta} R(\theta, \delta) d\Lambda'(\theta) \\ &\leq \int R(\theta, \delta_\Lambda) d\Lambda'(\theta) \leq \underbrace{\sup_{\theta} R(\theta, \delta_\Lambda)}_{\text{cond.}} = r_\Lambda \\ &\Rightarrow \Lambda \text{ is } \underset{\Lambda'}{\text{least favourable prior.}} \end{aligned}$$

[Corollary] If a Bayes estimator  $\delta_\Lambda$  has constant risk:  $R(\theta, \delta_\Lambda) = R(\theta', \delta_\Lambda)$  for  $\theta, \theta' \in \Omega$ , then  $\delta_\Lambda$  is minimax.

\* This is sufficient but not necessary cond.

Strategy: Find a support set  $\omega$  such that  $A(\omega) = 1$  and for which  $R(\theta, \delta_A)$  is maximum for all  $\theta \in \omega$ .

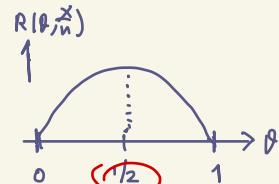
[Corollary] Define  $\omega_A = \left\{ \theta : R(\theta, \delta_A) = \sup_{\theta'} R(\theta', \delta_A) \right\}$   
TPE 5.1.6.

then a Bayes estimator  $\delta_A$  is minimax if  $A(\omega) = 1$ .  
iff

(Example) Suppose  $X \sim \text{Binomial}(n, \theta)$  for some  $\theta \in (0, 1)$  and we adopt the squared error loss function. Is the sample portion  $\frac{X}{n}$  minimax?

$$R(\theta, \frac{X}{n}) = E \left( \left( \frac{X}{n} - \theta \right)^2 \right) = \frac{\theta(1-\theta)}{n}$$

$$\sup_{\theta} R(\theta, \frac{X}{n}) = R\left(\frac{1}{2}, \frac{X}{n}\right) = \frac{1}{4n}.$$



TPE 5.1.5 / 5.1.6 is not helpful here as  $A\left\{\frac{1}{2}\right\} = 1 \Rightarrow \delta_A(X) = \frac{1}{2} (\neq \frac{X}{n})$

With Beta( $a, b$ ) as the prior,  $\delta_{a,b}(X) = \frac{X+a}{n+a+b}$  (see Lecture \_).

For any  $a$  and  $b$ ,

$$\begin{aligned} R(\theta, \delta_{a,b}(X)) &= E_{\theta} \left( \left( \frac{X+a}{n+a+b} - \theta \right)^2 \right) \\ &= \frac{1}{(n+a+b)^2} E_{\theta} \left( (X+a - \theta(n+a+b))^2 \right) \\ &= \frac{1}{(n+a+b)^2} \left[ n\theta(1-\theta) + \{a(\theta-1) + \theta b\}^2 \right], \quad \checkmark \end{aligned}$$

which is a quadratic function of  $\theta$ . To eliminate  $\theta$ , we need to

$$\text{set } \begin{cases} -n + (a+b)^2 = n - 2a(a+b) = 0 \quad [\text{coeff. of } \theta^2] \\ \dots \dots \dots \\ n - 2a(a+b) = 0 \dots \quad [\text{coeff. of } \theta] \end{cases}$$

$\Rightarrow$  Note that  $a, b > 0$

$$a+b = \sqrt{n}, \quad n - 2a\sqrt{n} = 0 \dots$$

$$\begin{aligned} \Rightarrow a = \sqrt{n} - b \Rightarrow n - 2(\sqrt{n}-b)\sqrt{n} &= n - 2n + 2b\sqrt{n} \\ \Rightarrow \dots \quad n = 2b\sqrt{n} \Rightarrow \begin{pmatrix} b = \sqrt{n}/2 \\ a = \sqrt{n}/2 \end{pmatrix} \end{aligned}$$

$\Rightarrow \text{Beta}\left(\frac{\sqrt{n}}{2}, \frac{\sqrt{n}}{2}\right)$  is a least favorable prior with constant risk. Then our Bayes estimator  $\delta_{\frac{\sqrt{n}}{2}, \frac{\sqrt{n}}{2}}(X) = \frac{X + \sqrt{n}/2}{n + \sqrt{n}}$  is minimax with constant risk of  $\frac{1}{4(\sqrt{n}+1)^2}$ . (check).

We conclude that  $\frac{X}{n}$  is NOT minimax.  $\Leftarrow \frac{1}{4n + (8\sqrt{n} + 4)} < \frac{1}{4}$ . //

(Example) Minimax for iid Normal r.v.s. with unknown mean  $\theta$ :  $X_1, \dots, X_n \sim N(\theta, \frac{\sigma^2}{n})$

Observe that  $\bar{X}_n$  has constant risk of  $\frac{\sigma^2}{n}$ . However,  $\bar{X}_n$  is not Bayes for any prior. [Recall TPE 4.2.3: unbiased estimators are Bayes only in the degenerate situations of zero risk]

We cannot conclude via TPE 5.1.5 that  $\bar{X}_n$  is minimax.

Consider instead:  $\delta_{\alpha, \mu_0}(x) = \alpha \bar{X}_n + (1-\alpha) \mu_0$ , for  $\alpha \in (0, 1)$ ,  $\mu_0 \in \mathbb{R}$

The corresponding worst-case risk:

$$\begin{aligned} \sup_{\theta} E_{\theta} \left( (\theta - \delta(x))^2 \right) &= \sup_{\theta} \left( \alpha^2 \text{Var}_{\theta}(\bar{X}_n) + (1-\alpha)^2 (\theta - \mu_0)^2 \right) \\ &= \frac{\alpha^2 \sigma^2}{n} + (1-\alpha)^2 \sup_{\theta} (\theta - \mu_0)^2 = +\infty. \end{aligned}$$



→ Extend our minimax results to the limits of Bayes estimators rather than restricting our attention to "regular" Bayes estimator only.

Definition (Least favourable sequences of priors)

$$\Lambda \sim N(\mu, \tau^2)$$

Let  $\{\Lambda_m\}$  be a sequence of priors with minimal average risk  $r_{\Lambda_m} = \inf_{\delta} \int R(\theta, \delta) d\Lambda_m(\theta)$ . Then,  $\{\Lambda_m\}$  is a least favourable sequence of priors if there is a real number  $r$  such that  $r_{\Lambda_m} \rightarrow r < \infty$  and  $r \geq r_{\Lambda'}$ , for any other prior  $\Lambda'$ .

[THEOREM]

5.1.12

(Extend 5.1.4)

Suppose there is a real number  $r$  such that  $\{\Lambda_m\}$  is a sequence of priors with  $r_m \rightarrow r < \infty$ . Let  $\delta$  be any estimator such that  $\sup_{\theta} R(\theta, \delta) = r$ . Then i)  $\delta$  is minimax and ii)  $\{\Lambda_m\}$  is a least favourable sequence of priors.

- Proof. i) Let  $\delta'$  be any other estimator. Then for any  $m$ ,

$$\sup_{\theta} R(\theta, \delta') \geq \int R(\theta, \delta') d\Lambda_m(\theta) \geq r_{\Lambda_m}.$$

Sending  $m \rightarrow \infty$  yields

$$\sup_{\theta} R(\theta, \delta') \geq r = \sup_{\theta} R(\theta, \delta)$$

$\Rightarrow \delta$  is minimax.

ii) Let  $\Lambda'$  be any prior, then

$$r_{\Lambda'} = \int R(\theta, \delta_{\Lambda'}) d\Lambda'(\theta) \leq \int R(\theta, \delta) d\Lambda'(\theta) \leq \sup_{\theta} R(\theta, \delta) = r.$$

$\Rightarrow \{\Lambda_m\}$  is least favourable,

■ No uniqueness statement is guaranteed here because of the limit.

■ Generalise the scope so that Bayes estimators whose Bayes risks converge to the maximum risk can serve as a possible candidate.

(Example) (Minimax for iid Normal r.v.s. cont'd)

By the result above, it suffices to find a sequence  $\{\Lambda_m\}$  such that

$$r_{\Lambda_m} \rightarrow \frac{\sigma^2}{n} = r.$$

Let  $\{\Lambda_m\}$  be conjugate priors  $\{N(0, m^2)\}$  in which case  $\Lambda_m$  tends to the uniform prior on  $\mathbb{R}$  (also improper with  $\pi(\theta) = 1, \forall \theta \in \mathbb{R}$ ).

Recall that

$$\theta | (X_1, \dots, X_n) \sim N \left( \frac{n\bar{X}_n}{\frac{n}{\sigma^2} + \frac{1}{m^2}}, \underbrace{\frac{1}{\frac{n}{\sigma^2} + \frac{1}{m^2}}}_{\perp\!\!\!\perp (X_1, \dots, X_n)} \right)$$

$$r_{\Lambda_m} = \frac{1}{\frac{n}{\sigma^2} + \frac{1}{m^2}} \xrightarrow{m \rightarrow \infty} \frac{\sigma^2}{n} = \sup_{\theta} R(\theta, \bar{X}_n).$$

$\Rightarrow \bar{X}_n$  is minimax and  $\{\Lambda_m\}$  is least favourable.

- Minimality via submodel restriction.

Another technique of deriving a minimax estimator for a general family of models is to restrict the problem/discussion to a subset of that family.

(Example) Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, \sigma^2)$ , where neither  $\theta$  nor  $\sigma^2$  is known.

Note that  $\sup_{(\theta, \sigma^2)} R(\theta, \sigma^2), \bar{X}_n = \sup_{\sigma^2} \frac{\sigma^2}{n} = \infty$ .

