

# STA261 - Module 1

Statistics

Rob Zimmerman

University of Toronto

July 5-7, 2022

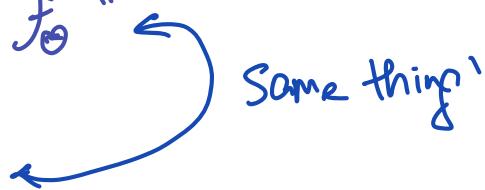
# Data and samples

- *Data* is factual information collected for the purposes of inference (Merriam-Webster)
- *Inference* is the act of passing from statistical sample data to generalizations (as of the value of population parameters) usually with calculated degrees of certainty (also Merriam-Webster)
- We collect a *sample* of data from a *population* associated with some probability distribution, and we would like to infer unknown properties of that distribution
- Example 1.1:
  - height of typical UofT student  $\sim N(\mu, \sigma^2)$  for some  $\mu \in \mathbb{R}$   
assume  $\sigma^2 > 0$
  - # of courses STA261 students are taking  $\sim \text{Poisson}(\lambda)$ ? for some  $\lambda > 0$ .

# Random variables versus observed data (this is really important)

- Our data sample goes through two phases of life: first as a *random sample*, and then as *observed data*
- A random sample is a set of *random variables*; observed data is a set of *constants*; the same goes for functions thereof  
 $Z \sim N(0, 1) \Rightarrow P(Z > 0) = \frac{1}{2}$   
 $Z \in \mathbb{R} \Rightarrow P(Z > 0) \in [0, 1]$   
 $Z = 2.61 \Rightarrow P(Z > 0) = 1.$
- We denote random variables using uppercase letters, and constants using lowercase letters:
- Example 1.2:  
 $\vec{X}_n = (X_1, X_2, \dots, X_n)$  random sample  
 $\vec{x}_n = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$  observed data
- It is **very** important to clearly distinguish between the two quantities. But why?

# iid-ness

- “iid” stands for “**independent and identically distributed**”
- This term is used everywhere in statistics, because it saves a lot of time
- Instead of “Let  $X_1, X_2, \dots, X_n$  be a random sample from some distribution with pdf/pmf  $f_0$ ”  
we write “let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f_0$ ”  
  
Ex: “Let  $Z_1, Z_2, \dots, Z_n \stackrel{\text{iid}}{\sim} N(0,1)$ ”

~~$\vec{X} \stackrel{iid}{\sim} f_\theta$~~       Ignore! Not worth the confusion.  
 $\Rightarrow X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta$   
 $\Rightarrow$  Density of  $\vec{X}$  is  $f(\vec{x}) = \prod_{i=1}^n f_\theta(x_i)$

A slight abuse of notation: we'll usually write  $f_\theta(\vec{x})$  for the joint pdf/pmf of  $\vec{X} = (X_1, \dots, X_n)$  and  $f_\theta(x_i)$  for the pdf/pmf of any  $X_i$ , in the iid case.

So, for example, if  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$ , then

$$f_\lambda(\vec{x}) = \prod_{i=1}^n f_\lambda(x_i) = \frac{e^{-n\lambda} \lambda^{\sum x_i}}{\prod_{i=1}^n x_i!}$$

# Statistics

$\bar{X}$

- **Definition 1.1:** A **statistic** is a function of the (random) data sample which is free of any unknown constants

$$T(\vec{X}) = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i; \quad T(\vec{X}) = 261$$

- **Example 1.3:**  $T(\vec{X}) = 24 \cdot X_3^2$   $T(\vec{X}) = X_{\text{max}}$  (sample maximum)

Say  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, 1)$ ,  $\mu \in \mathbb{R}$ . Let  $T(\vec{X}) = \frac{1}{n} \sum_{i=1}^n X_i - \mu$ .

- A statistic is useful when it allows us to summarize the data sample in ways that helps us with inference

NOT A STATISTIC!!  
 $\mu$  is unknown.

- Different statistics are useful for different models

- **Example 1.4:**  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ ,  $\mu \in \mathbb{R}$ ,  $\sigma^2 > 0$ .

Intuitively,  $T(\vec{X}) = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  is useful if we're interested in  $\mu$ .

- Ex:  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$ ,  $p \in (0, 1)$ .

Then  $T(\vec{X}) = \bar{X}_n$  is also useful for understanding  $p$ .

# Parameters and Statistical Models

- Many classical probability distributions have *parameters* associated with them

- Example 1.5:  $N(\mu, \sigma^2)$        $\text{Exp}(x)$        $\text{Bin}(n, p)$   


- Definition 1.2: A **statistical model** is a set of pdfs/pdfs  $\{f_\theta(\cdot) : \theta \in \Theta\}$  defined on the same sample space, where each  $\theta$  is a fixed **parameter** in a known **parameter space**  $\Theta$ . When  $\Theta \subseteq \mathbb{R}^k$  for some  $k \in \mathbb{N}$ , the set is also called a **parametric model** (or **parametric family**).  


- Example 1.6:  $\left\{ \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2}\right) : \mu \in \mathbb{R} \right\} = \{ N(x, 1) : \mu \in \mathbb{R} \}$   
 $\left\{ \lambda e^{-\lambda x} : \lambda > 0 \right\} = \{ \text{Exp}(\lambda) : \lambda > 0 \}$   
  
 $\Theta = \{\theta_1, \theta_2\}$   
 $\Rightarrow \{f_\theta : \Theta \in \Theta\} = \{f_{\theta_1}, f_{\theta_2}\}$

- Statistical inference is classically concerned with figuring out which one of those distributions generated the data, based on the data sample we have available
- This amounts to inferring the particular parameter  $\theta$

# Parameters and Statistical Models: More Examples

- Example 1.7:

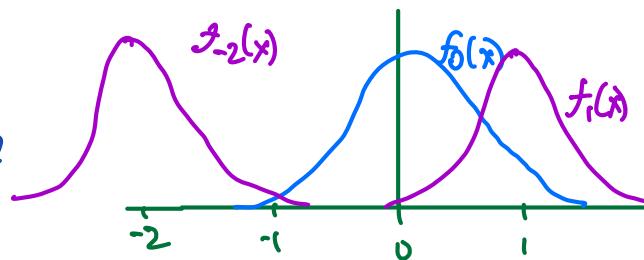
$$\begin{aligned} & \left\{ N(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0 \right\} \\ = & \left\{ N(\mu, \sigma^2) : (\mu, \sigma^2) \in \mathbb{R} \times (0, \infty) \right\} \end{aligned}$$

Maybe instead,  $\left\{ N(\mu, \sigma^2) : \mu \in (0, \infty), \sigma^2 \in (0, \infty) \right\}$

if we know in advance that  $\mu > 0$ .

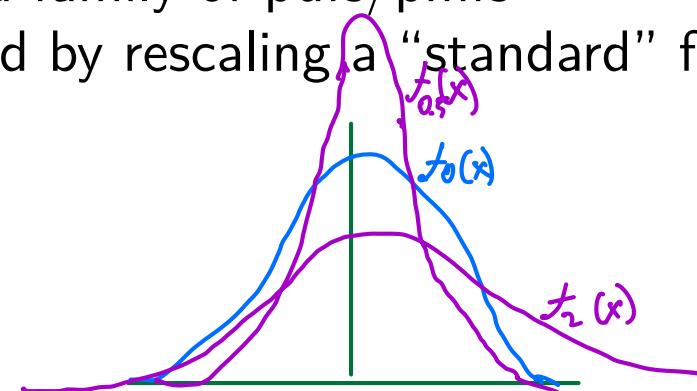
# Important Parametric Families: Location-Scale Families

- **Definition 1.3:** A **location family** is a family of pdfs/pdfs  $\{f_\mu(\cdot) = f_0(\cdot - \mu) : \mu \in \mathbb{R}\}$  formed by translating a “standard” family member  $f_0(\cdot)$ .



- **Example 1.8:**  $\{N(\mu, 1) : \mu \in \mathbb{R}\}$

- **Definition 1.4:** A **scale family** is a family of pdfs/pdfs  $\{f_\sigma(\cdot) = f_1(\cdot/\sigma)/\sigma : \sigma > 0\}$  formed by rescaling a “standard” family member  $f_1(\cdot)$ .

