

Risk, sufficiency, completeness, and Ancillarity

①

- X — RV. $X \sim P_\theta$ $\theta \in \mathbb{R}^p$

model: a set of distributions for X : $\mathcal{P} = \{P_\theta: \theta \in \Omega\}$

- Statistic: a function of the data X .

ex.
$$g(X) = \bar{X} = \frac{1}{n}(X_1 + \dots + X_n)$$

- Inference: Estimation hypothesis testing

- estimation $g(X) \rightarrow g(\theta)$

ex. $X \sim \text{Bin}(100, \theta)$ Flip a coin 100 times

$\theta \in [0, 1] = \Omega$ $P_\theta = \text{Bin}(100, \theta)$

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

An estimator of θ is $g(X) = \frac{X}{100}$

- good or bad? decision theory

Loss function: $L(\theta, d) \geq 0$ estimate $g(\theta)$ by a value d

Risk function: $R(\theta, g) = E_\theta L(\theta, g(X))$

ex. $X \sim \text{Bin}(100, \theta)$ $g(x) = \frac{x}{100}$ $g(\theta) = \theta$

$L(\theta, d) = (\theta - d)^2$

$R(\theta, \delta) = E_{\theta} \left(\theta - \frac{X}{100} \right)^2 = \frac{\theta(1-\theta)}{100}$ $\theta \in [0, 1]$

- compare different estimators:

$\delta_0(x) = \frac{x}{100}$ $\delta_1(x) = \frac{x+3}{100}$ $\delta_2(x) = \frac{x+3}{106}$

$R(\theta, \delta_0) = \frac{\theta(1-\theta)}{100}$ $R(\theta, \delta_1) = \frac{9 + 100\theta(1-\theta)}{100^2}$ $R(\theta, \delta_2) = \frac{(980)(1+\theta)}{106^2}$

when θ is near $\frac{1}{2}$, δ_2 is preferred.

minimax estimator

Sufficient Statistics

ex. X, Y — independent $f_{\theta}(x) = \theta e^{-\theta x}$ $x > 0$

U independent of X and Y $U \sim \text{unif}(0, 1)$

$T = X + Y$ $\left. \begin{array}{l} \tilde{X} = uT \\ \tilde{Y} = (1-u)T \end{array} \right\}$

(3)

density of T :

$$P(T \leq t | Y=y) = P(X+Y \leq t | Y=y) = F_X(t-y)$$

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

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

$$P_T(t) = t\theta^2 e^{-\theta t} \quad t \geq 0$$

joint density of (T, U) :

$$P_\theta(t, u) = \begin{cases} t\theta^2 e^{-\theta t} & t > 0, u \in (0, 1) \\ 0 & \text{o.w.} \end{cases}$$

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

$$= \iint \mathbb{1}_B(x, y) \frac{1}{x+y} P_\theta(x+y, \frac{x}{x+y}) dy dx$$

joint density (\tilde{X}, \tilde{Y}) :

$$\frac{P_\theta(x+y, \frac{x}{x+y})}{x+y} = \theta^2 e^{-\theta(x+y)} \quad x \geq 0, y \geq 0$$

→ same as the joint density of (X, Y) .

T provides as much information about θ as the pair (X, Y) .

$$T = X+Y \longrightarrow \text{sufficient statistic.}$$

conditional distributions Q_t for X and Y given $T=t$ (4)

$$Q_t(B) = P\left[\begin{pmatrix} X \\ Y \end{pmatrix} \in B \mid T=t\right] = P\left[\begin{pmatrix} U_t \\ (1-U_t) \end{pmatrix} \in B\right]$$

does not depend on θ .

Def.: $X \rightarrow \mathcal{P} = \{P_\theta : \theta \in \mathcal{R}\}$.

$\Rightarrow T=T(X)$ is a sufficient statistic for \mathcal{P} if for every t and θ , the conditional distribution of X under P_θ given $T=t$ does not depend on θ .

