# Inference Problem.

(i) You are given a collection of prob. measures

$$\{P_\theta : \theta \in \Theta\} \text{ on a sample space } (X, \mathcal{F})$$

where $X$ is a set and $\mathcal{F}$ is a $\sigma$-field on $X$.

(ii) Observe $X \sim P_\theta$ for some $\theta \in \Theta$

(iii) Infer $\theta$ from $X$.

Let $L(\theta, \boxed{\delta(X)})$ be the loss in estimating $\theta$ by $\delta(X)$, an estimator.

Define $\boxed{R(\theta, \delta) = E_{X \sim P_\theta} \left\{ L(\theta, \boxed{\delta(X)}) \right\}}$ to be the risk function of the estimator $\delta$.

e.g. $X_1, \ldots, X_n \overset{iid}{\sim} N(\theta, 1), \theta \in \mathbb{R}, X \sim \mathbb{R}^n$

$$\boxed{X = (X_1, \ldots, X_n)} \qquad \underset{\underset{P_\theta(X \in A)}{\parallel}}{P_\theta(A)} = \frac{1}{(\sqrt{2\pi})^n} \int_A e^{-\sum_{i=1}^{n} \frac{(X_i - \theta)^2}{2}} \, dx_1, \ldots dx_n.$$

$$L(\theta, \delta(X)) = \left\{ \theta - \delta(X) \right\}^2$$

Two proposed estimator $\begin{cases} \delta_1(X) = \bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i \\ \delta_2(X) = 0 \end{cases}$

Correspondingly,

$$\boxed{R(\theta, \delta_1) = E_\theta (\bar{X} - \theta)^2 = \frac{1}{n}} \text{ and } R(\theta, \delta_2) = E_\theta(\theta^2) = \theta^2$$

## Strategies

### Strategy 1

**Def** We say $\delta(X)$ is unbiased for $\theta$ if $E_{X \sim P_\theta} \delta(X) = \theta, \forall \theta \in \Theta$

For the previous example,

$$E_\theta \left\{ \delta_1(X) \right\} = E_\theta \left\{ n^{-1} \sum_{i=1}^{n} X_i \right\} = \theta \qquad E_\theta \left\{ \delta_2(X) \right\} = 0$$

Later on, we shall show that $\delta_1$ is the "best" amongst the class of all unbiased estimators in the problem.

### Strategy 2 (Minimax)

We shall look at $\sup_{\theta \in \Theta} R(\theta, \delta)$ for compassion.

In our example, $\sup_{\theta \in \Theta} R(\theta, \delta_1) = \frac{1}{n}, \quad \sup_{\theta \in \Theta} R(\theta, \delta_2) = +\infty$

### Strategy 3 (Bayes)

Assume $\theta$ is random and has a distribution. $\pi$.

We may compare the __Bayes risk__, which is $\boxed{E_{\theta \sim \pi} \{R(\theta, \delta)\}}$

In our case, let $\pi \sim N(\mu, \tau)$

Bayes risk of $\delta_1$, is $E_{\theta \sim \pi} \{R(\theta, \delta_1)\} = E_{\theta \sim \pi} \{\frac{1}{n}\} = \frac{1}{n}$

$\qquad\qquad \delta_2, \qquad \cdots \qquad \delta_2 = E_{\theta \sim \pi} \{\theta^2\} = \mu^2 + \tau$

$\boxed{\text{Strategy 4}}$  What happens when $n$ is large ?

In this case, by WLLN, $\delta_1(X) = \bar{X}_n \xrightarrow{P} \theta$.

$\qquad\qquad\qquad \delta_2(X) \xrightarrow{P} 0$.

__Asymptotic optimality__  (to be learnt)


__Chap 1__  __Probability & Measure__

C & B.

__Def__  The set $S$ of all possible outcomes of a particular experiment is called the __sample space__
for the experiment.

e.g.  $S = \{H, T\}$ (coin flip)  ;  $S = [0, \infty)$  (stock price)

__Def__  Event $=$ any subset of $S$ including $S$ itself.

equality :  $A \subset B$,  $B \subset A$  $\Leftrightarrow$  $A = B$

union : $\cup$    intersection : $\cap$

De Morgan's Laws :  $(A \cup B)^c = A^c \cap B^c$

$\qquad\qquad\qquad\qquad (A \cap B)^c = A^c \cup B^c$

__Def__  Two events $A$ and $B$ are __disjoint__ (or mutually exclusive) if $A \cap B = \phi$.

For the disjoint events $A_i$ $(i = 1, \ldots, n)$,  $A_i \cap A_j = \phi$ for $i \neq j$.

Furthermore, if $\bigcup_{i=1}^{\infty} A_i = S$, the collection of $A_i$ $(i \in \mathbb{Z}^+)$ forms a partition of $S$

__Prob. Theory__    Remark  $(\bigcup_{i=1}^{\infty} A_i^c)^c = \bigcap_{i=1}^{n} A_i \in \mathcal{B}$

__Def__  A measure $\mu$ and on a set $S$ assigns a non-negative value $\mu(A)$ to a subset of $S$

__Def__  A collection of subsets of $S$ is called a $\sigma$-algebra, denote by by $\mathcal{B}$, or $\sigma(S)$, if

a.  $\phi \in \mathcal{B}$

b.  $A \in \mathcal{B} \Rightarrow A^c \in \mathcal{B}$      (closure under complementation)

c.  If $A_i \in \mathcal{B}$ for $i = 1, 2, \ldots$ , then $\bigcup_{i=1}^{\infty} A_i \in \mathcal{B}$.   (closure under countable unions)

Def A function $\mu$ on a $\sigma$-field of $\mathcal{A}$ is a **measure** if

a. For every $A \in \mathcal{A}$, $0 \le \mu(A) \le \infty$, $\mu: \mathcal{A} \to [0, \infty]$

b. If $A_i$ are disjoint, $A_i \in \mathcal{A}$, then

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mu(A_i)$$

c. $\mu(\Omega) = 1$ (prob. measure)

## 10 ways looking at a random variable

1. A function $X: \Omega \to \mathbb{R}$ such that images $X^{-1}(B)$ for any Borel set are elements of $\mathcal{A}$ is called a <u>random variable</u>. A p-tuple of r.v.'s is called a random vector.

2. Associated with a random vector of $X$ on $(\Omega, \mathcal{A}, P)$ is a distribution function d.f.

$$F_{\underline{X}}(\underline{X}) = F_{X_1, X_2, \dots X_p}(x_1, x_2, \dots x_p) = P(w: X_1(w) \le x_1, \dots X_p(w) \le x_p)$$

Remark: $F$ is right continuous. càdlàg.

3. For any scalar r.v. $X$ with d.f. $F$, the quantity $Q(u) = F^{-1}(u) = \inf\{x: F(x) \ge u\}$

is called the $u^{th}$ quantile of $X$ (as of $F$), $u \in (0,1)$

$Q(1/2) = $ median. $Q(1/4)$, $Q(3/4) = $ Q1. lower quantile. Q3. upper quantile.

4. If the df $F$ is absolutely continuous with respect to the measure $\mu$, then $F$ has a density $f$ w.r.t. $\mu$. $\quad F(x) = \int_{-\infty}^{x} f(u) du. \quad f(x) = F'(x)$

5. The expectation of a r.v. $X$ is $\quad E(X) = \int_{\Omega} X(w) dP(w) = \int_{-\infty}^{\infty} x \, dF(x) = \int_{-\infty}^{\infty} x \underset{\underset{\text{R.N. derivative}}{\downarrow}}{f(x)} dx$

Likewise, define the expectations of functions of $X$

$$E\{g(X)\} = \int_{\Omega} g(X(w)) dP(w) = \int g(x) dF(x)$$

e.g. $g(x) = I(x \in B) = \begin{cases} 1 & \text{if } x \in B \\ 0 & \text{ow.} \end{cases} \quad E\{I(x \in B)\} = \int_B dP(w) = P(B).$

6. Moments. Moments of higher powers of $X - \mu$ are often used to describe the basic characteristics of the distributions of r.v.'s. In particular, we denote

$$\mu_k = E(X - EX)^k, \quad k = 1, 2, \dots \qquad \text{as the } k\text{th central moment of } X.$$

e.g. $k = 2$, $\mu_2$: $Var(X) = E(X - EX)^2 = \sigma^2$ [variance].

$k = 3$, $\dfrac{\mu_3}{\sigma^3}$: skewness $(X)$ [asymmetry]

$k = 4$ $\dfrac{\mu_4}{\sigma^4}$: kurtosis $(X)$

# 7. Moment Generating Function mgf.

To compute moments, it is often convenient to use the mgf, which is defined as.

$$M_X(t) = E\{\exp(tX)\} = \int e^{tx} dF(x) \qquad \text{Laplace transform}$$

where $M_X(t)$ exists and its derivative exists in some neighborhood of $0$.

Essentially, we have $\quad V_k = m_X^{(k)}(0) = E(X^k) \quad , \quad k = 0, 1, 2, \dots$

The property hold :

(a) For constant $(\mu, \sigma)$ $\quad : \quad m_{\mu + \sigma X}(t) = \exp(\mu t) \, m_X(\sigma t)$

(b) For indep. $X, Y$ $\quad : \quad m_{X+Y}(t) = m_X(t) \, m_Y(t)$.

Example   Suppose we have a discrete r.v. on $\{0, 1, 2, \dots\}$ with

$$P(X = j) = a_j$$

We define the "generating function" of $X$ as $\quad g(z) = \sum_{j \geq 0} a_j z^j = \sum_{j \geq 0} P(X=j) z^j$

Since $\sum_{j \geq 0} a_j = 1$, it is clear that

$$|g(z)| \leq \sum_j |a_j| |z|^j \leq \sum_j a_j = 1 \quad \text{for } |z| \leq 1. \checkmark$$

Consider the derivatives.

$$g'(z) = a_1 + 2a_2 z + 3a_3 z^2 + \cdots = \sum_{n=1}^{\infty} n a_n z^{n-1}$$

$$g''(z) = 2a_2 + 6a_3 z + \cdots = \sum_{n=1}^{\infty} n(n-1) a_n z^{n-2}$$

$$\vdots$$

$$g^{(j)}(z) = \sum_{n=j}^{\infty} (n(n-1)\cdots(n-j+1)) a_n^{n-j} = \sum_{n-j} \binom{n}{j}(j!) a_n z^{n-j}$$

Thus, $\quad g^{(j)}(0) = j! \, a_j \quad$ or $\quad a_j = (j!)^{-1} g^{(j)}(0)$.

So all the information about the $a_j$'s are contained within the function $g$ and is made accessible, by differentiating and evaluating $g^{(k)}$ at $0$.

Suppose the moments exist, then

$$g'(1) = \sum_{n=0}^{\infty} n a_n = \sum_{n=0}^{\infty} n P(X=n) = EX.$$

$$g''(1) = \sum_{n=0}^{\infty} n^2 a_n - \sum_{n=0}^{\infty} n a_n = E X^2 - (EX)^2$$

The distribution of a non-negative integer valued r.v. is uniquely determined by its generating function. $\quad a_j = (j!)^{-1} g^{(j)}(0)$.

We can write in a slightly fancy notation:

$$g(X) = E(z^X) = E(e^{-\lambda X})$$

if $X$ takes arbitrary real values, and consider $0 < j \leq 1$, any such

$z$ can be written as $e^{-\lambda}$, for $0 \leq \lambda < \infty$

$$E(e^{-\lambda X}) = \sum_{j=0}^{\infty} p_j e^{-j\lambda}$$

$\downarrow$       $\longrightarrow$ prob. that $X$ takes the value of $x_j$.

$$E(e^{-\lambda X}) = \int e^{-\lambda u} f(u)\, du.$$

8. characteristic function:

$$\phi_X(\lambda) \triangleq E[e^{i\lambda X}] = \int e^{i\lambda u} f(u)\, du \qquad \leftarrow \text{characteristic function.}$$

$$|E(e^{itX})| \leq E|e^{itX}| = E|\cos t X + i \sin t X|$$

$$= E\left(\cos^2 t X + \sin^2 t X\right)$$

$$= 1.$$

see e.g. 2.3.10 _Non unique moments_

9. Cumulants.

$$K_X(t) = \log M_X(t)$$

10. Conditional prob.

The conditional prob. of an event $B$ given that an event $A$ has occurred

is $P(B|A) = \dfrac{P(A \cap B)}{P(A)}$

If $(X,Y)$ has a jonit density of $f_{X,Y}(x,y)$ and $X$ with marginal $f_X(x)$

then the conditional density

$$f_{Y|X}(y) = \frac{f_{X,Y}(x,y)}{f_X(x)}$$

Let $X$ and $Y$ be random $k$-vectors

(a) If $\phi_X(t) = \phi_Y(t)$ for all $t \in \mathbb{R}^k$, then $F_X = F_Y$.   ($\phi_X(t) = E(\exp(it^T X))$

(b) If $m_X(t) = m_Y(t) < \infty$ for all $t$ in the neighborhood of $0$, then $F_X = F_Y$.

       $E\{\exp(t^T Y)\}$.

Pf: (a) See Billingsley (1968. P.395) [Inversion formula].

For any $a = (a_1, \ldots a_k) \in \mathbb{R}^k$, $b = (b_1, \ldots b_k) \in \mathbb{R}^k$ and $(a,b] = (a_1, b_1] \times (a_2, b_2] \times \cdots (a_k, b_k]$

satisfying $F_X = 0$, $F_X((a,b]) = \lim_{c \to \infty} \int_{-c}^{c} \cdots \int_{-c}^{c} \phi_X(t_1, \ldots t_k) \prod^k \frac{e^{-it_j a_j} - e^{it_j b_j}}{\ldots} d t_j.$

③

$\phi_X(t)$ , $M_X(t)$

Characteristic Function & Moment Generating Function

[Thm] Let $X$ and $Y$ be random $k$-vector

(a) If $\phi_X(t) = \phi_Y(t)$ for all $t \in \mathbb{R}^k$, then $F_X(x) = P(X \le x) = F_Y(x)$

(b) If $M_X(t) = M_Y(t) < \infty$ for all $t$ in the neighborhood of $0$, then $F_X = F_Y$.

Pf for (b) = First consider the case $k=1$.

From $e^{s|x|} \le e^{sx} + e^{-sx}$, we can conclude that $|X|$ has an mgf that is finite in the neighborhood $(-c, c)$ for some $c > 0$ and a constant $s$.

Observe that

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

we obtain

$$\left| \phi_X(t+a) - \sum_{j=0}^{n} \frac{a^j}{j!} E\left[ (iX)^j e^{itX} \right] \right| \le \frac{|a| E|X|^{n+1}}{(n+1)!} \qquad (*)$$

Since

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

} tuto

and

$$\frac{\partial^r \phi_X(t)}{\partial t_1^{r_1} \cdots \partial t_k^{r_k}} = (-1)^{\frac{r}{2}} E\left( X_1^{r_1} X_2^{r_2} \cdots X_k^{r_k} e^{it^T X} \right), \quad r = \sum_{i=1}^{k} r_i \qquad (ii)$$

We can write $\phi_k(t+a) = \sum_{j=0}^{\infty} \frac{\phi_k^{(j)}(t)}{j!} a^j$ , $|a| < c$ $\qquad (\dagger)$ ,

which also holds when $\phi_X$ is replaced by $\phi_Y$.

Under the assumption that $m_X = m_Y < \infty$ in the neighborhood of $0$, $X$ and $Y$ has the same moments of all orders. By (ii), $\phi_X^{(j)}(0) = \phi_Y^{(j)}(0)$ for all $j = 1, 2, \dots$ which, and $(\dagger)$ with $t=0$ imply that $\phi_X$ and $\phi_Y$ are the same on the integral $(-c, c)$ and hence have identical derivatives there.

Choose $t = c - \varepsilon$ and $-c + \varepsilon$ and an arbitrary $\text{soma}$ small $\varepsilon > 0$ in $(\dagger)$, we can show that $\phi_X$ and $\phi_Y$ also agree on $(-2c + \varepsilon, 2c + \varepsilon)$ and hence $[-2c, 2c]$. Likewise, by the same arguement, $\phi_X$ and $\phi_Y$ are the same on $(-3c, 3c)$. Hence $\phi_X(t)$ and $\phi_Y(t)$ for all $(t)$ and by (a) $F_X = F_Y$.

e.g. tuto.

For $k \geq 2$, suppose $F_X \neq F_Y$, then by (a), 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 $F_{t^T X} \neq F_{t^T Y}$. But $m_X = m_Y < \infty$ in a neighborhood of $0 \in \mathbb{R}^k$. This implies that $m_{t^T X} = m_{t^T Y} < \infty$ in a nbh of $0 \in \mathbb{R}$, which leads to the conclusion that $F_{t^T X} = F_{t^T Y}$. contradiction.

Some useful inequalities
→ Shorack and Wellerer (1986)

CB Chapter 3.5: ineq, $\to$ convergence. Keener

[Thm] Let $Z$ be a real r.v. and $g$ a non-negative, non-decreasing function on the support of $Z$, i.e. a set $B$ such that $P(Z \in B) = 1$, then

$$P(Z \geq a) \leq \frac{Eg(Z)}{g(a)}$$

Pf: observe that $g(a) I(Z \geq a) \leq g(Z) I(Z \geq a) \leq g(Z)$

Taking expectation. done.

Examples:

(a) Markov: $Z = |X|$, $g(t) = t^+ \overset{\text{max}(0,t)}{\underset{\|}{}} \Rightarrow P(|X| \geq a) \leq \frac{E|X|}{a}$

(b) Chebyshev: $Z = |X|$, $g(t) = t^2 \Rightarrow P(|X| \geq a) \leq \frac{E(X^2)}{a^2}$

(c) Bernstein: $Z = X$, $g(t) = e^{st} \Rightarrow P(X \geq a) \leq \frac{E(e^{sX})}{e^{sa}}$

One can construct examples for which they are actually sharp. For example, in the case of Markov inequality, suppose

$$X = \begin{cases} a & \mu/a \\ 0 & 1 - \mu/a \end{cases},$$

then $E(X) = \mu$ and obviously $P(X \geq a) = \frac{E(X)}{a}$.

[Thm] (Cauchy-Schwarts)

Let $X = (X_1, \ldots, X_p)$ be a $p$-vector of real r.v.'s and $U = E(XX^T)$. The matrix $U$ is symmetric, non-negative definite with singularity ($|U| = 0$) iff there exists a $p$-vector $\alpha \neq 0$ such that $\underline{E(\alpha^T X)^2 = 0}$ (*)

· Proof: Since Expectation(E) is applied componentwise (and multiplication commutes) symmetry is immediate. Non-negative definiteness follows from

$$\alpha^T U \alpha = E(\alpha^T X)^2 \geq 0.$$

If equality holds, then clearly (*) holds and $U$ is singular since $U\alpha = 0$. On the other hand, if $U$ is singular, there must exist $\alpha \neq 0$ s.t. $U\alpha = 0$ as the equality holds

④

[Corollary] For r.v. 's $X_1$ and $X_2$

$$\{E(X_1 X_2)\}^2 \leq E(X_1^2) E(X_2^2)$$

and centering $X_1$ and $X_2$ at there respective measure yields.

$\widetilde{X_j} = X_j - E(X_j) \quad j = 1, 2, \cdots$

$$\{cov(X_1, X_2)\}^2 \leq Var(X_1) Var(X_2) . \Rightarrow corr(X_1, X_2) \in [-1, 1]$$

$$\overset{\shortparallel}{\underset{\sqrt{Var(X_1) Var(X_2)}}{cov(X_1, X_2)}}$$

Proof: Consider the previous thm with $p=2$, and recalling that $|u|$ may be expressed

as the product of its eigenvalue which are non-negative, we have

$$0 \leq |u| = \frac{E(X_1^2) E(X_2^2)}{- \{E(X_1 X_2)\}^2}$$

[Thm] (Jensen) If $X$ and $g(X)$ are integrable r.v. 's and $g(\cdot)$ is convex, then

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

Proof: Convexity of $g$ implies that for any $\xi$, there exists a line $L$ through the point

$(\xi, g(\xi))$ such that the graph $g$. is above the line, i.e.

$$g(x) \geq g(\xi) + \lambda(x - \xi)$$

In particular, we let $\xi = E(X)$, then for all $x$,

$$g(x) \geq g(E(X)) + \lambda \{x - E(X)\}.$$

Note that $\lambda$ depends on $\xi$ but not on $X$. Now, let $x = X$ and

$$g(X) \geq g(E(X)) + \lambda \{X - EX\}.$$

Taking expectation, done.

[Corollary] (Liapounrv) $\{E|X|^r\}^{1/r}$ is $\nearrow$ in $r$ for $r \geq 0$
↳ Hölder. Minkowski ...

Proof: By Jensen's ineq, since $|X|^r$ is convex in $|X|$ for $r \geq 1$,

Remark: we have $(E|X|)^r \leq E|X|^r$ in which case $E|X| \leq (E|X|^r)^{1/r}$

Moment inequality:
Now, replace $|X|$ by $|X|^q$ for $0 < q < r$,

If $X$ has a $r$th

moment, it also as $(E|X|^q)^{1/q} \leq (E|X|^{rq})^{1/(rq)} \triangleq (E|X|^s)^{1/s}$ where $s = rq$.

have the $q$th for $0 < q < s < \infty$, since $r \geq 1$, so $q \leq rq = s$.

moment for $q < r$.

· Convergence Results.

We are interested in sequences $X_1, X_2, \ldots$ of r.v.'s on a p-space $(\Omega, A, P)$

## I. Convergence in probability

Let $\{X_i\}_{i \geq 1}$ and $X$ be real-valued r.v.'s on $(\Omega, A, P)$.