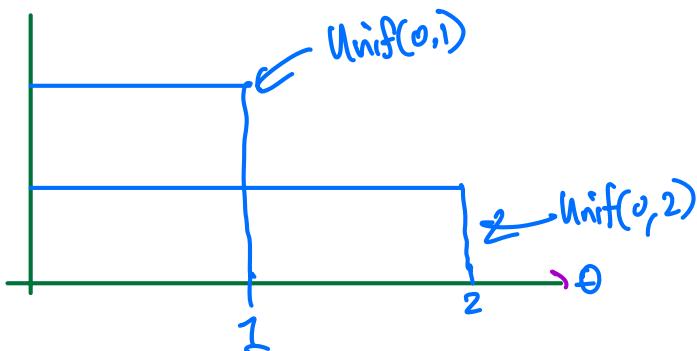


- **Example 1.9:**  $\{N(0, \sigma^2) : \sigma^2 > 0\}$

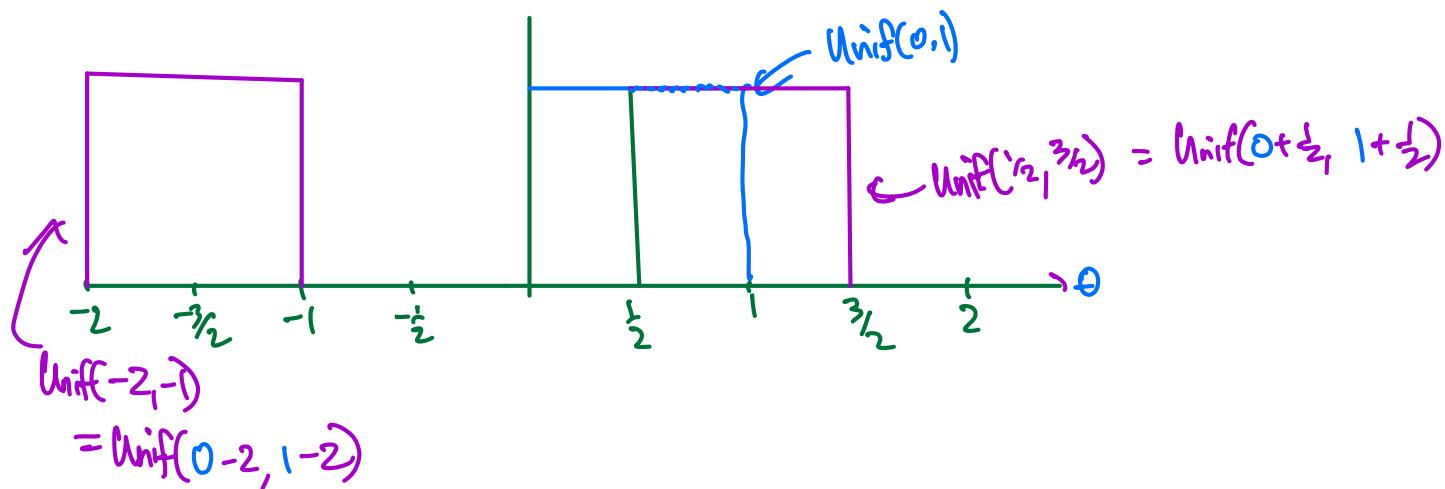
- **Definition 1.5:** A **location-scale family** is a family of pdfs/pdfs  $\{f_{\mu, \sigma}(\cdot) = f_{0,1}\left(\frac{\cdot - \mu}{\sigma}\right) / \sigma : \mu \in \mathbb{R}, \sigma > 0\}$  formed by translating and rescaling a “standard” family member  $f_{0,1}(\cdot)$ .

- **Example 1.10:**  $\{N(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0\}$

# Poll Time!



Not a location family!



# Important Parametric Families: Exponential Families

- **Definition 1.6:** An **exponential family** is a parametric family of pdfs/pdfs of the form

$$f_{\theta}(x) = h(x) \cdot g(\theta) \cdot \exp \left( \sum_{j=1}^k w_j(\theta) \cdot T_j(x) \right),$$

for some  $k \in \mathbb{N}$ , where all functions of  $x$  and  $\theta$  are known.

Usually  $k=1$  when  $\theta \in \mathbb{R}$ ,  
whence  $f_{\theta}(x) = h(x) \cdot g(\theta) \cdot \exp(w(\theta) \cdot T(x))$

- Lots of theory simplifies considerably if we assume our random sample comes from an exponential family

- Many of your favourite distributions are included

- Example 1.11:

$$X \sim \text{Exp}(\lambda)$$

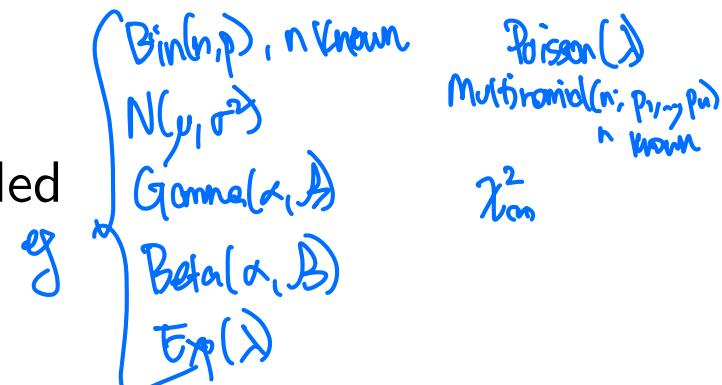
$$f_{\lambda}(x) = \lambda e^{-\lambda x}$$

$$= 1 \cdot \lambda \cdot \exp(-\lambda \cdot x)$$

$\uparrow h(x)$        $\uparrow g(\lambda)$        $\uparrow w(\lambda)$        $\uparrow T(x)$

$$X \sim \text{Bernoulli}(\theta), \theta \in (0,1)$$

$$\begin{aligned} f_{\theta}(x) &= \theta^x (1-\theta)^{1-x} \\ &= (1-\theta) \left( \frac{\theta}{1-\theta} \right)^x \\ &= (1-\theta) \cdot \exp \left( \lg \left( \frac{\theta}{1-\theta} \right)^x \right) \end{aligned}$$



$$x = \exp(\lg(\theta)) \quad \forall x > 0$$

$$X \sim \text{Bernoulli}(\theta), \theta \in (0,1)$$

$$\begin{aligned} f_{\theta}(x) &= \theta^x (1-\theta)^{1-x} \\ &= (1-\theta) \left( \frac{\theta}{1-\theta} \right)^x \\ &= (1-\theta) \cdot \exp \left( \lg \left( \frac{\theta}{1-\theta} \right)^x \right) \end{aligned}$$

$\uparrow h(x)$        $\uparrow g(\theta)$        $\uparrow T(x)$        $\uparrow w(\theta)$

# A Quick Review of Conditional Distributions

- $X|Y$  is a random variable, which has its own distribution called a conditional distribution

- Remember Bayes' rule:  $P(A|B) = \frac{P(A \cap B)}{P(B)}$ ,  $P(B) \neq 0$

- $X | Y = y \quad f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$

- 
- $X | X = x$  degenerate/constant at  $x$

- Example 1.12:  $\mathbb{E}[X|Y]$  is a random variable (a function of  $Y$ )

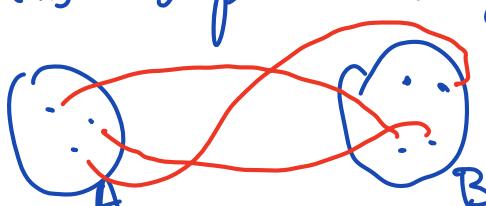
- Example 1.13:  $\mathbb{E}[X|X] = X$ . If  $X \perp\!\!\!\perp Y$ , then  $\mathbb{E}[X|Y] = \mathbb{E}[X]$

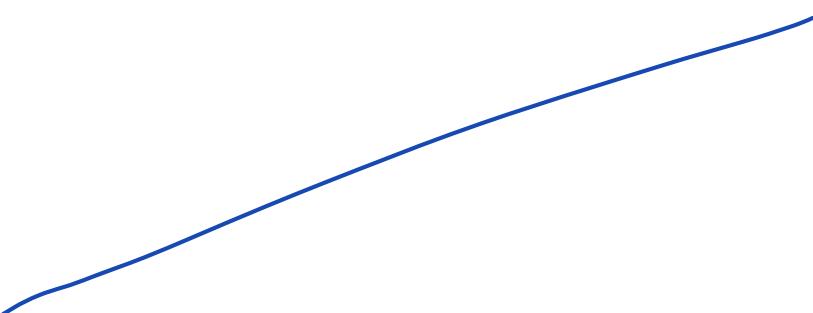
"Tower property":  $\mathbb{E}[\mathbb{E}[X|Y]] = \mathbb{E}[X]$

"Conditional variance":  $\text{Var}(X) = \mathbb{E}[\text{Var}(X|Y)] + \text{Var}(\mathbb{E}[X|Y])$

EXERCISE: Prove these,  
if you haven't already!