Suppose T is sufficient. $Q_t(B) = P_\theta(X \in B \mid T=t)$

$$P_\theta(X \in B \mid T) = Q_T(B). \quad P_\theta(X \in B) = E_T Q_T(B).$$

Construct "fake" data \tilde{X} from T taking $\tilde{X} \sim Q_t$ when $T=t$. then $\tilde{X} \mid T=t \sim Q_t$.

$$P_\theta(\tilde{X} \in B) = E_\theta P_\theta(\tilde{X} \in B \mid T) = E_\theta Q_T(B).$$

X and \tilde{X} have the same distribution.

Thm. $X \rightarrow \mathcal{P} = \{P_\theta : \theta \in \mathcal{R}\}$. $T=T(X)$ suff.

\Rightarrow For any $f(X)$ of $g(\theta)$ there exists a randomized estimator based on T that has the same risk function as $f(X)$.

Factorization Theorem.

- A sufficient statistic $T(x)$ conveys all of the information about θ from data X .

Def. A family of distributions $\mathcal{P} = \{P_\theta : \theta \in \Omega\}$ is dominated if there exists a measure μ with P_θ absolutely continuous w.r.t. μ , for all $\theta \in \Omega$.

Thm. $\mathcal{P} = \{P_\theta : \theta \in \Omega\}$ dominated by μ . A necessary and sufficient condition for T to be sufficient is that

$$P_\theta(x) = g_\theta(T(x)) h(x) \quad \text{for a.e. } x \text{ under } \mu.$$

ex. $X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta(x) = (\theta+1)x^\theta, \quad x \in (0,1) \quad \theta > -1$

$$P_\theta(x) = \prod_{i=1}^n f_\theta(x_i) = \prod_{i=1}^n (\theta+1)x_i^\theta = (\theta+1)^n \left(\prod_{i=1}^n x_i \right)^\theta, \quad x \in (0,1)^n$$

$$T = \prod_{i=1}^n x_i \quad \text{--- suff.}$$

ex. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(\theta, \theta+1) \quad f_\theta(x) = I_{(\theta, \theta+1)}(x)$

$$\prod_{i=1}^n I_{(\theta, \theta+1)}(x_i) = I_{(\theta, \theta+r)}(\min x_i) I_{(-\infty, \theta+1)}(\max x_i)$$

$$T = (\min x_i, \max x_i) \quad \text{--- suff.}$$

Minimal sufficiency

(6)

- If T is suff. if $T = f(\tilde{T}) \Rightarrow \tilde{T}$ is also suff.

Def. T is minimal suff. if T is suff. and for every suff. \tilde{T} there exists a function f s.t. $T = f(\tilde{T})$.

ex. $X_1, \dots, X_n \stackrel{iid}{\sim} N(0, 1)$

$\tilde{T} = \begin{pmatrix} \sum_{i=1}^n X_i \\ \sum_{i=n+1}^{2n} X_i \end{pmatrix}$ is suff. but not minimal suff.

$T = \sum_{i=1}^{2n} X_i$ - minimal suff. $T = f(\tilde{T})$. $f(t) = t_1 + t_2$.

- $P_\theta(x)$: viewed as a function of θ — likelihood
any suff. statistic must provide enough information to graph the shape of the likelihood.

Thm: $P_\theta(x) = g_\theta(T(x))h(x)$. if $P_\theta(x) \propto P_\theta(y)$ implies

$T(x) = T(y)$, then T is minimal sufficient

$P_\theta(x) = c(x, y)P_\theta(y)$

ex. \mathcal{P} — s-parameter exponential family.

$$p_\theta(x) = e^{\eta(\theta)T(x) - B(\theta)} h(x) \quad \theta \in \mathcal{R}$$

T — suff. Suppose $p_\theta(x) \propto_\theta p_\theta(y)$

$$e^{\eta(\theta)T(x)} \propto_\theta e^{\eta(\theta)T(y)} \quad \eta(\theta)T(x) = \eta(\theta)T(y) + c(x,y)$$

if θ_0 and θ_1 are any two points in \mathcal{R} .