We say that $X_n$ converges in prob. to $X$ if

$$\lim_{n \to \infty} P(|X_n - X| < \varepsilon) = 1, \quad \text{for any } \varepsilon > 0.$$

and we write usually $X_n \xrightarrow{P} X$.

Remark: Often $X$ will be a degenerate r.v. e.g. $X_n = \bar{X}_n = \dfrac{\sum_{i=1}^{n} Z_i}{n}$

where $Z_i \overset{iid}{\sim} N(\mu, \sigma^2)$, then $X_n \xrightarrow{P} \mu$, we can think of $\mu$ as the

degenerate r.v. $X$, which takes the value $\mu$ with prob 1.

## II Almost sure convergence / convergence with prob. 1

We say $X_n$ converges almost surely (a.s.), or converges with prob. 1

if $\quad P(\lim_{n \to \infty} X_n = X) = 1$.

or equivalently, for any $\varepsilon > 0$,

$$\lim_{n \to \infty} P(|X_m - X| < \varepsilon \text{ for all } m > n) = 1$$

$\}$ tuto

One can show that a.s. convergence $\Rightarrow$ convergence in prob.

(and we have counter examples to show that the converse is wrong)

## III. Convergence in the $q$th mean

$X_n$ Converges in the $q$th mean to $X$ if $\boxed{\lim_{n \to \infty} E|X_n - X|^q = 0}$

By the moment inequality introduced earlier,

$X_n \xrightarrow{q\text{th}} X \Rightarrow X_n \xrightarrow{p\text{th}} X$ for any $p < q$. Often, $q = 2$ in practice.

As an example of extreme behavior, suppose that $\boxed{X_n = \begin{cases} 0 & 1 - n^{-3} \\ n & n^{-3} \end{cases}}$

then taking $X = 0$, we have $\lim_{n \to \infty} E|X_n - X|^q = 0$ for $q = 1, 2$, but $E|X_n - X|^3 = 1$.

⑤

# IV. Convergence in distribution (in law)

$X_n$ converges in distribution to $X$ if for their respective distribution functions $\lim_{n\to\infty} F_n(x) = F(x)$ at each point of continuity of $F$.

We write $X_n \xrightarrow{d/D/L} X$ or as $F_n \Rightarrow F$ (For converges weakly to $F$)

Often, we write $X_n \rightsquigarrow X$    e.g. $X_n \rightsquigarrow N(0,1)$

$$ \mathrm{II} \Rightarrow \mathrm{I} \Rightarrow \mathrm{IV} \Rightarrow \mathrm{III} $$

## Big / Small $O$ notation

For positive deterministic sequences $\{a_n\}, \{b_n\}$

a) If there is a $\Delta < \infty$ s.t $a_n/b_n \le \Delta$ for sufficiently large $n$, we say $a_n = O(b_n)$

b) if $a_n/b_n \to 0$, we say $a_n = o(b_n)$

Clearly, if $a_n = O(n^r)$ and $b_n = O(n^s)$ then $a_n b_n = O(n^{r+s})$
and $a_n + b_n = O(n^{\max\{r,s\}})$

For sequences $\{X_n\}$ and $\{Y_n\}$ of r.v.'s on $(\Omega, A, P)$ and any $\varepsilon > 0$.

a*) If there exists $\Delta < \infty$ s.t. $P(|X_n| \ge \Delta |Y_n|) < \varepsilon$ for sufficiently large $n$, we write $X_n = O_p(Y_n)$

b*) If $P(|X_n| > \varepsilon |Y_n|) \xrightarrow{n\to\infty} 0$, then we write $X_n = o_p(Y_n)$

In many cases, $Y_n$ will be deterministic, we write correspondingly

$X_n = O_p(1)$ : "bounded in prob."

$X_n = o_p(1)$ : "tending $0$ in prob."

[Thm] (Slusky) ... $\delta$-method.

Let $\boxed{X_n \rightsquigarrow X}$ and $\boxed{Y_n \xrightarrow{P} y}$, a real constant. Then.
$X_n \xrightarrow{d} X$

(a) $X_n + Y_n \rightsquigarrow X + y$   (exercise).

(b) $X_n Y_n \rightsquigarrow y X$

Proof for (b)

Suppose $y=0$, and let $B>0$ be a real constant, and denote

$$X_n^B = X_n I(|X_n| \le B)$$

Then $\{|Y_n X_n| \ge \varepsilon\} = \{|Y_n| |X_n^B| \ge \varepsilon\} \cup \{|Y_n| |X_n - X_n^B| \ge \varepsilon\}$ ⊛

$\{|Y_n| |X_n^B| \ge \varepsilon\} \subseteq \{|Y_n| > \frac{\varepsilon}{B}\}$

and $P\{|Y_n| |X_n^B| \ge \varepsilon\} \le P\{|Y_n| \ge \frac{\varepsilon}{B}\} \to 0$.

By the hypothesis that $X_n = O_p(1)$, there exists $\delta > 0$ and $B_j < \infty$ s.t.

for $n$ sufficiently large, $P(|X_n - X_n^{B\delta}| > 0) < \delta$.

Since $\{|Y_n| |X_n - X_n^B| \ge \varepsilon\} \subseteq \{|X_n - X_n^B| > 0\}$

(*) and additivity implies that $\lim\limits_{n\to\infty} P\{\underline{|X_n| |Y_n|} \ge \varepsilon\} < \delta$. $\quad |X_n Y_n - 0|$

Since $\varepsilon$ and $\delta$ are arbitrary, we have shown that $X_n Y_n \xrightarrow{p} 0$.

The result follows by noticing that $Y_n$ can be replaced by $Y_n - y$.

[Thm] (Continuous mapping)

If $X_n \leadsto X$ and $g$ is continuous, $g(X_n) \leadsto g(X)$. Proof skipped.

[Thm] ($\delta$-method)

Suppose $a_n(X_n - b) \leadsto X$, where $a_n$ is a sequence of constants tending to $\infty$ and

$b$ is a fixed number. Let $g: \mathbb{R} \to \mathbb{R}$ be differentiable with continuous derivative $g'$

at $b$. Then $a_n\{g(X_n) - g(b)\} \leadsto \underbrace{g'(b)}_{\text{slope}} X$.
$\qquad\qquad\qquad \uparrow X \qquad\qquad \downarrow \text{Var}(X)$

Proof: By Slusky's thm,

$$X_n - b = a^{-1}\{a_n(X_n - b)\} \longrightarrow 0$$

and therefore $X_n \to b$. Now apply mean value thm to $g(X_n) - g(b)$,

we have $\qquad g(X_n) - g(b) = g'(X_n^*)(X_n - b)$

where $|X_n^* - b| \le |X_n - b|$ where $X_n^* \to b$, so by the conuity of $g'$ and

cont. mapping thm (ACMT). $g'(X_n^*) \to g'(b)$. Multiplying $a_n$ and again

applying Slusky, we have the result. The above argument generalizes to

$X_n, X \in \mathbb{R}^2$.

LLN, CLT (dependent)

Lecture 3    References: 1) C&B. Chapter 5 & 6↑   2) Keener Chapter 8. 2-3↑.   3) Ferguson (convergence results)

Asymptotic behavior of sample mean.

*(i) Law of Large Numbers (LLNs) : $\hat{\mu} \to \mu$   $\hat{\mu} = \bar{X} = n^{-1} \sum_{i=1}^{n} X_i$

*(ii) CLT   : $\sqrt{n} (\hat{\mu} - \mu) \rightsquigarrow N(0, \sigma^2)$

(i) weak LLN : Let $Z_1 \ldots$ be indep. r.v.'s with means $\mu_1 \cdots$ and variance $\sigma_1^2, \ldots$

$Z_n$   $n^{-2} \sum_{i=1}^{n} \sigma_i^2 \to 0$   as $n \to \infty$,   then $\bar{Z} = \hat{\mu} \to \bar{\mu}$

· Pf:   $P(|\bar{Z} - \bar{\mu}| \geq \varepsilon) \leq \dfrac{E(\bar{Z} - \bar{\mu})^2}{\varepsilon^2} = \dfrac{\frac{1}{n^2} \sum_{i=1}^{n} \sigma_i^2}{\varepsilon^2} \to 0$

Remark:   If $\sigma_i^2 = \sigma^2$, then $\frac{1}{n^2} \sum_{i=1}^{n} \sigma_i^2 = \sigma^2/n$.   $\bar{Z} - \mu = o_p(\frac{1}{n})$
(for iid)

SLLN (Kolmogrov)   For $\{Z_i, \mu_i, \sigma_i^2\}$ as above, if $\sum_{i=1}^{n} \dfrac{\sigma_i^2}{j^2} < \infty$, then $\bar{Z} \to \bar{\mu}$ a.s.

In particular, if $\{Z_i\}_{i=1}$'s are iid, with $E(Z_i) = \mu$, then $\bar{Z} \to \mu$ a.s.

An important special case if the SLLN involves taking

$Z_i = I_{[X_i, \infty)}(x) = I(X_i \leq x)$

for iid r.v.'s $X_i \sim F$ and fixed $x$. Observe that

$E(Z_i) = P(X_i \leq x) = F(x)$

We can infer that   $F_n(x) \triangleq n^{-1} \sum_{j=1}^{n} I(X_j \leq x) \xrightarrow{a.s.} F(x)$   as $n \to \infty$.
This can be strengthened.

[Thm] (Glivenko- Cantelli)

$P(\sup_x |F_n(x) - F(x)| \to 0) = 1.$

· Pf: Let $\varepsilon > 0$ and find an integer $k > 1/\varepsilon$ and numbers

$-\infty = x_0 < x_1 \leq \cdots \leq x_{k-1} < x_k = \infty$ such that. $F(x_j^-) \leq \dfrac{j}{k} \leq F(x_j)$, $j = 1, \ldots k-1$.

where $F(x^-) = P(X < x)$.

Note that if $x_{j-1} < x_j$, then $F(x_j^-) - F(x_{j-1}) \leq \varepsilon$.

From SLLN,   $F_n(x_j) \xrightarrow{a.s} F(x_j)$ and $F_n(x_j^-) \xrightarrow{a.s} F(x_j^-)$ for $j = 1, 2, \ldots k-1$.

Hence   $\Delta_n = \max\{|F_n(x_j) - F(x_j)|, |F_n(x_j^-) - F(x_j^-)|, j = 1, \ldots, k-1\} \to 0.$

Now, let $x$ be arbitrary and find $j$ s.t. $x_{j-1} < x \leq x_j$.

Then $F_n(x) - F(x) \leq F_n(x_j^-) - F(x_{j-1}) \leq F_n(x_j^-) - F(x_j^-) + \varepsilon$.

and $F_n(x) - F(x) \geq F_n(x_{j-1}) - F(x_{j*}) \geq F_n(x_{j-1}) - F(x_{j-1}) - \varepsilon$

Thus $\sup_x |F_n(x) - F(x)| \leq \Delta_n + \varepsilon \overset{a.s}{\to} \varepsilon$ as $n \to \infty$ 

(Van dei Vaant, 98. Chap 79)

Vanpnik (99).

Since this holds for all $\varepsilon > 0$ and the results follow $\square$.

(ii) CLT.

Suppose $Z_1, \dots, Z_n$ are iid $N(0,1)$ or $Z \sim N(0, I_n)$. we know that $\alpha^T Z \sim N(0, \alpha^T \alpha)$

So, for instance, if we take $\alpha = n^{-1/2} I_n$, we have $n^{-1/2} \sum_{i=1}^{n} Z_i \sim N(0,1)$.

or equivalently $\sqrt{n}\bar{Z} \sim N(0,1)$.

[Thm] Suppose $X_1, \dots$ are iid with $E(X_i) = \mu$, $Var(X_1) = \sigma^2$, then $\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \rightsquigarrow N(0,1)$.

Pf: Existence of $\mu$ and $\sigma^2$ implies that the moment expension of the cf of $X_1$ can be written as $\phi_{X_1}(t) = \exp\{i\mu t - \frac{1}{2}\sigma^2 t^2 + o(t^3)\}$.

Define $S_n = X_1 + \dots + X_n$, which has $\phi_{S_n}(t) = \phi_{X_1}^n(t)$

and let $U_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma^2} = \frac{\sqrt{n}(S_n/n - \mu)}{\sigma}$

$\phi_{U_n}(t) = E(e^{+itU_n}) = \phi_{X_1}^n(\frac{t}{\sigma\sqrt{n}}) \exp(-\frac{i\mu t\sqrt{n}}{\sigma})$

$= \left[\exp\{\frac{i\mu t}{\sigma\sqrt{n}} - \frac{\frac{1}{2}\sigma^2 t^2}{\sigma^2 n} + o(\frac{t^2}{\sigma^2 n})\}\right]^n \exp(-\frac{i\mu t\sqrt{n}}{\sigma})$

$= \exp\{-\frac{1}{2}t^2 + n\,o(\frac{t^2}{n})\} \longrightarrow \exp(-\frac{1}{2}t^2)$ as $n \to \infty$

which is the cf of $N(0,1)$.

Why normal?

If $\sqrt{n}(\bar{X}_n - \mu) \rightsquigarrow X$, what does $X$ look like?

Consider $Z_{2n} = \frac{X_1 + \dots + X_n + X_{n+1} + \dots + X_{2n}}{\sqrt{2n}}$

Clearly, $Z_{2n} \rightsquigarrow X$ but $Z_{2n} = \frac{X_1 + \dots + X_n}{\sqrt{2n}} + \frac{X_{n+1} + \dots + X_{2n}}{\sqrt{2n}} \overset{\Delta}{=} Z_{an} + Z_{bn}$

only normal distr. ↑ fufills the setup. (stable laws).

where $Z_{an}$ and $Z_{bn} \rightsquigarrow \frac{X}{\sqrt{2}}$

Chow and Teicher (3rd)

[Thm] (Lyapunov) Let $X_1, \dots$ be indep. with $EX_i = 0$ and $EX_i^2 < \sigma_i^2 < \infty$. $E|X_i|^3 < \infty$

and $S_n^2 = \sum_{i=1}^{n} \sigma_i^2$. If $\lim_{n \to \infty} \sum_{i=1}^{n} E|X_i|^3 = 0$, then $S_n^{-1} \sum_{i=1}^{n} X_i \rightsquigarrow N(0,1)$.

Lindeberg Feller cond. Martingale CLT ...

①

# CLT for dependent cases / sequences.

One solution : <u>$\alpha$-mixing</u>

Given a sequence $X_1, X_2, \ldots$ and sets $A \in \sigma(X_1, \ldots X_k)$.

and $B \in \sigma(X_{k+n}, X_{k+n+1}, \ldots)$ for $k \geq 1$ and $n \geq 1$, then if there exists

a sequence of real numbers $\alpha_n \to 0$ s.t.

$$|P(A \cap B) - P(A)P(B)| \leq \alpha_n$$

then $\{X_n\}$ is $\alpha$-mixing.

<u>Special case</u> If $\alpha_n = 0$ for $n > m$, then the sequence is said to be m-dependent.

## <u>CLT for $\alpha$-mixing sequences</u>

Suppose $X_1, X_2, \ldots$ is stationary and $\alpha$-mixing with $\alpha_n = O(n^{-5})$. $E(X_n) = 0$ and

$E(X_n'^2) < \infty$. Set $S_n = X_1 + \cdots + X_n$. If $n^{-1} \text{Var}(S_n) \to \sigma^2 = E(X_1^2) + \sum_{k=1}^{\infty} E(X_1 X_{k+1})$

converges absolutely with $\sigma^2 > 0$, then $S_n/\sigma\sqrt{n} \rightsquigarrow N(0,1)$.

<u>Extra examples</u> With assumptions in CLT, if $f$ is differentiable at $\mu$, then

$$\sqrt{n}\{f(\bar{X}_n) - f(\mu)\} \rightsquigarrow N(0, [f'(\mu)]^2 \sigma^2).$$

Pf : Write $f(\bar{X}_n) = f(\mu) + f'(\mu_n)(\bar{X}_n - \mu)$, where $\mu_n$ is an immediate point

between $\bar{X}_n$ and $\mu$. Since $|\mu_n - \mu| \leq |\bar{X}_n - \mu|$ and $\bar{X}_n \xrightarrow{P} \mu$ (LLNs) and since $f'$

is continuous, $f'(\mu_n) \xrightarrow{P} f'(\mu)$. If $Z \sim N(0, \sigma^2)$, then $\sqrt{n}(\bar{X}_n - \mu) \rightsquigarrow Z \sim N(0, \sigma^2)$

by CLT. Thus by Slusky's thm, $\boxed{\sqrt{n}\{f(\bar{X}_n) - f(\mu)\}} = f'(\mu_n)\{\sqrt{n}(\bar{X}_n - \mu)\} \to f'(\mu)Z$

$\sim N(0, [f'(\mu)]^2 \sigma^2)$.

## <u>Asymptotics of medians and percentiles</u>

For regularity, assume $F$ has a unique median $\theta$, so $F(\theta) = \frac{1}{2}$, and that $F'(\theta)$ exists,

which is finite and positive. We want to study the asymptotic distribution of

$\sqrt{n}(M_n - \theta)$, where $M_n$ denotes the sample median of $(X_1, \ldots X_n) \overset{iid}{\sim} f$. (with d.f. $F$)

$$P(\sqrt{n}(M_n - \theta) \leq a) = P(M_n \leq \theta + a/\sqrt{n})$$

Define $S_n = \#\{i \leq n : X_i \leq \theta + \frac{a}{\sqrt{n}}\} = \sum_{i=1}^{n} I(X_i \leq \theta + \frac{a}{\sqrt{n}})$

m: the middle integer.

Note that $\{M_n \leq \theta + \frac{a}{\sqrt{n}}\}$ if $\{S_n \geq m\}$. It is evident that, if we treat the

observation $i$ as a success if $X_i \leq \theta + \frac{a}{\sqrt{n}}$, then $S_n \sim$ Binomial $(n, F(\theta + \frac{a}{\sqrt{n}}))$

Let $Y_n \sim$ binomial $(n, p)$, then by CLT,

$$\sqrt{n}\left(\frac{Y_n}{n} - p\right) = \frac{Y_n - np}{\sqrt{n}} \rightsquigarrow N(0, p(1-p))$$

in which case $P\left(\frac{Y_n - np}{\sqrt{n}} > y\right) \to 1 - \Phi\left(\frac{y}{\sqrt{p(1-p)}}\right) = \Phi\left(\frac{-y}{\sqrt{p(1-p)}}\right)$ as $n \to \infty$.

where $\Phi(\cdot)$ denotes the cdf of $Z \sim N(0,1)$

Hence, the normal approximation for the binomial distribution gives

$$P(\sqrt{n}(M_n - \theta) \le a) = P(S_n > m-1)$$

$$= P\left\{\frac{S_n - nF(\theta + a/\sqrt{n})}{\sqrt{n}} > \frac{m-1-nF(\theta + a/\sqrt{n})}{\sqrt{n}}\right\}$$

$$= \Phi\left(\frac{[nF(\theta + a/\sqrt{n}) - m+1]/\sqrt{n}}{\sqrt{F(\theta + a/\sqrt{n})(1 - F(\theta + a/\sqrt{n}))}}\right) + o(1) \qquad (*) \qquad \text{Keener (P138)}$$

Since $F$ is continuous at $\theta$.

$$\left[F(\theta + a/\sqrt{n})\left\{1 - F(\theta + \frac{a}{\sqrt{n}})\right\}\right]^{1/2} \to 1/2 \quad \text{as } n \to \infty.$$

$\dfrac{\text{TPE}}{\text{TSH}} > \text{lehmann}.$

Because $F$ is differentiable at $\theta$,

$$\frac{nF(\theta + \frac{a}{\sqrt{n}}) - m+1}{\sqrt{n}} = \frac{aF(\theta + \frac{a}{\sqrt{n}}) - F(\theta)}{a/\sqrt{n}} - \frac{nF(\theta) - m + 1}{\sqrt{n}}$$

$$= \frac{aF(\theta + a/\sqrt{n}) - F(\theta)}{a/\sqrt{n}} + \frac{1}{2\sqrt{n}} \to aF'(\theta).$$

Since the numerator and the denominator of the argument of $\Phi$ in $(*)$ both converge, we can write

$$P(\sqrt{n}(M_n - \theta) \le a) \to \Phi(2aF'(\theta))$$

The limit here is the cdf of a normal r.v. with mean $0$ and variance $[4\{F'(\theta)\}^2]^{-1}$ evaluates at $a$ and so

$$\sqrt{n}(M_n - \theta) \to N\left(0, \frac{1}{4[F'(\theta)]^2}\right)$$

[Thm] Let $X_1, \ldots, X_n$ be iid with common cdf $F$. and let $\tau \in (0,1)$ and let $\tilde{\theta}_n$ be the $\lfloor \tau n \rfloor^{th}$ order statistic for $X_1, \ldots, X_n$. where $\lfloor X \rfloor$, floor of $X$. If $F(\theta) = \tau$ and if $F'(\theta)$ exists and is finite and positive. then $\sqrt{n}(\tilde{\theta}_n - \theta) \rightsquigarrow N\left(0, \frac{\tau(1-\tau)}{[F'(\theta)]^2}\right)$

<u>Data Reduction</u> : Sufficiency ...

$= T(\underset{\sim}{X}) = T(X_1, \ldots, X_n)$

<u>Definition</u> (Statistic) A statistic $T$ is a function of the data.

**Def** (Sufficient statistic) A statistic is sufficient for a mode $P = \{P_\theta : \theta \in \Theta\}$ if for all $t$, the conditional distribution $X \mid T(X) = t$ does not depend on $\theta$.