# A Quick Review of Functions

- Let  $f : A \rightarrow B$  be a function
- If  $f$  is one-to-one, then "injective"
$$f(x) = f(y) \Leftrightarrow x = y$$

- If  $f$  is onto, then "surjective"
$$\forall b \in B, \exists a \in A \text{ s.t. } f(a) = b$$
- If  $f$  is a bijection, then  $f$  is one-to-one AND onto.
- Example 1.14:



# Freedom From $\theta$

- Most of the functions  $f_\theta(x)$  we will deal with have parameters involved in addition to the “independent variable”
- If the parameter  $\theta$  can vary too, then  $f_\theta(x)$  is really a function of both  $x$  and  $\theta$   
*i.e., there exists  $g: \Theta \times \mathbb{R} \rightarrow \mathbb{R}$  s.t.  $f_\theta(x) = g(\theta, x)$*
- If  $f_\theta(x)$  is actually *not* a function of  $\theta$  (i.e., it's constant with respect to  $\theta$ ), we might also say that it's “free of  $\theta$ ” or that it “does not depend on  $\theta$ ”
- Example 1.15:  $f_\theta(x) = x^2$  is free of  $\theta$   
 $f_\theta(x) = \theta x^2$  is not free of  $\theta$  Say  $X \sim N(\mu, 1)$ . Then  $P_\mu(X - \mu \leq x) = \Phi(x)$  is free of  $\mu$ .
- So if we say that the distribution of  $X$  is free of  $\theta$ , we mean that the cdf of  $X$  (and hence the pdf/pmf) is the same for all  $\theta \in \Theta$
- Example 1.16: If  $X \sim \text{Exp}(\lambda)$ , then the distribution of  $\lambda X$  is free of  $\lambda$   
*CHECK THIS!*

# Data Reduction: A Thought Experiment

- Is there a such thing as “more data than necessary”?
- Suppose that field researchers collect a sample  $\mathbf{X} = (X_1, X_2, \dots, X_n) \stackrel{iid}{\sim} f_\theta$ , where  $n$  is astronomically large; they want us statisticians to do inference on  $\theta$ , but sending us  $\mathbf{X}$  would take weeks
- Wouldn’t it be great if we didn’t need the entire sample  $\mathbf{X}$  to make inferences about  $\theta$ , but rather a much smaller statistic  $T(\mathbf{X})$  – perhaps just a single number – that still contained as much information about  $\theta$  as  $\mathbf{X}$  itself did?
- The researchers observe  $\mathbf{X} = \mathbf{x}$ , calculate  $T(\mathbf{x}) = t$  on their end, and then text  $t$  over to us
- Example 1.17:  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, 1)$ ,  $\mu \in \mathbb{R}$ .

Instead of sending us  $X_1, \dots, X_n$ , what if we only got  $\bar{x}_n = \frac{1}{n} \sum x_i$ ?  
Can we still make inferences about  $\mu$ ?

# Sufficiency

- How do we “encode” this idea?
- If we know that  $T(\mathbf{X}) = t$ , then there should be nothing else to glean from the data about  $\theta$
- Definition 1.7: A statistic  $T(\mathbf{X})$  is a **sufficient statistic** for a parameter  $\theta$  if the conditional distribution of  $\mathbf{X} | T(\mathbf{X}) = t$  does not depend on  $\theta$ .
- An interpretation: if the conditional distribution

$$\mathbb{P}(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})) = \frac{\mathbb{P}_\theta (\mathbf{X} = \mathbf{x} \text{ and } T(\mathbf{X}) = T(\mathbf{x}))}{\mathbb{P}_\theta (T(\mathbf{X}) = T(\mathbf{x}))}$$

is really free of  $\theta$ , then the information about  $\theta$  in  $\mathbf{X}$  and the information about  $\theta$  in  $T(\mathbf{X})$  are “equal”

- Example 1.18:  $T(\vec{X}) = \vec{X}$  is sufficient:

$$\mathbb{P}(\vec{X} = \vec{x} | T(\vec{X}) = \vec{x}) = \mathbb{P}_\theta(\vec{X} = \vec{x} | \vec{X} = \vec{x}) = 1$$

But no data reduction!

$$= \frac{\mathbb{P}_\theta(\vec{X} = \vec{x})}{\mathbb{P}_\theta(T(\vec{X}) = T(\vec{x}))}$$

# Sufficiency

- **Example 1.19:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$ , where  $\theta \in (0, 1)$ . Show that  $T(\mathbf{X}) = \sum_{i=1}^n X_i$  is sufficient for  $\theta$ .

$$T(\vec{X}) \sim \text{Bin}(n, \theta) \Rightarrow P(T(\vec{X}) = t) = \binom{n}{t} \theta^t (1-\theta)^{n-t}.$$

Then  $P_\theta(\vec{X} = \vec{x} | T(\vec{X}) = t)$

$$\begin{aligned}&= P_\theta(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n, \sum_i X_i = t) \\&= P_\theta(X_1 = x_1, \dots, X_n = x_n, \sum_{i=1}^n X_i = t) \\&= P_\theta(X_1 = x_1, \dots, X_n = t - \sum_{i=1}^{n-t} x_i) \\&= P_\theta(X_1 = x_1) \cdot \dots \cdot P_\theta(X_n = t - \sum_{i=1}^{n-t} x_i) \quad \text{by independence} \\&= \theta^{x_1} (1-\theta)^{1-x_1} \cdots \theta^{t - \sum_{i=1}^{n-t} x_i} (1-\theta)^{1-t + \sum_{i=1}^{n-t} x_i} \\&= \theta^t (1-\theta)^{n-t}\end{aligned}$$

$$\text{Therefore, } P_\theta(\vec{X} = \vec{x} | T(\vec{X}) = t) = \frac{\theta^t (1-\theta)^{n-t}}{\binom{n}{t} \theta^t (1-\theta)^{n-t}} = \frac{1}{\binom{n}{t}} \quad \text{which is free of } \theta.$$

$\therefore T(\vec{X})$  is sufficient for  $\theta$ .

# Sufficiency

- **Example 1.20:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$ , where  $\mu \in \mathbb{R}$  and  $\sigma^2$  is known. Show that the sample mean  $T(\mathbf{X}) = \bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$  is sufficient for  $\mu$ .

$$T(\bar{X}) \sim N(\mu, \sigma^2/n) \text{ so has density } f_T(t) = \frac{\sqrt{n}}{\sqrt{2\pi}\sigma} \exp\left(-\frac{n(t-\mu)^2}{2\sigma^2}\right).$$

For the numerator, observe that for  $t = \frac{1}{n} \sum_i x_i$ :

$$\begin{aligned}\sum(x_i - \mu)^2 &= \sum(x_i - t + t - \mu)^2 \\ &= \sum[(x_i - t)^2 + 2(x_i - t)(t - \mu) + (t - \mu)^2] \\ &= \sum(x_i - t)^2 + n(t - \mu)^2\end{aligned}$$

$$\begin{aligned}\text{Therefore, } f_{\bar{X}}(\bar{x}) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-n/2} \cdot \exp\left(-\frac{\sum(x_i - \mu)^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-n/2} \cdot \exp\left(-\frac{\sum(x_i - t)^2 + n(t - \mu)^2}{2\sigma^2}\right) \quad \text{from above} \\ &= f_{\bar{X}, \Pi}(\bar{x}, t)\end{aligned}$$

Therefore,

$$\begin{aligned}f_{\bar{X}, \Pi}(\bar{x}|t) &= \frac{(2\pi\sigma^2)^{-n/2} \cdot \exp\left(-\frac{\sum(x_i - t)^2 + n(t - \mu)^2}{2\sigma^2}\right)}{\frac{\sqrt{n}}{\sqrt{2\pi}\sigma} \exp\left(-\frac{n(t - \mu)^2}{2\sigma^2}\right)} \\ &= \frac{1}{\sqrt{n}(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{\sum(x_i - t)^2}{2\sigma^2}\right)\end{aligned}$$

which is free of  $\mu$ .

