

sufficiency (ct'd).

How to find sufficient stats?

Thm (factorization).

$\{P_\theta : \theta \in \Theta\}$, assuming there exists $P_\theta = \frac{dP_\theta}{d\mu}$ (for some base m.s.r μ).

then T is sufficient for the class

$\exists g_\theta, h$ s.t.

$$\stackrel{\text{def}}{=} P_\theta(x) = \underbrace{g_\theta(T(x))}_{\text{generating } T(x)} \cdot \underbrace{h(x)}_{\text{Conditional on } T, \text{ generate } X \text{ w.r.t. } \theta}.$$

(Proof assumes \exists joint density based on θ).

Conditionally on T ,
generate X
w.r.t. θ

Proof. " \Leftarrow ".

$$P_\theta(x | T(X)=t) = \frac{g_\theta(t) h(x) \mathbb{1}_{T(X)=t}}{\int_{T(z)=t} g_\theta(z) h(z) d\mu(z)}.$$

$$= \frac{h(x) \mathbb{1}_{T(X)=t}}{\int_{T(z)=t} h(z) d\mu(z)}$$

(Indp of θ).

" \Rightarrow ". Construction:

$$g_{\theta}(t) = P_{\theta}(T(X)=t)$$

marginal density

$$h(x) = P_{\theta_0}(X=x \mid T(X)=t) \quad \text{conditional density.}$$

(for arbitrary $\theta_0 \in \Theta$).

$$P_{\theta}(x) = P_{\theta}(T(X)=t, X=x) = P_{\theta}(T(X)=t) \cdot \underbrace{P_{\theta}(X=x \mid T(X)=t)}_{\text{indp of } \theta.}$$

(for $t=T(x)$.)

(see Keener textbook for a rigorous proof). So $g_{\theta} = h(x)$.

e.g. Exponential family.

$$P_{\theta}(x) = \exp\left(\eta(\theta)^T T(x) - B(\theta)\right) \cdot h(x).$$

By factorization, T is sufficient

, $N(\mu, \Sigma)$. in \mathbb{R}^d

$$P_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{d/2} \cdot \det(\Sigma)^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$$\propto \exp\left(-\frac{1}{2} x^T \Sigma^{-1} x + \mu^T \Sigma^{-1} x\right). \quad \eta(\mu, \Sigma) = \begin{bmatrix} \Sigma^T \mu \\ \text{vec}(\Sigma^{-1}) \end{bmatrix}$$

$$T(x) = \begin{bmatrix} x \\ \text{vec}(xx^T) \end{bmatrix} \in \mathbb{R}^{d+d^2}$$

(actually, we only need $d + \frac{d(d+1)}{2}$ dimensions).

- $X \sim N(\mu, \text{Id})$. $T(X) = X \in \mathbb{R}^d$ gives an exp family
- $X \sim N(0, \Sigma)$ $T(X) = \text{vec}(XX^\top)$ gives an exp family
- $X \sim \text{Binom}(n, p)$ with known n .

$$p_p(x) = \binom{n}{x} \cdot p^x (1-p)^{n-x}$$

$$= \binom{n}{x} \cdot \exp\left(x \cdot \log \frac{p}{1-p} + n \log(1-p)\right).$$
 $T(X) = X$. $\eta(p) = \log \frac{p}{1-p}$.

$\eta(\theta)$: natural parametrization

Fact. $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} P_\theta$ where $\{P_\theta : \theta \in \Theta\}$ is an exp family.

$$P_\theta(x) = \exp(\eta(\theta)^\top T(x) - B(\theta)) \cdot h(x)$$

then $P_\theta(x_1, \dots, x_n) = \exp(\eta(\theta)^\top \sum_{i=1}^n T(x_i) - n B(\theta)) \prod_{i=1}^n h(x_i)$.

is also an exp family w/ $\sum_{i=1}^n T(x_i)$ sufficient.

Corollary. $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \Sigma)$

then (X_1, \dots, X_n) follows an exp family,

with $T(X) = \begin{bmatrix} \sum_{i=1}^n X_i \\ \sum_{i=1}^n \text{vec}(X_i X_i^\top) \end{bmatrix}$

e.g. $\text{Unif}([\theta-1, \theta+1]) \quad \theta \in \mathbb{R}$.