**Example** (Weighted coin flips)

Let $X_1, X_2, \ldots, X_n$ be iid. according to Bernoulli $(\theta)$. is the number of heads. i.e. $\sum_{i=1}^{n} X_i$ sufficient? To check this, let's show the conditional distribution of $\underline{X}$ given $\sum_{i=1}^{n} X_i$.

We have $\quad P_\theta(X) = \prod_{i=1}^{n} \theta^{X_i} (1-\theta)^{1-X_i} = \theta^{\sum_{i=1}^{n} X_i} (1-\theta)^{n - \sum_{i=1}^{n} X_i}$.

So the conditional distribution is

$$P_\theta(X = x \mid T(x)) = \frac{P_\theta(X=x, T(x)=t)}{P(T(x)=t)} = \frac{I(\sum_{i=1}^{n} X_i = t) \, \theta^t (1-\theta)^{n-t}}{\binom{n}{t} \theta^t (1-\theta)^{n-t}} = \frac{I(\sum_{i=1}^{n} X_i = t)}{\binom{n}{t}},$$

which does not depend on $\theta$, so the sum of heads is a sufficient stat.

**Example** (Max of uniform)

Let $X_1, \ldots, X_n$ be iid uniform $(0, \theta)$. Then $T(X) = \max(X_1, \ldots, X_n)$ is sufficient.

To see the intuition, think of $X_1, \ldots, X_n$ as $n$ numbers on the real line, then the remaining $n-1$ numbers, given the maximum is fixed at $t$, behaves like $n-1$ iid random samples drawn from $\mathcal{U}(0, t)$.

**Example** (Order statistic)

Let $X_1, \ldots, X_n$ be iid with any model. Then the order statistic

$$T = \{ X_{(1)} \leq X_{(2)} \leq \cdots \leq X_{(n)} \} \text{ are sufficient.}$$

**[Thm]** (TPE 1.6. Thm 6.1)      Point estimation.    P33.

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

**[Thm]** Neyman-Fisher Factorization Criterion (NFFC) $\kappa$ TSH. P.19.

Suppose each $P_\theta \in P$ has density $P(X; \theta)$ w.r.t a common $\sigma$-finite measure $\mu$, i.e. $\frac{dP_\theta}{d\mu} = p(x; \theta)$. Then $T(x)$ is sufficient iff $p(x; \theta) = g_\theta(T(x)) h(x)$ for some $g_\theta$ and $h$.

· **Pf:** (Discrete)

Suppose $p(x; \theta) = g_\theta(T(x)) h(x)$. Since $P_\theta(X = x \mid T(X) = t) = 0$ whenever $T \neq T(x)$. so we may focus our attention to the case where $P_\theta(X = x \mid T(X) = T(x))$.

We can write $P_\theta(X=x \mid T(X) = T(x)) = \dfrac{P_\theta(X=x, T(X)=T(x))}{P_\theta(T(X)=T(x))} = \dfrac{P_\theta(X=x)}{P_\theta(T(X)=T(x))}$

$= \dfrac{g_\theta(T(x)) \, h(x)}{\sum\limits_{x'} P(x'; \theta) I(T(x') = T(x))} = \dfrac{g_\theta(T(x)) \, h(x)}{\sum\limits_{x'} g_\theta(T(x)) h(x') I(T(x')=T(x))}$

$= \dfrac{h(x)}{\sum\limits_{x'} h(x') I(T(x')=T(x))}$ , which is ind. of $\theta$ and hence

$T$ is sufficient.

Conversely, suppose $P_\theta(X=x \mid T(X)=T(x))$ is indep. of $\theta$.

Then, defining $h(x) = P_\theta(X=x \mid T(X)=T(x))$, we have

$P(x; \theta) = P_\theta(X=x) = P_\theta(X=x, T(X)=T(x))$

$= P_\theta(X=x \mid T(X)=T(x)) \, P_\theta(T(X)=T(x)) = h(x) \, g_\theta(T(x))$

which establishes the criterion.

## Example (normal)

Let $X_i$ be iid $N(\mu, \sigma^2)$ and $\theta = (\mu, \sigma^2)$. The joint density is ~~$p(x;\theta)$~~

$P(x; \theta) = \prod\limits_{i=1}^{n} \dfrac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} = \left(\dfrac{1}{\sqrt{2\pi\sigma^2}}\right)^n e^{\frac{1}{2\sigma^2}\left(-\sum\limits_{i=1}^{n} x_i^2 + 2\mu \sum\limits_{i=1}^{n} x_i - n\mu^2\right)}$

$= g_\theta(T(x))$ where $T(x) = \left(\sum\limits_{i=1}^{n} x_i^2, \sum\limits_{i=1}^{n} x_i\right)$

## Example

Suppose $X$ and $Y$ are indep. with common Lebesgue density $f_\theta(x) = \theta e^{-\theta x} I(x \geq 0)$. Let $U$ be indep. of $X$ and $Y$ and uniformly distributed on $(0,1)$. Take $T = X+Y$ and define $\widetilde{X} = U T$ and $\widetilde{Y} = (1-U)T$.

To find $f_{\widetilde{X}, \widetilde{Y}}$, observe that $P(T \leq t \mid Y=y) = P(X+Y \leq t \mid Y=y)$

$= E\{I(X+Y \leq t) \mid Y=y\}$

$= F_X(t-y)$

$\Rightarrow F_T(t) = P(T \leq t) = E\{F_X(t-Y)\} = \int_0^\infty 1 - e^{-\theta(t-y)} \theta e^{-\theta y} \, dy$

$= 1 - e^{-\theta t} - t\theta e^{-\theta t}$

Hence $f_T(t) = \dfrac{\partial F(t)}{\partial t} = t\theta^2 e^{-\theta t}$, $t \geq 0$.

Also, $f_{T,U}(t, u; \theta) = t\theta^2 e^{-\theta t} I(t \geq 0, u \in [0,1])$.

From which, we have $P\left(\binom{\widetilde{X}}{\widetilde{Y}} \in B\right) = \iint I_B(tu, t(1-u)) f_{T,U}(t, u; \theta) \, du \, dt$

$\overset{\text{simplification}}{=} \iint I_B(x, y)(x+y)^{-1} f_{T,U}\left(x+y, \dfrac{x}{x+y}\right) dy \, dx.$

Thus $(\widetilde{X}, \widetilde{Y})$ has the density

$$\frac{f_{X,Y}\left(x+y, \frac{x}{x+y}\right)}{x+y} = \begin{cases} \theta^2 e^{-\theta(x+y)} & , \ x \geqslant 0, \ y \geqslant 0 \\ 0 & , \ ow \end{cases}$$

Lecture 4    Keener : Ch2&3    C&B: Ch6 .

References: 1. C&B ☆
2. keener ☆☆
3. TPE ☆☆☆
4. TSH ☆☆☆
5. Bickel & Doksum ☆
6. Serrich (? CMU) ☆☆
7. Shao Jun ☆☆
8. Van der Vaart ☆☆☆
(asymptotic).
9. Ferguson (?).

Recall : $X, Y$ .  $u$ = unif $(0,1)$ , $uT$, $(1-u)T$ , $T = X + Y$ .

NFFC:    $P(X; \theta) = g_\theta (T(X)) h(x)$

Def  Exponential families : A dominated family

$\{P_\theta : \theta \in \Theta \}$ is said to form a $k$-dimensional exponential family if the corresponding density function $\{P_\theta (x)\}_{\theta \in \Theta}$ are of the form  $P_\theta (x) = \exp\{ \sum_{i=1}^{k} \eta_i (\theta) T_i(x) - B(\theta) \} h(x)$ , where $h, T_1, ..., T_k : X \to \mathbb{R}$.
$B, \eta_1, ... \eta_k : \Theta \to \mathbb{R}$.

By NFFC , we can see that $(T_1, ..., T_k)$ is sufficient .

E.g.  $X_1, ..., X_n \overset{iid}{\sim} N(\theta, \sigma^2)$ , $\Theta = \mathbb{R} \times (0, \infty)$

$$P_\theta (X) = \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp\left( -\sum_{i=1}^{n} \frac{(X_i - \mu)^2}{2\sigma^2} \right) = \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp\left\{ -\sum_{i=1}^{n} \frac{X_i^2}{\sigma^2} + \frac{\mu \sum_{i=1}^{n} X_i}{\sigma^2} - \frac{n\mu^2}{2\sigma^2} \right\}$$

Hence .   $T_1 (X) = \sum_{i=1}^{n} X_i^2$ , $\eta_1 (\theta, \sigma^2) = -\frac{1}{2\sigma^2}$

$T_2 (X) = \sum_{i=1}^{n} X_i$ , $\eta_2 (\theta, \sigma^2) = \frac{\mu}{\sigma^2}$  natural parameters.

$B(\mu, \sigma^2) = \frac{n\mu^2}{2\sigma^2} - \frac{n}{2} \log (2\pi\sigma^2)$

$h(X) = I( X \in (-\infty, \infty)$

*

A measure $\nu$ is dominated by the measure $\mu$ if $\nu \ll \mu$ , which means that for ~~some~~ all measurable $A$ , $\mu(A) = 0$ implies $\nu(A) = 0$ . A family of prob. measure $(P_\theta)_{\theta \in \Theta}$ is dominated by $\mu$ iff for each $\theta \in \Theta$ , the measure $P_\theta$ is dominated by $\mu$ .

E.g. $X_1, ..., X_n \overset{iid}{\sim}$ Cauchy, i.e. $P_\theta (X) = \frac{1}{\pi} \frac{1}{1 + (X - \theta)^2}$ is the density of $P_\theta$ w.r.t Lebesgue measure. In this case , $T(X) = (X_{(1)}, ..., X_{(n)})$ is sufficient , where $(X_{(1)} \leq \cdots \leq X_{(n)})$ are the order statistics.

Indeed $T(X)$ is minimal sufficient .

[Thm] ( Pitman - Koopman - Darmois ) (1936).

Suppose $(X_1, ..., X_n)$ are iid with density $\{P_\theta : \theta \in \Theta \}$ w.r.t Lebesgue measure, which are continuous in $x$ for $\theta$ fixed and support on an interval $I \subseteq \mathbb{R}$. Suppose there exists a sufficient statistic $(T_1, ..., T_k)$ which are continuous:

(i) If $k=1$, then $P_\theta (X) = e^{\eta(\theta) T(X) - B(\theta)} h(X)$ .

(ii) If $n > k > 1$ , and the function $X \mapsto P_\theta (x)$ are continuous differentiable.
then  $P_\theta (x) = e^{\sum_{i=1}^{k} \eta_i(\theta) T(X) - B(\theta)} h(X)$ .

**Def** An exponential family is in <u>canonical form</u> when the density has the form

$$P_\eta(x) = \exp\left\{ \sum_{i=1}^{k} \eta_i T_i(x) - A(\eta) \right\} h(x)$$

This parametrizes the density in terms of the <u>natural parameters</u> $\eta$ instead of $\theta$.

**Def** The set of all valid natural parameters $\theta$ is called the <u>natural parameter space</u> = for each $\eta \in \Theta$, there exists a normalizing constant $A(\eta)$ s.t. $\int P_\eta(x) dx = 1$.

     Equivalently, $\quad \Theta = \left\{ \eta : 0 < \int \exp\left( \sum_{i=1}^{k} \eta_i T_i(x) \, h(x) \right) d\mu(x) < \infty \right\}$

Thus, for any canonical exponential family, $P = \{ P_\eta, \eta \in H \}$, we have $\eta \in \Theta$.

## Reducing the dimension

There are two cases when the superficial dimension of a $k$-dim exponential family $P = \{ P_\eta : \eta \in H \}$ can be reduced.

<u>Case 1</u> The $T_i(x)$'s satisfy an affine equality constraint $\forall x \in X$.

E.g. $X \sim \text{Exp}(\eta_1, \eta_2)$ i.e. $p(x; \eta_1, \eta_2) = \exp\left\{ -\eta_1 x - \eta_2 x + \log(\eta_1 + \eta_2) \right\}$.

$\exp\left\{ -(\eta_1 + \eta_2) x + \log(\eta_1 + \eta_2) \right\} I(x > 0)$.

     Hence $T_1(x) = T_2(x) = x$. i.e. they are linearly dependent. $\Rightarrow$ unidentifiable

**Def** If $P = \{ P_\theta, \theta \in \Theta \}$, then $\theta$ is unidentifiable if for two parameters $\theta_1 \neq \theta_2$, $P_{\theta_1} = P_{\theta_2}$.

     In the ~~piecewise~~ previous example, $p(x; \eta_1 + a, \eta_2 - a) = p(x; \eta_1, \eta_2)$ for any $a < \eta_2$.

<u>Case 2</u> The $\eta_i$'s satisfy an affine equality constraint for all $\eta \in H$.

E.g. $p(x; \eta) \propto \exp(\eta_1 x + \eta_2 x^2)$ for all $(\eta_1, \eta_2)$ satisfying $\eta_1 + \eta_2 = 1$.

$$= \exp\left\{ \eta_1 (x - x^2) + x^2 \right\}.$$

**Def**

A canonical exponential family $P = \{ P_\eta : \eta \in H \}$ is <u>minimal</u> if    *no linear combination constraints*.

- $\sum_{i=1}^{k} \lambda_i T_i(x) = \lambda_0 \; \forall x \in X \Rightarrow \lambda_i = 0 \; \forall i \in \{0, \dots, k\}$. (no affine $T_i$ equality).

- $\sum_{i=1}^{k} \lambda_i \eta_i = \lambda_0 \quad \forall \eta \in H \Rightarrow \lambda_i = 0 \; \forall i \in \{0, \dots, k\}$ (no affine $\eta_i$ equality).
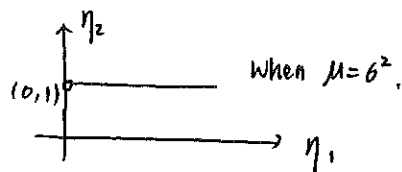
Keener Ch.5.

**Def**

Suppose $P = \{ P_\eta : \eta \in H \}$ is a $k$-dimensional minimal exponential family. If $H$ contains an open $k$-dim rectangle, then $P$ is called <u>full rank</u>, otherwise $P$ is <u>curved</u>.

We illustrate three types of exponential families via normal dis. $N(\mu, \sigma^2)$, where $\eta_1 = \frac{1}{2\sigma^2}$, $\eta_2 = \frac{\mu}{\sigma^2}$, $T_1(x) = x^2$ and $T_2(x) = x$.

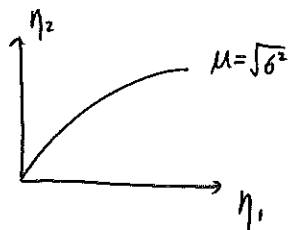## I) <u>Non-minimal</u> (so the dimension can be reduced)

when $\mu = \sigma^2$, $\eta_1 = \frac{1}{2\sigma^2}$, $\eta_2 = 1$.



when $\mu = \sigma^2$.

## II) <u>Minimal & Curved</u>
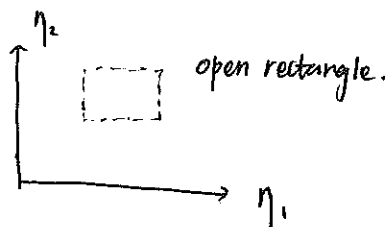
e.g. $\mu = \sqrt{\sigma^2}$, so $\eta_1 = \frac{1}{2\sigma^2}$, $\eta_2 = \frac{1}{\sqrt{\sigma^2}}$

$\eta_2^2 = \eta_1$



$\mu = \sqrt{\sigma^2}$

## III) <u>Minimal & Full rank</u>

e.g. no extra constraint on $(\mu, \sigma^2)$

where the natural parameter space

is $(0, +\infty) \times \mathbb{R}$.



open rectangle.

## <u>Properties of exponential family</u>

<u>Property 1</u>: If $X_1, \ldots, X_n \overset{\text{iid}}{\sim} p(x; \theta) = \exp\{ \sum_{i=1}^{k} \eta_i^{(\theta)} T_i(x) - B(\theta)\} h(x)$,

then $p(x_1, \ldots x_n; \theta) = \exp\{ \sum_{i=1}^{k} \eta_i(\theta) \sum_{j=1}^{n} T_i(x_j) - nB(\theta)\} \prod_{j=1}^{n} h(x_j)$.

By NFFC, $(\sum_{j=1}^{n} T_1(x_j), \ldots, \sum_{j=1}^{n} T_k(x_j))$ is therefore a sufficient statistic.

Hence, exponential family data is highly compassible. (Pitman- koopman- Darmois).

<u>Property 2</u>: If $f$ is integrable and $\eta \in \Theta$, then $G(f, \eta) = \int f(x) \exp\{ \sum_{i=1}^{k} \eta_i T_i(x)\} d\mu(x)$ is

infinitely differentiable w.r.t. $\eta$ and the derivatives can be obtained by differentiating under

the integral sign. (see. TSH, 2.7.1).

<u>Property 3</u>: Moments of $T_i$'s.

Take, in particular, $f(x) = 1$, then

$G(f, \eta) \triangleq \int \exp\{ \sum_{i=1}^{k} \eta_i T_i(x)\} h(x) d\mu(x) = \exp\{A(\eta)\}$.
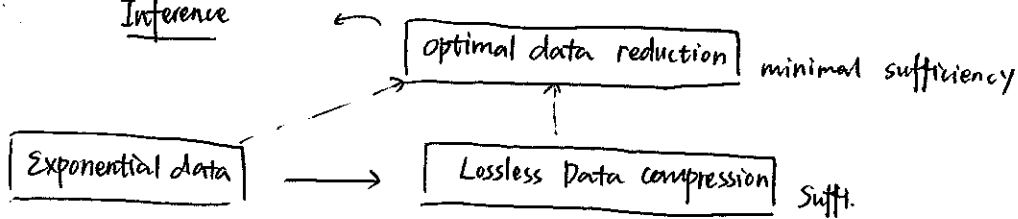
$\frac{\partial G(f, \eta)}{\partial \eta_i} = \int T_i(x) \exp\{ \sum_{i=1}^{k} \eta_i T_i(x)\} h(x) d\mu(x) = \frac{\partial A(\eta)}{\partial \eta_i} \exp\{A(\eta)\}$.

$\frac{\partial A(\eta)}{\partial \eta_i} = \int T_i(x) \exp\{ \sum_{i=1}^{k} \eta_i T_i(x) - A(\eta)\} h(x) d\mu(x) = E_\eta\{T_i(x)\}$.

# Minimal Sufficiency

**Def** A sufficient statistic $T$ is minimal ~~sufficien~~ if for every statistic $T'$. $T$ is a function of $T'$.

Equivalently, $T$ is minimal if for every sufficient statistic $T$, $T(x) = T(y)$ whenever $T'(x) = T'(y)$.



**[Thm]** Let $\{P_\theta(x)\}_{\theta \in \Theta}$ be a family of densities w.r.t some measure $\mu$ (usually Lebesgue). Suppose that there exists a ~~stochastic~~ statistic $S.t.$ for every $x, y \in X$.

$$P_\theta(x) = C_{x,y}\, P_\theta(y) \quad \Leftrightarrow \quad T(x) = T(y)$$

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

**Pf :** **[T is sufficient].**

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

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

and select an arbitrary representative $x_t$ from each $A_t$.

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 $A_t$, this implies that $T(y) = T(x_{T(y)})$

From the assumption of the thm,

$$P_\theta(y) = C_{y, x_{T(y)}} \, P_\theta(x_{T(y)})$$

$$= h(y)\, g_\theta(T(y))$$

which yields sufficiency of $T$ by NFFC.

**[T is minimal]**

Consider another sufficient statistic $T'$. By NFFC, $P_\theta(x) = \tilde{g}_\theta(T'(x))\, \tilde{h}(x)$

Take any $x, y$ s.t. $T'(x) = T'(y)$. Then

$$P_\theta(x) = \tilde{g}_\theta(T'(x))\, \tilde{h}(x)$$

$$= \tilde{g}_\theta(T'(y))\, \tilde{h}(y) \cdot \frac{\tilde{h}(x)}{\tilde{h}(y)}$$

$$= P_\theta(y) \cdot C_{x,y}.$$

Hence, $T(x) = T(y)$ by the assumption of the thm. So, $T'(x) = T'(y)$ implies that $T(x) = T(y)$ for any sufficient statistic $T'$ and $x, y$. As a result, $T$ is a minimal sufficient statistic.

Remark For any underline{minimal} $k$-dim exponential family, the statistic $(\sum_{j=1}^{n} T_1(x_j), \ldots \sum_{j=1}^{n} T_k(x_j))$ is a minimal sufficient statistic. (Keener Ex 3.12).

Remark The support of $X$ should be indep. of $\theta$. $\{x \in X : p_\theta(x) > 0\}$

e.g. $U(0, \theta)$, Binomial $(n, \theta)$.

## Ancillarity and completeness

$U(\theta, \theta+1)$ qualifying may appear.

E.g. Consider $X_1, \ldots, X_n \overset{iid}{\sim}$ Cauchy $(\theta)$ where distribution is given by
$$p_\theta(x) = \frac{1}{\pi} \cdot \frac{1}{1+(x-\theta)^2} = f(x-\theta).$$

then $(X_{(1)}, \ldots X_{(n)})$ is minimal sufficient    (TPE $\S$ 1.5)

Def A statistic $A$ is ancillary for $X \sim P_\theta \in P$ if the distribution of $A(X)$ does not dep. on $\theta$.

e.g. Consider the previous e.g $*$, then $A(X) = X_{(n)} - X_{(1)}$ is ancillary even though $(X_{(1)}, \ldots X_{(n)})$ is minimal sufficient. To see this, note that $X_i = Z_i + \theta$ for $Z_i \overset{iid}{\sim}$ Cauchy $(0)$, $X_{(i)} = Z_{(i)} + \theta$ and $A(X) = A(Z)$, which does not dep. on $\theta$.

