

# Sufficiency

## 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  $\theta$

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  $\sum x_i = t$

# Interpretations of Sufficiency

Recall we only care about  $X$  in the first place because it is (indirectly) informative about  $\theta$

Sufficiency means only  $T(x)$  is informative

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

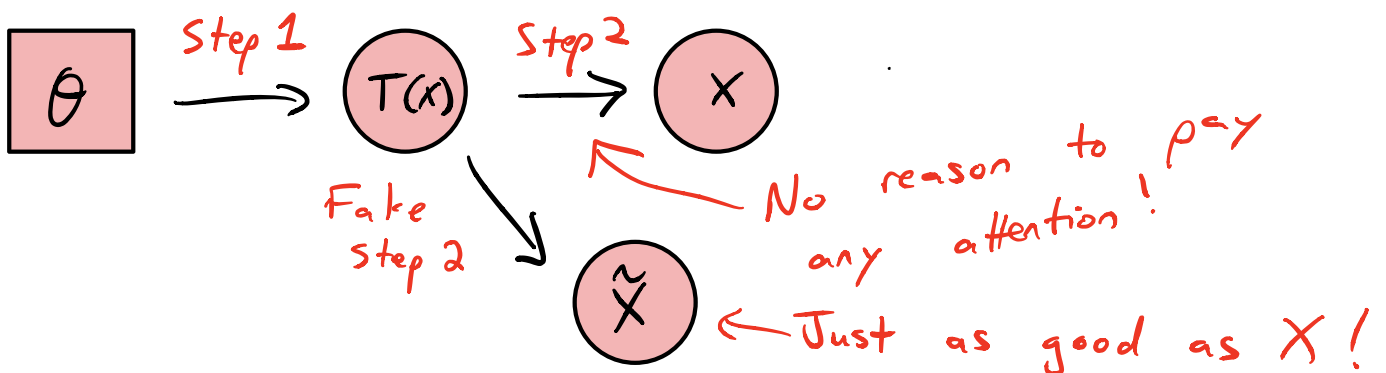
- 1) Generate  $T$ : distribution dep. on  $\theta$   
*"Pick a slice of  $\mathcal{X}$ "*
- 2) Generate  $X|T$ : does not dep on  $\theta$   
*"Generate within the slice"*  
*Dist. over slices depends on  $\theta$ , not dist. within each slice.*

## 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  $X$   
*no  $\theta$*

In graphical model form:



# Factorization Theorem

There is a very convenient way to verify sufficiency of a statistic based only on the density:

## Theorem (Factorization Theorem)

Let  $\mathcal{P} = \{p_\theta : \theta \in \Theta\}$  be a family of distributions dominated by  $\mu$  ( $p_\theta \ll \mu \forall \theta$ ), densities  $p_\theta$ .

$T$  is sufficient for  $\mathcal{P}$  iff there exist non-neg. functions  $g_\theta, h$  such that

$$p_\theta(x) = g_\theta(T(x)) h(x) \quad \text{for a.e. } x \text{ under } \mu$$
$$[\mu(\{x : p_\theta \neq g_\theta(T(x)) \cdot h(x)\}) = 0]$$

Rigorous proof in Keener 6.4

"Physics proof": (rigorous for discrete  $\mathcal{X}$ )

$$(\Leftarrow) p_\theta(x | T=t) = 1_{\{T(x)=t\}} \cdot \frac{g_\theta(t) h(x)}{\sum_{T(z)=t} g_\theta(t) h(z) d_\mu(z)}$$

$$(\Rightarrow) \text{ Take } g_\theta(t) = \int_{T(x)=t} p_\theta(x) d\mu(x) = P_\theta(T(X)=t)$$

