

Sufficiency

9/5/2023

Outline

- 1) Review
- 2) Sufficiency
- 3) Factorization Theorem

Sufficiency

Motivation: Coin flipping

Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\theta)$

$$\Rightarrow X \sim \prod_i \theta^{x_i} (1-\theta)^{1-x_i} \quad \text{on } \{0, 1\}^n$$

$$\begin{aligned} \text{Then } T(X) = \sum X_i &\sim \text{Binom}(n, \theta) \\ &= \theta^t (1-\theta)^{n-t} \binom{n}{t} \quad \text{on } \{0, \dots, n\} \end{aligned}$$

$(X_1, \dots, X_n) \rightarrow T(X)$ is throwing away data. How do we justify this?

In exp. fam. lingo, $T(X)$ is the "sufficient statistic" for X . Today we'll see why we call it that.

Definition Let $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ be a statistical model for data X . $T(X)$ is sufficient for \mathcal{P} if $P_\theta(X|T)$ does not depend on θ

Example (Cont'd)

$$\begin{aligned} P_\theta(X=x | T=t) &= \frac{P_\theta(X=x, T=t)}{P_\theta(T=t)} \\ &= \frac{\cancel{\theta^{\sum x_i}} \cancel{(1-\theta)^{n-\sum x_i}} \mathbf{1}\{\sum x_i = t\}}{\cancel{\theta^t} \cancel{(1-\theta)^{n-t}} \binom{n}{t}} \\ &= \mathbf{1}\{\sum x_i = t\} / \binom{n}{t} \end{aligned}$$

So given $T(X)=t$, X is uniform on all seq.s with $\sum x_i = t$

Factorization Theorem

Often, we can identify sufficient stats by inspecting the density.

Theorem (Factorization Theorem)

Let $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ be a model with densities $p_\theta(x)$ wrt common measure μ .

$T(x)$ is sufficient iff there exist $g_\theta(t)$, $h(x)$ with

$$p_\theta(x) = g_\theta(T(x)) h(x)$$

for μ -almost-every x : $\mu(\{x : p_\theta(x) \neq g_\theta(T(x)) \cdot h(x)\}) = 0$

[Avoids counterexamples from changing $p_{\theta_0}(x_0)$ some θ_0, x_0]

Rigorous proof in Keener 6.4

Proof (discrete \mathcal{X}): Assume $w \log \mu = \#$ on \mathcal{X}

$$\begin{aligned} (\Leftarrow) \quad P_{\theta}(X=x | T=t) &= \frac{P_{\theta}(X=x, T(x)=t)}{P_{\theta}(T(x)=t)} \\ &= \frac{\cancel{g_{\theta}(t)} h(x) 1\{T(x)=t\}}{\sum_{T(z)=t} \cancel{g_{\theta}(t)} h(z)} \end{aligned}$$

(\Rightarrow) Assume $T(x)$ sufficient.

$$\begin{aligned} \text{Take } g_{\theta}(t) &= \sum_{T(x)=t} p_{\theta}(x) \\ &= P_{\theta}(T(X)=t) \end{aligned}$$

For any $\theta_0 \in \Theta$, let

$$\begin{aligned} h(x) &= p_{\theta_0}(x) / \sum_{T(z)=T(x)} p_{\theta_0}(z) \\ &= \cancel{P_{\theta_0}}(X=x | T(X)=T(x)) \end{aligned}$$

Then,

\nwarrow no dep. on θ

$$\begin{aligned} g_{\theta}(T(x)) h(x) &= P_{\theta}(T=T(x)) P(X=x | T=T(x)) \\ &= P_{\theta}(X=x) \end{aligned}$$

□

Interpretations of Sufficiency

X is informative about θ only because its distribution depends on θ .