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

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

Remark: For many important situations, completeness is the needed condition for the minimal sufficient and ancillary statistics to be independent.

Remark Complete, sufficient statistics give "optimal" unbiased estimator.

## Lecture 5 UMVUE, Cramér Rao Lower Bound.

Def A statistic $T$ is complete for $\longrightarrow$ See above.

Remark (i) If $T$ is complete sufficient, then $T$ is minimal sufficient (Bahader's Thm)    (TPE)
(ii) Complete sufficient statistics yield optimal unbiased estimators.

Example Let $X_1, \ldots, X_n \overset{iid}{\sim}$ Bernoulli $(\theta)$, $\theta \in (0,1)$, Then $T(X) = \sum_{i=1}^{n} X_i$ is sufficient.
Suppose $E_\theta \{f(T(X))\} = 0$ for all $\theta \in (0,1)$, this means $\sum_{j=0}^{n} f(j) \binom{n}{j} \theta^j (1-\theta)^{n-j} = 0$ $\forall \theta \in (0,1)$
Dividing both sides by $\theta^n$, and using $\beta = \frac{\theta}{1-\theta}$, we can write
$$\sum_{j=1}^{n} f(j) \binom{n}{j} \beta^j = 0 \qquad \forall \beta > 0.$$

(12)

If $f$ are non-zero, then the polynomial on the left is a ploynomial of degree at most $n$ which can only have $n$ roots at most. Hence, it is impossible to have LHS equal $0$ for writing $\beta$ unless $f=0$, in which case $T$ is complete.

## Example

Let $X_1, \ldots, X_n \overset{iid}{\sim} N(\theta, 6^2)$ with an unkown $\theta \in \mathbb{R}$ and a known $6^2 > 0$. We'd like to examine if $\bar{X}_n = n^{-1} \sum_{i=1}^{n} X_i$ is complete. To simplify our calculation, we consider the case with $n=1$ and $6=1$ so that $T(X) = X \sim N(\theta, 1)$.

Suppose $E_\theta \{f(X)\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(x) e^{-\frac{(x-\theta)^2}{2}} dx = 0 \quad \forall \theta \in \mathbb{R}$.

By multiplying $\sqrt{2\pi} e^{-\frac{\theta^2}{2}}$ on both sides, we obtain

$$\int_{-\infty}^{\infty} f(x) e^{-\frac{x^2}{2}} e^{\theta x} dx = 0 \qquad \forall \theta \in \mathbb{R}. \qquad (*)$$

Now, we decompose $f$ into its positive part and negative part as $f(x) = f_+(x) - f_-(x)$, where $f_+(x) = \max\{f(x), 0\}$, $f_-(x) = \max\{-f(x), 0\}$. Then $f_+(x) \geq 0$ and $f_-(x) \geq 0$ for all $x \in \mathbb{R}$. and $f_+(x) = f_-(x)$ iff $f_+(x) = f_-(x) = 0$.

If $f(x) \geq 0$ a.s. or $f(x) \leq 0$ a.s. , then $(*)$ implies that $f(x) = 0$ a.s. because setting $\theta = 0$ gives us and integral of a non-negative (or non-positive) function of being $0$. This is completeness.

In other words, if $f_+$ and $f_-$ have non-zero components, and we may write

$$\frac{\int_{-\infty}^{\infty} f_+(x) e^{-\frac{x^2}{2}} e^{\theta x} dx}{\int_{-\infty}^{\infty} f_+(x) e^{-\frac{x^2}{2}} dx} = \frac{\int_{-\infty}^{\infty} f_-(x) e^{-\frac{x^2}{2}} e^{\theta x} dx}{\int_{-\infty}^{\infty} f_-(x) e^{-\frac{x^2}{2}} dx} \qquad (**)$$

with the equality of the numerators follow from $(*)$ and equality of the denominator follows from $(*)$ with setting $\theta = 0$. The quantity $\frac{f_+(x) e^{-\frac{x^2}{2}}}{\int_{-\infty}^{\infty} f_+(x) e^{-\frac{x^2}{2}} dx}$ defines a prob. density and the LHS of $(**)$ is the MGF of this density.

Likewise, the RHS of $(**)$ is the MGF of $\frac{f_-(x) e^{-\frac{x^2}{2}}}{\int_{-\infty}^{\infty} f_-(x) e^{-\frac{x^2}{2}} dx}$

$\Rightarrow \quad f_+(x) = f_-(x)$ a.s. $\quad \Rightarrow \quad f_+(x) = f_-(x) = 0$ a.s.

$\qquad \Rightarrow \quad f(x) = 0 \qquad \Rightarrow \quad T$ is complete

## [Thm] (Basu's Thm)

If $T$ is complete and sufficient for $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ and $A$ is ancillary, then $T(X) \perp\!\!\!\perp A(X)$.

Pf: N.T.S. $P_\theta(A \in \mathcal{A} \mid T) = P_\theta(A \in \mathcal{A})$ a.s. $P_\theta$.

Observations:
- LHS is free of $\theta$ by sufficiency of $T$ $\quad\}$ free of $\theta$.
- RHS is free of $\theta$ since $A$ is ancillary.

$E_\theta (LHS) = E_\theta (RHS)$ by tower property.

$\Rightarrow$ LHS = RHS with prob. 1 by completeness of $T$. //

Example $X_1, \ldots, X_n \overset{iid}{\sim} N(\mu, \sigma^2)$ $(\mu, \sigma^2$ both unkown$)$

Claim: $\bar{X}_n \perp\!\!\!\perp (n-1)^{-1} \sum_{i=1}^{n} (X_i - \bar{X})^2 \underbrace{\quad}$ Sample variance.

Pf: Fix any $\sigma > 0$ and consider a submodel $P_\theta = \{ N(\mu, \sigma^2), \mu \in \mathbb{R} \}$

In each submodel, $\bar{X}_n$ is complete and sufficient, and $n^{-1} \sum_{i=1}^{n} (X_i - \bar{X})^2$ is ancillary.

By Basu's thm, $\bar{X}_n \perp\!\!\!\perp n^{-1} \sum_{i=1}^{n} (X_i - \bar{X})^2$ under $N(\mu, \sigma^2)$ for any $\mu$.

Since $\sigma$ is an arbitrary, we have $\bar{X}_n \perp\!\!\!\perp n^{-1} \sum_{i=1}^{n} (X_i - \bar{X})^2$ for the full model. //

From Data compression to Risk Reduction / "optimal" estimation

[Thm] (Rao-Blackwell Thm) Keener 3.28.

Suppose that $T$ is sufficient for $P = \{ P_\theta : \theta \in \Theta \}$, that $\delta(X)$ is an estimation for $g(\theta)$

for which $E\{\delta(X)\}$ exists, and that $R(\theta, \delta) = E_\theta \{ L(\theta, \delta(X)) \} < \infty$.

If $L(\theta, \cdot)$ is convex, then

Jensen's ineq:
$$E\{ \phi(X) \} \geqslant \phi \{ E(X) \}$$
$$\phi(\cdot) \text{ convex}$$

$$R(\theta, \eta) \leqslant R(\theta, \delta)$$

where $\eta(T(X)) = \underline{E\{ \delta(X) | T(X) \}}$.

If $L(\theta, \cdot)$ is strictly convex, then $R(\theta, \eta) < R(\theta, \delta)$ for any $\theta$

unless $\eta(T'(X)) = \delta$ with prob. 1.

Pf: By Jensen's inequality,
$$E_\theta \{ L(g(\theta), \delta(X)) | T \} \geqslant L[ g(\theta, E_\theta \{ \delta(X) | T \}) ]$$
$$= L(g(\theta), \eta(T))$$

Taking another expectation, we have

$$E_\theta \{ L(g(\theta), \delta(X)) \} \geqslant E_\theta \{ L(g(\theta), \eta(T)) \}$$

$$\Rightarrow R(g(\theta), \delta) \geqslant R(g(\theta), \eta) \qquad //.$$

Example Let $X_1, X_2, \ldots, X_n \overset{iid}{\sim}$ Bernoulli $(\theta)$, $\theta \in (0, 1)$

Consider the loss function $L(\theta, d) = (\theta - d)^2$ [ squared loss function ].

Suppose we consider first an unreasonable estimator $\delta(X) = X_1$.

We have shown that $T(X) = \bar{X}_n$ is sufficient, so we may apply

Rao-Blackwell thm to improve $\delta$. In particular,

$$\eta(T(X)) = E\{\delta(X) | T(X)\}$$
$$= \frac{1}{n} \sum_{i=1}^{n} E(X_i | \bar{X}_n)$$
$$\overset{iid}{=} \frac{1}{n} \sum_{i=1}^{n} E(X_i | \bar{X}_n)$$
$$= E(\bar{X}_n | \bar{X}_n)$$
$$= \bar{X}_n.$$

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

R.B. gives a strict improvement. //

Remark : Rao-Blackwell thm, however, does not necessarily lead to a uniformly optimal estimator. For example, consider $\delta$ naive $(X) = \frac{1}{2}$, then

$$\eta(T(X)) = E\{\delta \text{ naive } (X) | \bar{X}\} = \frac{1}{2} \text{ as well.}$$

Since $R(\theta, \eta) = (\frac{1}{2} - \theta)^2$, neither R.B ised outcome is uniformly better across all $\theta$.

## Unbiased Estimation

An estimator is _unbiased_ if $E_\theta\{\delta(X)\} = g(\theta)$. We attempt to find an unbiased estimator with uniformly minimum risk. i.e. un biased $\delta$ satisfying $R(\theta, \delta) \le R(\theta, \delta')$. for all $\theta \in \theta$. and an unbiased estimator $\delta'$. Such an estimator is called uniformly minimum risk unbiased estimator (UMRUE). If, in particular, $L(\theta, \delta) = (\theta - \delta)^2$ is the choosen loss function, then an UMRUE becomes UMVUE.

$$E_\theta\{g(\theta) - \delta(X)\}^2 = \underbrace{[E_\theta\{\delta(X) - g(\theta)\}]^2}_{\underset{\longrightarrow 0 \text{ for } \delta(X) \text{ is an unbiased est. for } g(\theta).}{Bias^2}} + \underbrace{E_\theta[\delta(X) - E_\theta\{\delta(X)\}]^2}_{Variance}$$

[Thm] (Lehmann- Scheffé Thm)

If $T$ is a complete and sufficient statistics, and $E_\theta\{h(T(X))\} = g(\theta)$. i.e. $h(T(X))$ is unbiased for $g(\theta)$, then $h(T(X))$ is

    (1) the only function of $T(X)$ that is unbiased for $g(\theta)$

    (2) an UMRUE under any convex loss function

    (3) the unique UMRUE (up to a P-null set) under any strictly convex loss function

    (4) the unique UMVUE (up to a P-null set)

Pf: (1) Suppose $E_\theta\{\tilde{h}(T(X))\} = g(\theta)$, then $E_\theta\{\tilde{h}(T(X)) - h(T(X))\} = 0 \quad \forall \theta \in \Theta$.

thus $\tilde{h}(T(X)) = h(T(X))$ a.s. for each $\theta$ due to completeness.

(2) Consider an unbiased estimator $\delta(X)$, and let $\tilde{h}(T(X)) = E_\theta\{\delta(X) | T(X)\}$.

Then $E_\theta\{\tilde{h}(T(X))\} = E_\theta\{\delta(X)\} = g(\theta)$ by tower property of conditional expectation.

By (1). $\tilde{h}(T(X)) = h(T(X))$ and by the Rao-Blackwell thm,

$R(g(\theta), h(T(\cdot))) \leq R(g(\theta), \tilde{h}(T(\cdot)))$ for all $\theta$ if the loss function is convex.

Therefore, $h(T(X))$ is an UMRUE under any convex loss function.

(3) If the loss function is strictly convex, $R(g(\theta), h(T(\cdot))) < R(g(\theta), \delta)$ unless $\delta(X) = h(T(X))$ a.s. Thus, $h(T(X))$ is the unique UMRUE.

(4) Done by (3). $\quad //$ .

Strategies for obtaining UMVUEs (k)

A — Rao-Blackwellisation / Conditioning $E(\cdot | T)$

B — Solve for $\delta$ satisfying $E_\theta\{\delta(T(x))\} = g(\theta), \forall \theta \in \Theta$.

Example

Suppose $X_1, \ldots, X_n \overset{iid}{\sim}$ Bernoulli $(\theta)$

$T(X) = \overset{\sum_{i=1}^{n} X_n}{X_n}$ is complete and suff. and

$E\{T(x)/n\} = \theta$. Therefore $\bar{X}_n$ is an UMRUE for $\theta$ under any convex function.

Suppose now we are interested in estimating $g(\theta) = \theta^2$.

If we choose $\delta(X) = I(X_1 = X_2 = 1) = X_1 X_2$, then $E_\theta\{\delta(X)\} = \theta^2$ is unbiased.

Apply strategy A to obtain:

$$E\{\delta(X) | T(X) = t\} = P(X_1 = X_2 = 1 | T(X) = t)$$
$$= \frac{P(X_1 = X_2 = 1, \sum_{i=3}^{n} X_i = t-2)}{P(T(X) = t)}$$
$$= \frac{\theta^2 \binom{n-2}{t-2} \theta^{t-2} (1-\theta)^{n-t} I(t \geq 2)}{\binom{n}{2} \theta^t (1-\theta)^{n-t}}$$
$$= \frac{t(t-1) I(t \geq 2)}{n(n-1)} = \frac{t(t-1)}{n(n-1)}$$

Hence, $\dfrac{T(X)\{T(X) - 1\}}{n(n-1)}$ is the UMVUE.

Example Suppose $X_1, \ldots, X_n \overset{iid}{\sim}$ uniform $(0, \theta)$. In this case, $T(X) = X_{(n)} = \max_{1 \leq i \leq n} X_i$ is a complete and sufficient statistic and $\delta(X) = 2X_1$ is an unbiased estimator of $\theta$.

Given the knowledge of $X_{(n)}$, $X_1$ is equal to $X_{(n)}$ with prob. $1/n$ and distributed according to uniform $(0, X_{(n)})$ with prob. $1 - 1/n$.

So $\quad P(X_1 = x_1 \mid T(X)) = \frac{1}{n} I(T(X) = x_1) + \frac{(1 - \frac{1}{n}) I(0 < x_1 < T(X))}{T(X)}$

Hence, our UMVUE,

$$E\{\delta(X) \mid T(X)\} = 2 E\{X_1 \mid T(X)\}$$

$$= 2 \left\{ \frac{1}{n} T(X) + (1 - \frac{1}{n}) \int_0^{T(X)} \frac{x_1 dx_1}{T(X)} \right\}$$

$$= 2 \left\{ \frac{T(X)}{n} + (1 - \frac{1}{n}) \frac{T(X)}{2} \right\}$$

$$= \left( \frac{n+1}{n} \right) T(X) \qquad // .$$

## Example

Let $X_1, X_2, \ldots, X_n$ iid Poisson $(\theta)$. Since this is a one-dimensional full-rank exponential family, $X$ is a complete sufficient statistic. $X$ is furthermore unbiased and therefore UMV for $\theta$. Suppose that our goal is to estimate $g(\theta) = e^{-a\theta}$ for some given $a \in \mathbb{R}$. We need to find an estimator $\delta$ such that $E\{\delta(X)\} = g(\theta)$ for all $\theta$. Under our model, we may reexpress this system of equations as:

$$\sum_{X=0}^{\infty} \delta(X) \frac{e^{-\theta} \theta^X}{X!} = e^{-a\theta} \quad \text{for all } \theta$$

$$\Rightarrow \sum_{X=0}^{\infty} \frac{\delta(X) \theta^X}{X!} = e^{(1-a)\theta} = \sum_{X=0}^{\infty} \frac{(1-a)^X \theta^X}{X!}$$

$$\Rightarrow \delta(X) = (1-a)^X \quad \text{is the UMVUE of } g(\theta).$$

## Remark:

This estimator is not satisfying. If $a = 2$, for example, it will change its sign according to $X$ even though we realize that our estimand $e^{-\theta a}$ is non-negative. The estimator is in fact in admissible when $a > 1$ and dominated by $\max[\delta(X), 0]$.

Ch5 of TPE.

Suppose we have $\delta_i$ UMVU for $g(\theta)$ for $i \in \{1, 2\}$. Is $\delta_1 + \delta_2$ then UMVU for $\underline{g_1(\theta) + g_2(\theta)}$?

[Thm] (Characterization of UMVUEs, see TPE 2.1.7).

Let $\Delta = \{\delta : E_\theta(\delta^2) < \infty\}$. Then $\delta_0 \in \Delta$ is UMVUE for $g(\theta) = E(\delta_0)$ iff $E\{\delta_0(\theta) u\} = 0$ for every $u \in \mathcal{U}$, where $\mathcal{U} = \{$unbiased estimator of $0\}$

$$= \{u : X \to \mathbb{R} \text{ s.t. } E_\theta(u(X)) = 0, \ E_\theta(u(X)^2) = \infty\}.$$

Pf: If $\delta_0$ is UMVUE, let us consider $\delta_\lambda = \delta_0 + \lambda u$ for $\lambda \in \mathbb{R}$, $u \in \mathcal{U}$.

Since $\delta_0$ has the minimal variance,

$$\text{Var}(\delta_\lambda) = \text{Var}(\delta_0) + \lambda^2 \text{Var}(u) + 2\lambda \text{cov}(\delta_0, u) \geqslant \text{Var}(\delta_0) \quad (\text{UMVU}). \quad (\#)$$

Consider the quadratic form $q(\lambda) = \lambda^2 \text{Var}(u) + 2\lambda \text{cov}(\delta_0, u)$,

then $q$ has the roots $0$ and $-2\text{cov}(\delta_0, u) / \text{Var}(u)$,

If the roots are distinct, then the form must be negative at some point, which would violet the inequality (#). Hence $-2\text{cov}(\delta_0, u)/\text{Var}(u) = 0$, and thus $E(u\delta_0) = \text{cov}(\delta_0, u) = 0$.

For the converse result, we assume $E(\delta_0 u) = 0 \ \forall u \in \mathcal{U}$. and consider any unbiased estimator $\delta$ for $g(\theta)$. Then $\delta - \delta_0 \in \mathcal{U}$ so $E\{\delta_0(\delta - \delta_0)\} = 0$.

This implies that $E(\delta_0 \delta) = E(\delta_0^2)$ and substracting $E(\delta_0)E(\delta)$ on both sides, we have

$$\text{Var}(\delta_0) = \text{cov}(\delta_0, \delta) \overset{cs}{\leq} \sqrt{\text{Var}(\delta_0)\text{Var}(\delta)}$$

Hence $\text{Var}(\delta_0) \leq \text{Var}(\delta)$ for any arbitrary estimator $\delta$ and $\delta_0$ is UMVUE.

$$\forall u \in \mathcal{U}, \quad E((\delta_1 + \delta_2)u) = E(\delta_1 u) + E(\delta_2 u) = 0$$
$$\Rightarrow \delta_1 + \delta_2 \text{ is UMVU for } g_1(\theta) + g_2(\theta)$$

## Lecture 6

· Rao-Blackwell, LS, Basu, ...
Cramer-Rao Lower Bound · TPE §2.5 & 2.6.   Keener Ch.3?   C&B Ch.7.

Assume the following:

(a) $\Theta \subseteq \mathbb{R}$ is an open interval

(b) $\{P_\theta : \theta \in \Theta\}$ have common support $A$

(c) $P_\theta'(X) = \dfrac{\partial P_\theta(X)}{\partial X}$ exists and is finite for all $x \in A$.

Define $I(\theta) = E_\theta\left\{\dfrac{\partial}{\partial\theta}\log P_\theta(x)\right\}^2 = \int_A \left\{\dfrac{P_\theta'(x)}{P_\theta(x)}\right\}^2 P_\theta(x)\,d\mu = \int_A \dfrac{\{P_\theta'(x)\}^2}{P_\theta(x)}\,d\mu$ to be the information function.

**Lemma** (i) Assume (a)-(c) hold, and $\dfrac{\partial}{\partial\theta}\int_A P_\theta(x)\,d\mu = \int_A \dfrac{\partial}{\partial\theta}P_\theta(x)\,d\mu$.

then $I(\theta) = \text{Var}_\theta\left\{\dfrac{\partial}{\partial\theta}\log P_\theta(x)\right\}$

(ii) In addition, $P_\theta''(x)$ exists $\forall \theta \in \Theta$, $x \in A$ and

(e) $\dfrac{\partial^2}{\partial\theta^2}\int_A P_\theta(x)\,d\mu = \int_A \dfrac{\partial^2}{\partial\theta^2}P_\theta(x)\,d\mu$, then $I(\theta) = -E_\theta\left\{\dfrac{\partial^2}{\partial\theta^2}\log P_\theta(x)\right\}$.

· Pf: (i) We need to show $I(\theta) = \text{Var}\left(\dfrac{\partial}{\partial\theta}\log P_\theta(x)\right)$. Assume that $I_\theta = E_\theta\left(\dfrac{\partial}{\partial\theta}\log P_\theta(x)\right)^2 < \infty$.