$$h(x) = p_\theta(x) / \int_{T(z)=t} p_\theta(z) d\mu(z)$$

$$= P_\theta(X=x | T(X)=T(x))$$

$$g_\theta(T(x)) h(x) = P_\theta(T=T(x)) P(X=x | T=T(x)) \quad \square$$

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

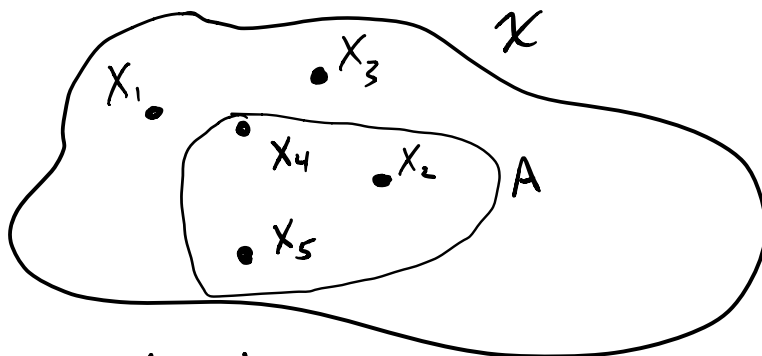
$$\mathcal{P}^{(1)} = \{P_{\theta}^{(1)} : \theta \in \Theta\} \text{ on } \mathcal{X} \subseteq \mathbb{R}$$

$P_{\theta}$  is invariant to perm.s of  $X = (X_1, \dots, X_n)$

$\Rightarrow$  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  $\delta_{X_i}(A) = \mathbb{1}\{X_i \in A\}$



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

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

$$\underline{E}_X. \quad X_1, \dots, X_n \stackrel{iid}{\sim} U[\theta, \theta+1]$$

$$= 1\{\theta \leq x \leq \theta+1\}$$

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

$$= 1\{\theta \leq X_{(1)}\} 1\{X_{(n)} \leq \theta+1\}$$

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

### Minimal Sufficiency

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

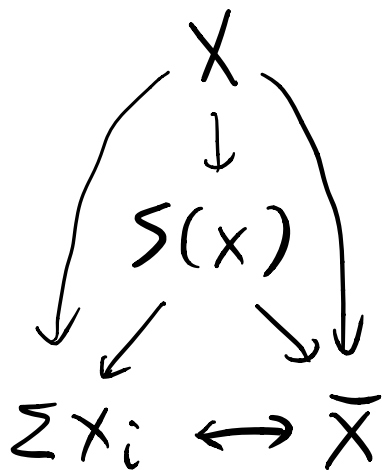
$$T(X) = \sum X_i \quad \text{sufficient}$$

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

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

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

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  $\Sigma X_i$

How to check minimal sufficiency? Basically,  
"equivalent to knowing likelihood ratios"

Definition Assume  $\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}$  has densities  
 $p_{\theta}(x)$  wrt common  $\mu$ , data  $X$ . The (log) likelihood  
function is the (log) density, reframed as a  
random function of  $\theta$ :

$$\text{Lik}(\theta; X) = p_{\theta}(x), \quad \ell(\theta; X) = \log \text{Lik}(\theta; X)$$

Note if  $T(X)$  is sufficient then

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

Thm 3.11 Assume  $\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}$ , densities  $P_{\theta}$   
 $T(X)$  sufficient for  $X$ .

If  $\text{Lik}(\theta; x) \propto_{\theta} \text{Lik}(\theta; y) \Rightarrow T(x) = T(y)$   
then  $T(X)$  is minimal sufficient

" $T$  determines the likelihood shape in a one-to-one fashion"

Proof Suppose  $S$  is sufficient and  $\nexists f$   
s.t.  $f(S(x)) = T(x)$

Then  $\exists x, y$  with  $S(x) = S(y)$ ,  $T(x) \neq T(y)$

$$\text{Lik}(\theta; x) = g_{\theta}(S(x)) h(x)$$

$$\propto_{\theta} g_{\theta}(S(y)) h(y)$$

$$= \text{Lik}(\theta; y)$$

But that implies  $T(x) = T(y)$  by assumption.



$$\underline{\text{Ex.}} \quad p_{\theta}(x) = e^{\eta(\theta)'T(x) - B(\theta)} h(x)$$

Is  $T(x)$  minimal?

