

# Sufficiency

## Outline

- 1) Sufficiency
- 2) Factorization Theorem
- 3) Examples
- 4) Minimal sufficiency

## Three models for coin flipping

Model 3  $X_{i,j} \stackrel{\text{ind.}}{\sim} \text{Bernoulli}(\theta_{i,t})$   $i=1, \dots, 48$   $j=1, \dots, n_i$   $\theta_{i,j} \downarrow \text{in } j$

Model 2  $X_{i+} \stackrel{\text{ind.}}{\sim} \text{Binom}(n_i, \theta_i)$   $X_{i+} = \sum_{j=1}^{n_i} X_{i,j}$

Model 1  $X_{++} \stackrel{\text{ind.}}{\sim} \text{Binom}(n, \theta)$   $X_{++} = \sum_{i=1}^{48} \sum_{j=1}^{n_i} X_{i,j}$

most assumptions  $\swarrow$   $\searrow$  fewest assumptions

These models are nested:  $\mathcal{P}_1 \subseteq \mathcal{P}_2 \subseteq \mathcal{P}_3$

Data keeps getting compressed too...  
are we losing anything by doing this?

Answer No.  $X_{++}$  is a sufficient statistic for  $\mathcal{P}_1$ , and  $(X_{1+}, \dots, X_{48+})$  is also sufficient for  $\mathcal{P}_2$

Def A statistic  $T(X)$  is any function of data  $X$

# Sufficiency

Def A statistic  $T(x)$  is sufficient for model  $\mathcal{P}$  if the conditional distribution of  $X | T(x)$  is the same for all  $P \in \mathcal{P}$

Check definition for  $T(x) = X_{++}$  in  $\mathcal{P}_1$ :

$$\begin{aligned} P_\theta(x) &= \prod_{i=1}^{48} \prod_{j=1}^{n_i} \theta^{x_{ij}} (1-\theta)^{1-x_{ij}} \\ &= \theta^{X_{++}} (1-\theta)^{n-X_{++}} \quad (\text{why no } \binom{n}{X_{++}}?) \end{aligned}$$

$$\begin{aligned} P_\theta(X=x | X_{++}=t) &= \frac{P_\theta(X=x, X_{++}=t)}{P_\theta(X_{++}=t)} \\ &= \frac{\mathbb{1}\{X_{++}=t\} \cdot \theta^t (1-\theta)^{n-t}}{\binom{n}{t} \theta^t (1-\theta)^{n-t}} \\ &= \mathbb{1}\{X_{++}=t\} / \binom{n}{t} \end{aligned}$$

Intuition Suppose we believe Model 1.

Big/small  $X_{++}$  more likely with big/small  $\theta$

But once we know  $X_{++} = 178,079$ ,  
all data sets  $X$  with that many same-side  
flips are equally likely, regardless of  $\theta$

Not true in Models 2 & 3  $\Rightarrow X_{++}$  no longer sufficient

# Factorization Theorem

Usually, we can recognize 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  $\mu$ .

$T(x)$  is sufficient iff there exist  $g_\theta(t)$ ,  $h(x) \geq 0$  with

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

Note we could absorb  $h$  into  $\mu$  as density

(define new base measure  $\nu$ ,  $\nu(A) = \int_A h(x) d\mu(x)$ )

$\Rightarrow \mathcal{P}$  has densities  $p_\theta(x) = g_\theta(T(x))$  wrt  $\nu$

Interp: after changing base measure,

density depends on  $x$  only through  $T(x)$

(Can't absorb  $g_\theta(T(x))$  into  $\mu$ : depends on  $\theta$ )

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

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

( $\Rightarrow$ ) Assume  $T(x)$  sufficient, let

$$g_\theta(t) = \mathbb{P}_\theta(T(X)=t)$$

$$h(x) = \mathbb{P}(X=x | T(X)=T(x))$$

$\leftarrow$  no dep. on  $\theta$

$$\begin{aligned} \Rightarrow g_\theta(T(x)) h(x) &= \mathbb{P}_\theta(T(X)=T(x) \text{ and } X=x) \\ &= \mathbb{P}_\theta(X=x) = p_\theta(x) \quad \square \end{aligned}$$