It suffices to show that $E_\theta\left(\dfrac{\partial}{\partial\theta}\log P_\theta(x)\right) = 0$. This is true because

$$E_\theta\left(\dfrac{\partial}{\partial\theta}\log P_\theta(x)\right) = \int_A \dfrac{P_\theta'(x)}{P_\theta(x)} P_\theta(x)\,d\mu = \int_A \dfrac{\partial}{\partial\theta}P_\theta(x)\,d\mu = \dfrac{\partial}{\partial\theta}\underset{\to 1.}{\underbrace{\int_A P_\theta(x)\,d\mu}} = 0$$

(ii) Note that $\dfrac{\partial^2}{\partial\theta^2}\log P_\theta(x) = \dfrac{\partial}{\partial\theta}\left(\dfrac{P_\theta'(x)}{P_\theta(x)}\right) = \dfrac{P_\theta''(x)}{P_\theta(x)} - \left(\dfrac{P_\theta'(x)}{P_\theta(x)}\right)^2$

thus $E\left\{\dfrac{\partial^2}{\partial\theta^2}\log P_\theta(x)\right\} = \int_A P_\theta''(x)\,d\mu - E_\theta\left\{\dfrac{\partial}{\partial\theta}\log P_\theta(x)\right\}^2$

$$= \int_A \dfrac{\partial^2}{\partial\theta^2}P_\theta(x)\,d\mu - I(\theta) \overset{(a)}{=} \dfrac{\partial^2}{\partial\theta^2}\int_A P_\theta(x)\,d\mu - I(\theta) = -I(\theta)$$

⑮

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

**Remark**  Information depends on parametrization. For example, if $\eta = \tau(\theta)$, where $\tau \in C^2$, & $\tau'(\theta) \neq 0$, then $I(\tau(\theta)) = \frac{I(\theta)}{\{\tau'(\theta)\}^2}$ because

$$E_\tau \left\{ \frac{\partial}{\partial \tau} \log p_{\theta(\tau)}(x) \right\}^2 = E_\theta \left\{ \frac{\partial}{\partial \theta} \log p_\theta(x) \cdot \frac{\partial \theta}{\partial \tau} \right\}^2 = \frac{E_\theta \left\{ \frac{\partial}{\partial \theta} \log p_\theta(x) \right\}^2}{[\tau'(\theta)]^2} = \frac{I(\theta)}{\{\tau'(\theta)\}^2}$$

**Example**  Suppose $X \sim N(\theta, 1)$, $\theta > 0$

$$p_\theta(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\theta)^2}{2}} \Rightarrow \log p_\theta(x) = -\log \sqrt{2\pi} - \frac{(x-\theta)^2}{2} \quad \checkmark$$

$$\Rightarrow \frac{\partial}{\partial \theta} \log p_\theta(x) = x - \theta, \quad I(\theta) = E(x-\theta)^2 = \text{Var}(X) = 1$$

$$\left( \text{or} \quad \frac{\partial^2}{\partial \theta^2} \log p_\theta(x) = -1 \Rightarrow I(\theta) = -(-1) = 1 \right)$$

Let $\eta = \theta^2$, $X \sim N(\sqrt{\eta}, 1)$. $\quad \log p_{\theta(\eta)}(x) = -\log(\sqrt{2\pi}) - \frac{(x-\sqrt{\eta})^2}{2}$

$$\Rightarrow \frac{\partial}{\partial \eta} \log p_{\theta(\eta)}(x) = \frac{x-\sqrt{\eta}}{2\sqrt{\eta}} \Rightarrow I(\eta) = \frac{E(x-\sqrt{\eta})^2}{4\eta} = \frac{1}{4\eta} = \frac{1}{4\theta^2} = \frac{1}{(\tau'(\theta))^k}$$

# Multi-parameter Cramér-Rao Lower Bound

Suppose the following conditions hold:

(a) $\Theta \in \mathbb{R}^k$ is an open set

(b) $\{P_\theta(x) : \theta \in \Theta\}$ have common support $I$

(c) $\frac{\partial P_\theta(x)}{\partial \theta_i}$ exist $\forall i = 1, \ldots, k$, $x \in I$ and is finite

(d) $\frac{\partial}{\partial \theta_i} \int_X P_\theta(x) \, d\mu = \int_X \frac{\partial}{\partial \theta_i} P_\theta(x) \, d\mu \quad \forall i = 1, \ldots, k$

(e) $\frac{\partial}{\partial \theta_i} \int_X \delta(x) P_\theta(x) \, d\mu = \int_X \frac{\partial}{\partial \theta_i} \delta(x) P_\theta(x) \, d\mu \quad \forall i = 1, \ldots, k$.

Define the $k \times k$ information matrix $I(\underline{\theta})$ by $I_{ij}(\underline{\theta}) = E_\theta \left\{ \left( \frac{\partial}{\partial \theta_i} \log P_\theta(x) \right) \left( \frac{\partial}{\partial \theta_j} \log P_\theta(x) \right) \right\}$

In particular, if $k = 1$, $I(\theta) = E_\theta \left( \frac{\partial}{\partial \theta} \log P_\theta(x) \right)^2$.

Assume $I(\underline{\theta})$ is finite and positive definite,

then $\text{Var}_\theta(\delta(x)) \geq \alpha^T I(\underline{\theta})^{-1} \alpha$, where $\alpha = (\alpha_1, \ldots, \alpha_k)^T$

$$= \left( \frac{\partial}{\partial \theta_1} E_\theta(\delta(x)), \ldots, \frac{\partial}{\partial \theta_k} E_\theta(\delta(x)) \right)^T.$$

In particular, if $\delta(x)$ is unbiased for $g(\theta)$,

then $\text{Var}_\theta(\delta(x)) \geq \alpha^T I(\theta)^{-1} \alpha$, $\alpha_i = \frac{\partial}{\partial \theta_i} \{ g(\theta) \}$. $i = 1, 2, \ldots, k$.

Pf: Let $\Psi_i(x) = \frac{\partial}{\partial \theta_i} \log P_\theta(x)$. then $E_\theta(\Psi_i(x)) = \int_X \left\{ \frac{\partial}{\partial \theta_i} \log P_\theta(x) \right\} P_\theta(x_i) \, d\mu$

$$= \int_X \frac{\frac{\partial}{\partial \theta_i} P_\theta(x)}{P_\theta(x)} P_\theta(x) \, d\mu = \int_X \frac{\partial}{\partial \theta_i} P_\theta(x) \, d\mu = \frac{\partial}{\partial \theta_i} \int P_\theta(x) \, d\mu = 0.$$

Fix a non-zero vector $(a_1, \ldots, a_k)$. Then $E_\theta \left\{ \sum_{i=1}^k a_i \Psi_i(x) \right\} = 0$

Claim: $\text{Var} \left( \sum_{i=1}^k a_i \Psi_i(x) \right) = a^T I(\underline{\theta}) a$.

Observe that $\text{Var} \left( \sum_{i=1}^k a_i \Psi_i(x) \right) = \sum_{i,j} a_i a_j \text{Cov}(\Psi_i(x), \Psi_j(x))$

$$= \sum_{i,j} a_i a_j E(\Psi_i(x) \Psi_j(y))$$

$$= \sum_{i,j} a_i a_j I_{ij}(\underline{\theta}) = a^T I(\underline{\theta}) a.$$

Finally, $\text{Cov} \left( \delta(x), \sum_{i=1}^k a_i \Psi_i(x) \right) = \sum_{i=1}^k a_i \text{Cov}(\delta(x), \Psi_i(x))$

$$= \sum_{i=1}^k a_i E(\delta(x) \Psi_i(x)) = \sum_{i=1}^k a_i \int_X \delta(x) \frac{\partial}{\partial \theta_i} \log P_\theta(x) \cdot P_\theta(x) \, d\mu.$$

$$= \sum_{i=1}^k a_i \int_X \delta(x) \frac{\partial}{\partial \theta} P_\theta(x) \, d\mu = \sum_{i=1}^k a_i \frac{\partial}{\partial \theta_i} \int_X \delta(x) P_\theta(x) \, d\mu = \sum_{i=1}^k a_i \alpha_i(\underline{\theta})$$

By Cauchy-Schwarz inequality.

$$\text{Var}(\delta(X)) \, \text{Var}\left(\sum_{i=1}^{n} a_i \psi_i(X)\right) \geq \text{Cov}\left(\sum_{i=1}^{k} a_i \psi_i(X), \delta(X)\right)$$

$$\Rightarrow \text{Var}(\delta(X)) \geq \sup_{a \neq 0} \frac{\left(\sum_{i=1}^{k} a_i \alpha_i(\theta)\right)^2}{a^T I(\theta) a} = \alpha^T I(\theta)^{-1} \alpha.$$

E.g. $X_1, \ldots, X_n \overset{iid}{\sim} N(\mu, \sigma^2)$, $\mu > 0$, $\sigma^2 > 0$.

Problem 1: We want to ~~simulate~~ estimate $g_1(\mu, \sigma^2) = \mu$

Problem 2: $\qquad\qquad\qquad\qquad\qquad\qquad\qquad g_2(\mu, \sigma^2) = \sigma^2$

Consider unbiased estimator only.

Claim 1. $\frac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X})^2$ is UMVUE for $\sigma^2$.

why? ① $\left(\sum_{i=1}^{n} X_i, \sum_{i=1}^{n} X_i^2\right)$ is complete sufficient. ② $\frac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X})^2 = \frac{1}{n-1}\left\{\sum_{i=1}^{n} X_i^2 - n\bar{X}^2\right\}$.

③ $\frac{1}{\sigma^2} \sum_{i=1}^{n} (X_i - \bar{X})^2 \sim \chi^2_{n-1} \Rightarrow E\left\{\frac{1}{n-1}\sum_{i=1}^{n}(X_i-\bar{X})^2\right\} = \frac{\sigma^2}{n-1} \cdot (n-1) = \sigma^2$ (unbiased).

Claim 2: $\bar{X}$ is UMVUE for $\mu$. $E(\bar{X}-\mu)^2 = \frac{\sigma^2}{n}$

$\qquad$ Note that $E\left\{\sum_{i=1}^{n}\frac{(X_i-\bar{X})^2}{n-1} \cdot \sigma^2\right\}^2 = \text{Var}\left(\sum_{i=1}^{n}\frac{(X_i-\bar{X})^2}{n-1}\right) = \frac{2(n-1)\sigma^4}{(n-1)^2} = \frac{2\sigma^4}{(n-1)}$

$\log p_{\mu,\sigma^2}(x) = -\frac{n}{2}\log\sigma^2 - \sum_{i=1}^{n}\frac{(X_i-\mu)^2}{2\sigma^2} + Cn$.

$\frac{\partial}{\partial\mu}\log p_{\mu,\sigma^2}(x) = \sum_{i=1}^{n}\frac{(X_i-\mu)}{\sigma^2}$ $\boxed{\frac{\partial^2}{\partial\mu^2}(\log p_{\mu,\sigma}(x)) = -\frac{\mu}{\sigma^2}}$ $I_{11}$

$\frac{\partial}{\partial\sigma^2}\log p_{\mu,\sigma^2}(x) = -\frac{n}{2\sigma^2} + \sum_{i=1}^{n}\frac{(X_i-\mu)^2}{2\sigma^4}$

$\frac{\partial^2}{\partial\mu\partial\sigma^2}\log p_{\mu,\sigma^2}(x) = -\sum_{i=1}^{n}\frac{(X_i-\mu)}{\sigma^4}$

$\frac{\partial^2}{\partial(\sigma^2)^2}\log p_{\mu,\sigma^2}(x) = \frac{n}{2\sigma^4} - \sum_{i=1}^{n}\frac{(X_i-\mu)^2}{\sigma^6}$

$\qquad I = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix}$

$\qquad I_{22} = -\frac{n}{2\sigma^4} + E\sum_{i=1}^{n}\frac{(X_i-\mu)^2}{\sigma^6} = -\frac{n}{2\sigma^4} + \frac{n\sigma^2}{\sigma^6} = \boxed{\frac{n}{2\sigma^4}}$

$\Rightarrow$ CRLB for $\mu = [1 \ 0] I^{-1}\begin{bmatrix}1\\0\end{bmatrix} = \frac{1}{I_{11}} = \frac{\sigma^2}{n}$ $\qquad\left.\begin{array}{l}\text{CRLB attained by } \bar{X} \\ \text{but not } S^2.\end{array}\right.$

$\Rightarrow$ CRLB for $\sigma^2 = [0 \ 1] I^{-1}\begin{bmatrix}0\\1\end{bmatrix} = \frac{1}{I_{22}} = \frac{2\sigma^4}{n}$

# Cramer-Rao lower bound  (Information inequality)

Suppose (a) – (d) hold, let $\delta(X)$ be an estimator s.t. $E_\theta\{\delta(X)\}^2 < \infty$. $I(\theta) \in (0, \infty)$

Assume further that $\int_A \delta(X) \frac{\partial}{\partial\theta} P_\theta(x)\, d\mu = \frac{\partial}{\partial\theta} E_\theta\{\delta(X)\}$. Then

$$Var_\theta(\delta(X)) \geq \frac{\left[\frac{\partial}{\partial\theta} E_\theta\{\delta(X)\}\right]^2}{I(\theta)}$$

Remark: If we want to estimate $g(\theta)$ using an unbiased estimator $\delta$,

then $Var_\theta(\delta(X)) \geq \frac{\{g'(\theta)\}^2}{I(\theta)}$.

Pf: Let $V = \frac{\partial}{\partial\theta} \log P_\theta(x)$, so $E_\theta(V^2) = I(\theta) = Var_\theta(V)$. Also $E_\theta(V) = 0$.

By Cauchy Schwarz inequality, $Var_\theta(V) \, Var_\theta(\delta(X)) \geq \{cov(V, \delta(X))\}^2$

$\Rightarrow \quad Var(\delta(X)) \geq \frac{\{cov(V, \delta(X))\}^2}{I(\theta)}$.

It suffices to show that $cov(V, \delta(X)) = \frac{\partial}{\partial\theta} E_\theta(\delta(X))$

Observe that $cov(V, \delta(X)) = E_\theta(V(\delta(X))) = \int_A \frac{\partial}{\partial\theta} \log P_\theta(x)\, \delta(x)\, P_\theta(x)\, d\mu$

$$= \int_A \frac{\partial P_\theta(x)}{\partial\theta} \delta(x)\, d\mu = \frac{\partial}{\partial\theta} \int_A P_\theta(x)\, \delta(x)\, d\mu$$

$$= \frac{\partial}{\partial\theta} E_\theta(\delta(X)), \text{ where } A = \{x \in X : P_\theta(x) > 0\}. \quad //$$

Example  Let $X_1, \ldots, X_n \overset{iid}{\sim} N(\theta, 1)$. We are interested in estimating $g(\theta) = \theta$.

Take $\delta(X) = \bar{X}_n$. Then $Var(\bar{X}_n) = \frac{1}{n}$.

For any unbiased estimator $\delta$, we have

$$Var(\delta(X)) \geq \frac{1}{I_n(\theta)} = \frac{1}{n},$$

where $I_n(\theta) = E\left\{\frac{\partial}{\partial\theta} \log P_\theta(x_1, \ldots, x_n)\right\}^2 = n$,

~~then~~ Hence $\bar{X}_n$ is UMVUE.

$I_n(\theta) = E\left\{\frac{\partial}{\partial\theta} \log P_\theta(X_1) + \cdots + \frac{\partial}{\partial\theta} \log P_\theta(X_n)\right\}^2$

$$= \sum_{i=1}^{n} E\left\{\frac{\partial}{\partial\theta} \log P_\theta(X_i)\right\} + \sum_{i \neq j} E\left\{\frac{\partial}{\partial\theta} \log P_\theta(X_i) \cdot \frac{\partial}{\partial\theta} \log P_\theta(X_j)\right\}^{\nearrow 0}$$

$$= \sum_{i=1}^{n} 1$$

$$= n.$$

<u>Example</u>  $X_1,\ldots,X_n \overset{iid}{\sim}$ Poisson $(\lambda)$, $\lambda > 0$, We want to verify that $\bar{X}_n$ is UMVUE for $\lambda$.

Observe that  $Var(\bar{X}_n) = \dfrac{Var(X_1)}{n} = \dfrac{\lambda}{n}$.   $E(\bar{X}_n) = \lambda$.

Also, $Var_\lambda(\delta(X)) \geqslant \dfrac{1}{n I_1(\lambda)}$.

It suffices to show that $I_1(\lambda) = \dfrac{1}{\lambda}$.

Recall that  $p_\lambda(x) = \dfrac{e^{-\lambda}\lambda^x}{x!}$

$\log p_\lambda(x) = -\lambda + x\log\lambda - \log(x!)$

$\dfrac{\partial}{\partial\lambda}\log p_\lambda(x) = -1 + \dfrac{x}{\lambda}$  $\Rightarrow$  $I_1(\lambda) = E\left(\dfrac{x}{\lambda}-1\right)^2 = \dfrac{Var(x)}{\lambda^2} = \dfrac{\lambda}{\lambda^2} = \dfrac{1}{\lambda}$  ,.

Remark: Assume (a)-(d) hold, then $Var_\theta(\delta(x)) \geqslant \dfrac{\{\frac{\partial}{\partial\theta}E_\theta(\delta(x))\}^2}{I(\theta)}$

If the equality holds, then  $\dfrac{\partial}{\partial\theta}\log p_\theta(x) = a(\theta)\delta(x) + b(\theta)$   a.s. $*P_\theta$ (measure)  $^{w.r.t.}$

Assume that $\frac{\partial}{\partial\theta}\log p_\theta(x)$ is continuous in $\theta$, then $\theta\mapsto a(\theta)$ and $\theta\mapsto b(\theta)$ are also continuous, provided $\delta$ is not a degenerate r.v.. It follows that for fixed $\theta_0 \in \Theta$,

$$p_\theta(x) = p_{\theta_0}(x)\, e^{\{\int_{\theta_0}^\theta a(t)dt\}\delta(x) + \{\int_{\theta_0}^\theta b(t)dt\}}$$

Thus, $p_\theta$ is a 1-parameter exponential family and $\delta(x)$ is the natural sufficient statistic.

Let  $A = \{x : \frac{\partial}{\partial\theta}\log p_\theta(x) = a(\theta)\delta(x) + b(\theta)\}$.  then $\exists\ x_1 \neq x_2$ s.t. $\delta(x_1) \neq \delta(x_2)$

For the $x_1, x_2$ equality holds in $\frac{\partial}{\partial\theta}\log p_\theta(x) = a(\theta) \cdot \delta(x) + b(\theta) = h(x,\theta)$

$\Rightarrow h(x_1, \theta) = a(\theta)\delta(x_1) + b(\theta)$

$\quad h(x_2, \theta) = a(\theta)\delta(x_2) + b(\theta)$

$\Rightarrow \dfrac{h(x_1,\theta) - h(x_2,\theta)}{\delta(x_1) - \delta(x_2)} = a(\theta)$

$\Rightarrow a$ is cont. $\Rightarrow b$ is cont.

[Thm] Let $p_\theta(x) = e^{\eta(\theta)T(x) - B(\theta)}\tilde{h}(x)$, $\theta \in \Theta$ open interval. Let $\tau(\theta) = E_\theta(T)$

Assume $T$ is not a constant r.v., then

(a) $\tau'(\theta) \neq 0$ and $I(\tau(\theta)) = \dfrac{1}{Var_\theta(T)}$

(b) $I(h(\theta)) = \left(\dfrac{\eta'(\theta)}{h'(\theta)}\right)^2 Var_\theta(T)$.

- Average Risk Optimality (Bayes Estimator) TPE Ch.4.

  Suppose $\{P_\theta : \theta \in \Theta\}$ is a collection of prob. measures on $X$ dominated by $\sigma$-finite measure $\mu$. Assume that now $\theta$ is a random variable on $\Theta$ with dist. $\pi$, which is regarded as the _prior dist._

Suppose we want to estimate $g(\theta)$, where $g: \Theta \to \mathbb{R}$. For an estimator $\delta(X)$, let the loss incurred be $L(g(\theta), \delta(X))$. Then the risk function, as defined before, is

$$R(g(\theta), \delta) = E_{X \sim P_\theta}\{L(g(\theta), \delta(X))\} = E\{L(g(\theta), \delta(X)) \mid \theta\}.$$

Define the Bayes risk of $\delta$ by $r(\pi, \delta) = E_{\theta \sim \pi}\{R(g(\theta), \delta(X))\}$. An estimator $\delta_0$ is said to be a _Bayes estimator_ if it minimizes the Bayes risk, i.e. for any other estimator, we have $r(\pi, \delta_0) \leq r(\pi, \delta)$.

The conditional distribution of $(\theta|X)$ is called the _posterior distribution_.