doesn't const. in  $\theta$   
change with  $x$

Assume  $\text{Lik}(\theta; x) \propto_{\theta} \text{Lik}(\theta; y)$ . WTS  $T(x) = T(y)$

$$\text{Lik}(\theta; x) \propto_{\theta} \text{Lik}(\theta; y)$$

$$\Leftrightarrow e^{\eta(\theta)'T(x)} = e^{\eta(\theta)'T(y)} c(x, y) \quad \forall \theta$$

$$\Leftrightarrow \eta(\theta)'T(x) = \eta(\theta)'T(y) + a(x, y) \quad \forall \theta$$

$$\Leftrightarrow (\eta(\theta_1) - \eta(\theta_2))'T(x) = (\eta(\theta_1) - \eta(\theta_2))'T(y) \quad \forall \theta_1, \theta_2$$

$$\Leftrightarrow \eta(\theta_1) - \eta(\theta_2) \perp T(x) - T(y) \quad \forall \theta_1, \theta_2$$

$$\Leftrightarrow T(x) - T(y) \perp \text{Span}\{\eta(\theta_1) - \eta(\theta_2) : \theta_1, \theta_2 \in \Theta\}$$

We were trying to show  $T(x) - T(y) = 0$ , not quite there yet.

If  $\text{Span}\{\eta(\theta_1) - \eta(\theta_2) : \theta_1, \theta_2 \in \Theta\} = \mathbb{R}^S$ ,  
we are done.

Otherwise,  $T(x)$  really might not be minimal!

e.g.,  $\eta(\theta) = \begin{pmatrix} \theta \\ 0 \end{pmatrix}$  : then  $T_1(x)$  sufficient

[could  $T(x)$  still be minimal?]

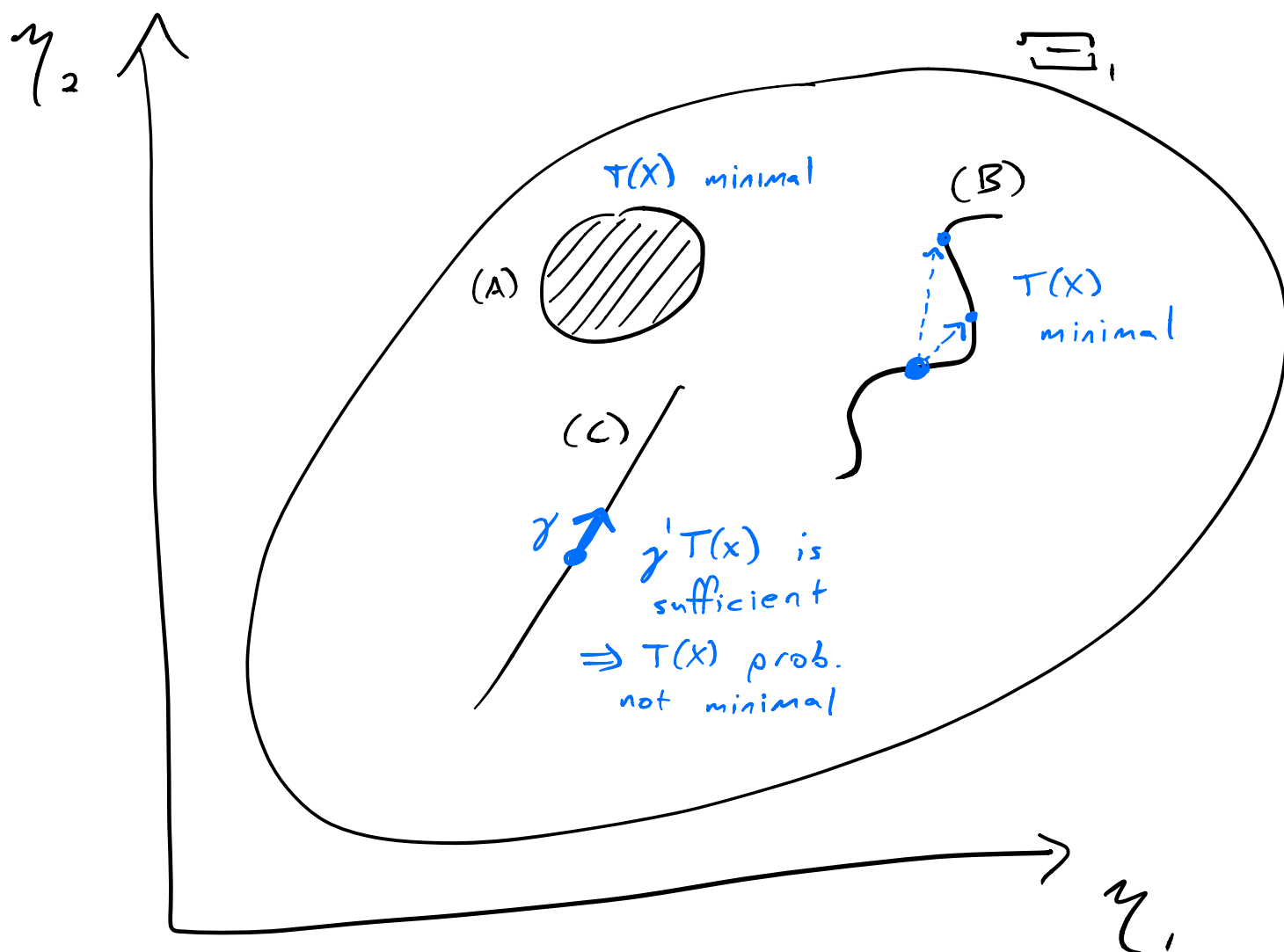
Ex.  $X \sim N_2(\mu(\theta), I_2)$   $\theta \in \mathbb{R}$