Proof similar for general densities

- careful about conditioning in cts spaces

## Examples

Ex. Normal location family

$$X_1, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\theta, 1) = \frac{1}{\sqrt{2\pi}} e^{-(x-\theta)^2/2}$$

$$p_{\theta}(x) = (2\pi)^{-n/2} \prod_{i=1}^n e^{-(x_i-\theta)^2/2}$$

$$= e^{\theta \sum x_i - n\theta^2/2} \cdot \frac{e^{-\sum x_i^2/2}}{(2\pi)^{n/2}}$$

(collect factors with  
no dep. on  $\theta$ )

$\Rightarrow \sum X_i$  is sufficient

Ex. Poisson family

$$X_1, \dots, X_n \stackrel{iid}{\sim} \text{Pois}(\theta) = \frac{\theta^x e^{-\theta}}{x!} \quad \text{for } x = 0, 1, 2, \dots$$

$$p_{\theta}(x) = \prod_{i=1}^n \frac{\theta^{x_i} e^{-\theta}}{x_i!} = \theta^{\sum x_i} e^{-n\theta} \cdot \frac{1}{\prod x_i!}$$

$\Rightarrow \sum X_i$  is sufficient

These two examples have something important in common!  
(next lecture)

Ex. Uniform location family

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

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

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

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

# Interpretations of Sufficiency

$X$  is informative about  $\theta$  only because its distribution depends on  $\theta$ .

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

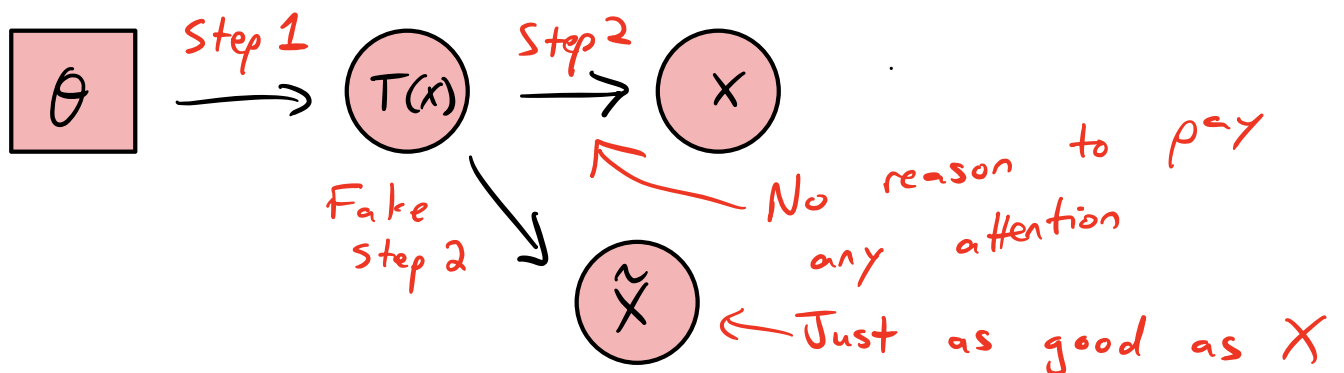
- 1) Generate  $T$  : distribution dep. on  $\theta$
- 2) Generate  $X|T$  : does not dep on  $\theta$

## 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  $\tilde{X} \sim P_\theta$

In graphical model form:



## Order Statistics

For  $x_1, \dots, x_n \in \mathbb{R}$ , define order statistics

$$\min_i x_i = X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)} = \max_i x_i$$

Ex (iid sampling on  $\mathbb{R}$ )  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P_\theta$ ,

any model  $\mathcal{P} = \{P_\theta^n : \theta \in \Theta\}$  on  $\mathcal{X} \subseteq \mathbb{R}$