Define the marginal distribution of $X$ as $M$ (which has the density $m$ w.r.t. $\mu$)

$$m(x) = \int_\Theta P_\theta(x) \, \pi(d\theta).$$

Example $X_1, \ldots, X_n \overset{iid}{\sim} N(\mu, \sigma^2)$, $\underset{\text{known}}{}$, $\Theta = \mathbb{R}$. Assume that $\theta \sim N(\mu, \tau^2)$ (prior dist.). $\underset{\text{hyperparameters}}{}$

$$P_\theta(\underset{\sim}{x}) = \left(\frac{1}{\sqrt{2\pi \sigma^2}}\right)^n e^{\sum_{i=1}^n \frac{(x_i - \theta)^2}{2\sigma^2}}, \quad \pi(\theta) = \frac{1}{\sqrt{2\pi\tau^2}} e^{-\frac{(\theta - \mu)^2}{2\tau^2}}$$

Joint density $= P_\theta(x)\pi(\theta) \propto e^{-\sum_{i=1}^n \frac{(x_i - \theta)^2}{2\sigma^2}} e^{-\frac{(\theta - \mu)^2}{2\tau^2}}$

posterior density $\propto e^{-\sum_{i=1}^n (x_i - \theta)^2 / 2\sigma^2} e^{-\frac{(\theta - \mu)^2}{2\tau^2}}$

$$\propto \cdots$$
$$= e^{-\frac{\theta^2}{2}\left[\frac{n}{\sigma^2} + \frac{1}{\tau^2}\right] + \theta\left[\frac{\sum x_i}{\sigma^2} + \frac{\mu}{\tau^2}\right]}$$

Now if $(\theta|X) \sim N(a, b)$, this has density proportional to

$$e^{-\frac{(\theta - a)^2}{2b^2}} \propto e^{-\frac{\theta^2}{2b^2} + \frac{\theta a}{b^2}}$$

$$\frac{1}{b^2} = \frac{1}{\sigma^2} + \frac{1}{\tau^2} \Rightarrow b^2 = \frac{\sigma^2 \tau^2}{n\tau^2 + \sigma^2}$$

$$\frac{a}{b^2} = \frac{\sum x_i}{\sigma^2} + \frac{\mu}{\tau^2} \Rightarrow a = \frac{\sum_{i=1}^n x_i / \sigma^2 + \mu/\tau^2}{n/\sigma^2 + 1/\tau^2} = \frac{\mu\sigma^2 + n\bar{x}\tau^2}{\sigma^2 + n\tau^2}$$

$\Rightarrow$ Posterior dist. is $N\left(\frac{\mu\sigma^2 + n\bar{x}\tau^2}{\sigma^2 + n\tau^2}, \frac{\sigma^2\tau^2}{\sigma^2 + n\tau^2}\right)$

**Remark:** If the prior has density $\pi$, the posterior has density $\pi(\theta|X)$ w.r.t. the same ~~dist.~~ dominated measure

$$\to \pi(\theta|X) m(X) = P_\theta(X) \pi(\theta)$$

$$\Rightarrow \pi(\theta|X) \propto P_\theta(X) \pi(\theta)$$

$$\propto g_\theta(T(X)) \pi(\theta) \text{ — posterior depends on X through sufficiency.}$$

**[Thm]** If $L(g(\theta), \delta(X)) = \{g(\theta) - \delta(X)\}^2$, then $\delta_0(X) = E\{g(\theta)|X\}$ is a Bayes estimate. with Bayes risk $E\{Var(g(\theta)|X)\}$

If $\delta(X)$ is another Bayes estimator, then $\delta_0(X) = \delta(X)$ w.p. 1.

**Pf:** Let $\delta$ be any estimator, then the risk of $\delta$ is

$$E\{\delta(X) - g(\theta))^2\} = E\{\delta(X) - \delta_0(X) + \delta_0(X) - g(\theta)\}^2$$

$$= E\{(\delta(X) - \delta_0(X))^2\} + E\{(\delta_0(X) - g(\theta))^2\} + 2E\{(\delta(X) - \delta_0(X))(\delta_0(X) - g(\theta))\}.$$

$$\geq E\{(\delta(X) - \delta_0(X))\}^2 + E\{\delta_0(X) - g(\theta)\}^2 \text{ if } \underline{E\{(\delta(X) - \delta_0(X))(\delta_0(X) - g(\theta))\} = 0} \text{ (*)}$$

$$\Rightarrow \delta_0 \text{ is a Bayes estimator. Furthermore,} \qquad \delta \text{ is a Bayes estimator}$$
$$\text{✓}$$

$$\text{iff } \delta(X) = \delta_0(X) \quad w.p. 1.$$

Finally, Bayes risk of $\delta_0 = E[g(\theta) - E\{g(\theta)|X\}]^2$

$$= E(E[g(\theta) - E\{g(\theta)|X\}]^2 | X)$$

$$= E\{Var(g(\theta)|X)\}$$

To see (*),

$$E\{(\delta(X) - \delta_0(X))(\delta_0(X) - g(\theta))\}$$

$$= E[E\{(\delta_0(X) - \delta(X))(\delta_0(X) - g(\theta))\} | X]$$

$$= E[(\delta_0(X) - \delta(X)) E\{\delta_0(X) - g(\theta)|X\}] = 0.$$

**Next:** Least favorable prior.

**Lecture 7** (Average risk optimality) Bayes. est. minimaxity, admissibility.

[THM] If $L(g(\theta), \delta(X)) = |g(\theta) - \delta(X)|^2$, then $\delta_0(X) = E(g(\theta)|X)$ is a Bayes estimate

with Bayes risk $E(\text{Var}(g(\theta)|X))$. If $\delta(X)$ is another Bayes estimator, then

$\delta_0(X) = \delta(X)$ w.p. 1. (Pf in L6)

## Remarks

(a) Here $\delta_0(X) = \delta(X)$ w.p. 1 refers to the <u>joint</u> probability when $X$ and $\theta$ are both random.

This also means that $\delta_0(X) = \delta(X)$ w.p. 1 under the marginal dist. of $X$.

(b) This does not imply $P(\delta(X) = \delta_0(X) | \theta) = 1 \quad \forall \theta$.

(c) If, however, the marginal dist. of $X$ dominates $P_\theta$, $\theta \in \Theta$, then we have $\delta_0(X)$ is

the unique Bayes estimate in the sense that $P_\theta(\delta(X) = \delta_0(X)) = 1 \quad \forall \theta$.

## Example

Suppose $X \sim$ Binomial $(n, \theta)$, $\theta \in [0,1]$. $\pi_1(\theta) = U(0,1)$, $\pi_2(0) = \pi_2(1) = \frac{1}{2}$

Case 1:                                                    Case 2

**Case 1:** $\quad P(X=x) = \int_0^1 \binom{n}{x} \theta^x (1-\theta)^{n-x} d\theta = \binom{n}{x} \int_0^1 \theta^x (1-\theta)^{n-x} d\theta = \dfrac{\binom{n}{x}}{B(x+1, n-x+1)} = \dfrac{1}{n+1}$

$\uparrow p(x|\theta)$

Marginal dominates conditional, i.e. $P(X \in A) = 0 \Rightarrow P(X \in A | \theta) = 0$.

<u>Bayes estimate</u> is unique.

$\pi(\theta|X) \propto \theta^x (1-\theta)^{n-x} = $ Beta $(x+1, n-x+1) \Rightarrow$ Bayes estimate $= \dfrac{x+1}{n+2}$

**Case 2:** $\quad P(X=x) = \frac{1}{2} P(X=x | \theta=0) + \frac{1}{2} P(X=x | \theta=1)$

$= \frac{1}{2} \{ I(x=0) + I(x=n) \}$

$\Rightarrow P(X=0) = P(X=n) = \frac{1}{2}$ and $P(X=x) = 0$ for $x = \{1, 2, \dots, n-1\}$.

$\Rightarrow$ Marginal does <u>not</u> dominate the conditional.

$\Rightarrow$ Bayes estimate is not unique. Correspondingly, the Bayes estimate is $E(\theta|X)$

$E(\theta|X=0) = P(\theta=1|X=0) = \dfrac{P(X=0|\theta=1) P(\theta=1)}{P(X=0|\theta=1) P(\theta=1) + P(X=0|\theta=0) P(\theta=0)} = \dfrac{0}{0 + \frac{1}{2}} = 0$

$E(\theta|X=n) = \cdots = 1$

Then the class of all Bayes estimators is given by $\delta_0(0) = 0$, $\delta_0(n) = 1$,

$\delta_0(X)$ : any arbitrary values for $x \in \{1, \dots, n-1\}$.

<u>Lemma</u>  A Bayes estimator (w.r.t squared error) can never be unbiased, unless $\delta(x) = g(\theta)$ w.p. 1.

Pf: Let $\delta_0(x) = E\{g(\theta)|X\}$ be the Bayes estimator. Assume that $E\{\delta_0(x)|\theta\} = g(\theta)$ [is unbiased].

We claim that $I \triangleq E\{\delta_0(x) - g(\theta)\}^2 = 0$.

$$I = E\{\delta_0(x)\}^2 + E\{g(\theta)\}^2 - 2E\{\delta_0(x) g(\theta)\}$$

where $E\{\delta_0(x)g(\theta)\} = E\{E\{\delta_0(x)g(\theta)|X\}\} = E\{\delta_0(x) E\{g(\theta)|X\}\} = E\{\delta_0(x)\}^2$

or $E\{\delta_0(x) g(\theta)\} = E\{E\{\delta_0(x)g(\theta)|\theta\}\} = E\{g(\theta)\}^2$

$\Rightarrow I = E\{g(\theta)\}^2 - E\{\delta_0(x)\}^2 = E\{\delta_0(x)\}^2 - E\{g(\theta)\}^2$

$\Rightarrow I = 0.$  //.

## Conjugate

A class of prob. distributions $F$ is said to be a conjugate family of priors for a model $\{P_\theta : \theta \in \Theta\}$ if the posterior distribution $\pi(\theta|X)$ also belongs to $F$.

<u>Example</u>  ① $X_1, ..., X_n \sim N(\theta, \delta^2)$, $\theta \sim N(\mu, \tau^2)$, $\pi(\theta|X) \sim N(\cdot, \cdot)$

② $X_1, ..., X_n \overset{iid}{\sim}$ Binomial $(1, p)$, $p \sim$ Beta $(\alpha, \beta)$.

This has density $\pi_{\alpha, \beta}(p) = \dfrac{p^{\alpha-1}(1-p)^{\beta-1}}{Beta(\alpha, \beta)}$, where Beta $(\alpha, \beta) = \int_0^1 p^{\alpha-1}(1-p)^{\beta-1}dp$.

Note Beta $(\alpha, \beta) = \dfrac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$, where $\Gamma(\alpha) = \int_0^\infty e^{-x}x^{\Gamma-1}dx$, $\Gamma(h+1) = h!$ for integer $h$.

$$f_p(\underset{\sim}{x}) = p^{\sum_{i=1}^n X_i}(1-p)^{n-\sum_{i=1}^n X_i} \Rightarrow \pi(p|\underset{\sim}{X}) \propto p^{\sum_{i=1}^n X_i + \alpha - 1}(1-p)^{n - \sum_{i=1}^n X_i + \beta - 1}$$

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

$$E(p|\underset{\sim}{X}) = \dfrac{\sum_{i=1}^n X_i + \alpha}{n + \alpha + \beta} \rightarrow \dfrac{\sum X_i}{n} = \bar{X}_n$$

× empirical Bayes

× Hierachical Bayes.

③ $X_1, ..., X_n \overset{iid}{\sim} P_o(\lambda)$, $\lambda \sim \Gamma(\alpha, \gamma)$ distribution

$$\pi_{\alpha, \gamma}(\lambda) = \dfrac{e^{-\alpha\lambda}\lambda^{\gamma-1}\alpha^\gamma}{\Gamma(\gamma)}, \ \lambda > 0.$$

④ $X_1, ..., X_n \overset{iid}{\sim} U(0, \theta)$, $\theta \sim$ Pareto $(a, c)$, $\pi_{a, c}(\theta) = \dfrac{ac^a}{\theta^{a+1}}, \ \theta > c$.

## Minimaxity   Ch5. TPE.

<u>Def</u>  The minimax risk of an estimator $\delta(x)$ for estimating $g(\theta)$ is $\sup\limits_{\theta \in \Theta} R(g(\theta), \delta)$. An estimator $\delta_0$ is said to be <u>minimax</u>, if, for any other estimator $\delta$, we have $\sup\limits_{\theta \in \Theta} R(g(\theta), \delta_0) \leq \sup\limits_{\theta \in \Theta} R(g(\theta), \delta)$.

**Def** Given a prob. dist. $\pi$ (prior) on $\Theta$, define the Bayes risk of the prior $\pi$ by $r(\pi) = r(\pi, \delta_\pi)$
where $\delta_\pi$ is Bayes estimate w.r.t. $\pi$.

**Def** A prior $\pi$ is said to be __least favorable__ if $r(\pi) \geqslant r(\pi')$ for all $\pi'$ (other prior dist. on $\Theta$)

**[THM]** Suppose $\pi$ is a distribution on $\Theta$ s.t. $r(\pi) = r(\pi, \delta_\pi) = \sup\limits_{\theta \in \Theta} R(g(\theta), \delta_\pi)$

Then  (a) $\delta_\pi$ is minimax

(b) If $\delta_\pi$ is unique Bayes w.r.t. $\pi$, then $\delta_\pi$ is unique minimax.

(c) $\pi$ is least favorable.   ⟨Pf⟩ ↲

**Corollary** A Bayes estimator with constant risk is minimax.

Pf: This means $R(g(\theta), \delta_\pi) = \alpha$ (free of $\theta$).

$\Rightarrow$  $r(\pi, \delta_\pi) = E_{\theta \sim \pi}\{R(g(\theta), \delta_\pi)\} = \alpha$

and  $\sup\limits_{\theta \in \Theta} R(g(\theta), \delta_\pi) = \alpha$ .

**THM Pf:** (a)  Let $\delta$ be arbitrary, then

$$\sup\limits_{\theta \in \Theta} R(g(\theta), \delta) \geqslant \int_\Theta R(g(\theta), \delta)\, \pi(d\theta) = r(\pi, \delta) \overset{(\ast)}{\geqslant} r(\pi, \delta_\pi) = \sup\limits_{\theta \in \Theta} R(g(\theta), \delta_\pi)$$

bayes estimate. ↑ (pointing to $\delta_\pi$)

(b) Let $\delta \neq \delta_\pi$, i.e. $\exists \theta$ st $P_\theta(\delta(x) \neq \delta_\pi(x)) > 0$

$\Rightarrow$ ($\ast$) is a strictly inequality as $\delta_\pi$ is unique Bayes.

(c) Let $\pi'$ be any distribution. N.T.S. $r(\pi') \leqslant r(\pi)$

But observe that  $r(\pi') = \int_\theta R(g(\theta), \delta_{\pi'})\, \pi'(d\theta) \leqslant \int_\theta R(g(\theta), \delta_\pi)\, \pi'(d\theta) \leqslant \sup\limits_{\theta \in \Theta} R(g(\theta), \delta_\pi) = r(\pi)$

$\Rightarrow \pi$ is least favorable.  //.

**Example** Let $X_1, \ldots, X_n \overset{iid}{\sim} B(1, p)$. Find a minimax estimator for $p$.

Let the prior on $p$ be Beta$(\alpha, \beta)$, i.e. $\pi(p) \propto p^{\alpha-1}(1-p)^{\beta-1}$

Then the Bayes estimator is $\delta_\pi(X) = \dfrac{\sum\limits_{i=1}^{n} X_i + \alpha}{n + \alpha + \beta}$ .

$R(p, \delta_\pi) = E\left(\dfrac{\sum\limits_{i=1}^{n} X_i + \alpha}{n + \alpha + \beta} - p\right)^2 = \dfrac{np - np^2 + \alpha^2 - 2\alpha(\alpha+\beta)p + (\alpha+\beta)^2 p^2}{(n+\alpha+\beta)^2}$

To make this free of $p$, $\begin{cases} n = 2\alpha(\alpha+\beta) \\ n = (\alpha+\beta)^2 \end{cases} \Rightarrow \begin{cases} \alpha+\beta = \sqrt{n} \\ 2\alpha\sqrt{n} = n \end{cases} \Rightarrow \begin{cases} \alpha = \sqrt{n}/2 \\ \beta = \sqrt{n}/2 \end{cases}$

$\delta_\pi(X) = \dfrac{\sum\limits_{i=1}^{n} X_i + \sqrt{n}/2}{n + \sqrt{n}}$  is the unique minimax estimator (marginal dominates conditional) .

**Def** A sequence of priors $\{\pi_n\}_{n\geq 1}$ is least favorable if $\lim_{n\to\infty} r(\pi_n) = \sup_\pi r(\pi)$

**[THM]** Suppose $\{\pi_n\}_{n\geq 1}$ is a sequence of priors such that $\lim_{n\to\infty} r(\pi_n) = \sup_{\theta\in\Theta} R(g(\theta),\delta_0)$, then

(a) $\delta_0$ is minimax

(b) $\{\pi_n\}$ is least favorable.   (Pf ↓)

**Example**   $X_1,\ldots,X_n \overset{iid}{\sim} N(\theta, \overset{known}{\delta^2})$. Find a minimax estimator for $\theta$ with the squared loss function.

(motivate the above THM)   Claim: $\bar{X}_n$ is minimax.

Let $\pi_{\mu,\tau^2}(\theta) = N(\mu,\tau^2)$. The Bayes estimator is $\delta_\pi = \dfrac{\frac{n\bar{X}_n}{\delta^2} + \frac{\mu}{\tau^2}}{\frac{n}{\delta^2} + \frac{1}{\tau^2}}$

Bayes risk: $r(\pi) = r(\pi, \delta_\pi) = \dfrac{1}{\frac{n}{\delta^2} + \frac{1}{\tau^2}}$.

Here $\delta_0 = \bar{X}_n$, $R(\theta, \bar{X}_n) = E(\bar{X}-\theta)^2 = \dfrac{\delta^2}{n} \Rightarrow \sup_{\theta\in\Theta} R(\theta, \bar{X}_n) = \dfrac{\delta^2}{n}$

Also, $\lim_{\tau\to\infty} r(\pi_\tau) = \dfrac{\delta^2}{n} = \sup_{\theta\in\Theta} R(\theta, \bar{X}_n)$

$\Rightarrow \bar{X}_n$ is minimax. Also, $\{\pi_\tau\}_{\tau\in N}$ is a least favorable distribution.

**[THM] 2**

**Pf:** (a) Let $\delta$ be any other estimator. Then

$$\sup_{\theta\in\Theta} R(g(\theta),\delta) \geq \int_\Theta R(g(\theta),\delta)\, \pi_n(d\theta) = r(\pi_n,\delta) \geq r(\pi_n)$$

Take limit to get $\sup_{\theta\in\Theta} R(g(\theta),\delta) = \lim_{n\to\infty} r(\pi_n)$

(b) N.T.S  $\sup_\pi r(\pi) = \lim_{n\to\infty} r(\pi_n)$.

Observe that $\sup_\pi r(\pi) \geq r(\pi_n) \Rightarrow \sup_\pi r(\pi) \geq \lim_{n\to\infty} r(\pi_n)$.

For any $\pi$, $r(\pi) = \inf_\delta r(\pi,\delta) \leq r(\pi,\delta_0) \leq \sup_{\theta\in\Theta} R(g(\theta),\delta_0) = \lim_{n\to\infty} r(\pi_n)$

Hence, $\sup_\pi r(\pi) \leq \lim_{n\to\infty} r(\pi_n)$.   //.

**[Lemma]** Suppose $\delta(X)$ is minimax for $g(\theta)$ on the parameter set $\theta\in\Theta_0$, where $\Theta_0 \subseteq \Theta$.
If $\sup_{\theta\in\Theta_0} R(g(\theta),\delta) = \sup_{\theta\in\Theta} R(g(\theta),\delta)$, then $\delta$ is minimax for $\theta\in\Theta$.

Pf  See TPE.

**Example** $X_1,\ldots,X_n \overset{iid}{\sim} N(\mu,\delta^2)$, $\mu\in\mathbb{R}$, $\delta^2 > 0$   (both unknown)

$\theta = (\mu,\delta^2) \in \mathbb{R}\times(0,\infty)$

For any estimator $\delta$, $\sup_{\theta\in\Theta} R(\mu,\delta) \geq \sup_{\theta\in\Theta; \delta=\delta_0} R(\mu,\delta) \geq \dfrac{\delta_0^2}{n}$

$\Rightarrow \sup_{\theta\in\Theta} R(\mu,\delta) \geq \sup_{\delta_0 > 0} \dfrac{\delta_0^2}{n} = +\infty$

cf.

Example : Assume $\mu \in \mathbb{R}$, $0 < 6 \leq M$, $\Theta = \mathbb{R} \times [0, M]$.

In this case $\bar{X}$ is again minimax. This is because

Let $\Theta_0 = \mathbb{R} \times \{M\}$. In this case, we know that $\bar{X}_n$ is minimax and $\sup_{\theta \in \Theta_0} R(\mu, \bar{X}_n) = \frac{M^2}{n}$

Also, $R(\mu, \bar{X}_n) = \frac{6^2}{n} \Rightarrow \sup_{\theta \in \Theta} R(\mu, \bar{X}) = \sup_{6 \in [0,M]} \frac{6^2}{n} = \frac{M^2}{n} = \sup_{\theta \in \Theta} R(\mu, \bar{X}_n)$

$\Rightarrow \bar{X}_n$ is minimax on $\Theta$.

## Admissibility

An estimator $\delta$ is said to be inadmissible if $\exists \delta'$ s.t. $R(g(\theta), \delta') \leq R(g(\theta), \delta)$ with strict inequality for some $\theta \in \Theta$.

An estimator $\delta$ is admissible if there is no such $\delta'$.

Remark  If loss function is strictly convex, any estimator which is not a function of the minimal sufficient statistic is inadmissible (Rao-Blackwell).

Lemma  Any unique Bayes estimator is admissible.

TPE 5.2.   Suppose $\delta$ is a unique Bayes estimator, which is not admissible.

$\Rightarrow \exists \delta'$ better than $\delta \Rightarrow \delta'$ is Bayes.    Contradiction.

Lemma  An admissible estimator with constant risk is minimax.

Lemma : If $\delta$ is unique minimax, then $\delta$ is admissible.

## Lecture 8

### Asymptotic Optimality ~ M-estimator        $\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N(0, \Sigma)$      C&B Ch 7.
K. Ch 9.

Let $\{X_1, X_2, ..., X_n\}$ be iid from $\{P_\theta : \theta \in \Theta\}$ with pdf $p_\theta(\cdot)$ w.r.t. some $6$-finite measure. Suppose we want to estimate $g(\theta)$ and a candidate estimator is $\delta_n(X_1, ..., X_n)$.