$\therefore T(\bar{X})$  is sufficient for  $\mu$ .

# The Factorization Theorem

- **Theorem 1.1 (Factorization theorem):** Let  $\mathbf{X} = (X_1, \dots, X_n) \sim f_\theta(\mathbf{x})$ , where  $f_\theta(\mathbf{x})$  is a joint pdf/pmf. A statistic  $T(\mathbf{X})$  is sufficient for  $\theta$  if and only if there exist functions  $g_\theta(t)$  and  $h(\mathbf{x})$  such that

$$f_\theta(\mathbf{x}) = h(\mathbf{x}) \cdot g_\theta(T(\mathbf{x})) \quad \text{for all } \theta \in \Theta,$$

where  $h(\mathbf{x})$  is free of  $\theta$  and  $g_\theta(T(\mathbf{x}))$  only depends on  $\mathbf{x}$  through  $T(\mathbf{x})$ .

- In other words,  $T(\mathbf{X})$  is sufficient whenever the “part” of  $f_\theta(\mathbf{x})$  that actually depends on  $\theta$  is a function of  $T(\mathbf{x})$ , rather than  $\mathbf{x}$  itself

*Proof. (Discrete case).* Want to show  $\frac{P_\theta(\vec{X}=\vec{x} \wedge T(\vec{X})=t)}{P_\theta(T(\vec{X})=t)}$  is free of  $\theta$

↑  
Continuous case  
requires measure theory

iff  $P_\theta(\vec{X}=\vec{x}) = h(\vec{x}) \cdot g_\theta(t)$ .

# The Factorization Theorem

( $\Rightarrow$ ) Assume  $T$  is sufficient for  $\Theta$ .

$$\begin{aligned} \text{Then } P_\theta(\vec{X} = \vec{x}) &= P_\theta(\vec{X} = \vec{x} \wedge T(\vec{X}) = t) \\ &= P_\theta(\vec{X} = \vec{x} \mid T(\vec{X}) = t) \cdot P_\theta(T(\vec{X}) = t) \quad = h(\vec{x}) \cdot g_\theta(t). \\ &\qquad\qquad\qquad \underbrace{\qquad\qquad}_{=: h(\vec{x})} \qquad\qquad\qquad \underbrace{\qquad\qquad}_{=: g_\theta(t)} \\ &\qquad\qquad\qquad \text{because } T \text{ is sufficient} \end{aligned}$$

( $\Leftarrow$ ) Assume  $P_\theta(\vec{X} = \vec{x}) = h(\vec{x}) \cdot g_\theta(t)$  for some  $h, g_\theta$ .

$$\begin{aligned} \text{Then the pmf of } T \text{ is } P(T(\vec{X}) = t) &= \sum_{\vec{x}: T(\vec{x})=t} P_\theta(\vec{X} = \vec{x}, T(\vec{X}) = t) \\ &= \sum_{\vec{x}: T(\vec{x})=t} P_\theta(\vec{X} = \vec{x}) \\ &= \sum_{\vec{x}: T(\vec{x})=t} h(\vec{x}) \cdot g_\theta(t) \text{ by ccs.} \\ &= \left( \sum_{\vec{x}: T(\vec{x})=t} h(\vec{x}) \right) g_\theta(t) \end{aligned}$$

$$\begin{aligned} \text{Therefore, } P_\theta(\vec{X} = \vec{x} \mid T(\vec{X}) = t) &= \frac{P_\theta(\vec{X} = \vec{x})}{P_\theta(T(\vec{X}) = t)} = \frac{h(\vec{x}) \cdot g_\theta(t)}{\left( \sum_{\vec{x}: T(\vec{x})=t} h(\vec{x}) \right) g_\theta(t)} = \frac{h(\vec{x})}{\sum_{\vec{x}: T(\vec{x})=t} h(\vec{x})} \text{ which is free of } \Theta. \quad \therefore T(\vec{X}) \text{ is sufficient} \\ &\qquad\qquad\qquad \text{for } \Theta. \quad \square \end{aligned}$$

# Poll Time!

From the proof of the FT,

$$f_\theta(x) = P_\theta(X=x) = P_\theta(X=x \mid T(X)=T(x)) \cdot P_\theta(T(X)=T(x))$$

$\underbrace{P_\theta(X=x \mid T(X)=T(x))}_{h(x)}$        $\underbrace{P_\theta(T(X)=T(x))}_{g_\theta(x)}$

$$\begin{aligned} \text{But } h(x) &= P(X=x \mid T(X)=T(x)) \\ &= P(X=x \mid X=x) = 1. \end{aligned}$$

So  $h(x)$  must be constant wrt  $x$ , but it can be any constant  $c \neq 0$   
as long as  $g_\theta(x) = \frac{1}{c} \cdot f_\theta(x)$ .

# The Factorization Theorem: Examples

- **Example 1.21:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$ , where  $\theta \in (0, 1)$ . Show that  $T(\mathbf{X}) = \sum_{i=1}^n X_i$  is sufficient for  $\theta$ . *Let  $t = \sum x_i$*

$$\begin{aligned} f_\theta(\vec{x}) &= \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} \\ &= \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} \\ &= 1 \cdot \theta^t (1-\theta)^{n-t} \end{aligned}$$

*↑*  
 $h(\vec{x})$   
 $g_\theta(t)$

. By the Factorization Theorem,  $T(\vec{x})$  is sufficient for  $\theta$ .

- **Example 1.22:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$ , where  $\mu \in \mathbb{R}$  and  $\sigma^2$  is known. Show that the sample mean  $T(\mathbf{X}) = \bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$  is sufficient for  $\mu$ . *Let  $t = \frac{1}{n} \sum x_i$*

$$\begin{aligned} \text{Then } f_\mu(\vec{x}) &= \prod_{i=1}^n f_\mu(x_i) \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{\sum(x_i - \mu)^2}{2\sigma^2}\right) \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{\sum(x_i - t)^2 + n(t - \mu)^2}{2\sigma^2}\right) = \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{\sum(x_i - t)^2}{2\sigma^2}\right) \cdot \exp\left(-\frac{n(t - \mu)^2}{2\sigma^2}\right) \end{aligned}$$

$\therefore h(\vec{x})$   
 $\therefore g_\mu(t)$

By the FT,  $T(\vec{x})$  is sufficient for  $\mu$ .

# The Factorization Theorem: Examples

- **Example 1.23:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$ , where  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ . The **sample variance** is the statistic  $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ . Show that  $T(\mathbf{X}) = (\bar{X}_n, S_n^2)$  is sufficient for  $(\mu, \sigma^2)$ .

$$f_{\mu, \sigma^2}(\vec{x}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{\sum(x_i - \bar{x})^2 + n(\bar{x} - \mu)^2}{2\sigma^2}\right)$$

$$= 1 \cdot \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{(n-1)t_2 - n(t_1 - \mu)^2}{2\sigma^2}\right)$$

$\uparrow h(\vec{x})$ 
 $\underbrace{\quad}_{S_{\mu, \sigma^2}(t_1, t_2)}$ 
By the FT,  $T(\vec{x})$  is sufficient for  $(\mu, \sigma^2)$ .

$$\stackrel{=} {\frac{1}{n-1} \sum (x_i - \bar{x})^2}$$

- **Example 1.24:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(0, \theta)$  where  $\theta > 0$ . Show that  $\bar{X}_n$  is not sufficient for  $\theta$ , and find a statistic that is.

$$f_\theta(\vec{x}) = \prod_{i=1}^n \frac{1}{\theta} \cdot \mathbb{1}_{0 \leq x_i \leq \theta}$$

$$= \theta^{-n} \cdot \mathbb{1}_{0 \leq x_1 \leq \theta, \dots, 0 \leq x_n \leq \theta}$$

$$= \underbrace{\mathbb{1}_{0 \leq x_m}}_{h(\vec{x})} \cdot \underbrace{\theta^{-n} \cdot \mathbb{1}_{x_m \leq \theta}}_{g_\theta(x_m)}$$

By the FT,  $T(\vec{x}) = X_m$  is sufficient for  $\theta$ .