$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unif}([\theta-1, \theta+1]).$

$$T(X) = (X_{(1)}, X_{(n)})$$

where $X_{(1)}, X_{(2)}, \dots, X_{(n)}$

are sorted version of X_1, \dots, X_n

$$P_\theta(x_1, \dots, x_n) = 2^{-n} \cdot \mathbb{1}_{\{x_{(1)} \geq \theta-1, x_{(n)} \leq \theta+1\}}.$$

Use of sufficient stats: unbiased estimation.

Estimate $g(\theta)$ from $X \sim P_\theta$.

Def. δ is unbiased for $g(\theta)$ if

$$\mathbb{E}_\theta[\delta(X)] = g(\theta) \quad (\forall \theta \in \Theta)$$

Def. (UMVU) For any other δ' unbiased.

we have $\text{var}_\theta(\delta(X)) \leq \text{var}_\theta(\delta'(X)) \quad (\forall \theta \in \Theta)$

Among all the possible unbiased estimators, δ is the optimal.

Def (complete stats)

We call $T(X)$ is complete for $f_{P_\theta}(\theta; \Theta)$ if

$\underbrace{\mathbb{E}_\theta[f(T(x))]}_{\text{A system of integral eq.}} = 0, \forall \theta \in \Theta$ implies $f(T) = 0$ a.s.
 only have trivial solutions.

Ideal: no redundancy of info in T .
 e.g. If $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unif}([\theta, \theta+1])$.

$$T(X_1, \dots, X_n) = (X_1, \dots, X_n)$$

$$f(T) = X_1 - X_2$$

$$\text{But for } T(X_1, \dots, X_n) = (X_{(1)}, X_{(n)}).$$

you cannot find such functions.

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}([0, \theta])$.

$$T(X) = \max_{1 \leq i \leq n} X_i. \quad P_\theta(T \leq t) = \left(\frac{t}{\theta}\right)^n \text{ for } t \in [0, \theta].$$

$$\mathbb{E}_\theta[f(T(X))] = \frac{n}{\theta^n} \int_0^\theta f(t) t^{n-1} dt = 0 \quad (\forall \theta).$$

$$\Rightarrow \int_0^\theta f(t) t^{n-1} dt = 0 \quad \forall \theta$$

$$\text{So. } f(t) = 0 \quad \text{almost everywhere.}$$

Def. Exp family $P_\theta(x) = \exp(\eta(\theta)^T T(x) - B(\theta)) h(x)$
 is called full-rank if $\eta(\theta)$ has an interior point.
 $(T_1(x), T_2(x), \dots, T_d(x))$
 (linearly indp.)

Fact. Full-rank exp family $\Rightarrow T$ is sufficient & complete.

Proof. Let $\eta(\theta_0)$ be an interior point of $\eta(\mathbb{R})$.

Let v be density of $T(X)$ under θ_0 .
 (ideal: local perturbation around θ_0).

$$0 = \int_{\mathbb{R}^d} f(t) \cdot \frac{P_\theta(T(x)=t)}{v(t)} \cdot v(t) dt$$

$$\Rightarrow 0 = \int_{\mathbb{R}^d} f(t) \cdot \exp((\eta(\theta) - \eta(\theta_0))^T t) v(t) dt \quad (\forall \theta \in \mathbb{R}).$$

$$\text{So. } \int_{\mathbb{R}^d} f(t) e^{\eta^T t} v(t) dt = 0$$

for $\eta \in \mathbb{B}(0, r_0)$

multivariate
Laplace transform
of $f \cdot v$.

By uniqueness of Laplace transform.
 $f \cdot v = 0$ (a.e.).

Thm (Lehmamn - Scheffé)

If T is sufficient & complete for $(\theta \in \Theta)$
 Suppose that \exists unbiased estimator δ for $g(\theta)$.
 then $\delta^* = E[\delta(X) | T]$ is UMVU.