Restrict our attention to  $\Omega = \{(\theta, \sigma^2) : \theta \in \mathbb{R}^2, \sigma^2 \leq B\}$ ,

where  $B$  is a known constant.

$$\bullet \sup_{\theta \in \mathbb{R}, \sigma^2 \leq B} R(\theta, \sigma^2, \bar{X}_n) = \frac{B}{n}$$

$$= \sup_{\theta \in \mathbb{R}, \sigma^2 = B} R(\theta, \sigma^2, \bar{X}_n) \leq \sup_{\theta \in \mathbb{R}, \sigma^2 = B} R(\theta, \sigma^2, \delta)$$

$$\leq \sup_{\theta \in \mathbb{R}, \sigma^2 \leq B} R(\theta, \sigma^2, \delta)$$

(Example). (Weighted squared error loss)

Let  $X \sim \text{Binomial}(n, \theta)$ , with the loss function  $L(\theta, d) = \frac{(d-\theta)^2}{\theta(1-\theta)}$  ←  
 (hence  $w(\theta) = \frac{1}{\theta(1-\theta)}$ ). Note that here for any  $\theta$ ,  $R(\theta, \frac{X}{n}) = \frac{1}{n}$ ; that  
 is to say, the risk is constant in  $\theta$ .

Recall that with the loss function  $L(d, \theta) = w(\theta)(d-\theta)^2$ , the associated  
 Bayes estimator is given by  $E_{\Theta|X}(\Theta | \omega(\Theta) | X) / E_{\Theta|X}(\omega(\Theta) | X)$ .

Applying this result, we find that the Bayes estimator has the form

$$\delta_A(X) = \frac{E_{\Theta|X}\left(\frac{1}{1-\Theta} | X\right)}{E_{\Theta|X}\left(\frac{1}{\Theta(1-\Theta)} | X\right)}.$$

We use a prior conjugate to the binomial likelihood:  $\Theta \sim \text{Beta}(a, b) \stackrel{d}{=} \text{Beta}(a, b)$

$\xrightarrow{\text{a, b, posterior}}$  for some  $a, b > 0$

Suppose we observe  $X=x$ . If  $a+x > 1$ , and  $b+n-x > 1$ , then we can  
 claim that the resulting estimator is given by

$$\rightarrow \boxed{\delta_{a,b}(x) = \frac{a+x-1}{a+b+n-2}} \quad (\text{to be proved in a min}).$$

and it minimises the posterior risk.

In particular, the estimator  $\delta_{1,1}(x) = \frac{x}{n}$  minimises the posterior risk  
 with respect to the uniform prior after observing  $0 < x < n$ . ← verify

If we can verify also that this form remains unchanged when  $x \in \{0, n\}$   
 $\rightarrow$  then we can claim that  $\delta_{1,1}(x) = \frac{x}{n}$  is Bayes with constant risk  
 and hence minimax.

$\rightarrow$  Recall that the Beta function can be evaluated as

$$\int_0^1 x^{k_1-1} (1-x)^{k_2-1} dx = \frac{\Gamma(k_1) \Gamma(k_2)}{\Gamma(k_1 + k_2)}, \text{ where } k_1, k_2 > 0.$$

Therefore, if  $Y \sim \text{Beta}(a, b)$ , where  $a, b > 0$ , we can evaluate the expectation

$$\begin{aligned} E\left(\frac{1}{1-y}\right) &= \int_0^1 \frac{1}{1-y} \left\{ \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} y^{a-1} (1-y)^{b-1} \right\} dy \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 \left\{ y^{a-1} (1-y)^{b-2} \right\} dy \end{aligned}$$

$$\begin{aligned}
&= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(a)\Gamma(b-1)}{\Gamma(a+b-1)} \\
&= \frac{\Gamma(a+b)}{\Gamma(a+b-1)} \cdot \frac{\Gamma(b-1)}{\Gamma(b)} \cdot \frac{\Gamma(a)}{\Gamma(a)} = \frac{a+b-1}{b-1}.
\end{aligned}$$

Similarly  $E\left(\frac{1}{Y(1-Y)}\right) = \dots = \frac{(a+b-2)(a+b-1)}{(a-1)(b-1)}$ , where  $a > 1$ .

Combining these results / identities, we have, for  $a, b > 1$ ,

$$\frac{E\left(\frac{1}{Y(1-Y)}\right)}{E\left(\frac{1}{Y}\right)} = \frac{a-1}{a+b-2}.$$

To see this is the case, note that the posterior risk under the prior  $\Delta_{1,1}$  after observing  $X=x$  and deciding  $\delta(x)=d$  is

$$\int_0^1 \underbrace{\frac{(d-\theta)^2}{\theta(1-\theta)}}_{\text{Risk}} \cdot \frac{\Gamma(x+1+n-x+1)}{\Gamma(x+1)\Gamma(n-x+1)} \cdot \theta^x (1-\theta)^{n-x} d\theta,$$

which, in the case of  $X=0$ , implies to  $\int_0^1 \frac{(d-\theta)^2}{\theta} (1-\theta)^{n-1} d\theta$ .

This integral converges for  $d=0$  and diverges otherwise, so the posterior risk is minimized by choosing  $\underline{\delta}(0)=D \in 0/n$ . Similarly in the case of  $X=n$ , the posterior risk is minimised by choosing  $\underline{\delta}(n)=1 = n/n$ . This confirms that  $\underline{\delta}_{1,1}(X) = \underline{X}/n$  minimises the posterior for any outcome  $X$  and is indeed Bayes. This estimator has constant risk, we conclude that  $\underline{X}/n$  is minimax. ||

### Randomised Minimax Estimators.

Randomised estimators are functions of the form  $\delta(X, U)$ , where  $U \sim \text{Uniform}(0,1)$

Non-random estimators usually achieve better or equal average risk and they are oftentimes available.

Observation: With convex losses, we can "ignore" randomised estimators because we can find a better or equally good non-randomised/deterministic estimator.

Recall that the data  $X$  is sufficient, so by Rao-Blackwell theorem, the non-random estimator  $\tilde{\delta}(X) = E(\delta(X, U) | X)$  is no worse than  $\delta(X, U)$ .

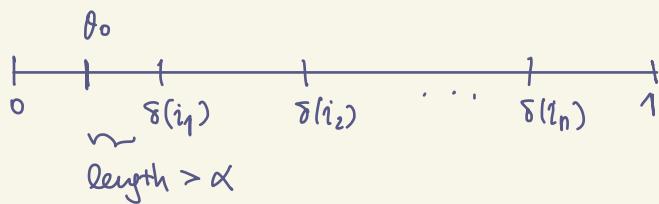
(Example) (Randomised minimax estimator).

Let  $X \sim \text{Binomial}(n, \theta)$ , where  $\theta \in (0, 1)$ , and consider estimation of  $\theta$  under the 0-1 loss:

$$L(\theta, d) = \begin{cases} 0 & \text{if } |d - \theta| < \alpha \\ 1 & \text{otherwise} \end{cases}$$

First consider an arbitrary non-random estimator  $\delta$ . Since  $X$  can only take the  $n+1$  values  $\{0, 1, \dots, n\}$ , the estimator  $\delta(X)$  can also only take  $n+1$  values  $\{\delta(0), \delta(1), \dots, \delta(n)\}$ .

If  $\alpha < \frac{1}{2(n+1)}$ , then we can always find  $\theta_0$  such that  $|\delta(x) - \theta_0| \geq \alpha$  for every  $x \in \{0, 1, \dots, n\}$ .



$R(\theta_0, \delta(X)) = 1$ , which corresponds to the maximum risk of any non-random  $\delta$ .

Consider instead the estimator  $\delta'(X, U) = U$ , which is completely random and independent of the data  $X$ . Then for any  $\theta \in (0, 1)$ ,

$$\begin{aligned} R(\theta, \delta') &= E\left(L(\theta, \underbrace{\delta'(X, U)}_U)\right) \\ &= P(|U - \theta| \geq \alpha) \\ &= 1 - P(\theta - \alpha < U < \alpha + \theta) \\ &\leq 1 - \alpha < 1 < \text{maximum risk of any non-random estimator } \delta. \end{aligned}$$

||

Next lecture:

?  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, \sigma^2)$ ,  $L(\theta, d) = (\theta - d)^2$ , when will  $\bar{X}_n$  be admissible?

known

## Lecture 10

14 Nov 2022

- \* Admissibility of minimax est.
- \* Simultaneous estimation: James-Stein estimator,
- \* Testing.

### Minimax Estimators and submodels.

Recall that an estimator  $\delta^M$  is minimax if its maximum risk is minimal:

$$\inf_{\delta} \sup_{\theta \in \Omega} R(\theta, \delta) = \sup_{\theta} R(\theta, \delta^M).$$

[Lemma] Suppose that  $\delta$  is minimax for a submodel  $\theta \in \Omega_0 \subset \Omega$  and  
TPE 5.1.15

$$\sup_{\theta \in \Omega_0} R(\theta, \delta) = \sup_{\theta \in \Omega} R(\theta, \delta)$$

Then  $\delta$  is minimax for the full model  $\theta \in \Omega$ .

→ Allows us to find a minimax est. for a particular tractable submodel, then we show that the worst-case risk doesn't rise as we expand the model.