# The Factorization Theorem: Examples

- **Theorem 1.2:** Let  $X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta$  be a random sample from an exponential family, where

$$f_\theta(x) = h(x) \cdot g(\theta) \cdot \exp \left( \sum_{j=1}^k w_j(\theta) \cdot T_j(x) \right).$$

Then  $T(\mathbf{X}) = \left( \sum_{i=1}^n T_1(X_i), \dots, \sum_{i=1}^n T_k(X_i) \right)$  is sufficient for  $\theta$ .

*Proof.* Important EXERCISE!

# The Factorization Theorem: Examples

- **Example 1.25:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$ , where  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ . Show that  $T(\mathbf{X}) = (\sum_{i=1}^n X_i^2, \sum_{i=1}^n X_i)$  is sufficient for  $(\mu, \sigma^2) = \theta$

Exp family:  $f_\theta(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

$$\begin{aligned} &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2} + \frac{2\mu x}{2\sigma^2} - \frac{\mu^2}{2\sigma^2}\right) \\ &= 1 \cdot \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{1}{2\sigma^2} \cdot x^2 + \frac{\mu}{\sigma^2} x\right) \end{aligned}$$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$

$h(x) \quad g(\theta) \quad w_1(\theta) \quad w_2(\theta)$

$T_1(x) \quad T_2(x)$

By Theorem 1.2,  $T(\mathbf{x}) = (\sum_{i=1}^n X_i^2, \sum_{i=1}^n X_i)$  is sufficient for  $\theta = (\mu, \sigma^2)$ .

# The Factorization Theorem: Examples

- **Example 1.26:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(\{1, 2, \dots, \theta\})$ , where  $\theta \in \mathbb{N}$ . Show that  $T(\mathbf{X}) = X_{(n)}$  is sufficient for  $\theta$ .

$$\begin{aligned}f_\theta(\mathbf{x}) &= \prod_{i=1}^n f_\theta(x_i) \\&= \prod_{i=1}^n \frac{1}{\theta} \mathbb{1}_{x_i \in \{1, 2, \dots, \theta\}} \\&= \theta^{-n} \cdot \mathbb{1}_{x_i \in \{1, 2, \dots, \theta\} \forall i} \\&= \underbrace{\mathbb{1}_{x_i \in \mathbb{N} \forall i}}_{h(\mathbf{x})} \cdot \underbrace{\theta^{-n} \cdot \mathbb{1}_{x_{(n)} \leq \theta}}_{g_\theta(x_{(n)})}\end{aligned}$$

By the FT,  $T(\mathbf{x}) = X_{(n)}$  is sufficient for  $\theta$ .

# If There's One, There's More...

- If we have some sufficient statistic, we can always come up with (infinitely) many others...
- Theorem 1.3: Let  $T(\mathbf{X})$  be sufficient for  $\theta$  and suppose that  $r(\cdot)$  is a bijection. Then  $r(T(\mathbf{X}))$  is also sufficient for  $\theta$ .

Proof. If  $T(\vec{x})$  is sufficient for  $\Theta$ , then by the FT

$$\begin{aligned}f_{\theta}(\vec{x}) &= h(\vec{x}) \cdot g_{\theta}(T(\vec{x})) \text{ for some } h(\cdot), g_{\theta}(\cdot). \\&= h(\vec{x}) \cdot g_{\theta}(r^{-1}(r(T(\vec{x})))) \\&= h(\vec{x}) \cdot \tilde{g}_{\theta}(r(T(\vec{x}))) \text{ where } \tilde{g}_{\theta}(t) = g_{\theta}(r^{-1}(t)).\end{aligned}$$

By the FT,  $r(T(\vec{x}))$  is sufficient for  $\Theta$ .  $\square$

# Too Many Sufficient Statistics

- So there are lots of sufficient statistics out there
- We saw that  $T(\mathbf{X}) = \mathbf{X}$  is always sufficient – it's also pretty useless as far as data reduction goes
- There are usually “better” ones out there – how do we get the best bang for our buck?
- Another issue: the factorization theorem makes it easy to show that a statistic is sufficient (if it actually is), but less so to show that a statistic is *not* sufficient
- We will develop theory that takes care of both of these issues at once

# Minimal Sufficiency

- **Definition 1.8:** A sufficient statistic  $T(\mathbf{X})$  is called a **minimal sufficient statistic** if, for any other sufficient statistic  $U(\mathbf{X})$ , there exists a function  $h$  such that  $T(\mathbf{X}) = h(U(\mathbf{X}))$ .
- In other words, a minimal sufficient statistic is some function of *any other sufficient statistic*  $N(\mu, 1)$ :  $\vec{X}$  is sufficient for  $\mu$ .  
So is  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  which is a function of  $\vec{X}$  (but not the other way around)
- A minimal sufficient statistic achieves the greatest reduction of data possible (while still maintaining sufficiency)
- **Example 1.27:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$ , where  $\mu \in \mathbb{R}$  and  $\sigma^2$  is known. Show that  $T(\mathbf{X}) = (\bar{X}_n, S_n^2)$  is not minimal sufficient for  $\mu$ .  
We saw that  $\bar{X}_n$  is sufficient for  $\mu$ . But  $T(\bar{X})$  is not a function of  $\bar{X}_n$ , so it can't be minimal sufficient for  $\mu$ .

# Poll Time!

# A Criterion For Minimal Sufficiency

- It's usually not that hard to show that a statistic is not minimal sufficient
  - But how can we possibly show that a statistic *is* minimal?
  - **Theorem 1.4:** Let  $f_\theta(\mathbf{x})$  be the pdf/pmf of a sample  $\mathbf{X}$ . Suppose there exists a function  $T(\cdot)$  such that for any  $\mathbf{x}, \mathbf{y} \in \mathcal{X}^n$ ,  $T(\mathbf{x}) = T(\mathbf{y})$  if and only if the ratio  $f_\theta(\mathbf{x})/f_\theta(\mathbf{y})$  is free of  $\theta$ . Then  $T(\mathbf{X})$  is minimal sufficient for  $\theta$ .
- No proof!*
- This criterion is easier to apply than it looks
  - **Example 1.28:**  $X_1, \dots, X_n$  are Bernoulli( $\theta$ ),  $\theta \in (0, 1)$ . Then  $T(\vec{x}) = \sum_i x_i$  is minimal sufficient for  $\theta$ . Let  $\vec{x}, \vec{y} \in \mathcal{X}^n$ . Then
- $$\frac{f_\theta(\vec{x})}{f_\theta(\vec{y})} = \frac{\theta^{\sum x_i} (1-\theta)^{n-\sum x_i}}{\theta^{\sum y_i} (1-\theta)^{n-\sum y_i}} = \theta^{\sum x_i - \sum y_i} (1-\theta)^{\sum y_i - \sum x_i}, \text{ which is free of } \theta \text{ iff } \sum x_i = \sum y_i, \text{ ie } T(\vec{x}) = T(\vec{y}).$$
- $\therefore T(\vec{x})$  is minimal sufficient for  $\theta$ .

# Minimal Sufficiency: Examples

- **Example 1.29:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$ , where  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ . Show that  $T(\mathbf{X}) = (\bar{X}, S^2)$  is minimal sufficient for  $(\mu, \sigma^2) = \theta$ .

Let  $s_x^2 = \frac{1}{n-1} \sum_{i=1}^{n-1} (x_i - \bar{x})^2$ .

From Ex. 1.23,  $f_\theta(\vec{x}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{\sum(x_i - \bar{x})^2 + n(\bar{x} - \mu)^2}{2\sigma^2}\right)$   
 $= \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{(n-1)s_x^2 + n(\bar{x} - \mu)^2}{2\sigma^2}\right).$

Let  $\vec{x}, \vec{y} \in \mathcal{X}^n$ . Then

$$\begin{aligned} \frac{f_\theta(\vec{x})}{f_\theta(\vec{y})} &= \frac{\frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{(n-1)s_x^2 + n(\bar{x} - \mu)^2}{2\sigma^2}\right)}{\frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{(n-1)s_y^2 + n(\bar{y} - \mu)^2}{2\sigma^2}\right)} \\ &= \exp\left(\frac{(n-1)(s_y^2 - s_x^2) + n(\bar{x} + \bar{y} - 2\mu)(\bar{y} - \bar{x})}{2\sigma^2}\right) \end{aligned}$$

is free of  $(\mu, \sigma^2)$  iff  $s_x^2 = s_y^2$  and  $\bar{x} = \bar{y}$ .

