# Sufficiency

## Outline

1) Sufficiency

2) Factorization Theorem

3) Examples

4) Minimal sufficiency

# Three models for coin flipping

**Model 3** $\quad X_{i,j} \overset{ind}{\sim} \text{Bernoulli}(\theta_{i,t})$ $\quad \begin{matrix} i = 1, \dots, 48 \\ j = 1, \dots, n_i \end{matrix}$ $\quad \theta_{i,j} \searrow \text{in } j$

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

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

<span style="color:red">most assumptions ↙</span> $\qquad$ <span style="color:red">fewest assumptions ↙</span>

These models are <u>nested</u>: $\quad \mathcal{P}_1 \subseteq \mathcal{P}_2 \subseteq \mathcal{P}_3$

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

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

<u>Def</u> A <u>statistic</u> $T(x)$ is any function of data $X$

# Sufficiency

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

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

$$P_\theta(x) = \prod_{i=1}^{48} \prod_{j=1}^{n_i} \theta^{X_{ij}}(1-\theta)^{1-X_{ij}}$$

$$= \theta^{X_{++}}(1-\theta)^{n-X_{++}} \qquad \color{red}{\left(\text{why no } \binom{n}{X_{++}}?\right)}$$

$$\mathbb{P}_\theta(X = x \mid X_{++} = t) = \frac{\mathbb{P}_\theta(X = x, X_{++} = t)}{\mathbb{P}_\theta(X_{++} = t)}$$

$$= \frac{1\{x_{++} = t\} \cdot \theta^t(1-\theta)^{n-t}}{\binom{n}{t}\theta^t(1-\theta)^{n-t}}$$

$$= 1\{x_{++} = t\} / \binom{n}{t}$$

__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 $\mathcal{S} = \{P_\theta : \theta \in \Theta\}$ 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

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

Note we could absorb $h$ into $\mu$ as density
(define new base measure $\nu$, $\nu(A) = \int_A h(x)\,d\mu(x)$)

$\Rightarrow \mathcal{S}$ 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 $\mathcal{X}$

$(\Longleftarrow)$ $\mathbb{P}_\theta(X = x \mid T = t) = \dfrac{\mathbb{P}_\theta(X = x, T(x) = t)}{\mathbb{P}_\theta(T(x) = t)}$

$$= \frac{\cancel{g_\theta(t)}\, h(x)\, \mathbb{1}\{T(x) = t\}}{\displaystyle\sum_{T(z) = t} \cancel{g_\theta(t)}\, h(z)}$$

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

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

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

$\underset{\color{red}{\text{no dep. on } \theta}}{\color{red}\nwarrow}$

$\Longrightarrow g_\theta(T(x))\, h(x) = \mathbb{P}_\theta\big(T(X) = T(x) \text{ and } X = x\big)$

$$= \mathbb{P}_\theta(X = x) = p_\theta(x) \qquad \boxtimes$$

Proof similar for general densities
requires care about conditioning

# Examples

**Ex.** Uniform location family

$$X_1, \ldots, X_n \overset{iid}{\sim} U[\theta, \theta+1]$$
$$= 1\{\theta \le x \le \theta+1\}$$

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

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

$$\implies (X_{(1)}, X_{(n)}) \quad \text{is} \quad \text{sufficient.}$$

**Ex.** Normal location family

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

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

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

<span style="color:red">(collect factors with no dep. on $\theta$)</span>

$$\implies \sum X_i \quad \text{is} \quad \text{sufficient}$$

**Ex.** Poisson family

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

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

$$\implies \sum X_i \quad \text{is} \quad \text{sufficient}$$

# 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 $\Theta$
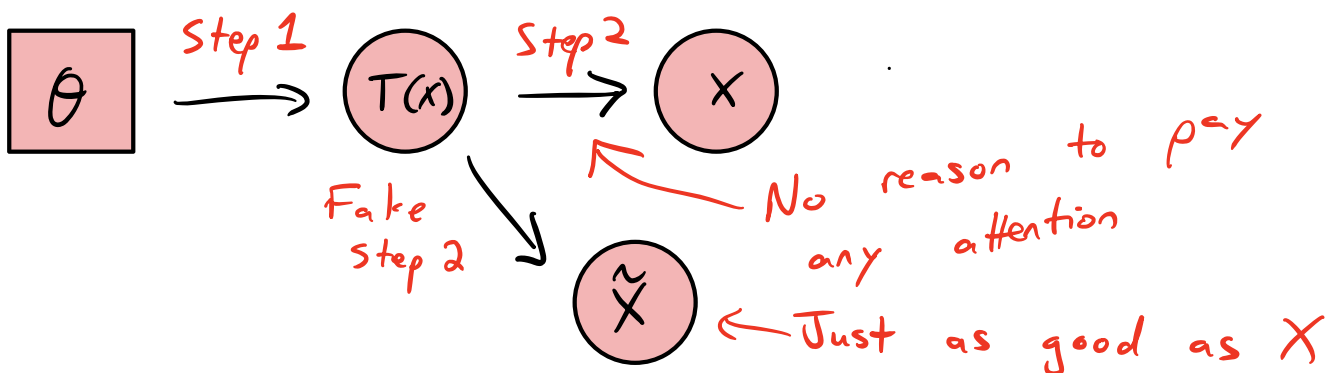2) Generate $X | T$ : does __not__ dep on $\Theta$