(Example) Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , where both  $\mu$  and  $\sigma^2$  are unknown. Our parameter vector  $\theta = (\mu, \sigma^2)$  and our parameter space  $\Omega = \mathbb{R} \times \mathbb{R}^+$ . Our goal is to estimate  $\mu$ . Our loss function is the rel. squared error loss, which is given by

$$L((\mu, \sigma^2), d) = \frac{(d - \mu)^2}{\sigma^2} \quad \leftarrow \text{(c.f. } \underbrace{\tilde{L}((\mu, \sigma^2), d) = (\mu - d)^2}_{\text{unbounded risk.}} \right)$$

We consider the submodel where  $\sigma^2 = 1$ , i.e.  $\Omega_0 = \mathbb{R} \times \{1\}$  and our loss function reduces to  $L(\mu, 1), d) = (d - \mu)^2$ .

Recall that  $\bar{X}_n$  is minimax for  $\Omega_0$  because  $R((\mu, 1), \bar{X}) = \frac{1}{n} \quad \forall (\mu, 1) \in \Omega$

Thus the risk does not depend on  $\sigma^2$ .

$R((\mu, 1), \bar{X}_n) = R((\mu, \sigma^2), \bar{X}_n)$ , we have that the maximum risks are equal. [i.e.  $\sup_{\theta \in \Omega_0} R(\theta, \delta) = \sup_{\theta \in \Omega} R(\theta, \delta)$ ]. It follows from Lemma 5.1.15 that  $\bar{X}_n$  is minimax on  $\Omega$ .

## Admissibility of minimax estimators.

- 1) Admissibility can give rise to minimaxity: If  $\delta$  is admissible with constant risk, then  $\delta$  is also minimax.
- 2) Minimaxity does not guarantee admissibility: It only ensures that the worst case risk is optimal.  
(see also TPE 5.2.4)

(Example) Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, \sigma^2)$  with  $\sigma^2$  known, and  $\theta$  the estimand. Then the minimax estimator is  $\bar{X}_n$  under the squared error loss and we would like to determine whether or not  $\bar{X}_n$  is admissible.

Q: Let's try to answer a more general question: When is  $a\bar{X}_n + b$ , for  $a, b \in \mathbb{R}$  (i.e. affine function of  $\bar{X}_n$ ) admissible?

Case 1  $0 < a < 1$ : In this case  $a\bar{X}_n + b$  is a convex combination of  $\bar{X}_n$  and  $b$ . By our previous results, it is a Bayes estimator with respect to some Gaussian prior on  $\theta$ . Since we adopt the squared error loss function, which is strictly convex, the Bayes estimator is unique. By Theorem 5.2.4 of TPE (a unique Bayes est. is always admissible),  $a\bar{X}_n + b$  is admissible.

Case 2  $a=0$ : In this case  $b$  is also a unique Bayes estimator w.r.t. a degenerate prior distribution with the unit mass at  $\theta = b$ . So, by Theorem 5.2.4,  $b$  is also admissible.

Case 3  $a=1, b \neq 0$ : In this case  $\bar{X}_n + b$  is not admissible because it is dominated by  $\bar{X}_n$ . Note that  $\bar{X}_n + b$  and  $\bar{X}_n$  have the same variance, but  $\bar{X}_n$  has a strictly smaller bias.

In general, the risk of  $a\bar{X}_n + b$  is

$$\begin{aligned} E((a\bar{X}_n + b - \theta)^2) &= E\left(\{a(\bar{X}_n - \theta) + b + \theta(a-1)\}^2\right) \\ &= \dots = \frac{a^2 \sigma^2}{n} + \{b + \theta(a-1)\}^2 \end{aligned} \quad (*)$$

Case 4  $a > 1$ : We have by (\*)

$$E((a\bar{X}_n + b - \theta)^2) \geq \frac{a^2 \sigma^2}{n} > \frac{\sigma^2}{n} = R(\theta, \bar{X}_n)$$

Hence, we conclude that  $\bar{X}_n$  dominates  $a\bar{X}_n + b$  when  $a > 1$  in which case  $a\bar{X}_n + b$  is inadmissible.

$$\begin{aligned}\text{Case 5 } a < 0: \quad E((a\bar{X}_n + b - \theta)^2) &> \left\{ b + \theta(a-1) \right\}^2 \\ &= (a-1)^2 \left( \theta + \frac{b}{a-1} \right)^2 \\ &> \left( \theta + \frac{b}{a-1} \right)^2, \quad \left( \theta - \frac{b}{a-1} \right)^2 \xrightarrow{\text{est.}}\end{aligned}$$

which is the risk of predicting with the constant  $\frac{-b}{a-1}$ . So,  $\frac{-b}{a-1}$  dominates  $a\bar{X}_n + b$  and therefore  $a\bar{X}_n + b$  is again inadmissible.

Case 6  $a = 1, b = 0$ : We shall proceed with a limiting Bayes argument.

Suppose  $\bar{X}_n$  is inadmissible, then, WLOG, we let  $\sigma^2 = 1$ , we have

$$R(\theta, \bar{X}_n) = \frac{1}{n}.$$

By our hypothesis, there must exist an estimator  $\delta'$  such that  $R(\theta, \delta') \leq 1/n$  for all  $\theta$  and  $\underline{R}(\theta, \delta') \leq 1/n$  for at least one  $\theta \in \Omega$ .

Because  $R(\theta, \delta)$  is continuous in  $\theta$ , there must exist  $\epsilon > 0$  and an interval  $(\theta_0, \theta_1)$  containing  $\theta'$  such that

$$\rightarrow R(\theta, \delta') < \frac{1}{n} - \epsilon \quad \forall \theta \in (\theta_0, \theta_1). \quad (\#)$$

Let  $r'_T$  be the average risk of  $\delta'$  w.r.t. the prior distribution  $N(0, T^2)$  on  $\theta$ .

(We used the same prior to prove that  $\bar{X}_n$  was the limit of a Bayes estimator, and hence minimax. We did this by letting  $T \rightarrow \infty$  (improper prior:  $\pi(\theta) = 1 \forall \theta$ ).

Let  $r_T$  be the average risk of a Bayes estimator  $\delta_T$  under the same prior.

Note that  $\delta_T \neq \delta'$  because  $R(\theta, \delta_T) \rightarrow \infty$  as  $\theta \rightarrow \infty$ , which is not consistent with  $R(\theta, \delta') \leq 1/n \quad \forall \theta \in \mathbb{R}$ . So  $r_T < r'_T$  because the Bayes estimator is unique almost surely with respect to the marginal distribution of  $\theta$ .

We shall look into the following ratio, which we shall show to become arbitrarily large and form a contradiction that

$$r_T < \underline{r}_T'$$

Using the form of the Bayes risk  $\underline{r}_T$ , we can write

$$\frac{\frac{1}{n} - \underline{r}_T'}{\frac{1}{n} - r_T} = \frac{\frac{1}{\sqrt{2\pi T}} \int_{-\infty}^{\infty} \left\{ \frac{1}{n} - R(\theta, \delta') \right\} \exp\left(-\frac{\theta^2}{2T^2}\right) d\theta}{\frac{1}{n} - \frac{1}{n+1/T^2}}$$

pmf  $N(0, T^2)$

Applying (#), we find:

$$\begin{aligned} \frac{\frac{1}{n} - \underline{r}_T'}{\frac{1}{n} - r_T} &\geq \frac{\frac{1}{\sqrt{2\pi T}} \int_{\theta_0}^{\theta_1} \varepsilon e^{-\frac{\theta^2}{2T^2}} d\theta}{\frac{1}{n(nT^2+1)}} \\ &= \frac{n(1+nT^2)}{T\sqrt{2\pi}} \varepsilon \int_{\theta_0}^{\theta_1} e^{-\frac{\theta^2}{2T^2}} d\theta \end{aligned}$$

$\theta \nearrow \infty$   $(\theta_1, -\theta_0)$

As  $T \rightarrow \infty$ , the first expression  $\frac{n(1+nT^2)\varepsilon}{T\sqrt{2\pi}} \rightarrow \infty$  and since the integrand converges non-tangentially to 1, Lebesgue's MCT ensures that the integral approaches the positive quantity  $\theta_1 - \theta_0$ . So, for sufficiently large  $T$ , we must have

$$\frac{\frac{1}{n} - \underline{r}_T'}{\frac{1}{n} - r_T} > 1 \Leftrightarrow \underline{r}_T' < \underline{r}_T$$

However, this is a contradiction because  $\underline{r}_T$  is the optimal avg. Mtk (since it is the Bayes mtk). So our assumption that there was a dominating estimator  $\delta$  is wrong and in this case  $a\bar{X}_n + b = \bar{X}_n$  is admissible.  $\square$

### Simultaneous estimation

So far, we consider only one parameter.

(Example). Let  $X_1, X_2, \dots, X_p$  be independent with  $X_i \sim N(\theta_i, \sigma^2)$  for  $1 \leq i \leq p$ . For the sake of simplicity / laziness, let  $\sigma^2 = 1$ . Our goal is to estimate  $\underline{\theta} = (\theta_1, \dots, \theta_p)^T$  under the loss function  $L(\underline{\theta}, \underline{d}) = \sum_{i=1}^p (d_i - \theta_i)^2$ .

A natural estimator for  $\theta$  is  $\hat{X} = (X_1, X_2, \dots, X_p)^T$ .

It can be shown that  $\hat{X}$  is the UMVUE, the MLE, a generalized Bayes and a minimax estimator for  $\underline{\theta}$ .

So it would be natural for us to think that  $X$  is admissible. However, it turns out that this is not the case when  $p \geq 3$ .