Def We say $\delta_n(\underline{X})$ is consistent for $g(\theta)$ if
$$\delta_n(\underline{X}) \xrightarrow[P_\theta]{P} g(\theta) \quad \forall \theta \in \Theta, \text{ i.e.}$$

$$\forall \theta \in \Theta, \forall \varepsilon > 0, \quad P_\theta\left(|\delta_n(\underline{X}) - g(\theta)| > \varepsilon\right) \to 0 \quad \text{as} \quad n \to \infty.$$

Example. $X_1, ..., X_n \overset{iid}{\sim} \text{Bin}(1, \theta)$, UMVUE for $g(\theta) = \theta$ is $\bar{X}_n$.

$\bar{X}_n \xrightarrow[P_\theta]{P} \theta$ by WLLN $\Rightarrow \bar{X}_n$ is consistent for $\theta$.

Remarks  For $X_1, ..., X_n \overset{iid}{\sim} F$

a) Assume $E_F |X_1| < \infty$, then $\frac{1}{n} \sum_{i=1}^{n} X_i \xrightarrow{P} E_F X_1$   (WLLN)

b) Assume $E_F X_1^2 < \infty$, then $W_n \triangleq \frac{\sum_{i=1}^{n} X_i - n E_F X_1}{\sqrt{n \, \text{Var}_F(X_1)}} \xrightarrow{d} N(0,1)$   (CLT)

i.e. $\lim_{n \to \infty} P(W_n \leq t) \to \Phi(t) \quad \forall t \in \mathbb{R}.$

**Def** Let $L(\theta|X_1,\ldots,X_n) = \prod_{i=1}^{n} P_\theta(X_i)$ be the likelihood function, and $l(\theta|X_1,\ldots,X_n) = \log L(\theta|X_1,\ldots,X_n)$ be the log-likelihood function. If there exists a unique $\hat\theta_n$, which is a global maximizer of $\theta \mapsto L(\theta|\underline{X})$ or $\theta \mapsto l(\theta|\underline{X})$, then define $\hat\theta_n$ as the MLE for $\theta$.

**Example** Suppose $X_1,\ldots,X_n \overset{iid}{\sim} \text{Bin}(1,\theta)$. $P_\theta(x) = \theta^x(1-\theta)^{1-x}$, $\theta \in (0,1)$.

$$L_n(\theta|\underline{X}) = \prod_{i=1}^{n} P_\theta(X_i) = \theta^{\sum_{i=1}^{n}X_i}(1-\theta)^{n-\sum_{i=1}^{n}X_i}$$

$$l_n(\theta|\underline{X}) = \sum_{i=1}^{n} X_i \log\theta + (n-\sum_{i=1}^{n}X_i)\log(1-\theta)$$

$$l_n'(\theta|\underline{X}) = \frac{\sum_{i=1}^{n}X_i}{\theta} - \frac{n-\sum_{i=1}^{n}X_i}{1-\theta}$$

regularity conditions

$$l_n''(\theta|\underline{X}) = -\frac{\sum_{i=1}^{n}X_i}{\theta^2} - \frac{n-\sum_{i=1}^{n}X_i}{(1-\theta)^2} < 0 \Rightarrow l_n(\cdot|\underline{X}) \text{ is strictly concave.}$$

Observe that $l_n'(\theta|\underline{X})\big|_{\theta=\hat\theta_n} = 0 \Rightarrow \frac{\sum_{i=1}^{n}X_i}{\theta} = \frac{n-\sum_{i=1}^{n}X_i}{1-\theta}$

$$\Rightarrow \hat\theta_n = \frac{\sum_{i=1}^{n}X_i}{n} = \bar{X}_n$$

$\Rightarrow$ MLE exists and equals $\bar{X}_n$.

Also, $\bar{X}_n \xrightarrow[P_\theta]{P} \theta \quad \forall \theta \in (0,1) \quad$ [CONSISTENCY]

and $\sqrt{n}(\bar{X}_n - \theta) \xrightarrow[P_\theta]{d} N(0, \frac{1}{\theta(1-\theta)})$ [Asy. normality via CLT]

slow down the motion

[Thm] Suppose $X_1,\ldots,X_n$ are iid from $P_\theta$ for some $\theta \in \Theta$, with pdf $p_\theta(\cdot)$.

A0. $P_{\theta_1} \neq P_{\theta_2}$ whenever $\theta_1 \neq \theta_2$ (identifiability)

A1. $\{P_\theta, \theta \in \Theta\}$ have common support.

Then, $P_{\theta_0}(l_n(\theta_0|\underline{X}) > l_n(\theta|\underline{X})) \xrightarrow{n\to\infty} 1 \quad \forall \theta \neq \theta_0$

Pf: Let $T_n = \frac{1}{n}\sum_{i=1}^{n}\log\frac{P_\theta(X_i)}{P_{\theta_0}(X_i)}$, then $T_n \xrightarrow[P_{\theta_0}]{P} E_{\theta_0}\{\log\frac{P_\theta(X_1)}{P_{\theta_0}(X_1)}\}$

Now $E_{\theta_0}\{\log\frac{P_\theta(X_1)}{P_{\theta_0}(X_1)}\} = \int \log\{\frac{P_\theta(X)}{P_{\theta_0}(X)}\}P_{\theta_0}(x)d\mu = -\underline{D(\theta_0\|\theta)} < 0$ for $\theta \neq \theta_0$

entropy

$\Rightarrow P_{\theta_0}(T_n < 0) \xrightarrow{n\to\infty} 1 \quad$ But $T_n < 0 \iff \frac{1}{n}\sum_{i=1}^{n}\log\frac{P_\theta(X_i)}{P_{\theta_0}(X_i)} < 0$

$\iff \log\prod_{i=1}^{n}P_\theta(X_i) < \log\prod_{i=1}^{n}P_{\theta_0}(X_i)$

$\iff l_n(\theta|\underline{X}) < l_n(\theta_0|\underline{X})$

[Corollary] Suppose (A0) and (A1) hold. If $\Theta$ is finite, then the MLE $\hat\theta_n$ exists with high prob. (prob $\to 1$) and $P_{\theta_0}(\hat\theta_n = \theta_0) \xrightarrow{n\to\infty} 1$.

Pf: Let $\Theta = \{\theta_0, \theta_1,\ldots,\theta_k\} \Rightarrow P_{\theta_0}(l_n(\theta_0|\underline{X}) > l_n(\theta_j|\underline{X})) \xrightarrow{n\to\infty} 1, \quad 1 \leq j \leq k.$

$\Rightarrow P_{\theta_0}(l_n(\theta_0|\underline{X}) > \max_{1\leq j\leq k} l_n(\theta_j|\underline{X})) \xrightarrow{n\to\infty} 1$

Let $A_n = \{ X : l_n(\theta_n | X) > \max_{1 \leq j \leq k} l_n(\theta_j | X) \}$

If $X \in A_n$, then $\hat{\theta}_n(X) = \theta_0$ and $P_{\theta_0}(A_n) \to 1$.


## MLE expansion $\quad \sqrt{n}(\hat{\theta}_n - \theta_0) \to 0$

$$0 = l_n'(\hat{\theta}_n) = l_n'(\theta_0) + \boxed{(\hat{\theta}_n - \theta_0)} \, l_n''(\theta_0) + \tfrac{1}{2}(\hat{\theta}_n - \theta_0)^2 l_n'''(\xi_n) \quad \xi_n \in [\theta_0, \hat{\theta}_n]$$

（interest）

$$\Rightarrow (\hat{\theta}_n - \theta_0)\left\{ l_n''(\theta_0) + \tfrac{1}{2}(\hat{\theta}_n - \theta_0)^2 l_n'''(\xi_n) \right\} = -l_n'(\theta_0)$$

$$\Rightarrow \sqrt{n}(\hat{\theta}_n - \theta_0) = \frac{-l_n'(\theta_0)/\sqrt{n}}{-l_n''(\theta_0)/n - \tfrac{1}{2}(\hat{\theta}_n - \theta_0)\, l'''(\xi_n)/n}.$$

It suffices to show

$$\frac{1}{\sqrt{n}} l_n'(\theta_0) \xrightarrow{d} N(0, I(\theta_0))$$

$$-\frac{1}{n} l_n''(\theta_n) \xrightarrow{P} I(\theta_0)$$

and $\frac{1}{n}(\hat{\theta}_n - \theta_0)\, l'''(\xi_n) \xrightarrow{P} 0$,

then $\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} \dfrac{N(0, I(\theta_0))}{I(\theta_0) + o} \overset{d}{=} N(0, I(\theta_0)^{-1})$.

Observe that $\quad u(\hat{\theta}_n) = 0 \Rightarrow u'(\theta_0) + u''(\theta_0)(\hat{\theta}_n - \theta_0) + \cdots$

$$\frac{1}{\sqrt{n}} l_n'(\theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \frac{\partial}{\partial \theta} \log p_\theta(X_i) \Big|_{\theta = \theta_0} \xrightarrow[\theta_0]{d} N(0, I(\theta_0)) \quad \text{by CLT}.$$

$$\frac{1}{n} l_n''(\theta_0) = \frac{1}{n} \sum_{i=1}^{n} \frac{\partial^2}{\partial \theta^2} \log p_\theta(X_i) \Big|_{\theta = \theta_0} \xrightarrow[\theta_0]{P} -I(\theta_0) \quad \text{by WLLN}.$$

$$\left| \frac{1}{n}(\hat{\theta}_n - \theta_0) l_n''(\xi_n) \right| \leq |\hat{\theta}_n - \theta_0| \cdot \frac{1}{n} \sum_{i=1}^{n} M(X_i) \xrightarrow[\theta_0]{P} 0 \cdot E_{\theta_0} M(X_i) = 0 \quad \text{by consistency and WLLN}.$$

Van der Vaart (98) Ch.5.   M. Z estimator

estimating equations · semi-parametric · Cox model  non-para

$E\{N(t) = Y(t) d\Lambda(t)\} = 0$

proportional $\lambda(t|Z) = \lambda_0(t) \exp\{\beta^T Z\}$  para-

↓ 0 placebo  ↓ 1 drug

poisson    mean.   hazards model    baseline    covariate
counting process    hazard    effect
(unspecified)

## Main references: TSH (3rd)

## Setup   Let $\{ P_\theta : \theta \in \Theta \}$ be a collection of prob. measure

on $X$, dominated by a $\sigma$-finite measure $\mu$. Let

$p_\theta(\cdot) = \dfrac{d P_\theta}{d\mu}$. Let $\Theta_0$ and $\Theta_1$ be two disjoint subsets

of $\Theta$. Given $X \sim P_\theta$ for some $\theta \in \Theta$, we want to

decide whether $\theta \in \Theta_0$ or $\theta \in \Theta_1$.

$\dfrac{f(t|z)}{S(t|z)} = \dfrac{f(t|z)}{1 - F(t|z)}$   $F(t) = P(T \leq t)$

the semi-parametric model

$\sqrt{n}(\hat{\beta} - \beta_0) \to ?$

∴ partial likelihood → profile likelih. (Murphy & vdV, 2000 JRSSB)

(Cox, 72, 75)

JRSSB  Biometrika

E.g. $X \in \mathbb{R}^n$, $X = (X_1, \ldots, X_n) \overset{iid}{\sim} N(\theta, 1)$

$\Theta \in \mathbb{R}$, $\Theta_0 = \{\theta\}$, $\Theta_1 = \{1\}$.

Top1 cited paper
Kaplan-Meier 337.

Simple vs simple / composite

Theorectical Statistics  $\Theta_1 = \{\theta : \theta > 1\}$.

Cox & Hinkley (74)
↓ strongly recommended by Ying Zhiliang.

(22)

**Def** A function $\phi : X \to \{0,1\}$ is called a non-randomized test function.

Types of error



| decision | $\theta \in \Theta_1$ | $\theta \in \Theta_0$ |
|---|---|---|
| $\phi = 1$ | $\vee$ | Type I |
| $\phi = 0$ | Type II | $\vee$ |

Type I → need to be controlled, more important

$$P_\theta(\phi = 1), \ \theta \in \Theta_0 \quad \text{Type I}$$

$$P_\theta(\phi = 0), \ \theta \in \Theta_1 \quad \text{Type II}$$

$*$ Power function of $\phi$ : $1 - $ prob. of type II error $= P_\theta(\phi = 1), \ \theta \in \Theta_1$  prob. of correctly reject $H_0$

$*$ size of a test $\phi$ : $\sup\limits_{\theta \in \Theta_0} P_\theta(\phi = 1)$   Type I

Let $\alpha \in (0,1)$, a test $\phi$ is called <u>level $\alpha$</u> if $\sup\limits_{\theta \in \Theta_0} P_\theta(\phi = 1) \leq \alpha$ .

**Def** A test $\phi$ is called <u>uniformly most powerful</u> level $\alpha$ test if given any other $\overset{\text{level}}{\vee} \alpha$ test $\psi$, we have $P_\theta(\phi = 1) \geq P_\theta(\psi = 1) \ \forall \ \theta \in \Theta_1$. (UMP)

**Def** A function $\phi : X \to [0,1]$ is called a randomized test function or just a test function. If $\phi(x) = p$, toss a coin with prob. of heads $p$. If heads choose $\Theta_1$, if tails choose $\Theta_0$. In all previous definitions, replace $P_\theta(\phi = 1)$ by $E_\theta \phi$.

[Thm] (Neyman - Pearson)

Suppose we want to test $H_0 : \theta = \theta_0$ versus $H_1 : \theta = \theta_1$ at level $\alpha$

(a) There exists a test $\phi$ satisfying

(i) $E_{\theta_0} \phi = \alpha$

(ii) $\exists \ k \in [0, \infty)$ s.t. $\phi = \begin{cases} 1 & \text{if } \frac{P_{\theta_1}(x)}{P_{\theta_0}(x)} > k \\ 0 & \text{o/w.} \end{cases}$

(b) If a test satisfies (i) and (ii) above, then $\phi$ is a Most Powerful test for testing $\theta = \theta_0$ vs $\theta = \theta_1$ at level $\alpha$.

(c) If $\phi$ is a Most Powerful at level $\alpha$, it must satisfy (ii), for the same $k$ as in (a). It also satisfies (i) unless $E_{\theta_1}(\phi) = 1$ (power $= 1$).

Pf : (a) If $\alpha = 0$, take $k = \infty$, $\phi = 0$. If $\alpha = 1$, take $k = 0$, $\phi = 1$.

For $\alpha \in (0,1)$, let $\alpha(c) = P_{\theta_0}(P_{\theta_1}(x) > c P_{\theta_0}(x))$, $c > 0$

$\qquad = P_{\theta_0}\left(\frac{P_{\theta_1}(x)}{P_{\theta_0}(x)} > c\right) = 1 - P_{\theta_0}\left(\frac{P_{\theta_1}(x)}{P_{\theta_0}(x)} \leq c\right)$

$\Rightarrow \alpha(\cdot)$ is a non-decreasing and right-continuous.

Also, $\alpha(c-) - \alpha(c) = P_{\theta_0}\left(\frac{P_{\theta_1}(x)}{P_{\theta_0}(x)} = c\right)$, $\alpha(\infty) = 0$, $\alpha(0-) = 1$.

$\Rightarrow \exists C_0$ s.t. $\alpha(C_0) \leq \alpha \leq \alpha(C_0-)$.

Let
$$\phi = \begin{cases} 1 & P_{\theta_1}(x) > C_0 P_{\theta_0}(x) \\ \dfrac{\alpha - \alpha(C_0)}{\alpha(C_0-) - \alpha(C_0)} & \text{if} \quad P_{\theta_1}(x) = C_0 P_{\theta_0}(x) \\ 0 & P_{\theta_1}(x) < C_0 P_{\theta_0}(x) \end{cases}$$

with if $\alpha(C_0-) = \alpha(C_0)$, set $\phi = 1$ on $P_{\theta_1}(x) = C_0 P_{\theta_0}(x)$

$\Rightarrow E_{\theta_0}\phi = P_{\theta_0}(P_{\theta_1}(x) > C_0 P_{\theta_0}(x)) + P_{\theta_0}(P_{\theta_1}(x) = C_0 P_{\theta_0}(x)) \cdot \dfrac{\alpha - \alpha(C_0)}{\alpha(C_0-) - \alpha(C_0)}$

$\quad\quad = \alpha(C_0) + [\alpha(C_0-) - \alpha(C_0)] \times \dfrac{\alpha - \alpha(C_0)}{\alpha(C_0-) - \alpha(C_0)}$

$\quad\quad = \alpha.$

(b) Let $\phi$ be of the MP form, i.e. $\exists k$ s.t. $E_{\theta_0}\phi = \alpha$, $\phi(x) = \begin{cases} 1 & \text{if } \frac{P_{\theta_1}(x)}{P_{\theta_0}(x)} > k \\ 0 & \text{o/w.} \end{cases}$

Let $\phi^*(x)$ be the test s.t. $E_{\theta_0}\{\phi^*(x)\} \leq \alpha$,

we need to show that $E_{\theta_1}\{\phi(x)\} - E_{\theta_1}\{\phi^*(x)\} \geq 0$.

Consider the integral:

$\int \{\phi(x) - \phi^*(x)\} \{P_{\theta_1}(x) - k P_{\theta_0}(x)\} d\mu$

$= \int_{\phi > \phi^*} \{\phi(x) - \phi^*(x)\} \{P_{\theta_1}(x) - k P_{\theta_0}(x)\} d\mu$

$\quad + \int_{\phi < \phi^*} \{\phi(x) - \phi^*(x)\} \{P_{\theta_1}(x) - k P_{\theta_0}(x)\} d\mu$

Observe if $\phi > \phi^*$, $\phi > 0 \Rightarrow P_{\theta_1}(x) \geq k P_{\theta_0}(x)$

$\quad\quad$ if $\phi < \phi^*$, $\phi < 1 \Rightarrow P_{\theta_1}(x) \leq k P_{\theta_0}(x)$

therefore, $0 \leq \int \{\phi(x) - \phi^*(x)\} \{P_{\theta_1}(x) - k P_{\theta_0}(x)\} d\mu$

$\quad\quad = E_{\theta_1}\phi(x) - E_{\theta_1}\phi^*(x) - k\{E_{\theta_0}\phi(x) - E_{\theta_0}\phi^*(x)\}$

$\Rightarrow E_{\theta_1}\phi(x) - E_{\theta_1}\phi^*(x) \geq k\{E_{\theta_0}\phi(x) - E_{\theta_0}\phi^*(x)\} \geq k(\alpha - \alpha) = 0$.

(c) Let $\phi^*$ be an MP test. Let $\phi$ be the test from (a). We have $E_{\theta_1}\phi(x) = E_{\theta_1}\phi^*(x) = \alpha$

$\Rightarrow \int (\phi(x) - \phi^*(x)) \{P_{\theta_1}(x) - k P_{\theta_0}(x)\} d\mu = 0$

$\Rightarrow \int_{\phi > \phi^*} \{\phi(x) - \phi^*(x)\} \{P_{\theta_1}(x) - k P_{\theta_0}(x)\} d\mu = \int_{\phi < \phi^*} \{\phi(x) - \phi^*(x)\} \{P_{\theta_1}(x) - k P_{\theta_0}(x)\} d\mu$

$\Rightarrow \phi = \phi^* \ \forall \ x$ s.t. $P_{\theta_1} \neq k P_{\theta_0}(\lambda)$, where $k$ is defined as in (a).

Also, we must have $E_{\theta_0}\{\phi^*(x)\} = \alpha$ unless $E_{\theta_1}\phi^*(x) = 1$

because $E_{\theta_0}(\phi^*(x)) = E_{\theta_0}(\phi(x)) = \alpha$ unless $k=0$. $Bn + k = 0 \iff E_{\theta_1}\phi^*(x) = 1$.

If $E_{\theta_0}(\phi^*(x)) < \alpha$, $E_{\theta_1}(\phi^*(x)) < 1$, then $\phi^*$ is not MP.

Remark: If $\{X: P_{\theta_1}(x) = k\, P_{\theta_0}(x)\}$ is of measure $0$, MP is unique.

E.g. $X_1, \ldots, X_n \overset{iid}{\sim} N(\theta, 1)$

test $H_0: \theta = 0$ vs $H_1: \theta = 1$ at level $\alpha$.

$$\frac{P_{\theta=1}(X_1, \ldots, X_n)}{P_{\theta=0}(X_1, \ldots, X_n)} = \frac{\left(\frac{1}{\sqrt{2\pi}}\right)^n e^{-\frac{1}{2}\sum_{i=1}^n (X_i - 1)^2}}{\left(\frac{1}{\sqrt{2\pi}}\right) e^{-\frac{1}{2}\sum_{i=1}^n (X_i^2)}} = e^{\sum_{i=1}^n X_i - \frac{n}{2}}$$

$\Rightarrow \phi = 1$ if $\frac{P_{\theta_1}(x)}{P_{\theta_0}(x)} > k$ $\iff \sum_{i=1}^n X_i - \frac{n}{2} > \log k \iff \sum_{i=1}^n X_i > k' = \log k + \frac{n}{2}$

$\Rightarrow \phi(x) = \begin{cases} 1 \\ 0 \end{cases}$ if $\begin{array}{l} \sum X_i > k' \\ \sum X_i < k' \end{array}$

where $\alpha = E_{\theta=0}\,\phi(x) = P_{\theta_0}\left(\sum_{i=1}^n X_i > k'\right) \Rightarrow k' = \sqrt{n}\, 3_{1-\alpha}$. $P(Z \leq 3_{1-\alpha}) = 1-\alpha$.

<u>Lecture 9</u>  UMP, MLR, least favorable dist, ...  TSH. Ch3  K. Ch12.  C&B Ch8.

