

# STA261 - Module 2

## Point Estimation

Rob Zimmerman

University of Toronto

July 13-15, 2021

# Extracting Information

- In Module 1, we learned about how a statistic can capture (or not capture) the information provided by our data sample  $\mathbf{X} = (X_1, \dots, X_n) \sim f_\theta$  about the unknown parameter  $\theta \in \Theta$
- For the remainder of the course, our focus will be on how to *extract* that information
- In Module 2, we have one goal: to estimate the parameter  $\theta$  or some function of the parameter  $\tau(\theta)$  as best we can (whatever that means)
- Example 2.1:
  - Heights of STA261 students  $\sim N(\mu, 3)$ .  
Want to estimate  $\mu$ . Maybe take  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i$ . Seems reasonable!
  - Event of voting for Candidate X in next election  $\sim \text{Bernoulli}(p)$ .  
Maybe want to estimate  $\tau(p) = \log\left(\frac{p}{1-p}\right)$  "log-odds of  $p$ "

# Point Estimation

- How do we estimate  $\theta$  from the observed data  $\mathbf{x}$ ?
- Ideally, we want some statistic  $T(\mathbf{X})$  such that  $T(\mathbf{x})$  will be close to  $\theta$
- **Definition 2.1:** Suppose  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$ . A **point estimator**  $\hat{\theta} = \hat{\theta}(\mathbf{X})$  is a statistic used to estimate  $\theta$ , which is not a function of  $\theta$  itself.
- How do we find good point estimators?

# Poll Time!

# Choosing “Good” Point Estimators

- A point estimator  $\hat{\theta}(\mathbf{X})$  is a random variable, so it has its own distribution (as does any statistic)
- Definition aside, it would seem that the best point estimator is the constant  $\hat{\theta}(\mathbf{X}) = \theta$ , but of course this is unattainable
- The constant  $\theta$  has  $\mathbb{E}_\theta [\theta] = \theta$  and  $\text{Var}_\theta (\theta) = 0$
- It would be nice if the distribution of  $\hat{\theta}(\mathbf{X})$  got close to these properties:  
$$\mathbb{E}_\theta [\hat{\theta}(\mathbf{X})] \approx \theta \text{ and } \text{Var}_\theta (\hat{\theta}(\mathbf{X})) \approx 0$$
- It would also be good if  $\text{Var}_\theta (\hat{\theta}(\mathbf{X}))$  got lower as the sample size  $n$  got bigger (if we're willing to pay good money for more samples, we should demand a higher precision in return)

In Module 5, we'll have further demands  
for when  $n \rightarrow \infty$

# Moments Are (Often) Functions of Parameters

- Here's one approach to choosing  $\hat{\theta}$
- In parametric families, it is often the case that the moments (i.e.,  $\mathbb{E}_\theta [X]$ ,  $\mathbb{E}_\theta [X^2]$ ,  $\mathbb{E} [X^3]$ , and so on) are functions of the parameters
- Example 2.2:  $X \sim N(\mu, \sigma^2) \Rightarrow \mathbb{E}[X] = \mu, \quad \mathbb{E}[X^2] = \mu^2 + \sigma^2$

$$X \sim \text{Bin}(n, p) \Rightarrow \mathbb{E}[X] = np, \quad \mathbb{E}[X^2] = np(1-p) + n^2 p^2$$

$$X \sim \text{Poisson}(\lambda) \Rightarrow \mathbb{E}[X] = \lambda, \quad \mathbb{E}[X^2] = \lambda + \lambda^2$$

$$X \sim \text{Exp}(\lambda) \Rightarrow \mathbb{E}[X] = \frac{1}{\lambda}, \quad \mathbb{E}[X^n] = \frac{n!}{\lambda^n} \quad (\text{exercise})$$