When  $p \geq 3$ ,  $X$  is dominated by the James-Stein estimator:  
(J-S)

$$\delta(X) = (\delta_1(X), \dots, \delta_p(X)),$$

$$\text{where } \delta_i(X) = \left(1 - \frac{p-2}{\|X\|_2^2}\right) \underline{x_i}.$$

→ Motivation for the J-S estimator:

View it from the empirical Bayes framework.

a) Suppose  $\theta_i \stackrel{iid}{\sim} N(0, A)$ , then the Bayes estimator for  $\theta_i$  is

$$\delta_{A,i}(X) = \frac{x_i}{1 + \frac{1}{A}} = \boxed{1 - \frac{1}{A+1}} \underline{x_i} \quad \checkmark$$

b) We need to choose  $A$ .

Marginalising over  $\theta$ , we can show that  $X$  has the distribution

$$x_i \stackrel{iid}{\sim} N(0, A+1) \quad \leftarrow \quad \int p(x; \theta) \pi(\theta; A) d\theta.$$

We may use  $X$  and the knowledge of this marginal distribution to find an estimate of  $1/(A+1)$ . Instead of using the MLE, we shall used an unbiased estimate:

$$\text{It can be shown that } E\left(\frac{1}{\|X\|_2^2}\right) = \frac{1}{(p-2)(A+1)} \quad (\text{Exercise})$$

( $\frac{1}{A+1} \|X\|_2^2$  follows a  $\chi_p^2$  distribution)

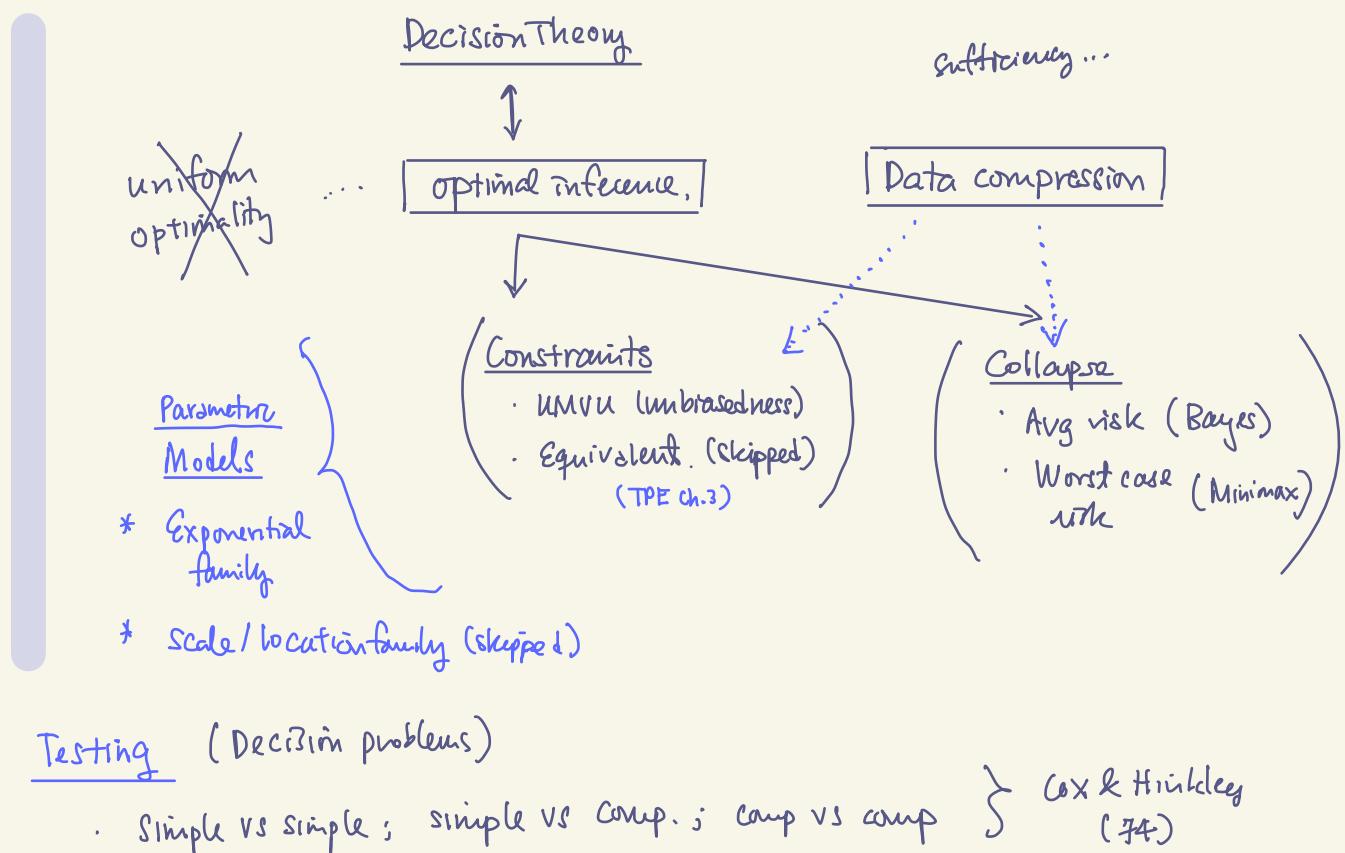
So,  $1 - \frac{p-2}{\|X\|_2^2}$  is UMVU for  $\boxed{1 - \frac{1}{A+1}}$ .

↳ J-S:  $\delta(x_i) = \left(1 - \frac{p-2}{\|X\|_2^2}\right) x_i$ .

Note that:  $E(\|X\|_2^2) = E\left(\sum_{i=1}^p X_i^2\right) = p + \sum_{i=1}^n \theta_i^2$   
 $= p + \|\theta\|_2^2 > p.$   $\frac{p\sigma^2}{p}$

→ We may view the J-S estimator as a method for correcting the bias in the size of  $X$ . (see TPE 5.5.1)

# Point Estimation



## Model Setup

We assume that the data are sampled from  $X \sim P_\theta$ , where  $P_\theta \in \mathcal{P} = \{P_\theta : \theta \in \Omega\}$ .

We divide the models in  $\mathcal{P}$  into two disjoint subclasses known as hypotheses:

$$H_0: \theta \in \Omega_0 \subset \Omega \quad (\text{null hypothesis})$$

$$H_1: \theta \in \Omega_1 \subset \Omega \setminus \Omega_0 \quad (\text{alternative hypothesis})$$

( $H_a$ )

Our goal is to infer (from the data) which hypothesis is "correct".

Our decision space is  $\mathcal{D} = \{\text{accept } H_0, \text{reject } H_0\}$ .  
↑  
do not reject

| Decision \ Truth                 | existing drug is better.<br>$\theta \in \Omega_0$ ( $H_0$ is correct) | new drug is better.<br>$\theta \in \Omega_1$ ( $H_1$ is correct) |
|----------------------------------|---|--|
| Clinical trials<br>(sample size) | Reject $H_0$  | Type I error   |
| Accept $H_0$                     |   | 1 - $\beta$ = Type II error.                                     |
| Loss function $L(\theta, d)$     |   |  |

Test function  $\phi(X) \in [0,1]$ . (critical function)

This test function  $\phi$  indicates that  $\delta_\phi(X, U)$  rejects  $H_0$  with probability  $\phi(X)$ . i.e.

$$\phi(X) = P\left(\delta_\phi(X, U) = \text{Reject } H_0 \mid X\right)$$

Definition The power function of a test  $\phi$  is  $\beta(\theta) = E_\theta(\phi(X))$

$$= P_{\theta}(\text{Reject } H_0)$$

If  $\theta_0 \in \Omega_0$  then  $\beta(\theta_0) = R(\theta_0, \delta_\phi) = \text{Type I error}$

If  $\theta_1 \in \Omega_1$ , then  $\beta(\theta_1) = 1 - R(\theta_1, \delta_\phi) = 1 - \text{Type II error}$ .

Our ideal goal is to minimize  $\beta(\theta_0)$  uniformly for all  $\theta \in \Omega_0$

and maximize  $\beta(\theta_1)$  uniformly for all  $\theta_1 \in \Omega_1$ .

## Lecture 11

21/Nov/2022.



### Neyman-Pearson Paradigm

Fix  $\alpha \in (0, 1)$  to control the Type I error. (Level of significance)

We require that our procedure can satisfy the following risk bound:

$$\sup_{\theta_0 \in \Omega_0} E_{\theta_0}(\phi(X)) = \sup_{\theta_0 \in \Omega_0} \beta(\theta_0) \leq \alpha$$

$\uparrow$   
size of the test.

Optimality goal: Find a level  $\alpha$  test that maximizes the power

$\beta(\theta_1) = E_{\theta_1}(\phi(X))$  for each  $\theta_1 \in \Omega_1$ . Such a test is called uniformly most powerful test (UMP).

MP test for the "Simple vs simple" case

Definition A hypothesis  $H_0$  is called simple if  $|\Omega_0| = 1$ , otherwise it is called composite. The same is true for  $H_1$ .

For the simple versus simple case, we will adopt the notation:

$$H_0: X \sim p_0 \quad (\text{density corr. to } P_{\theta_0}) \quad \theta = \theta_0$$

$$H_1: X \sim p_1 \quad (\text{density corr. to } P_{\theta_1}) \quad \theta = \theta_1$$

Our goal can be compactly described as:

$$\max_{\phi} E_{p_1}(\phi(X))$$

such that  $E_{p_0}(\phi(X)) \leq \alpha$ .

predefined value  
chosen by user

[Lemma] (Neyman-Pearson)

$(\theta = \theta_0)$

✓ (i) Existence. For testing  $H_0: p_0$  vs  $H_1: p_1$ , there is a test  $\phi(x)$  and a constant  $k_\alpha$  such that

$$\rightarrow \begin{cases} \text{a)} & E_{p_0} \phi(x) = \alpha \quad (\text{size} = \text{level}) \\ \text{b)} & \phi(x) = \begin{cases} 1 & \text{if } \frac{p_1(x)}{p_0(x)} > k_\alpha \\ 0 & \text{if } \frac{p_1(x)}{p_0(x)} < k_\alpha \end{cases} \end{cases}$$

rejection region  
critical

reject if  $LR > k_\alpha$

[This test is called a likelihood ratio test]

(ii) Sufficiency. If a test satisfies (a) and (b) for some constant  $k_\alpha$ , it is most powerful for testing  $H_0: p_0$  vs  $H_1: p_1$  at level  $\alpha$  [ $\phi$  is MP].