$P_\theta^n$  invariant to perm.s of  $X = (X_1, \dots, X_n)$

$\Rightarrow$  All permutations of  $x$  are equally likely

$\Rightarrow$  Order statistics  $S(X) = (X_{(i)})_{i=1}^n$  sufficient

$X \rightsquigarrow S(X)$  forgets orig. ordering of observations



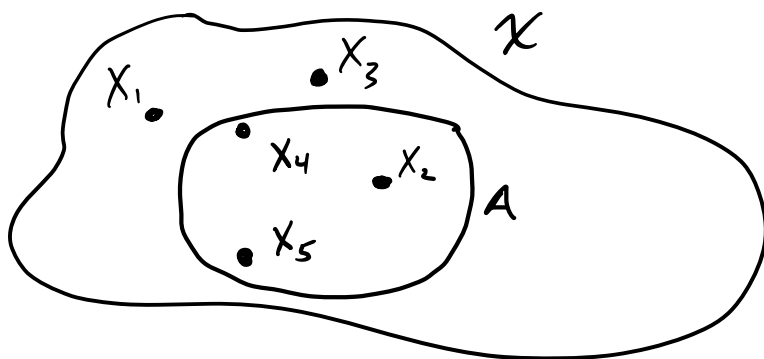
# Empirical Distribution

Order statistics depend on total ordering of  $\mathcal{X}$   
What about more general sample space?

Define Dirac measure  $\delta_x(A) = 1_{\{x \in A\}}$

Empirical distribution  $\hat{P}_n(\cdot) = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}(\cdot)$

random measure on  $\mathcal{X}$ , determined by sample



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

Ex (iid sampling)  $X_1, \dots, X_n \stackrel{iid}{\sim} P_\theta$

any model  $\mathcal{P} = \{P_\theta^n : \theta \in \Theta\}$  on any  $\mathcal{X}$

$\hat{P}_n$  is sufficient

$X \rightsquigarrow \hat{P}_n$  records which values observed,  
how many times

# Minimal Sufficiency

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

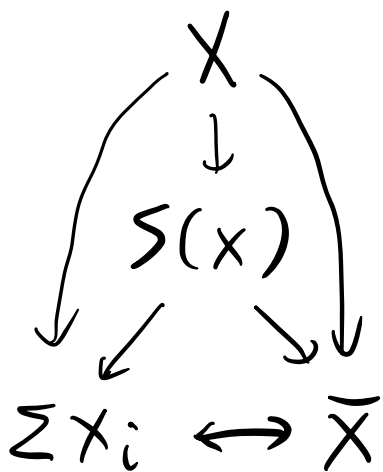
$$T(X) = \sum X_i \quad \text{sufficient}$$

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

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

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

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  $S(X)$ ,

$$T(X) = f(S(X)) \quad \text{for some } f \text{ (a.s. in } \mathcal{P})$$

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

## Recognizing minimal sufficiency

Assume  $\mathcal{P}$  has densities  $p_\theta$ , sample space  $\mathcal{X}$

Define equivalence relation on  $\mathcal{X}$ :

$$x \equiv_{\mathcal{P}} y \text{ if } \frac{p_\theta(x)}{p_\theta(y)} \text{ doesn't depend on } \theta$$

Note any sufficient statistic  $T$  can only collapse together equivalent values: if  $T(x) = T(y) = t$

$$\frac{p_\theta(x)}{p_\theta(y)} = \frac{\mathbb{P}_\theta(X=x, T(X)=t)}{\mathbb{P}_\theta(X=y, T(X)=t)} = \frac{\mathbb{P}(X=x \mid T(X)=t)}{\mathbb{P}(X=y \mid T(X)=t)}$$

So, for any sufficient stat  $T(x)$ ,  $T(x) = T(y) \Rightarrow x \equiv_{\mathcal{P}} y$

For minimal sufficient stats, the reverse implication also holds:

Theorem (Bahadur)  $T(X)$  is minimal sufficient if

$$X \equiv_{\mathcal{P}} Y \Leftrightarrow T(x) = T(y)$$

Interp: a minimal sufficient stat. collapses the sample space into exactly these equiv. classes.