## Sufficiency Principle

If $T(x)$ is sufficient for $P$ 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
$\underset{\color{red}\leftarrow \text{no } \Theta}{}$
be just as good as $X$ since $\tilde{X} \sim P_\Theta$

In graphical model form:



Step 1 $\quad$ Step 2

$\Theta \longrightarrow T(x) \longrightarrow X$

Fake Step 2

No reason to pay any attention

$\tilde{X}$ ← Just as good as $X$

# Order Statistics

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

$$\min_i x_i = x_{(1)} \leq x_{(2)} \leq \cdots \leq x_{(n)} = \max_i x_i$$

**Ex** ( iid sampling on $\mathbb{R}$ ) $X_1, \ldots, X_n \overset{iid}{\sim} P_\theta$,

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

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

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

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

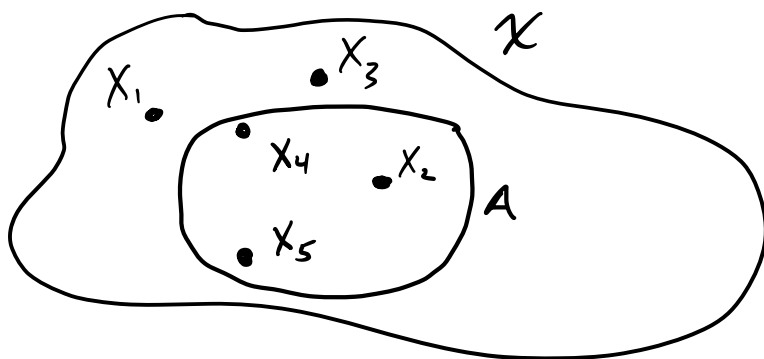$X \rightsquigarrow S(X)$ forgets orig. ordering of observations

# 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) $X_1, ..., X_n \overset{iid}{\sim} P_\theta$

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

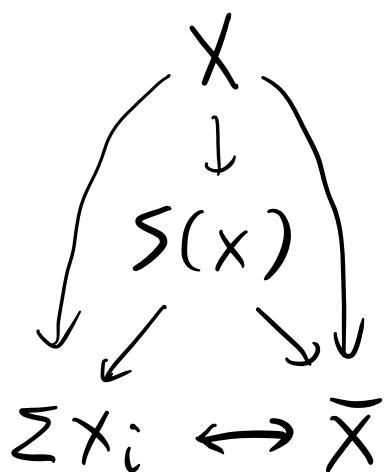Consider $X_1, \ldots, X_n \overset{iid}{\sim} N(\theta, 1)$

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

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

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

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

Which can be recovered from which others?

$$X$$
$$\downarrow$$
$$S(x)$$
$$\sum X_i \longleftrightarrow \bar{X}$$

← 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:** $p_\theta(x) = g_\theta(T(x)) h(x)$
$$= (g_\theta \circ f)(S(x)) h(x) \qquad \boxtimes$$

**Definition:** $T(X)$ is <u>minimal sufficient</u> if

1) $T(x)$ is sufficient

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

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

# 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

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

<span style="color:red">function of $x$ with parameter $\theta$</span>

<span style="color:red">function of $\theta$</span>

<span style="color:red">data $X$ determines which function</span>

The log-likelihood function is its log:

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

The likelihood up to scaling (or $\ell$ up to vertical shift) is a minimal sufficient statistic

If $T(X)$ is sufficient then

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

# 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

"$Lik(\cdot\,;x) \propto Lik(\cdot\,;y) \implies T(x) = T(y)$"

equivalently.

"$\ell(\cdot\,;x) - \ell(\cdot\,;y) = const(x,y) \implies T(x) = T(y)$"

# Ex    Laplace location family

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

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

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



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

Slope $= n - 2k$

$$\ell(\theta; x) = \ell(\theta; Y) + \text{const} \iff X, Y \text{ same order statistics}$$

$\implies$ order stats are minimal suff!

# Minimal sufficiency for exp. fam.s

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

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

Is $T(X)$ minimal? <span style="color:red">(always sufficient)</span>

Suppose $x$ and $y$ give same likelihood shape:

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

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

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

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

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

<span style="color:red">(That is, if $\Xi$ is not contained in a lower-dim affine space)</span>

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

<span style="color:red">[Can we conclude $T(X)$ is <u>not</u> minimal?]</span>

Other parameterizations:

$$\rho_\Theta(x) = e^{\eta(\Theta)'T(x) - \beta(\Theta)} h(x) \qquad \Theta \in \Theta$$

$T(x)$ minimal if $\text{span}\{\eta(\Theta_1) - \eta(\Theta_2) : \Theta_1, \Theta_2 \in \Theta\} = \mathbb{R}^S$



$\eta_2$

$T(x)$ minimal

(B)

(A)

$T(x)$ minimal

(C)

$\gamma$

$\gamma' T(x)$ is sufficient
$\Rightarrow T(x)$ prob. not minimal

$\eta_1$