$$[\eta(\theta_0) - \eta(\theta_1)] T(x) = [\eta(\theta_0) - \eta(\theta_1)] T(y)$$

$$[\eta(\theta_0) - \eta(\theta_1)] \cdot [T(x) - T(y)] = 0$$

$T(x) - T(y)$ is orthogonal to every vector in

$$\eta(\mathcal{R}) \ominus \eta(\mathcal{R}) \triangleq \{ \eta(\theta_0) - \eta(\theta_1) : \theta_0 \in \mathcal{R}, \theta_1 \in \mathcal{R} \}$$

it must lie in the orthogonal complement of the linear

span of $\eta(\mathcal{R}) \ominus \eta(\mathcal{R})$. If $\eta(\mathcal{R}) \ominus \eta(\mathcal{R}) = \mathbb{R}^s$. $T(x) = T(y)$

ex. X_1, \dots, X_n iid $f_\theta(x) = \frac{1}{2} e^{-|x-\theta|}$

$$\text{joint } p_\theta(x) = \frac{1}{2^n} e^{-\sum_{i=1}^n |x_i - \theta|}$$

$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ — order statistic — suff.

$$p_\theta(x) \propto p_\theta(y) \quad \sum_{i=1}^n |x_i - \theta| = \sum_{i=1}^n |y_i - \theta| + c$$

piecewise linear functions of θ with a slope that increases by two at each order statistic. The difference can only be constant in θ if x and y have the same order statistic.

Completeness

(8)

Def. A statistic T is complete for a family $\mathcal{P} = \{P_\theta : \theta \in \Omega\}$ if $E_\theta f(T) = c$ for all θ implies $f(T) = c$ (a.e. \mathcal{P}).

ex. $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{Unif}(0, \theta)$

joint density $I(\min X_i > 0) I(\max X_i < \theta) \frac{1}{\theta^n}$

$T = \max\{X_1, \dots, X_n\}$ sufficient.

$$P_\theta(T \leq t) = P_\theta(X_1 \leq t, \dots, X_n \leq t) = P_\theta(X_1 \leq t) \dots P_\theta(X_n \leq t) \\ = \left(\frac{t}{\theta}\right)^n \quad t \in (0, \theta)$$

T has density $\frac{n t^{n-1}}{\theta^n} \quad t \in (0, \theta)$

if $E_\theta f(T) = c$ for all $\theta > 0$

$$E_\theta(f(T) - c) = \frac{n}{\theta^n} \int_0^\theta (f(t) - c) t^{n-1} dt = 0$$

$f(T) = c \Rightarrow T$ is complete.

Thm: if T is complete and sufficient, then T is minimal sufficient.

pf. let \tilde{T} be a minimal suff. statistics. $\tilde{T} = f(T)$

Define $g(\tilde{T}) = E_{\theta}(T | \tilde{T})$ — independent of θ .

$E_{\theta} g(\tilde{T}) = E_{\theta}(T)$. $E_{\theta}(T - g(\tilde{T})) = 0$ for all θ .

$T - g(\tilde{T}) = T - g(f(T)) \Rightarrow T = g(\tilde{T})$ a.e. P .

Def. An exponential family with densities $f_{\theta}(x) = e^{\eta(\theta)T(x) - \beta(\theta)} h(x)$ $\theta \in \Omega$ is said to be of full rank if the interior of $\eta(\Omega)$ is not empty. and if T_1, \dots, T_s do not satisfy a linear constraint $v \cdot T = c$ (a.e. μ).

If $\Omega \subset \mathbb{R}^s$ and η is continuous and one-to-one, and the interior of Ω is nonempty, then the interior of $\eta(\Omega)$ cannot be empty.

$\eta(\Omega) \ominus \eta(\Omega)$ will be all of \mathbb{R}^s . $\Rightarrow T$ will be minimal suff.

Thm. In an exponential family of full rank, T is complete.