Proof: δ^* is unbiased

Consider another unbiased estimator δ' .

$$\tilde{\delta}(T) = E[\delta'(X) | T]$$

$$\forall \theta, \text{Var}_{\theta}(\tilde{\delta}) \leq \text{Var}(\delta').$$

On the other hand, by completeness.

$$E_{\theta}[\tilde{\delta}(T) - \delta^*(T)] = 0 \quad (\forall \theta \in \Theta).$$

$$\text{So } \tilde{\delta} = \delta^* \quad (\text{a.s.}).$$

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

$$\delta(X) = T(X) = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{is UMVU.}$$

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Ber}(\theta)$ $g(\theta) = \theta^k$

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

Unbiased estimator

$$f(X) = X_1 X_2 \cdots X_k.$$

$$\text{UMVU : } \mathbb{E}[\delta(X) | T(X)] = \frac{T(T-1) \cdots (T-k+1)}{n(n-1) \cdots (n-k+1)}.$$

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} N(0, \sigma^2)$ $g(\sigma^2) = \sigma^2$.

$$T(X) = \sum_{i=1}^n X_i^2 \quad \frac{T(X)}{n} \text{ is UMVU.}$$

Indeed, UMVU is not the "optimal" estimator

$$\delta_c(T) = cT.$$

$$\arg \min_c \text{MSE}(c) = \frac{1}{n+2}.$$

(sufficiency, unbiased estimation)

Additional remarks.

Least useful

most useful

UMVU.

sufficiency

debiasing

Cramér-Rao.

Completeness

Rao-Blackwell.

(approx unbl.).

Thm. (Cramér-Rao)

If $\hat{\theta}(X)$ is unbiased estimator for $g(\theta) \in \mathbb{R}$.

Assuming $(P_\theta : \theta \in \Theta)$ has shared support $(\forall \theta)$.

$\log P_\theta(X) \in C^2$, $\mathbb{E}|\nabla \log P_\theta(X)|^2 < +\infty$
twice \sim differentiable

then we have

$$\text{var}_\theta(\hat{\theta}(X)) \geq g(\theta)^T I(\theta)^{-1} g(\theta)$$

$$\text{where } I(\theta) := \mathbb{E}_\theta [\nabla \log P_\theta(X) \cdot \nabla \log P_\theta(X)^T].$$

(usually not achieved by unbiased estimation,
achievable (asymptotically) by MLE plug-in).

$$l(\theta; X) = \log P_\theta(X).$$

Proof. two simple facts.

$$\begin{aligned} \mathbb{E}_\theta[\nabla_\theta l(\theta; X)] &= \int (\nabla_\theta l(\theta; X) \cdot e^{l(\theta; X)}) d\mu(X) \\ &= \nabla_\theta \left(\int e^{l(\theta; X)} d\mu(X) \right) = 0. \end{aligned}$$

• (Fisher's identity).

Note that $\mathbb{E}_{\theta} \left[\frac{\partial}{\partial \theta_j} \ell(\theta; x) \right] = 0$.

take more derivative.

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta_i} \mathbb{E}_{\theta} \left[\frac{\partial}{\partial \theta_j} \ell(\theta; x) \right] \\ &= \frac{\partial}{\partial \theta_i} \int \frac{\partial}{\partial \theta_j} \ell(\theta; x) \cdot e^{\ell(\theta; x)} d\mu(x). \\ &= \int \left(\frac{\partial^2 \ell(\theta; x)}{\partial \theta_j \partial \theta_i} + \frac{\partial \ell(\theta; x)}{\partial \theta_i} \cdot \frac{\partial \ell(\theta; x)}{\partial \theta_j} \right) e^{\ell(\theta; x)} d\mu(x) \end{aligned}$$

So we have

$$\mathbb{E}_{\theta} \left[\nabla \ell(\theta; x) \nabla \ell(\theta; x)^T \right] = - \mathbb{E}_{\theta} \left[\nabla^2 \ell(\theta; x) \right].$$

(holding true only for well specified model at true θ).