# Towards the Method of Moments

- Suppose  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$  and we want to estimate  $\mu$
- We know that  $\mathbb{E}[X_1] = \mu$  and  $\mathbb{E}[X_1^2] - \mathbb{E}[X_1]^2 = \sigma^2$
- So if we took  $\hat{\mu}_1(\mathbf{X}) = X_1$ , then we'd have  $\mathbb{E}_{\mu}[\hat{\mu}_1(\vec{X})] = \mathbb{E}[X_1] = \mu$
- Can we do better? What about  $\hat{\mu}_2(\vec{X}) = \bar{X}_n$ ? Then  $\mathbb{E}_{\mu}[\hat{\mu}_2(\vec{X})] = \mu$   
and  $\text{Var}_{\mu}(\hat{\mu}_2(\vec{X})) = \frac{\sigma^2}{n} < \sigma^2 = \text{Var}_{\mu}(\hat{\mu}_1(\vec{X}))$
- Now suppose we want to estimate both  $\mu$  and  $\sigma^2$
- If we let  $m_1(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n X_i$  and  $m_2(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n X_i^2$ , then  
 $m_1(\mathbf{X}) \xrightarrow{d} \mu$  and  $m_2(\mathbf{X}) \xrightarrow{d} \mu^2 + \sigma^2$  by LLN
- Therefore  $m_2(\mathbf{X}) - m_1(\mathbf{X})^2 \xrightarrow{d} \sigma^2$  by the continuous mapping theorem  
So take  $\hat{\mu}(\vec{X}) = m_1(\vec{X})$  and  $\hat{\sigma}^2(\vec{X}) = m_2(\vec{X}) - m_1(\vec{X})^2$

# The Method of Moments

- Effectively, we're replacing the true moments with the sample moments
- Definition 2.2:** Suppose we have  $k$  parameters  $\theta_1, \theta_2, \dots, \theta_k$  to estimate in a parametric model, and each one is some function of the first  $k$  moments:

$$\theta_j = \psi_j \left( \mathbb{E}_\theta [X], \mathbb{E}_\theta [X^2], \dots, \mathbb{E}_\theta [X^k] \right), \quad 1 \leq j \leq k.$$

The **Method of Moments (MOM)** estimator for  $\theta_j$  is defined by choosing

$$\hat{\theta}_j(\mathbf{X}) = \psi_j \left( m_1(\mathbf{X}), m_2(\mathbf{X}), \dots, m_k(\mathbf{X}) \right), \quad 1 \leq j \leq k.$$

$$\downarrow \\ \frac{1}{n} \sum_{i=1}^n X_i^K$$

Principle: LLN + Continuous Mapping

but continuity of  $\psi_j$ 's not required!

# Method of Moments: Examples

- **Example 2.3:** Suppose  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$ , where  $\lambda > 0$ . Find the MOM estimator for  $\lambda$ .

$$\lambda = \mathbb{E}[X].$$

$\Rightarrow$  The MOM estimator is  $\hat{\lambda}(\vec{x}) = \bar{X}_n$ .

# Method of Moments: Examples

- **Example 2.4:** Suppose  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bin}(k, \theta)$ , where  $k \in \mathbb{N}$  and  $\theta$  is known. Find the MOM estimator for  $k$ .

$$\mathbb{E}[X] = k\theta \quad \Rightarrow \quad k = \frac{\mathbb{E}[X]}{\theta} = \psi_1(\mathbb{E}[X])$$

$$\Rightarrow \text{Our MOM estimator is } \hat{k}(\vec{x}) = \frac{\bar{X}_n}{\theta}.$$

- Could this be a problem?  
Yes! There's no reason for  $\hat{k}$  to be a natural number! Even though  $\mathbb{N} = \mathbb{N}$ .