Def. A statistic V is called ancillary if its distribution does not depend on θ . \rightarrow no information about θ

Thm (Basu) if T is complete and suff. for $P = \{P_\theta, \theta \in \Omega\}$ ⁽¹⁰⁾
 and if V is ancillary, then T and V are independent under P_θ for any $\theta \in \Omega$.

pf: $g_A(t) = P_\theta(V \in A | T=t)$ $g_A(T) = P_\theta(V \in A | T)$

$P_A = P_\theta(V \in A)$. By sufficiency and ancillarity, neither P_A nor $g_A(t)$ depend on θ .

$$P_A = P_\theta(V \in A) = E_\theta g_A(T)$$

By completeness, $g_A(T) = P_A$.

$$P_\theta(T \in B, V \in A) = E_\theta 1_B(T) 1_A(V) = E_\theta E_\theta(1_B(T) 1_A(V) | T)$$

$$= E_\theta 1_B(T) E_\theta(1_A(V) | T) = E_\theta 1_B(T) g_A(T)$$

$$= E_\theta 1_B(T) P_A = P_\theta(T \in B) P_\theta(V \in A).$$

ex. $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ $P = P_\sigma = \{N(\mu, \sigma^2)^n, \mu \in \mathbb{R}\}$

$$\bar{X} = \frac{1}{n}(X_1 + \dots + X_n).$$

Joint density: $\frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left[\frac{n\mu}{\sigma^2}\bar{X} - \frac{n\mu^2}{2\sigma^2} - \frac{1}{2\sigma^2}\sum_{i=1}^n X_i^2\right]$

full rank exponential family.

\bar{X} is complete.

$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ — sample variance. ancillary.

$Y_i = x_i - \mu \sim N(0, \sigma^2)$. $\bar{Y} = \bar{x} - \mu$ $x_i - \bar{x} = Y_i - \bar{Y}$.

$S^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$. joint distribution of Y_1, \dots, Y_n depends on

σ but not μ . $\Rightarrow \bar{X}$ and S^2 are independent.

Convex Loss and the Rao-Blackwell Theorem.

Def. $f: C \rightarrow \mathbb{R}$. $C \subset \mathbb{R}^p$ convex

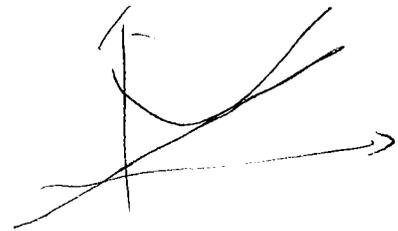
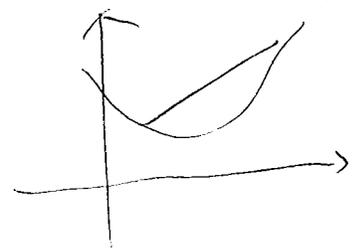
f convex if $f(rx + (1-r)y) \leq r f(x) + (1-r) f(y)$.
 $r \in (0,1)$.

strictly convex

$f'' \geq 0$ $\Rightarrow f$ convex

Thm. f — convex on an open interval C and if t is an arbitrary point in C , then there exists a constant

$c = c_t$ st. $f(t) + c(x-t) \leq f(x)$ $\forall x \in C$.



Thm: (Jensen's Inequality) C — an open interval. (12)
 f — convex on C . $P(X \in C) = 1$. $EX < \infty$.

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

pf. with $z = EX$. $f(EX) + c(X - EX) \leq f(x) \quad \forall x \in C$

$$f(EX) + c(X - EX) \leq f(X) \quad (\text{a.e. } P)$$

ex. $\frac{1}{x}$ $-\log x$ strictly convex $(0, +\infty)$

$$x > 0. \quad \frac{1}{EX} \leq E \frac{1}{x}. \quad \log EX \geq E \log x$$

Thm (Rao-Blackwell). T — suff. for $\mathcal{P} = \{P_\theta: \theta \in \Omega\}$.

δ — an estimator of $g(\theta)$ define $\eta(T) = E(\delta(X) | T)$.

$\forall \theta \in \Omega$, $R(\theta, \delta) < \infty$, $L(\theta, \cdot)$ is convex. then

$$R(\theta, \eta) \leq R(\theta, \delta).$$

$\forall L(\theta, \cdot)$ is strictly convex, the inequality will be strict unless $\delta(X) = \eta(T)$. (a.e. P_θ).

pf. Jensen's Inequality:

(13)

$$L(\theta, \eta(T)) \leq E_{\theta} (L(\theta, S(X)) | T).$$

Taking expectations, $R(\theta, \eta) \leq R(\theta, S)$.

- If estimators are judged by their risk, it should be a function of T , but not X .