We can think of the data as being generated in two stages:

- 1) Generate T : distribution dep. on θ
- 2) Generate $X|T$: does not dep on θ

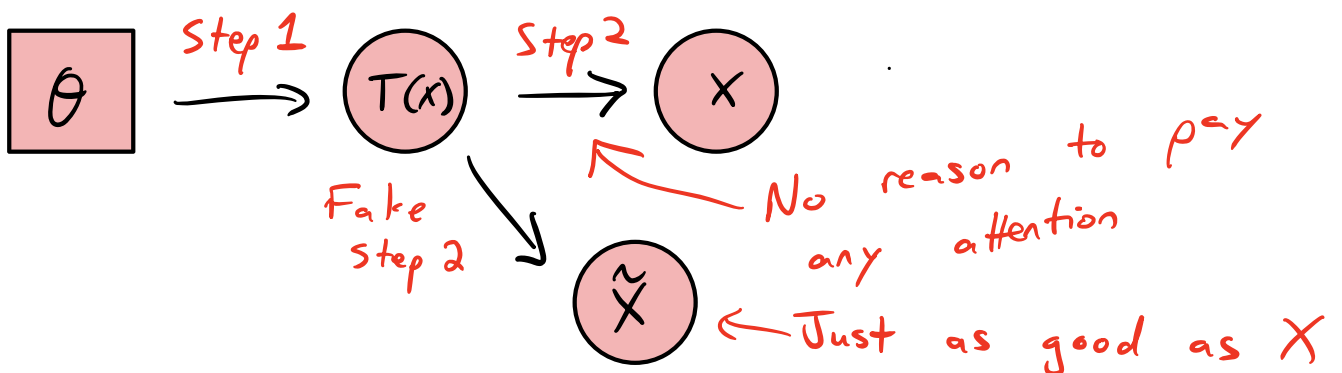
Sufficiency Principle

If $T(X)$ is sufficient for \mathcal{P} then any statistical procedure should depend on X only through $T(X)$

In fact, we could throw away X and generate a new $\tilde{X} \sim P(X|T)$ and it would be just as good as X since $\tilde{X} \sim P_\theta$

\leftarrow no θ

In graphical model form:



Examples

Ex. Exponential Families

$$p_{\theta}(x) = \underbrace{e^{\eta(\theta)'T(x) - B(\theta)}}_{g_{\theta}(T(x))} \underbrace{h(x)}_{h(x)}$$

Ex. Uniform location family

$$X_1, \dots, X_n \stackrel{iid}{\sim} U[\theta, \theta+1] \\ = 1\{\theta \leq x \leq \theta+1\}$$

$$p_{\theta}(x) = \prod_{i=1}^n 1\{\theta \leq x_i \leq \theta+1\} \\ = 1\{\theta \leq X_{(1)}\} 1\{X_{(n)} \leq \theta+1\}$$

$\Rightarrow (X_{(1)}, X_{(n)})$ is sufficient.

Order Statistics / Empirical Distribution

Ex. $X_1, \dots, X_n \stackrel{iid}{\sim} P_\theta^{(1)}$ for any model

$$\mathcal{P}^{(1)} = \{P_\theta^{(1)} : \theta \in \Theta\} \text{ on } \mathcal{X} \subseteq \mathbb{R}$$

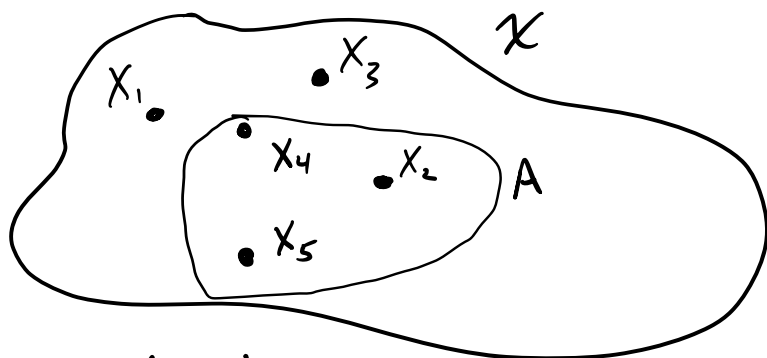
P_θ is invariant to perm.s of $X = (X_1, \dots, X_n)$

\Rightarrow All permutations of x are equally likely

\Rightarrow order statistics $(X_{(i)})_{i=1}^n$ ($X_{(k)} = k^{\text{th}}$ smallest) are sufficient. [Note $(X_i)_{i=1}^n \rightsquigarrow (X_{(i)})_{i=1}^n$ loses information, specifically the orig. ordering]