By Theorem 1.4,  $T(\vec{x}) = (\bar{x}, S^2)$  is minimal sufficient for  $\theta = (\mu, \sigma^2)$ .

# Minimal Sufficiency: Examples

- **Example 1.30:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$ , where  $\lambda > 0$ . Find a minimal sufficient statistic for  $\lambda$ .

EXERCISE !

# Minimal Sufficiency: Examples

- A minimal sufficient statistic isn't always as minimal as you would expect...
- Example 1.31: Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unif}([\theta, \theta + 1])$ , where  $\theta \in \mathbb{R}$ . Show that  $T(\mathbf{X}) = (X_{(1)}, X_{(n)})$  is minimal sufficient for  $\theta$ .

$$\begin{aligned}f_\theta(\vec{x}) &= \prod_{i=1}^n f_\theta(x_i) = \prod_{i=1}^n \mathbb{1}_{\theta \leq x_i \leq \theta+1} \\&= \mathbb{1}_{\theta \leq x_1, x_2, \dots, x_n \leq \theta+1} \\&= \mathbb{1}_{\theta \leq x_{(1)} \wedge x_{(n)} \leq \theta+1} \\&= \mathbb{1}_{x_{(n)} - 1 \leq \theta \leq x_{(n)}}\end{aligned}$$

Let  $\vec{x}, \vec{y} \in \mathcal{X}^n$ . Then  $\frac{f_\theta(\vec{x})}{f_\theta(\vec{y})} = \frac{\mathbb{1}_{x_{(n)} - 1 \leq \theta \leq x_{(n)}}}{\mathbb{1}_{y_{(n)} - 1 \leq \theta \leq y_{(n)}}}$  is free of  $\theta$  iff  $x_{(1)} = y_{(1)}$  and  $x_{(n)} = y_{(n)}$ .

Hence, by Theorem 1.4,  $T(\vec{x}) = (X_{(1)}, X_{(n)})$  is minimal sufficient for  $\theta$ .

# Poll Time!

(Q: how do you feel about sufficiency?)

# The “Opposite” of Sufficiency?

- We know that a sufficient statistic contains all the information about  $\theta$  that the original sample has
- What about a statistic that contains *no* information about  $\theta$ ?
- Why would such a thing be useful?
- **Definition 1.9:** A statistic  $D(\mathbf{X})$  is an **ancillary statistic** for a parameter  $\theta$  if the distribution of  $D(\mathbf{X})$  does not depend on  $\theta$

Eg:  $X_1, X_2 \stackrel{iid}{\sim} N(\mu, 1), \mu \in \mathbb{R}$   
 $T(\vec{X}) = (X_1 - \frac{1}{2}\sum X_i, X_2 - \frac{1}{2}\sum X_i)$   
right not depend on  $\mu$

# Ancillarity

- **Definition 1.10:** A statistic  $D(\mathbf{X})$  is an **ancillary statistic** for a parameter  $\theta$  if the distribution of  $D(\mathbf{X})$  does not depend on  $\theta$
- **Example 1.32:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unif}([\theta, \theta + 1])$ , where  $\theta \in \mathbb{R}$ . Show that the range statistic  $R(\mathbf{X}) := X_{(n)} - X_{(1)}$  is ancillary for  $\theta$ .  
Let  $Y_i = X_i - \theta$ . Then  $Y_i \stackrel{iid}{\sim} \text{Unif}(0, 1)$  and  $Y_{(j)} = X_{(j)} - \theta$  is also free of  $\theta$ .

Then  $P_\theta(R(\bar{X}) \leq r)$

$$= P_\theta(X_m - X_m \leq r)$$

$$= P_\theta([X_m - \theta] - [X_m - \theta] \leq r)$$

$$= P_\theta(Y_m - Y_m \leq r)$$

$$= P_\theta(\text{Beta}(n, 1) - \text{Beta}(1, n) \leq r) \text{ does not depend on } \theta.$$

Hence  $R(\bar{X})$  is ancillary for  $\theta$ .

# Ancillarity: Examples

- Did we actually use the uniform distribution anywhere in the previous example?
- **Theorem 1.5:** Let  $X_1, \dots, X_n$  be a random sample from a location family with cdf  $F(\cdot - \theta)$ , for  $\theta \in \mathbb{R}$ . Then the range statistic is ancillary for  $\theta$ .

*Proof.* Let  $Y_i = X_i - \theta \sim F(\cdot)$ .

Then  $P_\theta(R(\bar{x}) \leq r)$

$$= P_\theta(X_m - X_{m1} \leq r)$$

$$= P_\theta((X_m - \theta) - (X_{m1} - \theta) \leq r)$$

$$= P_\theta(Y_{m1} - Y_{m1} \leq r), \text{ which is free of } \theta \text{ since the distribution of } Y_{m1}, Y_{m1} \text{ is free of } \theta.$$

□

# Ancillarity: Examples

- **Example 1.33:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ . Show that  $D(\mathbf{X}) = \frac{X_1 + \dots + X_{n-1}}{X_n}$  is ancillary for  $\sigma^2$ .

Let  $Z_i = \frac{X_i}{\sigma} \sim N(0, 1)$ .

Then  $P_{\sigma^2}(D(\bar{X}) \leq x)$

$$= P_{\sigma^2}\left(\frac{X_1}{X_n} + \frac{X_2}{X_n} + \dots + \frac{X_{n-1}}{X_n} \leq x\right)$$

$$= P_{\sigma^2}\left(\frac{X_1/\sigma}{X_n/\sigma} + \frac{X_2/\sigma}{X_n/\sigma} + \dots + \frac{X_{n-1}/\sigma}{X_n/\sigma} \leq x\right)$$

$$= P_{\sigma^2}\left(\frac{Z_1}{Z_n} + \dots + \frac{Z_{n-1}}{Z_n} \leq x\right) \text{ is free of } \sigma^2.$$

Hence  $D(\bar{X})$  is ancillary for  $\sigma^2$

- **Theorem 1.6:** Let  $X_1, \dots, X_n$  be a random sample from a scale family with cdf  $F(\cdot/\sigma)$ , for  $\sigma > 0$ . Then any statistic which is a function of the ratios  $X_1/X_n, \dots, X_{n-1}/X_n$  is ancillary for  $\sigma$ .

EXERCISE!

# Ancillarity: Examples

- Recall that if  $Z_1, \dots, Z_n \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ , then the distribution of  $Y = \sum_{i=1}^n Z_i^2$  is called a **chi-squared distribution with  $n$  degrees of freedom**, which we write as  $Y \sim \chi_{(n)}^2$ .
- Theorem 1.7:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$  with  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ . Then  $\frac{n-1}{\sigma^2} S^2 \sim \chi_{(n-1)}^2$ .

*Proof ( $n = 2$ )*.

$$\begin{aligned} X_1 - X_2 &\sim N(0, 2\sigma^2) \\ &\stackrel{d}{=} \sqrt{2}\sigma \cdot N(0, 1) \\ \Rightarrow (X_1 - X_2)^2 &\stackrel{d}{=} 2\sigma^2 \cdot N(0, 1)^2 \\ &= 2\sigma^2 \cdot \chi_{(1)}^2 \end{aligned}$$

$$\begin{aligned} (n-1)S^2 &= \sum_{i=1}^2 (X_i - \bar{X})^2 = (X_1 - \frac{1}{2}(X_1 + X_2))^2 + (X_2 - \frac{1}{2}(X_1 + X_2))^2 \\ &= (\frac{1}{2}X_1 - \frac{1}{2}X_2)^2 + (\frac{1}{2}X_2 - \frac{1}{2}X_1)^2 \\ &= \frac{1}{2}(X_1 - X_2)^2 \end{aligned}$$

$$\begin{aligned} &\stackrel{d}{=} \frac{1}{2} \cdot 2\sigma^2 \cdot \chi_{(1)}^2 \\ &\stackrel{d}{=} \sigma^2 \cdot \chi_{(1)}^2 \end{aligned}$$

(ie,  $(n-1)S^2 \stackrel{d}{=} \sigma^2 \cdot \chi_{(1)}^2$  where  $\chi_{(1)}^2 \sim \chi_{(1)}^2$ )

General  $n$ : exercise (use induction).

$$\Rightarrow \frac{n-1}{\sigma^2} S^2 \sim \chi_{(n-1)}^2. \quad \square$$

- Example 1.34:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$  with  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ . Show that the sample variance  $S^2$  is ancillary for  $\mu$ .