(iii) Necessity. If a test  $\phi$  is MP at level  $\alpha$ , then it satisfies (b) for some  $k_\alpha$  and it also satisfies (a) unless there exists a test of size  $< \alpha$  with power 1.

(Example). Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , with  $\sigma^2$  known. Consider the following two hypotheses:

$$H_0: \mu = 0 \quad \text{vs} \quad H_1: \mu = \mu_1 \quad (\text{known}), \quad \mu_1 > 0$$

$$L(x) = \frac{p_1(x)}{p_0(x)} = \frac{\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_i - \mu)^2}{2\sigma^2}\right\}}{\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{x_i^2}{2\sigma^2}\right\}} = \frac{\exp\left\{-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}\right\}}{\exp\left\{-\frac{\sum_{i=1}^n x_i^2}{2\sigma^2}\right\}}$$

$$= \exp\left(\frac{1}{\sigma^2} \mu_1 \underbrace{\sum_{i=1}^n x_i}_{\text{suff. stat.}} - \frac{n\mu_1^2}{2\sigma^2}\right).$$

$\uparrow \text{const.}$

We have the following observation:

$$L(x) > k \Leftrightarrow \mu_1 \frac{\sum_{i=1}^n x_i}{\sigma^2} - \frac{n\mu_1^2}{2\sigma^2} > \log k$$

$$\Leftrightarrow \mu_1 \sum_{i=1}^n x_i > k' \quad \checkmark$$

$$\Leftrightarrow \begin{cases} \sum_{i=1}^n x_i > k'' & \text{if } \mu_1 > 0 \\ \sum_{i=1}^n x_i < k'' & \text{if } \mu_1 < 0. \end{cases}$$

Let's focus on the first case where  $\mu_1 > 0$ . We can rewrite our test in a different form as follows:

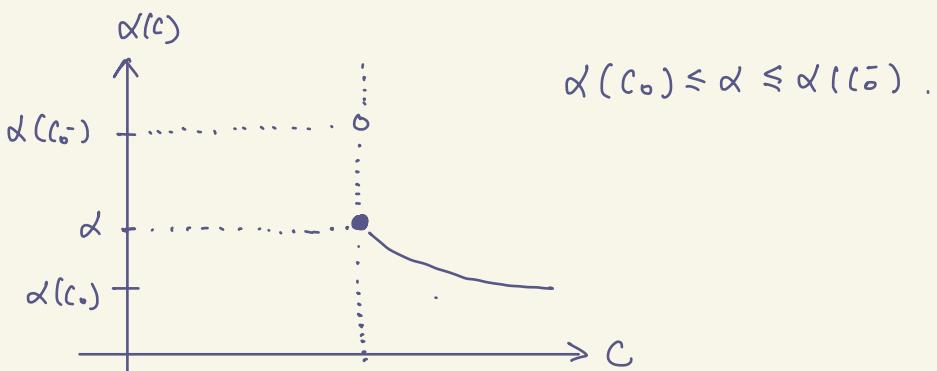
$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} > k'''.$$

By NP Lemma, we reject  $H_0$  iff  $\sqrt{n}(\bar{X}_n - \mu) / \sigma > k_\alpha$  is MP, where  $k''' = k_\alpha = \gamma_{1-\alpha}$  is the  $(1-\alpha)$  quantile of  $N(0,1)$ .

$$P_{\mu=0} \phi(X) = \alpha = P_{\mu=0} \left( \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} > k''' \right) < \gamma_{1-\alpha}.$$

Proof. Let  $r(x) = \frac{p_1(x)}{p_0(x)}$  be the likelihood ratio and denote the cumulative distribution function of  $r(X)$  under  $H_0$  by  $F_0$ .

(i) Existence. Let  $\alpha(c) = P_0(r(X) > c) = 1 - F_0(c)$ . Then  $\alpha(c)$  is a non-increasing, right continuous function of  $c$  (i.e.  $\lim_{\varepsilon \downarrow 0} \alpha(c+\varepsilon) = \alpha(c)$ ). Note that  $\alpha(c)$  is not necessarily left-continuous at every value of  $c$ , but the left-hand limits exist; denote them as  $\alpha(c^-) = \lim_{\varepsilon \downarrow 0} \alpha(c-\varepsilon)$



We define our test function to be

$$\phi(x) = \begin{cases} 1 & \text{if } r(x) > c_0 \\ \gamma & \text{if } r(x) = c_0 \\ 0 & \text{if } r(x) < c_0, \end{cases}$$

for some constant  $\gamma \in (0, 1)$ . We note that  $\phi$  only depends on  $x$  through  $r(x)$ . The test  $\phi$  always rejects if the likelihood ratio strictly exceeds the threshold  $c_0$  and never rejects it if the ratio falls strictly below  $c_0$ .

By construction, if we take  $R = c_0$ , then  $\phi$  satisfies condition (b) of part (i). We choose  $\gamma$  such that  $\phi$  also satisfies part (a). The size of  $\phi$  is given by

Type I error:

$$E_0(\phi(X)) = P_0(r(X) > c_0) + \gamma P_0(r(X) = c_0)$$

$$= \alpha(c_0) + \gamma \left\{ \frac{\alpha(c_0^-)}{\alpha(c_0^-) - \alpha(c_0)} - \frac{\alpha(c_0)}{\alpha(c_0^-) - \alpha(c_0)} \right\}$$

If  $\alpha(c_0^-) = \alpha(c_0)$ , then  $\alpha(c_0) = \alpha$  and we automatically have  $E_0(\phi(X)) = \alpha$  for any choice of  $\gamma$ .

otherwise, set  $\gamma = \frac{\alpha - \alpha(c_0)}{\alpha(c_0^-) - \alpha(c_0)}$

that gives  $E_0(\phi(X)) = \alpha$ .

- (ii) Let  $\phi$  satisfy (i)(a) & (b), and let  $\phi'$  be any other level  $\alpha$  test so that  $E_0(\phi'(X)) = \int \phi'(x) p_0(x) d\mu(x) \leq \alpha$ .

In order to show  $\phi$  is most powerful, we bound the power difference  $E_1(\phi(X)) - E_1(\phi'(X))$  from below by the size difference  $E_0(\phi(X)) - E_0(\phi'(X))$  up to a constant multiple.

To do this, we claim that the following holds:

$$\rightarrow \int \left\{ \begin{array}{c} \phi(x) - \phi'(x) \\ \text{tve} \end{array} \right\} \left\{ \begin{array}{c} p_1(x) - k p_0(x) \\ \text{tve} \end{array} \right\} d\mu(x) \geq 0. \quad (*)$$

To see this, we consider the following three cases:

- 1) If  $p_1(x) > k p_0(x)$ , then  $\phi(x) = 1$  because of our construction. Since  $\phi'(x) \leq 1$ , the integrand is non-negative.
- 2) If  $p_1(x) < k p_0(x)$ , then  $\phi(x) = 0$  by construction. Since  $\phi'(x) \geq 0$ , the integrand is non-negative too.
- 3) If  $p_1(x) = k p_0(x)$ , then the integrand is zero.

By exhaustion, the inequality  $(*)$  holds. Hence, we have

$$\begin{aligned} \int \{ \phi(x) - \phi'(x) \} p_1(x) d\mu(x) &\geq \underbrace{k \int \{ \phi(x) - \phi'(x) \} p_0(x) d\mu(x)}_{\text{red arrow}} \\ &= \underbrace{k \{ E_0(\phi(x)) - E_0(\phi'(x)) \}}_{\substack{\alpha \\ \leq \alpha}} \geq 0 \end{aligned}$$

We conclude that  $E_1(\phi(x)) \geq E_1(\phi'(x))$ ;  $\phi$  is the most powerful at level  $\alpha$ .

(iii) Necessity. Suppose  $\phi^*$  is most powerful at level  $\alpha$ , and let  $\phi$  be a likelihood ratio test satisfying (i)(a) & (b). We have to show that  $\phi^*(x) = \phi(x)$  except possibly when  $p_1(x)/p_0(x) = k$  for  $\mu$ -a.e.  $x$ . Define the sets

$$S^+ = \{x: \phi(x) > \phi^*(x)\}$$

$$S^- = \{x: \phi(x) < \phi^*(x)\}$$

$$S_0 = \{x: \phi(x) = \phi^*(x)\}$$

$$\text{and } S = (S^+ \cup S^-) \cap \{x: p_1(x) \neq k p_0(x)\}.$$

We want to show that  $\mu(S) = 0$ .

Suppose it is not the case, i.e.  $\underline{\mu(S)} > 0$ . As in part (ii), we have  $(\phi - \phi^*)(p_1 - k p_0) > 0$  on  $S$ . Therefore,

$$\begin{aligned} & \int_X (\phi - \phi^*)(p_1 - k p_0) d\mu(x) \\ &= \int_{S^+ \cup S^-} (\phi - \phi^*)(p_1 - k p_0) d\mu(x) \\ &= \int_S (\phi - \phi^*)(p_1 - k p_0) d\mu(x). \end{aligned}$$

By the hypothesis,  $E_0(\phi(X)) = \alpha$  and  $E_0(\phi^*(X)) \leq \alpha$ , so the previous inequality implies

$$E_1(\phi(X)) - E_1(\phi^*(X)) > k \{ E_0(\phi(X)) - E_0(\phi^*(X)) \} \geq 0$$

that is  $E_1(\phi(X)) > E_1(\phi^*(X))$ . Hence contradiction.