For more general \mathcal{X} we can say the empirical distribution $\hat{P}_n(\cdot) = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}(\cdot)$

is sufficient, where $\delta_{X_i}(A) = 1\{X_i \in A\}$



$$\hat{P}_n(A) = \frac{3}{5}$$

[Not important that it's a measure in this context; just keeps track of which values came up how many times]

Minimal Sufficiency

Consider $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\theta, 1)$

$$p_{\theta}^{(1)}(x) = \frac{1}{\sqrt{2\pi}} e^{\theta x - \theta^2/2 - x^2/2}$$

exponential family with $T(x) = x$

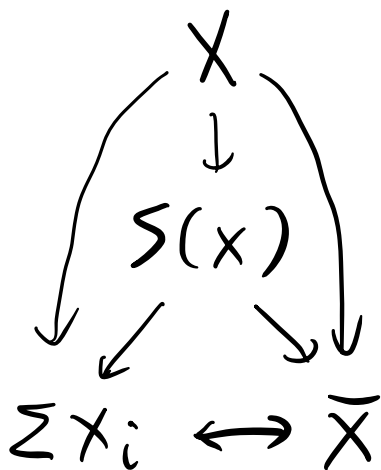
$$T(x) = \sum X_i \quad \text{sufficient}$$

$$\bar{X} = \frac{1}{n} \sum X_i \quad \text{also}$$

$$S(x) = (X_{(1)}, \dots, X_{(n)}) \quad \text{too}$$

$$X = (X_1, \dots, X_n) \quad \text{too}$$

Which can be recovered from which others?



these can be compressed further

These are the most compressed. Are they as compressed as possible?

Prop If $T(X)$ is sufficient and $T(X) = f(S(X))$
then $S(X)$ is sufficient

Proof : $p_{\theta}(x) = g_{\theta}(T(x)) h(x)$
 $= (g_{\theta} \circ f)(S(x)) h(x) \quad \square$

Definition: $T(X)$ is minimal sufficient if

- 1) $T(X)$ is sufficient
- 2) For any other sufficient $S(X)$,
 $T(X) = f(S(X))$ for some f
(a.s. in \mathcal{P})

So, no matter how many more suff. stats we add
to our diagram, they will all have arrows
pointing to ΣX_i

Likelihood Shape is Minimal

Definition

Assume $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ has densities $p_\theta(x)$

The likelihood function is the (random) function

$$\text{Lik}(\theta; X) = p_\theta(x)$$

function of θ data x determines which function function of x with parameter θ

The log-likelihood function is its log:

$$l(\theta; X) = \log \text{Lik}(\theta; X)$$

The likelihood up to scaling (or l up to vertical shift) is a minimal sufficient statistic

If $T(X)$ is sufficient then

$$\text{Lik}(\theta; x) = \underbrace{g_\theta(T(x))}_{T \text{ determines the "shape"}} \underbrace{h(x)}_{\text{scaling}}$$

HW 2: Likelihood ratios $\left(\frac{\text{Lik}(\theta_1; X)}{\text{Lik}(\theta_2; X)} \right)_{\theta_1, \theta_2 \in \Theta}$ minimal suff.

Recognizing Minimal Sufficient Statistics

$T(X)$ is minimal sufficient if

1) $T(X)$ is sufficient

(don't forget to check!)

2) $T(x)$ can be recovered from the likelihood shape

Keener Thm 3.11 formalizes condition 2

$$\text{"} \text{Lik}(\cdot; x) \propto \text{Lik}(\cdot; y) \Rightarrow T(x) = T(y) \text{"}$$