We showed that  $S^2 \sim \frac{\sigma^2}{n-1} \cdot \chi_{(n-1)}^2$  which is free of  $\mu$ . So  $S^2$  is ancillary for  $\mu$ .

# Poll Time!

$Z \cup f_2$  has no info about the unknown  $\Theta$  in our sample  $X_1, \dots, X_n$ .

$X_1 - X_n$  and  $X_1/X_n$  may be ancillary in some models, but not in general!

---

Support of a function:  $\text{Supp}(f) := \{x : f(x) \neq 0\}$

If  $f: \mathbb{R} \rightarrow \mathbb{R}$  is not supported on  $\mathbb{R}$ , multiply by an indicator for  $\text{Supp}(f)$  to be safe.

e.g. if  $X \sim \text{Unif}(0,1)$ , then  $f_X(x) = 1$  is FALSE!

$$\text{TRUE: } f_X(x) = \begin{cases} 1, & x \in (0,1) \\ 0, & \text{otherwise} \end{cases} = \mathbb{1}_{x \in [0,1]}$$

Indicators are useful!  $x \in A \cap B \iff \mathbb{1}_{A \cap B}(x) = 1 \iff \mathbb{1}_A(x) \cdot \mathbb{1}_B(x) = 1$

# Completeness: An Abstract Definition

- Everything so far has been about ways to reduce the amount of data we need while still retaining all information about  $\theta$
- We've seen that ancillary statistics are bad at it, sufficient statistics are good at it, and minimal sufficient statistics are very good at it
- We will study one more kind of statistic, but the definition isn't pretty
- **Definition 1.11:** A statistic  $U(\mathbf{X})$  is complete if any function  $h(\cdot)$  which satisfies  $\mathbb{E}_\theta [h(U(\mathbf{X}))] = 0$  for all  $\theta \in \Theta$  must also satisfy  $\mathbb{P}_\theta (h(U(\mathbf{X})) = 0) = 1$  for all  $\theta \in \Theta$ .

continuous:  $\int h(u) \cdot f_\theta(u) du = 0 \quad \forall \theta \in \mathbb{R}$  where  $f_\theta(\cdot)$  is the pdf  $U(\vec{x})$

discrete:  $\sum_{u \in \mathcal{U}} h(u) \cdot P_\theta(U(\vec{x})=u) = 0 \quad \forall \theta \in \mathbb{R}$ .

$U(\vec{x})$  complete means this implication is T:  $(\mathbb{E}_\theta [h(U(\vec{x}))] = 0 \quad \forall \theta \Rightarrow \mathbb{P}_\theta (h(U(\vec{x})) = 0) = 1 \quad \forall \theta)$

# Completeness: An Abstract Definition

- The concept of completeness is notoriously unintuitive – probably the most abstract one in our course – but it will pay off later
- For now, you can think about the finite case a bit like a finite-dimensional basis from linear algebra
- If  $\mathbf{v}_1, \dots, \mathbf{v}_n$  span  $\mathbb{R}^n$ , then  $\sum_{i=1}^n a_i \mathbf{v}_i = \mathbf{0}$  implies  $a_i = 0$  for all  $i$
- If  $U(\mathbf{X})$  is complete and supported on  $\{u_1, \dots, u_n\}$ , then  $\sum_{i=1}^n h(u_i) \cdot \mathbb{P}_{\theta}(U(\mathbf{X}) = u_i) = 0$  implies  $h(u_i) = 0$  for all  $i$
- The meaning will become clearer at the end of Module 2
- So why bring it up now?

# Showing Completeness is Very Difficult In General...

- **Example 1.35:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$  with  $\theta \in (0, 1)$ . Show that  $U(\mathbf{X}) = \sum_{i=1}^n X_i$  is complete.  $U(\vec{X}) \sim \text{Bin}(n, \theta)$ .

Suppose  $h(\cdot)$  is some function s.t.  $E_\theta[h(U(\vec{X}))] = 0 \quad \forall \theta \in (0, 1)$

$$\Rightarrow E_\theta[h(\sum_{i=1}^n X_i)] = 0 \quad \forall \theta \in (0, 1)$$

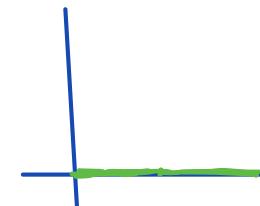
$$\text{Then } 0 = \sum_{j=0}^n h(j) \cdot \binom{n}{j} \theta^j (1-\theta)^{n-j}$$

$$= (1-\theta)^n \sum_{j=0}^n h(j) \cdot \left[ \binom{n}{j} \left( \frac{\theta}{1-\theta} \right)^j \right] \quad (1-\theta)^n \neq 0 \text{ so we can divide through by it}$$

$$= \sum_{j=0}^n h(j) \cdot \left[ \binom{n}{j} \left( \frac{\theta}{1-\theta} \right)^j \right] \quad \text{let } r = \frac{\theta}{1-\theta}$$

$$= \sum_{j=0}^n \hat{h}(j) \cdot r^j \quad \text{where } \hat{h}(j) = \binom{n}{j} h(j)$$

is a polynomial in  $r$  which is 0 for all  $r \in (0, \infty)$  (i.e., constant at 0).



Therefore  $\hat{h}(j)$  is 0  $\forall j$ .

$$\Rightarrow \binom{n}{j} h(j) = 0 \quad \forall j$$

$$\Rightarrow h(j) = 0 \quad \forall j$$

$$\Rightarrow h(\cdot) = 0 \text{ on } \{0, 1, \dots, n\}.$$

Hence,  $P_\theta(h(U(\vec{X})) = 0) = 1 \quad \forall \theta \in (0, 1)$ , so  $U(\vec{X})$  is complete.

# ...But for Exponential Families, There's Nothing To It

- **Theorem 1.8:** Let  $X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta$  be a random sample from an exponential family, where

$$f_\theta(x) = h(x) \cdot g(\theta) \cdot \exp \left( \sum_{j=1}^k w_j(\theta) \cdot T_j(x) \right).$$

Then  $T(\mathbf{X}) = \left( \sum_{i=1}^n T_1(X_i), \dots, \sum_{i=1}^n T_k(X_i) \right)$  is a complete statistic, as long as each component of  $\Theta$  contains an open interval in  $\mathbb{R}$ .<sup>1</sup>

$\Theta = \mathbb{R} \supseteq (0, 1) \quad \checkmark$   
 $\Theta = \mathbb{R} \times (0, \infty) \supseteq (0, 1) \times (0, 1) \quad \checkmark$   
 $\Theta = \{1, 2, 3, \dots, \theta\} \quad \times$   
DOES NOT CONTAIN AN OPEN INTERVAL

- Recall from Theorem 1.2 that in this case,  $T(\mathbf{X})$  is also sufficient for  $\theta$
- So it's really easy to find complete sufficient statistics for exponential families

---

<sup>1</sup>More generally,  $\Theta$  must contain an open set in  $\mathbb{R}^k$  – this requirement is sometimes called the “open set condition”

# Completeness: Examples

- **Example 1.36:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$ , where  $\mu \in \mathbb{R}$  and  $\sigma^2$  ~~is known~~. Show that  $\bar{X}_n$  is complete for  $\mu$ .

$$\begin{aligned}f_{\theta}(x) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x^2 - 2\mu x + \mu^2)}{2\sigma^2}\right) \\&= \underbrace{\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right)}_{h(x)} \cdot \underbrace{\exp\left(-\frac{\mu^2}{2\sigma^2}\right)}_{g(\theta)} \cdot \underbrace{\exp\left(\frac{\mu}{\sigma^2} \cdot x\right)}_{w(\theta)}\end{aligned}$$

Since  $\Theta = \mathbb{R}$  clearly contains an open interval,

$T(\vec{x}) = \bar{X}_n$  is complete for  $\mu$  by Theorem 1.8.