Proof of CRLB.

$$\begin{aligned} \nabla_{\theta} g(\theta) &\stackrel{(unbias)}{=} \nabla_{\theta} \left(\int f(x) \cdot e^{\ell(\theta; x)} d\mu(x) \right) \\ &= \int f(x) \cdot \nabla_{\theta} \ell(\theta; x) \cdot e^{\ell(\theta; x)} d\mu(x). \end{aligned}$$

How to insert $g(\theta)$?

$$\int \nabla_{\theta} \ell(\theta; x) e^{\ell(\theta; x)} d\mu(x) = 0.$$

$g(\theta)$ indep of x .

$$\begin{aligned}\nabla g(\theta) &= \int (\delta(x) - g(\theta)) \cdot \nabla_{\theta} l(\theta; x) \cdot e^{l(\theta; x)} d\mu(x) \\ &= \mathbb{E}_{\theta}[(\delta(x) - g(\theta)) \cdot \nabla_{\theta} l(\theta; x)].\end{aligned}$$

In 1D. Cauchy-Schwarz gives the result.

In general, consider quadratic form

$$u \in \mathbb{R}^d \mapsto \mathbb{E}_{\theta}[(u^T \nabla_{\theta} l(\theta; x) - (\delta(x) - g(\theta)))^2] \geq 0.$$

Expanding the quadratic form.

$$0 \leq u^T I(\theta) u - 2 u^T \mathbb{E}_{\theta}[(\delta(x) - g(\theta)) \cdot \nabla_{\theta} l(\theta; x)] + \text{var}_{\theta}(\delta(x)) = \nabla g(\theta).$$

This implies $\text{var}_{\theta}(\delta(x)) \geq \nabla g(\theta)^T I(\theta)^{-1} \nabla g(\theta)$.

Corollary. $g: \mathbb{R}^d \rightarrow \mathbb{R}^k$. δ unbiased.

$$\mathbb{E}[(\delta(x) - g(\theta))^{\otimes 2}] \underset{(PSD)}{\succ} (\nabla g(\theta))^T I(\theta)^{-1} \cdot (\nabla g(\theta)).$$

In particular, $g(\theta) = \theta$

$$\mathbb{E}[||\delta(x) - \theta||_2^2] \geq \text{Tr}(I(\theta)^{-1}).$$

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} p_\theta$
 $\ell_n(\theta; X_1, \dots, X_n) = \sum_{i=1}^n \ell(\theta; X_i)$.

$$I_n(\theta) = n \cdot I(\theta).$$

and CRLB

$$\frac{(\nabla g(\theta))^T I(\theta)^{-1} \nabla g(\theta)}{n}.$$

(asymptotically achieved by MLE, $n \rightarrow \infty$)

When $p_\theta(x)$ is singular / has discontinuities

CRLB is false.

e.g. $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unif}([0, \theta])$.

$$\text{UMVU. } f(x) = \max_i x_i \cdot \frac{n+1}{n}$$

$$\text{Var}_\theta(f(x)) = \frac{\theta^2}{n(n+2)} \text{ factor rate.}$$

Bayesian CRLB.

Idea. key step in CRLB

$$\nabla g(\theta) = \int (\delta(x) - g(\theta)) \nabla \ell(\theta; x) \cdot e^{\ell(\theta; x)} dx.$$

relying on unblased of f .

"look like" integration-by-parts formula.

use a prior distribution to integrate.

$$r_{\pi}(\delta) := \int_{\Theta} \mathbb{E}_{\theta} [(\delta(x) - g(\theta))^2] \pi(d\theta).$$

Theorem (van Trees inequality). for any estimator

$$r_{\pi}(\delta) \geq \left(\int_{\Theta} \nabla g(\theta) \pi(d\theta) \right)^T \left(\int_{\Theta} I(\theta) \pi(d\theta) + J(\pi) \right)^{-1} \cdot \left(\int_{\Theta} \nabla g(\theta) \pi(d\theta) \right)$$

where $J(\pi) := \int_{\Theta} \nabla_{\theta} \log \pi(\theta) \nabla_{\theta} \log \pi(\theta)^T \pi(d\theta)$

"Information theoretic's Fisher info".

Compared to CRLB:

- Pay a $J(\pi)$ term
- No need for unbiasedness.

Proof: see next lecture.

e.g. $v_0(x) = \cos^2\left(\frac{\pi x}{2}\right) I_{|x| \leq 1}$ supported on $[-1, 1]$.

$$J(v_0) = \pi^2.$$