Hence  $\mu(S) = 0$ .

It remains to show that the size of  $\phi^*$  is  $\alpha$  unless there exists a test of size strictly less than  $\alpha$  and power 1. For this, note that if size  $< \alpha$  and power  $< 1$ , we can add points to the rejection region until either the size is  $\alpha$  or power is 1.  $\square$

Definition For simple  $H_0: p_0$  vs  $H_1: p_1$ , we call  $\beta_{\phi}(p_1)$  =  $E_{p_1}(\phi(X))$  the power of  $\phi$ , i.e., the prob. of rejecting  $H_0$  under the alternative hypothesis.

[Corollary] Suppose  $\beta$  is the power of a most powerful level  $\alpha$  test of  $H_0: p_0$  vs  $H_1: p_1$  with  $\alpha \in (0, 1)$ , then  $\alpha < \beta$  (unless  $p_0 = p_1$ )

(This MP test rejects more often under  $H_1$  than  $H_0$ )

• Proof. Consider the test  $\phi_0(x) \equiv \alpha$ , which always rejects with probability  $\alpha$ . Since  $\phi_0$  is level  $\alpha$  and  $\beta$  is the max power, we have

$$\rightarrow \beta \geq E_{P_1}(\phi_0(x)) = \alpha.$$

Suppose  $\beta = \alpha$ . Then  $\phi_0(x) = \alpha$  is a most powerful level  $\alpha$  test. As a result,

$$\phi_0(x) = \begin{cases} 1 & \text{if } p_1(x)/p_0(x) > k \\ 0 & \text{if } p_1(x)/p_0(x) < k \end{cases} \quad \text{by NP(iii)}$$

Since  $\phi_0(x)$  never equals 0 or 1, it must be the case that  $p_1(x) = k p_0(x)$  with prob. 1.

Note that  $\int p_1(x) d\mu(x) = k \int p_0(x) d\mu(x) = 1$ , thus implies that  $k=1$  and hence  $P_0 = P_1$ .  $\square$

## Exponential Families and UMP one-sided tests

↑  
Uniformly most powerful

Goal: Extend the MP tests for simple alternatives up to UMP tests for composite alternatives.

(Example): (One-parameter exponential family). Consider  $X_1, \dots, X_n$  iid  $P_\theta$ , where  $p_\theta(x) \propto h(x) \exp(\theta T(x))$ , and we are interested in testing

$$H_0: \theta = \theta_0 \quad \text{vs} \quad H_1: \theta = \theta_1$$

We want to construct an MP test at level  $\alpha$ . The Crr. likelihood ratio is

$$\frac{\prod_{i=1}^n p_{\theta_1}(x_i)}{\prod_{i=1}^n p_{\theta_0}(x_i)} \propto \exp \left\{ \underline{\underline{(}\theta_1 - \theta_0)\sum_{i=1}^n T(x_i)} \right\}$$

Assuming that  $\theta_1 > \theta_0$ , we shall reject for large value of  $\sum_{i=1}^n T(x_i)$ .

That is, an UMP test has the following form:

$$\Phi(x) = \begin{cases} 1 & \text{if } \sum_{i=1}^n T(x_i) > k \\ \gamma & \text{if } \sum_{i=1}^n T(x_i) = k \\ 0 & \text{if } \sum_{i=1}^n T(x_i) < k \end{cases}$$

$H_0: \theta_1 < \theta_0$

where  $k$  and  $\gamma$  are chosen to satisfy the size constraint

$$\alpha = E_{\theta_0}(\Phi(X)) = P_{\theta_0}\left(\sum_{i=1}^n T(x_i) > k\right) + \gamma P_{\theta_0}\left(\sum_{i=1}^n T(x_i) = k\right)$$

Note that  $\sum_{i=1}^n T(x_i)$  has no explicit  $\theta$  dependence and that  $k$  and  $\gamma$  do not depend on  $\theta_1$  (with  $\theta_1 > \theta_0$ ). This means that  $\Phi$  is in fact a UMP test for testing

$$H_0: \theta = \theta_0 \quad \text{vs} \quad H_1: \theta > \theta_0. \quad \boxed{\theta \neq \theta_0}$$

Hence  $H_1$  is an example of a one-sided alternative, which arises when the parameter values of interest lie on only one side of the real-valued parameter  $\theta_0$ .

$\downarrow$   
Exponential family

Definition (Families with monotone likelihood ratio, MLR)

We say that the family of densities  $\{p_\theta: \theta \in \mathbb{R}\}$  has monotone likelihood ratio in  $T(X)$  if

(i)  $\theta \neq \theta'$  implies  $p_\theta \neq p_{\theta'}$  [identifiability].

(ii)  $\theta < \theta'$  implies  $p_{\theta'}(x)/p_\theta(x)$  is a non-decreasing function of  $T(x)$  [monotonicity]

(Example)

- ✓ 1)  $T(X) = \sum_{i=1}^n T(x_i)$  in the one-parameter exp. family.
- 2) (Double exponential). Let  $X \sim$  Double exponential ( $\theta$ ) with density  $p_\theta(x) = \frac{1}{2} \exp(-|x-\theta|)$ . It is easy to see that the model is identifiable, so we need to check only the 2nd condition.

Fix any  $\theta' > \theta$ , and consider the likelihood ratio

$$\frac{p_{\theta'}(x)}{p_\theta(x)} = \exp(|x-\theta| - |x-\theta'|).$$

Note that

$$|x-\theta| - |x-\theta'| = \begin{cases} \theta - \theta' & \text{if } x < \theta \\ 2x - \theta - \theta' & \text{if } \theta \leq x \leq \theta' \\ \theta' - \theta & \text{if } x > \theta' \end{cases},$$

which is non-decreasing in  $x$ . Therefore, the family has MLR in  $T(X) = X$ .

(Example) (Cauchy location model). Let  $X$  has density

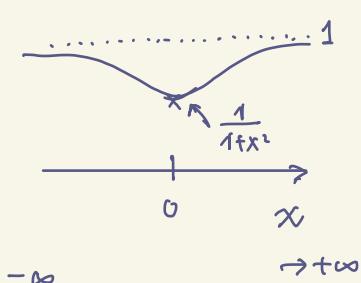
$$p_\theta = \frac{1}{\pi} \cdot \frac{1}{1+(x-\theta)^2}.$$

We find two points for which the MLR cond. fails.

For any fixed  $\theta > 0$ ,

$$\frac{p_\theta(x)}{p_0(x)} = \frac{1+x^2}{1+(\theta-x)^2} \rightarrow 1 \text{ as } x \rightarrow \pm\infty,$$

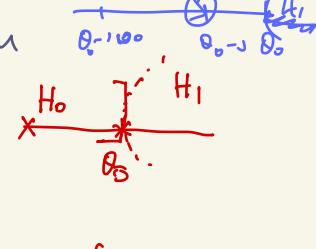
But  $p_\theta(0)/p_0(0) = \frac{1}{1+\theta^2} < 1$ . Thus the ratio must increase at some values of  $x$  and decrease at others. In particular, it is not monotone in  $x$ . Hence, we conclude that the LR in  $T(X) = X$  is not MLR.



[THEOREM] Suppose  $X \sim p_0$  has MLR in  $T(X)$  and we test  $H_0: \underline{\theta \leq \theta_0} \Rightarrow \theta = \theta_0 - \delta$  vs  $H_1: \theta > \theta_0$ , then

(ii) There exists a UMP test at level  $\alpha$  of the form

$$\phi(X) = \begin{cases} 1 & \text{if } T(X) > k \\ \gamma & \text{if } T(X) = k \\ 0 & \text{if } T(X) < k \end{cases}$$

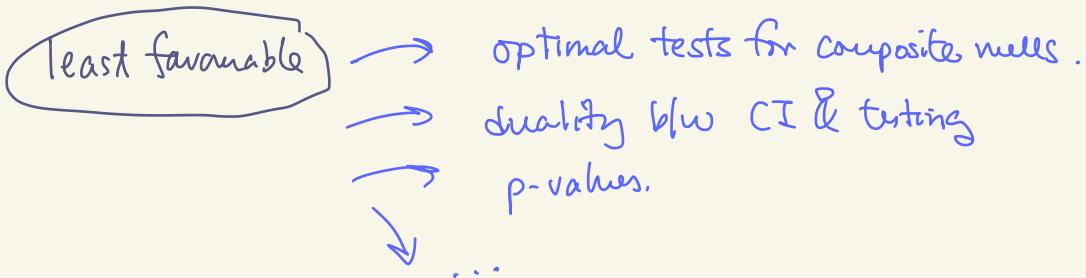


where  $k, \gamma$  are determined by the condition  $E_{\theta_0}(\phi(X)) = \alpha$

(ii) The power function  $\beta(\theta) = E_\theta(\phi(X))$  is strictly increasing when  $0 < \beta(\theta) < 1$ . Proof  $\rightarrow$  TSH.

Possible directions for constructing UMPs:

- 1) Reduce the composite alternative to a simple alternative
- \* 2) Collapse the composite null to a simple null.
- 3) Apply NP Lemma.



## Lecture 12

28/Nov/2022.