# Completeness: Examples

- **Example 1.37:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$ , where  $\lambda > 0$ . Show that  $\bar{X}_n$  is complete for  $\lambda$ .

$$f_\lambda(x) = \frac{\lambda^x e^{-\lambda}}{x!} = \frac{1}{x!} \cdot e^{-\lambda} \cdot \exp\left(\frac{1}{n} \sum_{i=1}^n x_i\right)$$

$w(\lambda)$   
I can multiply/divide  
by any constant I want!

$$\begin{matrix} h(x) \\ g(\lambda) \end{matrix}$$

$$T(x)$$

$\mathbb{R} = (0, \infty)$  contains an open interval, so  $T(\bar{X}) = \bar{X}_n$  is complete for  $\lambda$  by Theorem 1.8.

# Completeness: Examples

- **Example 1.38:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_{\mu, \sigma}$  where

$$f_{\mu, \sigma}(x) = \frac{1}{2\sigma} \exp\left(-\frac{|x - \mu|}{\sigma}\right), \quad x \in \mathbb{R},$$

where  $\sigma > 0$  and  $\mu$  is known. Find a complete statistic for  $\sigma$ .

$$T(\vec{x}) = \sum_{i=1}^n |X_i - \mu|$$

EXERCISE!

# Complete Statistics Are Minimal Sufficient!

- There is nothing resembling sufficiency in the definition of completeness; the two concepts seem completely unrelated
- And yet, Theorem 1.8 says that for exponential families, certain complete statistics are sufficient
- What about in general? The answer might surprise you...
- Theorem 1.9 (**Bahadur's theorem**): If a minimal sufficient statistic and a complete statistic both exist, then the complete statistic must also be minimal sufficient.  
*sufficient*      *sufficient*      *No proof...*
- That's *not* the same as saying that all minimal sufficient statistics are complete (which is unfortunately not true)

# Minimal Sufficient Statistics Are Not Always Complete

- But if a minimal sufficient statistic exists and it's not complete, then no complete statistic exists
- This is probably the simplest example of a minimal sufficient statistic that is not complete
- **Example 1.39:** Let  $X_1 \sim \text{Unif}(\theta, \theta + 1)$ , where  $\theta \in \mathbb{R}$ . Show that  $T(X_1) = X_1$  is minimal sufficient for  $\theta$ , but not complete.

$$f_\theta(x) = \mathbb{1}_{\theta \leq x \leq \theta+1}.$$

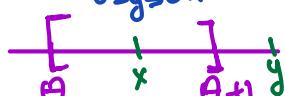
Let  $x, y \in X$ . Then

$$\frac{f_\theta(x)}{f_\theta(y)} = \frac{\mathbb{1}_{\theta \leq x \leq \theta+1}}{\mathbb{1}_{\theta \leq y \leq \theta+1}} \text{ is free of } \theta \text{ iff } x=y.$$

By Theorem 1.4,  $T(X) = X$  is minimal sufficient for  $\theta$ .

If  $x < y$ , then we could have  $x \leq \theta < y$ ,

so  $\mathbb{1}_{\theta \leq x \leq \theta+1} \neq \mathbb{1}_{\theta \leq y \leq \theta+1} \forall \theta$ :



However, consider  $h(x) = \sin(2\pi x)$ .

$$\text{We have } E_\theta[h(TX)]$$

$$= E_\theta[h(X)]$$

$$= \int_{\theta}^{\theta+1} \sin(2\pi x) dx$$

$$= 0 \quad \forall \theta \in \mathbb{R}.$$

But clearly,  $h(\cdot)$  is not identically 0.

$$\text{So } P_\theta(h(X)=0) \neq 1 \quad \forall \theta \in \mathbb{R}.$$

So  $T(X) = X$  is not complete.

# The Amazingly Useful Basu's Theorem

- Theorem 1.10 (**Basu's theorem**): Complete sufficient statistics are independent of *all* ancillary statistics.

Proof. (Discrete)

Let  $T = T(\vec{X})$  be a complete sufficient statistic.

Let  $S = S(\vec{X})$  be ancillary for  $\Theta$ .

It suffices to show  $P_\theta(S=s | T=t) = P(S=s)$ .

$$\text{By the Law of Total Probability, } P(S=s) = \sum_{t \in T} P(S=s | T=t) \cdot P_\theta(T=t) \quad (1)$$

$$\text{Moreover } 1 = \sum_{t \in T} P(T=t) \text{ so } P(S=s) = \left( \sum_{t \in T} P(T=t) \right) \cdot P(S=s) \quad (2)$$

$$\begin{aligned} \text{Therefore, } 0 &= (1) - (2) = \sum_{t \in T} \left[ P(S=s | T=t) - \underbrace{P(S=s)}_{=: h(t)} \right] \cdot P_\theta(T=t) \\ &= \sum_{t \in T} h(t) \cdot P_\theta(T=t) \\ &= E_\theta[h(T)] \quad \forall \theta \in \Theta. \end{aligned}$$

Since  $T$  is complete,  $h(t) = 0 \ \forall t \Rightarrow P(S=s | T=t) = P(S=s) \Rightarrow S \perp\!\!\!\perp T$ .  $\square$

# Poll Time!

# Basu's Theorem: Examples

- **Example 1.40:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$  where  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ . Show that the sample mean  $\bar{X}$  is independent of the sample variance  $S^2$ .

Example 1.36  $\Rightarrow \bar{X}_n$  is a complete sufficient statistic for  $\mu$ .

Example 1.34  $\Rightarrow S^2$  is ancillary for  $\mu$ .

By Basu's theorem,  $\bar{X}_n \perp\!\!\! \perp S^2$ .

-

- This is actually a characterizing property of the Normal distribution:  $\bar{X}_n \perp\!\!\! \perp S^2$  if and only if  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$  (interesting fact!)  
*(allowing for degenerate Normals with  $\sigma^2=0$ , since constants are independent of everything)*

# Basu's Theorem: Examples

- Example 1.41: Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Exp}(\theta)$ , where  $\theta > 0$ . Use Basu's theorem to find  $\mathbb{E}_\theta \left[ \frac{X_1}{X_1 + \dots + X_n} \right]$ .

$\{\text{Exp}(\theta) : \theta > 0\}$  is a scale family  $\Rightarrow$  by Theorem 1.6,  $\frac{X_1}{X_1 + \dots + X_n}$  is ancillary for  $\theta$ .

$\{\text{Exp}(\theta) : \theta > 0\}$  is an exponential family with  $T(x) = x \Rightarrow X_1 + \dots + X_n$  is complete sufficient for  $\theta$ .

By Basu's theorem,  $\frac{X_1}{X_1 + \dots + X_n} \perp\!\!\!\perp X_1 + \dots + X_n$ .

$$\begin{aligned} \text{Therefore, } \mathbb{E}(X_1) &= \mathbb{E}_\theta \left[ \frac{X_1}{X_1 + \dots + X_n} \cdot (X_1 + \dots + X_n) \right] \\ &= \mathbb{E} \left[ \frac{X_1}{X_1 + \dots + X_n} \right] \cdot \mathbb{E}_\theta [X_1 + \dots + X_n] \quad \text{by independence} \end{aligned}$$

$$\Rightarrow \frac{1}{\theta} = \mathbb{E} \left[ \frac{X_1}{X_1 + \dots + X_n} \right] \cdot \frac{n}{\theta}$$

$$\Rightarrow \mathbb{E} \left[ \frac{X_1}{X_1 + \dots + X_n} \right] = \frac{1}{n}.$$

$$f_\theta(x) = \theta e^{-\theta x} \Rightarrow \mathbb{E}[X] = \int_0^\infty x \cdot \theta e^{-\theta x} dx = \theta \cdot \frac{1}{\theta^2} = \frac{1}{\theta}.$$

# Basu's Theorem: Examples

- **Example 1.42:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_{\mu, \sigma}$  where

$$f_{\mu, \sigma}(x) = \frac{1}{2\sigma} \exp\left(-\frac{|x - \mu|}{\sigma}\right), \quad x \in \mathbb{R},$$

where  $\sigma > 0$  and  $\mu$  is known. Show that  $X_1/X_n$  is independent of  $\sum_{i=1}^n |X_i - \mu|$ .

This is a scale family, so  $\frac{X_1}{X_n}$  is ancillary for  $\sigma$  by Theorem 1.6.

Moreover,  $\sum_{i=1}^n |X_i - \mu|$  is a complete sufficient statistic by Example 1.38.

By Basu's theorem,  $\sum_{i=1}^n |X_i - \mu| \perp\!\!\!\perp \frac{X_1}{X_n}$ .