Proof: First show  $T(X)$  sufficient:

For any  $x$  with  $T(x) = t$ , we have

$$P_{\theta}(X = x \mid T(X) = t) = \frac{p_{\theta}(x)}{\sum_{z: T(z) = t} p_{\theta}(z)} = \frac{1}{\sum_{z: T(z) = t} p_{\theta}(z) / p_{\theta}(x)}$$

which doesn't depend on  $\theta$  because

$$T(z) = t = T(x) \Rightarrow p_{\theta}(z) / p_{\theta}(x) \text{ doesn't depend on } \theta$$

Next assume  $S(X)$  sufficient. If  $S(x) = S(y) = s$

then  $x \equiv_{\mathcal{P}} y$  so  $T(x) = T(y)$ . Set  $f(s) = T(x)$ .

Any other  $z$  with  $S(z) = s$  has

$$z \equiv_{\mathcal{P}} x \Rightarrow T(z) = T(x) = f(s) = f(S(z)).$$

## (log-) Likelihood functions

### 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  $\theta$       data  $x$  determines which function      function of  $x$  with parameter  $\theta$

The log-likelihood function is its log:

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

Note if  $x \equiv_p y$  is same as saying

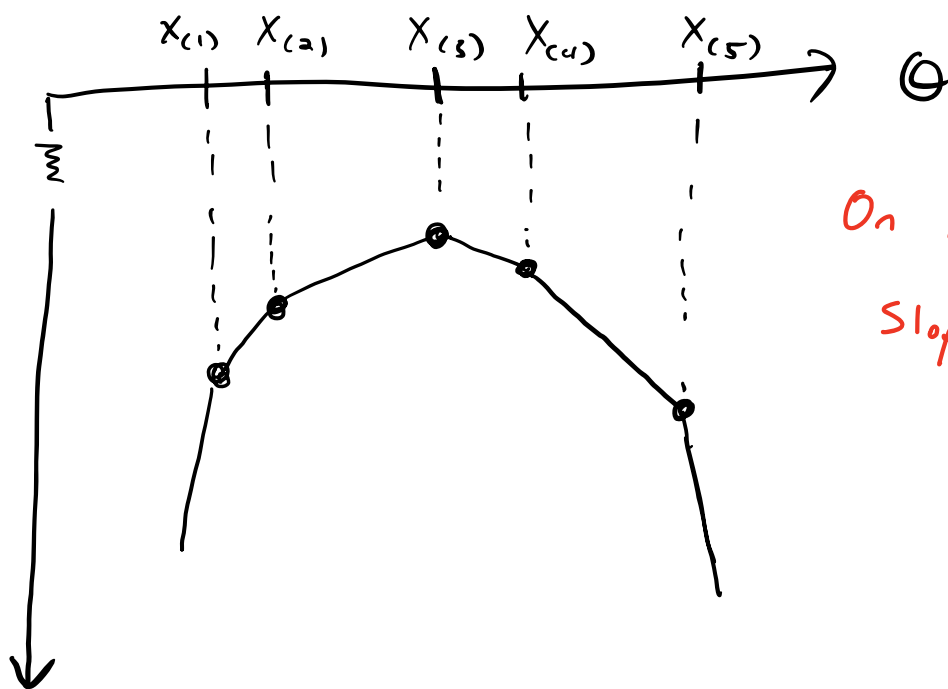
$$l(\theta; x) - l(\theta; y) = \frac{p_\theta(x)}{p_\theta(y)} \text{ is } \underline{\text{constant}}$$

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



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

$\Rightarrow$  order stats are minimal suff.