$$= e^{\mu(\theta)'x - B(\theta)} e^{-\frac{1}{2}x'x}$$

If  $\Theta = \mathbb{R}$ ,  $\mu(\theta) = a + \theta b$  for  $a, b \in \mathbb{R}^2$

$$\begin{aligned} p_\theta(x) &= e^{(a+\theta b)'x - B(\theta)} e^{-\frac{1}{2}x'x} \\ &= e^{\theta(b'x) - B(\theta)} e^{-\frac{1}{2}(x-2a)'x} \end{aligned}$$

$b'x$  is sufficient  $\Rightarrow X$  not minimal



Ex Laplace location family

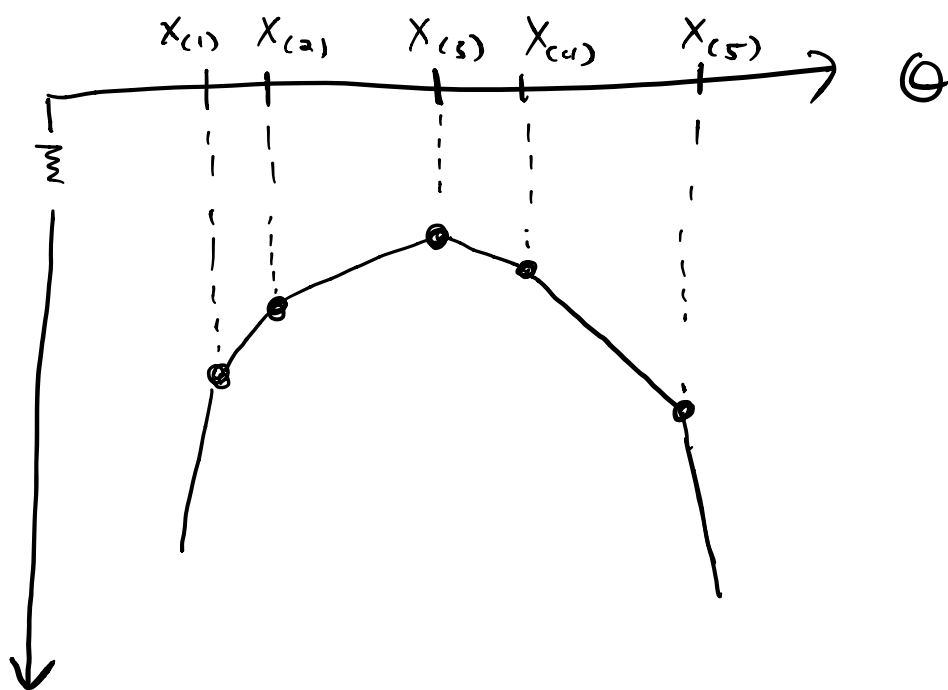
$$X_1, \dots, X_n \stackrel{iid}{\sim} \rho_\theta^{(1)}(x) = \frac{1}{2} e^{-|x-\theta|}$$

$$\rho_\theta(x) = \frac{1}{2^n} \exp \left\{ - \sum_{i=1}^n |x_i - \theta| \right\}$$

$$l(\theta; x) = \log \rho_\theta(x)$$

$$= - \sum_{i=1}^n |x_i - \theta| - n \log 2$$

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



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

Thm 3.11  $\Rightarrow$  order stats are minimal suff.