### Optimal Tests for Composite Nulls

$$H_0: X \sim f_\theta \quad , \quad \theta \in \Omega_0$$

$$H_1: X \sim g \quad , \quad \text{where } g \text{ is known.}$$

We may impose a prior distribution  $\Lambda$  on  $\Omega_0$ . So, we consider the new hypothesis

$$\rightarrow H_{\underline{\Lambda}}: X \sim h_{\underline{\Lambda}}(x) = \int_{\Omega_0} f_\theta(x) d\Lambda(\theta), \text{ "Simple"}$$

where  $h_{\underline{\Lambda}}(x)$  is the marginal distribution of  $X$  induced by  $\Lambda$ .

Let  $\beta_{\underline{\Lambda}}$  be the power of the MP level- $\alpha$  test  $\phi_{\underline{\Lambda}}$  for testing  $H_{\underline{\Lambda}}$  vs  $g$ .

(§3.8 of TSH)

Definition (Least favorable distribution)  $\Lambda$  is a least favorable distribution if  $\beta_{\underline{\Lambda}} \leq \beta_{\underline{\Lambda}'}$  for any prior  $\Lambda'$ .

[THEOREM] Suppose  $\phi_{\underline{\Lambda}}$  is an MP level  $\alpha$ -test for testing  $H_{\underline{\Lambda}}$  against  $g$ .

TSH 3.8.1

If  $\phi_{\underline{\Lambda}}$  is level- $\alpha$  for the original hypothesis  $H_0$

(i.e.  $F_{\theta_0}(\phi_{\underline{\Lambda}}) \leq \alpha \quad \forall \theta_0 \in \Omega_0$ ), then

- 1) The test  $\phi_{\underline{\Lambda}}$  is MP for the original test on  $H_0: \theta \in \Omega_0$  vs  $g$ .
- 2) The prior distribution  $\Lambda$  is least favorable.

• Proof. Let  $\phi^*$  be any other level  $\alpha$  test of  $H_0: \theta \in \Omega_0$  vs  $g$ .

Then  $\phi^*$  is also a level- $\alpha$  test for  $H_{\underline{\Lambda}}$  vs  $g$  because

$$E_\theta(\phi^*(x)) = \int \phi^*(x) f_\theta(x) d\mu(x) \leq \alpha \quad \forall \theta \in \Omega_0$$

which implies that

$$\begin{aligned} \int \phi^*(x) h_{\underline{\Lambda}}(x) d\mu(x) &= \int \phi^*(x) f_\theta(x) d\mu(x) \underbrace{d\Lambda(\theta)}_{\text{---}} \leq \int \alpha d\Lambda(\theta) \\ &= \alpha \end{aligned}$$

Since  $\phi_\Lambda$  is MP for  $H_1$  vs  $g$ , we have

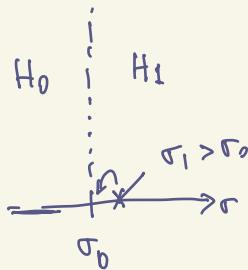
$$\int \phi^* g(x) d\mu(x) \leq \int \phi_\Lambda(x) g(x) d\mu(x)$$

Hence  $\phi_\Lambda$  is a MP test for  $H_0$  vs  $g$  because  $\phi_\Lambda$  is also level  $\alpha$ .

2. Let  $\Lambda'$  be any distribution on  $\mathcal{Q}_0$ . Since  $E_{\theta}(\phi_\Lambda(x)) \leq \alpha \forall \theta \in \mathcal{Q}_0$ , we know that  $\phi_{\Lambda'}$  must be Level- $\alpha$  for  $H_{\Lambda'}$  vs  $g$ . Thus  $\beta_{\Lambda'} \leq \beta_\Lambda$ . So  $\Lambda$  is the least favourable distribution.  $\square$

(Example) (Testing in the presence of nuisance parameters) let  $X_1, \dots, X_n$  be iid  $N(\theta, \sigma^2)$  where both  $\theta$  and  $\sigma^2$  are unknown. We consider testing  $H_0: \sigma \leq \sigma_0$  against  $H_1: \sigma > \sigma_0$ .

To find a UMP test, we follow the strategy discussed in the previous lecture:



- 1) First, we fix a simple alternative  $(\theta_1, \sigma_1)$  for some arbitrary  $\theta_1$  and  $\sigma_1 > \sigma_0$ .
- 2) We choose a prior distribution  $\Lambda$  to collapse our null hypothesis over. Intuitively, the least favourable prior should make the alternative hypothesis hard to distinguish.

One sensible choice is concentrating  $\Lambda$  on the boundary between  $H_1$  and  $H_0$  (i.e. the line  $\{\sigma = \sigma_0\}$ ). Thus  $\Lambda$  will be a probability distribution over  $\theta \in \mathbb{R}$  for the fixed  $\sigma = \sigma_0$ .

Given any test function  $\phi(x)$  and a sufficient statistic  $T$ , there exists a test function  $\eta$  that has the same power as  $\phi$  but depends on  $x$  only through  $T$ :

$$\eta(T(X)) = E(\phi(X) | T(X)).$$

Hence, we can restrict our attention to the sufficient statistic  $(Y, U)$ , where  $Y = \bar{X}$  and  $U = \sum_{i=1}^n (X_i - \bar{X})^2$ . We know that  $Y \sim N(\theta, \sigma^2/n)$ ,  $U \sim \sigma^2 \chi_{n-1}^2$ , and  $Y$  is independent of  $U$  by Basu's theorem.

Thus, for  $\Lambda$  supported on  $\sigma = \sigma_0$ , we have the joint density of  $(Y, U)$  under  $H_\Lambda$  as

$$C_0 U^{\frac{n-3}{2}} \exp\left(-\frac{U}{2\sigma_0^2}\right) \int \exp\left(-\frac{n}{2\sigma_0^2}(y-\theta)^2\right) d\Lambda(\theta) \quad (*)$$

and the joint density under the alternative hypothesis  $(\theta_1, \sigma_1)$  as

$$C_0 U^{\frac{n-3}{2}} \exp\left(-\frac{U}{2\sigma_1^2}\right) \exp\left(-\frac{n}{2\sigma_1^2}(y-\theta_1)^2\right). \quad (**)$$

We see from (\*) and (\*\*) that the choice of  $\Lambda$  only affects the distribution of  $Y$ . To achieve the maximum power against the alternative, we need to choose  $\Lambda$  such that the two distributions become as close as possible.

Under  $H_0$ :  $Y \sim N(\theta_0, \frac{\sigma_0^2}{n})$ . Under  $H_\Lambda$ , the distribution of  $Y$  is in a convolution form:  $Y = Z + \Theta$  for  $Z \sim N(0, \frac{\sigma_0^2}{n})$  and  $\Theta \sim \Lambda$ , where  $\Theta$  and  $Z$  are independent.

Hence, if we choose  $\Theta \sim N(\theta_1, \frac{\sigma_1^2 - \sigma_0^2}{n})$ ,  $Y$  will have the same distribution under the null and the alternative, which is  $N(\theta_1, \frac{\sigma_1^2}{n})$ . Under this choice of prior, the LRT rejects for large values  $\exp\left(-\frac{U}{2\sigma_1^2} + \frac{U}{2\sigma_0^2}\right)$ , i.e. rejects for large values of  $U$  (since  $\sigma_1 > \sigma_0$ ). So, the MP test rejects  $H_\Lambda$  if  $\sum_{i=1}^n (X_i - \bar{X})^2$  lies above some threshold determined by the size constraint. In particular, it rejects if  $\sum_{i=1}^n (X_i - \bar{X})^2 > \sigma_0^2 C_{n-1, 1-\alpha}$ , where  $C_{n-1, 1-\alpha}$  is the  $(1-\alpha)^{th}$  quantile of  $\chi_{n-1}^2$ .

- 3) Next, we check if the MP test is level  $\alpha$  for the composite null. For any  $(\theta, \sigma)$  with  $\sigma \leq \sigma_0$ , the prob. of rejection is

$$P_{\theta, \sigma} \left( \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} > \frac{\sigma_0^2 C_{n-1, 1-\alpha}}{\sigma^2} \right) = P \left( \chi_{n-1}^2 > \frac{\sigma_0^2}{\sigma^2} C_{n-1, 1-\alpha} \right)$$

while the equality holds iff  $\sigma = \sigma_0$

$$\leq \alpha$$

Hence from TSH 3.8.1 that our test is MP for testing the original null  $H_0$  vs  $N(\theta_0, \sigma_0)$ .

- 4) Finally, the MP level  $\alpha$  test for testing the composite null  $H_0$  versus an arbitrarily chosen  $(\theta_0, \sigma_0)$  does not depend on the choice of  $(\theta_0, \sigma_0)$ . Hence, it is UMP for testing the original composite null versus the composite alternative.

||

### Duality between Testing and Interval Estimation.

(interval)

Recall that a random set  $S(X)$  is a  $1-\alpha$  confidence region for a parameter  $\xi = \xi(\theta)$  if

$$P_\theta (\xi \in S(X)) \geq 1-\alpha \quad \forall \theta \in \Omega$$

For every  $\xi_0$ , let  $A(\xi_0)$  be the acceptance region for non-randomised level  $\alpha$  test of  $H_0: \xi(\theta) = \xi_0$  versus  $H_1: \xi(\theta) \neq \xi_0$  so that

$$P_\theta (X \in A(\xi_0)) \geq 1-\alpha, \quad \forall \theta \in \Omega$$

Define  $S(X) = \{ \xi : X \in A(\xi) \}$

then  $\xi(\theta) \in S(X)$  if and only if  $X \in A(\xi(\theta))$ , so

$$P_\theta (\xi(\theta) \in S(X)) = P_\theta (X \in A(\xi(\theta))) \geq 1-\alpha.$$

This shows that  $S(X)$  is a  $1-\alpha$  confidence interval for  $\xi$ .