\* Neyman - Pearson (Simple vs Simple)

Recap. the above example ↗

$\phi(x) = \begin{cases} 1 & > \\ \gamma & = \\ 0 & < \end{cases}$ \* (randomized test)

E.g. $X_1, X_2 \overset{iid}{\sim} Bernoulli(\theta)$

Test $H_0: \theta = \frac{1}{2}$ versus $H_1: \theta = \frac{2}{3}$ at level $\alpha = \frac{1}{2}$

| Sample | (0,0) | (0,1) | (1,0) | (1,1) |
|---|---|---|---|---|
| $P_{\theta_0}(X_1, X_2)$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{4}$ |
| $P_{\theta_1}(X_1, X_2)$ | $\frac{1}{9}$ | $\frac{2}{9}$ | $\frac{2}{9}$ | $\frac{4}{9}$ |
| $\frac{P_{\theta_1}}{P_{\theta_0}}$ | $\frac{4}{9}$ | $\frac{8}{9}$ | $\frac{8}{9}$ | $\frac{16}{9}$ |

$\Rightarrow k = \frac{8}{9}$

not reject $H_0$          reject $H_0$

Let $\phi = (X_1, X_2) = \begin{cases} 1 & (X_1, X_2) = (1,1) \\ 0 & (X_1, X_2) = (0,0) \\ \text{randomized} \end{cases}$

s.t. $E_{\theta_0}\phi(X_1, X_2) = \frac{1}{2}$

Test procedure: $\begin{array}{l} \phi(1,0) = 1, \quad \phi(0,1) = 0 \\ \phi(1,0) = 0, \quad \phi(0,1) = 1 \end{array} \Big\}$ non-randomized

$\phi(1,0) = \phi(0,1) = \frac{1}{2}$ randomized characterized by MP.

**[Corollary]** Let $\beta = \beta(\theta_1)$ denote the power of the MP test for testing $H_0: \theta = \theta_0$ vs $H_1: \theta = \theta_1$ at level $\alpha \in (0,1)$. Then $\beta \geq \alpha$. Further more, $\beta > \alpha$ unless $P_{\theta_1} = P_{\theta_0}$.

**Pf:** Let $\phi$ be the MP test from part (a) of NP lemma.

Let $\psi(\underset{\sim}{x}) \equiv \alpha \Rightarrow \beta = E_{\theta_1}\phi(\underset{\sim}{x}) \geq E_{\theta_1}\psi(\underset{\sim}{x}) = \alpha$.

Suppose $\beta = \alpha$, then $\psi$ is a MP test.

$\Rightarrow P_{\theta_1}(\underset{\sim}{x}) = k P_{\theta_0}(\underset{\sim}{x})$ a.s. $\mu$.

$\Rightarrow k = 1 \Rightarrow P_{\theta_1} = P_{\theta_0}$ //.

**E.g.** Suppose $X_1, \ldots, X_n \overset{iid}{\sim} N(\theta, 1)$ $\quad H_1: \theta \neq 0$ (no UMP test).

$\quad H_1: \theta < 0$

Test $\underline{H_0: \theta = 0}$ vs $\underline{H_1: \theta > 0}$ at level $\alpha$.

$\qquad$ simple $\qquad$ Composite

$\theta_1 < 0$

Fix $\theta_1 > 0$ $(\theta_1 \in \Theta_1)$

$\quad$ Test: $H_0: \theta = 0$ vs $H_1': \theta = \theta_1$ at level $\alpha$.

MP test for this problem is

$$\phi(\underset{\sim}{x}) = \begin{cases} 1 \\ 0 \end{cases} \quad \text{if} \quad \begin{array}{c} \sum X_i \overset{<}{>} \sqrt{n}\, z_{1-\alpha} \\ o/w. \end{array}$$

$\Rightarrow \phi$ is uniformly MP for testing $H_0: \theta = 0$ vs $H_1: \theta \overset{<}{>} 0$.

---

**Monotone Likelihood Ratio (MLR)**

**Def** Suppose $\Theta$ is an interval. We say that $\{P_\theta : \theta \in \Theta\}$ has the monotone likelihood ratio (MLR) property in a statistic $T(x)$ if $\forall\, \theta_1 < \theta_2 \in \Theta$, $\dfrac{P_{\theta_2}(\underset{\sim}{x})}{P_{\theta_1}(\underset{\sim}{x})}$ is a non-decreasing function of $T(x)$.

**E.g.** $P_\theta(x) = e^{\eta(\theta) T(x) - B(\theta)} h(x)$, $\theta \in (a,b)$, $\eta$ non-decreasing $\quad **$

$$\frac{P_{\theta_2}(x)}{P_{\theta_1}(x)} = e^{\{\eta(\theta_2) - \eta(\theta_1)\} T(x)} e^{-B(\theta_2) + B(\theta_1)} \triangleq g(T(x)),$$

where $g(t) = e^{\{\eta(\theta_2) - \eta(\theta_1)\} t} e^{-B(\theta_2) + B(\theta_1)}$

**E.g.** $X_1, \ldots, X_n \overset{iid}{\sim} \text{uniform}(0,\theta)$. $P_\theta(\underset{\sim}{x}) = \dfrac{1}{\theta^n} I(X_{(n)} < \theta)$

Let $T = X_{(n)}$, $\theta_1 < \theta_2$. If $0 < T < \theta_1$, $\dfrac{P_{\theta_2}(\underset{\sim}{x})}{P_{\theta_1}(\underset{\sim}{x})} = \left(\dfrac{\theta_1}{\theta_2}\right)^n$

$\quad$ if $\theta_1 \leq T < \theta_2$, $\dfrac{P_{\theta_2}(\underset{\sim}{x})}{P_{\theta_1}(\underset{\sim}{x})} = \infty$

$\quad$ if $\theta_2 \leq T$, $\dfrac{P_{\theta_2}(\underset{\sim}{x})}{P_{\theta_1}(\underset{\sim}{x})} = \dfrac{0}{0}$ (set to be $\infty$).

**Thm 3.4.1**

**[Thm]** Let $\{P_\theta(\cdot), \theta \in \Theta\}$ be MLR in $T(x)$ s.t. $P_{\theta_1} \neq P_{\theta_2}$ if $\theta_1 \neq \theta_2$ and $\Theta$ is an interval.

(a) For testing $H_0: \theta \leq \theta_0$ vs $H_1: \theta > \theta_0$ at level $\alpha \in (0,1)$, there exists a UMP test $\phi$ of the form

$$\phi(x) = \begin{cases} 1 & T(\underset{\sim}{x}) > c \\ \gamma & \text{if} \quad T(\underset{\sim}{x}) = c \quad \text{and} \quad E_{\theta_0}\phi(x) = \alpha. \\ 0 & T(\underset{\sim}{x}) < c \end{cases}$$

(b) The power function $\beta(\theta) = E_\theta \phi$ is strictly increasing on the set $\{\theta: 0 < \beta(\theta) < 1\}$,

   i.e. if $\theta_1 < \theta_2 \in \Theta$ s.t. $\beta(\theta_1), \beta(\theta_2) \in (0,1)$, then $\beta(\theta_1) < \beta(\theta_2)$.

(c) For all $\theta' \in \Theta$, the test of part (a) is UMP for testing $H_0: \theta \le \theta'$ vs $H_1: \theta > \theta'$

   at level $\alpha' = \beta(\theta')$

(d) For any $\theta < \theta_0$, $\phi$ minimizes $\beta(\theta)$ amongst all tests satisfying $E_{\theta_0} \psi(X) = \alpha$.

Pf: Let $f(c) = P_{\theta_0}(T(X) > c)$, $f(\infty) = 0$, $f(-\infty) = 1$,

$\exists \ c_0 \in [-\infty, \infty]$ s.t. $f(c_0 -) \ge \alpha \ge f(c_0)$

Let $\phi(X) = \begin{cases} 1 & T(X) > c_0 \\ \dfrac{\alpha - f(c_0)}{f(c_0-) - f(c_0)} & T(X) = c_0 \\ 0 & T(X) < c_0 \end{cases}$ , then we can check $E_{\theta_0} \phi(X) = \alpha$.

Fix $\theta_1 > \theta_0$, we need to show $\phi$ is MP for $\theta = \theta_0$ vs $\theta = \theta_1$.

Let $\dfrac{P_{\theta_1}(x)}{P_{\theta_0}(x)} = g_{\theta_0, \theta_1}(T(X))$ where $g_{\theta_0, \theta_1}(\cdot)$ is non-decreasing

Set $K = g_{\theta_0, \theta_1}(c_0)$. If $\dfrac{P_{\theta_1}(x)}{P_{\theta_0}(x)} > k \iff g_{\theta_0, \theta_1}(T(X)) > g_{\theta_0, \theta_1}(c_0)$

$$\Rightarrow T(X) > c_0 \Rightarrow \phi = 1.$$

$$\text{If } \dfrac{P_{\theta_1}(x)}{P_{\theta_0}(x)} < k \Rightarrow \phi = 0$$

$\Rightarrow \phi$ is of NP form $\Rightarrow \phi$ is MP for $\theta = \theta_0$ vs $\theta = \theta_1$ at level $\alpha$.

$\Rightarrow \phi$ is UMP for $\theta = \theta_0$ vs $\theta > \theta_0$ at level $\alpha$

N.T.S $\sup\limits_{\theta < \theta_0} E_\theta \phi \le \alpha$ s.t. $\phi$ is UMP for $\theta \le \theta_0$ vs $\theta > \theta_0$

Fix $\theta_0' \le \theta_0$, consider the test problem of $\theta = \theta_0'$ vs $\theta = \theta_0$ at level $\beta(\theta_0') = E_{\theta_0'} \phi$.

$\phi$ is MP for this problem.

$\Rightarrow \beta(\theta_0') = E_{\theta_0'} \phi \le \underbrace{E_{\theta_0} \phi}_{\text{power}} = \alpha$ (because size $\le$ power)

$\Rightarrow \phi$ is level $\alpha$ UMP test for $\theta \le \theta_0$ vs $\theta > \theta_0$

(b) Fix $\theta' \le \theta''$, assume $\beta(\theta'), \beta(\theta'') \in (0,1)$. N.T.S $\boxed{\beta(\theta') < \beta(\theta'')}$

Consider the problem of testing $\theta = \theta'$ vs $\theta = \theta''$ at level $\beta(\theta')$, $\phi$ is MP for this problem,

$\Rightarrow \beta(\theta') < \beta(\theta'')$ because of again "size $\le$ power"

(c) Repeat the proof [TSH 3.4.1].

(d) Fix $\theta' < \theta_0$, N.T.S $\phi$ minimizes $E_{\theta'} \tilde{\phi}$ for all tests with $E_{\theta_0} \tilde{\phi} \leq \alpha$.

$\Leftrightarrow \{ 1-\phi \text{ maximizes } 1-E_{\theta'} \tilde{\phi} \text{ subject to } 1-E_{\theta_0} \tilde{\phi} = 1-\alpha \}$

$\Leftrightarrow \{ \psi = 1-\phi \text{ maximizes } E_{\theta'} \tilde{\psi} \text{ subject to } E_{\theta_0} \tilde{\psi} = 1-\alpha \}$

$\theta \neq \theta' \Rightarrow P_\theta \neq P_{\theta'}$ (identifiability)

i.e. $\psi$ is MP for $\theta = \theta_0$ vs $\theta = \theta'$ at level $1-\alpha$

where $\psi = \begin{cases} 1 & T(x) < c_0 \\ \gamma & T(x) = c_0 \\ 0 & T(x) > c_0 \end{cases}$ and $E_{\theta_0} \psi = 1-\alpha$.

E.g. $X_1, \ldots, X_n \overset{iid}{\sim} U(0, \theta)$. Test $H_0 : \theta = 1$ vs $H_1 : \theta > 1$ at level $\alpha$.

By the thm, a UMP test is given by $\phi = \begin{cases} 1 \\ 0 \end{cases}$ if $\begin{matrix} X_{(n)} > K \\ X_{(n)} < K \end{matrix}$ and

$\alpha = E_{\theta=1} \phi = P_{\theta=1}(X_{(n)} > K) = 1 - P_{\theta=1}(X_{(n)} \leq K) = 1 - K^n \Rightarrow K = (1-\alpha)^{1/n}$.

E.g. (Cauchy location level)

Let $X$ have the density $P_\theta(x) = \dfrac{1}{1+(x-\theta)^2}$. We find two points at which the MLR condition fails.

For any fixed $\theta > 0$,

$\dfrac{P_\theta(x)}{P_0(x)} = \dfrac{1+x^2}{1+(x-\theta)^2} \to 1$ as $x \to \infty$ or $x \to -\infty$

but $P_\theta(0) / P_0(0) = \dfrac{1}{1+\theta^2}$, which is strictly less than 1.

Thus, the ratio must increase at some values of $x$ and decrease at others.

Hence, $P_\theta(x)$ is not monotone in $x$, or in other words, the likelihood ratio

$T(x) = x$ is not MLR.

Strategies for finding UMPs → existence not guaranteed.

1) Reduce the composite alternative to a simple alternative. If $H_1$ is composite, fix $\theta_1 \in \Theta_1$ and test $H_0$ against $H_1 : \theta = \theta_1$. (Hope that it does not depend on $\theta_1$) ✓

2) Collapse the composite null to a simple null $(\cdots)$ ☆.

3) Apply Neyman Pearson Lemma = Find the MP LRT test for simple vs simple case / use MLR trick. ✓

Least favorable distribution

Consider : $H_0 : X \sim f_\theta \quad \theta \in \Theta$
$\qquad\qquad H_1 : X \sim g$ (known).

We now impose a prior distribution $\pi$ on $\Theta_0$. So we consider a new set of Hypotheses:

$H_\pi : X \sim h_\pi(x) = \displaystyle\int_{\Theta_0} f_\theta(x) \, d\pi(\theta)$

vs $H_1 : X \sim g$

Let $\beta_\pi$ be the power of the MP level $\alpha$ test $\phi_\pi$ for testing $H_\pi$ vs. $g$.

**Def** $\pi$ is a least favorable distribution if $\beta_\pi \leq \beta_{\pi'}$ for any prior $\pi'$. (smallest power).

**[Thm]** (TSH 3.8.1) Suppose $\phi_\pi$ is a MP level $\alpha$ test for testing $H_\pi$ against $g$. If $\phi_\pi$ is level $\alpha$ for the original hypothesis $H_0$ (ie. $E_{\theta_0} \phi_\pi(X) \leq \alpha \ \forall \theta_0 \in \Theta_0$), then

  (a) The test $\phi_\pi$ is MP for the original $H_0 : \theta \in \Theta_0$ vs. $g$.

  (b) The distribution $\pi$ is least favorable.

**Pf:** (a) Let $\phi^*$ be any other level-$\alpha$ test of $H_0 : \theta \in \Theta_0$ vs $g$. Then $\phi^*$ is also a level $\alpha$ test for $H_\pi$ v.s. $g$ because

$$E_\theta \phi^*(X) = \int \phi^*(X) f_\theta(X) d\mu(X) \leq \alpha \quad \forall \theta \in \Theta_0.$$

which implies that

$$\int \phi^*(X) h_\pi(X) d\mu(X) = \iint \phi^*(X) f_\theta(X) d\mu(X) d\pi(\theta) \leq \int \alpha \, d\pi(\theta) = \alpha$$

Since $\phi_\pi$ is MP for $H_\pi$ vs $g$, we have

$$\int \phi^*(X) \underbrace{g(X) d\mu(X)}_{power} \leq \int \phi_\pi(X) g(X) d\mu(X)$$

Hence $\phi_\pi$ is a MP test for $H_0$ vs $g$ because $\phi_\pi$ is also level $\alpha$.

  (b) Let $\pi'$ be any distribution on $\Theta_0$. Since $E_\theta \phi_\pi(X) \leq \alpha \ \forall \theta \in \Theta_0$, we know that $\phi_\pi$ must be level $\alpha$ for $H_{\pi'}$ vs $g$. Thus $\beta_\pi \leq \beta_{\pi'}$ so $\pi$ is least favorable dist.

**Example** (Testing in the presence of nuisance parameters)

Let $X_1, \dots, X_n$ be iid $N(\theta, \sigma^2)$, where both $(\theta, \sigma^2)$ are unknown.

We consider the test $H_0 : \sigma \leq \sigma_0$ against $H_1 : \sigma > \sigma_0$. $\quad \Theta_0 = \{(\theta, \sigma), \theta \in \mathbb{R}, \sigma < \sigma_0\}$

1. Fix a simple alternative $(\theta_1, \sigma_1)$ for some arbitrary $\theta_1, \sigma_1 > \sigma_0$.

2.* Choose a prior $\pi$ to "collapse over null hypothesis". $\sigma = \sigma_0$

  Consider the boundary case between $H_0$ and $H_1$ : $\{\sigma = \sigma_0\}$

  $\pi$ will be a prob. dist. over $\theta \in \mathbb{R}$ for the fixed $\sigma = \sigma_0$.

  * observation : Given any test function $\phi(X)$ and a sufficient satistic $T$, there exists a test function $\eta$ that less than some power as $\phi$ but depends on $X$ only through $T$.

$$\eta(T(X)) = E\{\phi(X) | T(X)\}.$$

  Hence, we restrict our attention to sufficient statistics.

  $(Y, u)$ where $Y = \bar{X}_n$ and $u = \sum_{i=1}^{\infty} (X_i - \bar{X}_n)^2$. We know that $Y \sim N(\theta, \frac{\sigma^2}{n})$ and $u \sim \sigma^2 X_{n-1}$ and $Y$ is ind. of $u$ due to Basu's thm.

Thus, for $\pi$ supported on $\sigma = \sigma_0$, we obtain the joint density of $(Y, u)$ under $H_\pi$ as

$$C_0 u^{\frac{n-3}{2}} \exp\left(-\frac{u}{2\sigma^2}\right) \int \exp\left(-\frac{u}{2\sigma_0^2} (y-\theta)^2\right) d\pi(\theta)$$

and the joint density under the alternative hypothesis.

$$C_1 \, u^{\frac{n-3}{2}} \exp\left(-\frac{u}{2\delta_1^2}\right) \exp\left(-\frac{n}{2\delta_1^2}(y-\theta_1)^2\right).$$

To achieve the minimal maximum power against the alternative (i.e. to be least favorable) we need to choose $\pi$ s.t. the two distributions become as close as possible. Under $H_1$, $Y \sim N(\theta_1, \frac{\delta^2}{n})$. Under $H_\pi$, the distribution of $Y$ is in a convolution form, i.e. $Y = Z + \Theta$, for $Z \sim N(0, \frac{\delta_0^2}{n})$, $\Theta \sim \pi$, where $Z$ and $\Theta$ are indep.. Hence, if we choose $\Theta \sim N(\theta_1, \frac{\delta_1^2 - \delta_0^2}{n})$ $Y$ will become the same distribution under both $H_\pi$ and $H_1$, which is $N(\theta_1, \frac{\delta_1^2}{n})$. Under this prior, the LRT rejects for large values of $\exp\left\{-\frac{u}{2\delta_1^2} + \frac{u}{2\delta_0^2}\right\}$, i.e. large values of $u$.

So, the MP test rejects $H_\pi$ if $\sum_{i=1}^{n}(X_i - \bar{X})^2$ lies above the threshold determined by the size constraint. In particular, it rejects $H_\pi$ if $\sum_{i=1}^{n}(X_i - \bar{X}_n)^2 > \delta_0^2 \, C_{n-1, 1-\alpha}$, where $C_{n-1, 1-\alpha}$ is the $(1-\alpha)^{th}$ quantile of $\chi_{n-1}^2$.

3. Check if this MP test is of level $\alpha$ for the composite null. For any $(\theta, \delta)$ with $\delta \leq \delta_0$, the prob. of rejection is $P_{\theta, \delta}\left(\frac{\sum_{i=1}^{n}(X_i - \bar{X}_n)^2}{\delta^2} > \frac{\delta_0^2 \, C_{n-1, 1-\alpha}}{\delta^2}\right) = P\left(\chi_{n-1}^2 > \frac{\delta_0^2}{\delta^2} C_{n-1, 1-\alpha}\right) \leq \alpha$.

with equality holds iff $\delta = \delta_0$. Hence, it follows ~~that~~ from thm (TSH 3.8.1) that our test is MP for testing the original null vs $N(\theta_1, \delta_1)$

4. Finally, the MP level $\alpha$ test for testing the composite null $H_0$ vs an arbitrary choosen alternative $(\theta_1, \delta_1)$ ~~des~~ does not depend on $(\theta_1, \delta_1)$. Hence, it is UMP for testing $H_0$ against $H_1$.

$$KS: \sup_t |\hat{F}_0(t) - F(t)|$$

Emperical cdf

$$CrM: \int (\hat{F}_n(t) - F(t))^2 dF.$$

3 tests

* Likelihood ratio test   (Keener)

$X_1, \ldots, X_n \overset{iid}{\sim} P_\theta(\cdot)$. We want to test $H_0: \theta \in \Theta_0$ vs $H_1: \theta \in \Theta_1$.

$$LRT: \quad \Lambda(X_1, \ldots, X_n) = \frac{\sup_{\theta \in \Theta_0} P_\theta(X_1, \ldots, X_n)}{\sup_{\theta \in \Theta_0 \cup \Theta_1} P_\theta(X_1, \ldots, X_n)}$$

$$-2\log \Lambda(\underline{X}) \overset{d}{\to} \chi^2_{d(\Theta_0 \cup \Theta_1) - d(\Theta_0)}$$

* Wald test

$$\sqrt{n}(\hat{\theta}_0 - \theta_0) \overset{d}{\to} N(0, I^{-1}(\theta_0))$$

* Rao Score test

$$\frac{1}{\sqrt{n}}\sum_{i=1}^{n} U_{\theta_0}(X_i) \overset{d}{\underset{\theta_0}{\to}} N(0, I(\theta_0)) \qquad U_\theta(X_i) = \frac{\partial}{\partial \theta}\log P_\theta(X_i).$$