equivalently,

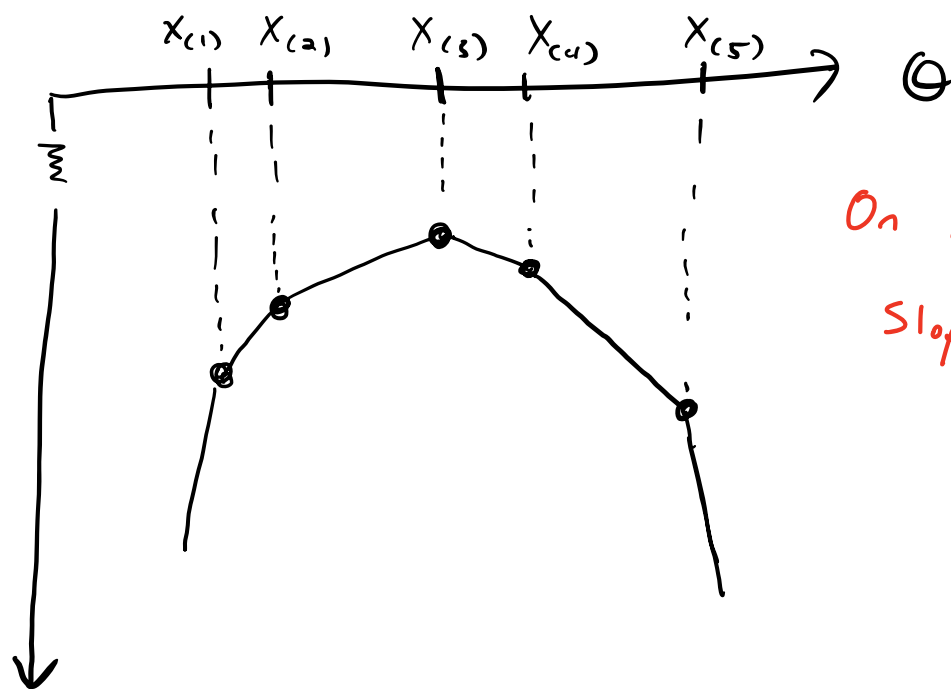
$$\text{"} \ell(\cdot; x) - \ell(\cdot; y) = \text{const}(x, y) \Rightarrow T(x) = T(y) \text{"}$$

Ex Laplace location family

$$X_1, \dots, X_n \stackrel{iid}{\sim} p_{\theta}^{(1)}(x) = \frac{1}{2} e^{-|x-\theta|}$$

$$l(\theta; x) = - \sum_{i=1}^n |x_i - \theta| - n \log 2$$

Piecewise linear in θ , knots at $x_{(i)}$



On $[x_{(k)}, x_{(k+1)}]$,

$$\text{Slope} = n - 2k$$

$$l(\theta; x) = l(\theta; y) + \text{const} \Leftrightarrow X, Y \text{ same order statistics}$$

\Rightarrow order stats are minimal suff.

Minimal sufficiency for exp. fam.s

$$\text{Suppose } p_{\eta}(x) = e^{\eta' T(x) - A(\eta)} h(x)$$

$$\ell(\eta; x) = \underbrace{T(x)' \eta}_{\text{random linear function of } \eta} - \underbrace{A(\eta)}_{\text{deterministic function of } \eta} + \underbrace{\log h(x)}_{\text{(random) const.}}$$

Is $T(x)$ minimal? (always sufficient)

Suppose x and y give same likelihood shape:

$$\ell(\eta; x) - \ell(\eta; y) = \text{const}(x, y)$$

$$\text{Then } (T(x) - T(y))' \eta = \text{const}(x, y) \quad \text{for } \eta \in \Xi$$

$$\Rightarrow T(x) = T(y) \quad \underline{\text{or}}$$

$$T(x) - T(y) \perp \text{Span}\{\eta_1, \eta_2 : \eta \in \Xi\}$$

If $\text{Span}\{\dots\} = \mathbb{R}^s$, $T(x)$ is minimal

(That is, if Ξ is not contained in a lower-dim affine space)

Otherwise might not be:

If $s=2$, $\Xi = \left\{ \begin{pmatrix} \theta \\ 0 \end{pmatrix} : \theta \in \mathbb{R} \right\}$ then $T_1(x)$ minimal

[Can we conclude $T(x)$ is not minimal?]

Other parameterizations:

$$p_{\theta}(x) = e^{\eta(\theta)'T(x) - \beta(\theta)} h(x)$$

$$\theta \in \Theta$$

$T(X)$ minimal if $\text{span}\{\eta(\theta_1) - \eta(\theta_2) : \theta_1, \theta_2 \in \Theta\} = \mathbb{R}^3$