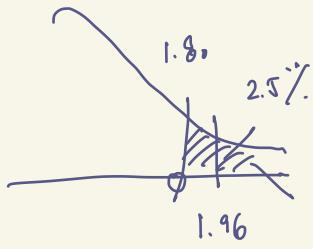
§3.3  
TSH.

### p-values

For varying  $\alpha$ , the resulting tests provide an example of the typical situation in which the rejection regions  $S_\alpha$  are nested in the sense that

the quantity

$$S_\alpha \subset S_{\alpha'} \text{ if } \alpha < \alpha'.$$



$$\hat{p} = \inf \{\alpha : X \in S_\alpha\}$$

↑  
Smallest significant level

is called a p-value.

It gives an idea of how strongly the data contradict the null hypothesis.

(Example) Let  $\Phi$  denote the standard normal cdf. Then the rejection region can be written as

$$S_\alpha = \left\{ X : X > \sigma z_{1-\alpha} \right\} = \left\{ X : \Phi\left(\frac{X}{\sigma}\right) > 1-\alpha \right\} = \left\{ X : 1 - \Phi\left(\frac{X}{\sigma}\right) < \alpha \right\}.$$

For a given observed value of  $X$ , the inf over all  $\alpha$  where the last inequality holds is  $\hat{p} = 1 - \Phi\left(\frac{X}{\sigma}\right)$ .

Alternatively, the p-value is  $P_0(X \geq x)$ , where  $x$  is the observed value of  $X$ . Note that under  $H_0$ , the distribution of  $\hat{p}$  is given by

$$P_0(\hat{p} \leq u) = P_0\left(1 - \Phi\left(\frac{X}{\sigma}\right) \leq u\right) = P_0\left(\Phi\left(\frac{X}{\sigma}\right) \geq 1-u\right) = u,$$

$\Rightarrow \hat{p}$  is uniformly distributed on  $(0, 1)$ .

||

(Lemma 3.3.1 of TSH + Example 3.3.2)

discrete

↑ EOF

Things to know..

### UMP $X$

- impose a reasonable restriction on the tests to be considered and the "optimal" tests within the class of tests under the restriction.
- Two typical strategies:  $\begin{cases} \textcircled{1} \text{ unbiasedness } \rightarrow \text{UMP unbiased test.} \\ \textcircled{2} \text{ invariance.} \end{cases}$

Recall that a UMP test  $\phi$  of size  $\alpha$  has the property that

$$E_\theta(\phi(X)) \leq \alpha \quad \theta \in \Omega_0 \quad \text{and} \quad E_\theta(\phi(X)) \geq \alpha \quad \theta \in \Omega_1. \quad (\dagger)$$

This means that  $\phi$  is at least as good as the naive test  $\phi = \chi$ .

Definition Let  $\alpha$  be a given level of significance. A test  $\phi$  for  $H_0: \theta \in \Omega_0$  versus  $H_1: \theta \in \Omega_1$  is said to be unbiased of level  $\alpha$  if and only if (t) holds. A test of size  $\alpha$  is called a <sup>and</sup> uniformly most powerful unbiased (UMPU) test if it is UMP within the class of unbiased tests of level  $\alpha$ .

(UMP) is UMPU

uniform.  $\theta = \theta_0$

$\theta \neq \theta_0$

$\left\{ \begin{array}{l} \theta > \theta_0 \\ \theta < \theta_0 \end{array} \right. \rightarrow \begin{array}{l} \rightarrow \\ \leftarrow \end{array}$

:

(See Keener § 12.7, Thm 12.26 & Example 12.19)

Exponential dist.

two-sided

$H_0: \theta = 1$  vs  $H_1: \theta \neq 1$

### Final Exam

12 / Dec. 12:30 - 15:00 (same classroom)

Coverage: Everything  
combining "2" cheat sheets

### Consistency / Asymptotic Normality of MLE

(Keener ch. 9) → EE, EM algo.

$\{\tilde{\theta}_n\}$  for  $\theta \in \Omega$ . weakly consistent :  $\tilde{\theta}_n \xrightarrow{P} \theta_0$   
 $\equiv$  strongly consistent :  $\tilde{\theta}_n \xrightarrow{a.s.} \theta_0$

Let  $X_1, \dots, X_n$  be iid with density  $f(x|\theta)$ . w.r.t. some  $\sigma$ -finite measure  $\nu$ , where  $\theta \in \Omega$ .

The likelihood function :  $L_n(\theta) = L_n(\theta | x_1, \dots, x_n) = \prod_{i=1}^n f(x_i | \theta)$

A MLE of  $\theta$  (denoted by  $\hat{\theta}$ ) is defined as  $\hat{\theta} = \underset{\theta \in \Omega}{\operatorname{argmax}} L_n(\theta)$ .

Kullback-Leibler information number is defined as

$$K(f_0, f_1) = E_0 \left( \log \frac{f_0(x)}{f_1(x)} \right) = \int \log \frac{f_0(x)}{f_1(x)} f_0(x) d\nu(x)$$

[Lemma] (Shannon-Kolmogorov Information Inequality). Let  $f_0(x)$  and  $f_1(x)$  be densities w.r.t.  $\nu$ . Then

$$K(f_0, f_1) = E_0 \left( \log \frac{f_0(x)}{f_1(x)} \right) \geq 0,$$

with equality holds if and only if  $f_1(x) = f_0(x)$  (a.e.  $d\nu$ ).

(Proof: Skipped)  $\rightarrow$  use Jensen's inequality.

Let  $\theta_0$  denote the true value of  $\theta$ . Then, the MLE is the value of  $\theta$  that maximises

$$\begin{aligned} l_n(\theta) - l_n(\theta_0) &= \log L_n(\theta) - \log L_n(\theta_0) \\ &= \sum_{j=1}^n \left( \log f(x_j | \theta) - \log f(x_j | \theta_0) \right). \end{aligned}$$

By the SLLN and the Lemma (sk)

$$\frac{1}{n} \log \frac{L_n(\theta)}{L_n(\theta_0)} = \frac{1}{n} \sum_{j=1}^n \log \frac{f(x_j | \theta)}{f(x_j | \theta_0)} \xrightarrow{\text{a.s.}} E_\theta \log \frac{f(x | \theta)}{f(x | \theta_0)} = \underline{-K(\theta_0, \theta)} \leq 0$$

unless  $f(x | \theta) = f(x | \theta_0)$ .

$\rightarrow \frac{L_n(\theta)}{L_n(\theta_0)}$  (LR) converges exp: fast to zero at the rate  $\exp(-nK(\theta_0, \theta))$

This implies already that if  $\Omega$  is finite, the MLE is strongly consistent.

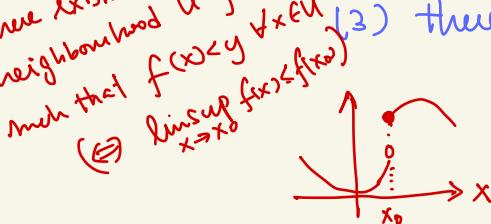
[THEOREM] let  $X_1, \dots, X_n$  be iid with density  $f(x | \theta)$ ,  $\theta \in \Omega$ , and let  $\theta_0$  denote the true value of  $\theta$ . If

& fn  $f: X \rightarrow \bar{\mathbb{R}}$  is  
called USC at  $x_0$  if  
for every real  $y > f(x_0)$ ,

(1)  $\Omega$  is compact

there exists a neighbourhood  $U$  of  $x_0$  such that  $f(x) < y \forall x \in U$   $\leftarrow$  (2)  $f(x | \theta)$  is upper semicontinuous in  $\theta$  for all  $x$

such that  $f(x) < y \forall x \in U$   $\leftarrow$  (3) there exists a function  $K$  such that  $E_\theta |K(X)| < \infty$



and  $U(x, \theta) = \log f(x|\theta) - \log f(x|\theta_0) \leq k(x)$ , for all  $x \in \mathcal{X}$

(4) for all  $\theta \in \Sigma$  and sufficiently small  $\epsilon > 0$

$\sup_{|\theta' - \theta| < \epsilon} f(x|\theta')$  is measurable in  $\mathcal{X}$

(5) (Identifiability)  $f(x|\theta) = f(x|\theta_0)$  (a.e. d $\mu$ )  $\Rightarrow \theta = \theta_0$ .

then,  $\hat{\theta}_n \xrightarrow{a.s.} \theta_0$ . (See proof. in Keener for example) ...

$$\begin{aligned} \text{[AN]} \quad & \Psi(x, \theta) = \frac{\partial}{\partial \theta} \log f(\theta|x) \\ & \text{In CLT.} \quad \text{id} \quad \text{in} \quad \text{MC} \quad \text{+} \quad \text{+} \\ & \hat{\theta}(x, \theta) = \theta_0 + \Psi(x, \theta_0)^T (\theta - \theta_0) \\ & + (\theta - \theta_0)^T \int_0^1 \int_0^1 \lambda \Psi(x, \theta_0 + \lambda \mu(\theta - \theta_0)) d\lambda d\mu(\theta - \theta_0) \\ & = \sqrt{n} (\hat{\theta} - \theta_0) \quad \text{remainder term...} \\ & \text{CLT. } N \quad \text{near } 0. \\ & \text{Normal.} \\ & \boxed{N(0, I(\theta_0)^{-1})} \quad (\text{Cramér}) \quad \boxed{\sqrt{n} (\hat{\theta}_n - \theta_0) \xrightarrow{d} N(0, I(\theta_0)^{-1})} \end{aligned}$$