# Poll Time!

$X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bin}(k, \theta)$ ,  $k$  known.

$$\mathbb{E}[X] = k\theta \Rightarrow \theta = \frac{\mathbb{E}[X]}{k}$$

$$\Rightarrow \underset{\text{mom}}{\hat{\theta}}(\vec{x}) = \frac{m_1(\vec{x})}{k} = \frac{1}{nk} \sum_{i=1}^n x_i.$$

# Method of Moments: Examples

- **Example 2.5:** The angle at which electrons are emitted in muon decay has a distribution with density  $f_\alpha(x) = (1 + \alpha x)/2$ , where  $x \in [-1, 1]$  and  $\alpha \in [-\frac{1}{3}, \frac{1}{3}]$ . Given a sample  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\alpha$ , find the MOM estimator for  $\alpha$ .

$$\mathbb{E}[x] = \int_{-1}^1 x \left( \frac{1+\alpha x}{2} \right) dx = \frac{1}{2} \left[ \frac{x^2}{2} + \frac{\alpha x^3}{3} \right]_{-1}^1 = \frac{\alpha}{3}$$

$$\Rightarrow \alpha = 3 \cdot \mathbb{E}[x]$$

So the MOM estimator of  $\alpha$  is  $\hat{\alpha}(\vec{x}) = 3 \bar{x}_n$ .

# Method of Moments: Examples

- **Example 2.6:** Suppose  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Gamma}(\alpha, \beta)$ , where  $\alpha, \beta > 0$ . Find the MOM estimators for  $\alpha$  and  $\beta$ .

$$\mu_1 = \mathbb{E}[X] = \frac{\alpha}{\beta} \quad \textcircled{1}$$

$$\mu_2 = \mathbb{E}[X^2] = \frac{\alpha}{\beta^2} + \frac{\alpha^2}{\beta^2} = \frac{\alpha + \alpha^2}{\beta^2} \quad \textcircled{2}$$

$$\textcircled{1} \Rightarrow \alpha = \mu_1 \cdot \beta$$

$$\Rightarrow \mu_2 = \frac{\mu_1 \cdot \beta + \mu_1^2 \cdot \beta^2}{\beta^2} = \frac{\mu_1}{\beta} + \mu_1^2$$

$$\Rightarrow \beta = \frac{\mu_1}{\mu_2 - \mu_1^2}$$

$$\Rightarrow \alpha = \frac{\mu_1^2}{\mu_2 - \mu_1^2}$$

Our MOM estimators are

$$\hat{\alpha}(\vec{x}) = \frac{m_1(\vec{x})^2}{m_2(\vec{x}) - m_1(\vec{x})^2}$$

$$\hat{\beta}(\vec{x}) = \frac{m_1(\vec{x})}{m_2(\vec{x}) - m_1(\vec{x})^2}$$

# Method of Moments: Advantages and Disadvantages

**GOOD:** Very simple; they always exist whenever the moments themselves do

**BAD:** If moments don't exist (e.g., Cauchy), this won't work

**BAD:** Observed values may not be in  $\mathbb{N}$

**BAD:** May not be sufficient (if we care)

**BAD:** May not have lowest variance possible (as we'll see....)

# The Likelihood Function

- **Definition 2.3:** Let  $\mathbf{X} \sim f_\theta$ , where  $f_\theta$  is a pdf or pmf in a parametric family. Given the observation  $\mathbf{X} = \mathbf{x}$ , the likelihood function for  $\theta$  is the function  $L(\cdot | \mathbf{x}) : \Theta \rightarrow [0, \infty)$  given by  $L(\theta | \mathbf{x}) = f_\theta(\mathbf{x})$ . ( $L(\theta | \mathbf{x})$  is a random function of  $\Theta$ )  
If  $\theta_1, \theta_2 \in \Theta$ , then  $L(\theta_1 | \mathbf{x}) > f_{\theta_1}(\mathbf{x})$  and  $L(\theta_2 | \mathbf{x}) = f_{\theta_2}(\mathbf{x})$
- Interpret this as the “probability” of observing the sample  $\mathbf{x}$ , given that the sample came from  $f_\theta$   
NOT  $P(\theta = \theta | \mathbf{x} = \mathbf{x})$  !!!
- So  $L(\theta_1 | \mathbf{x}) > L(\theta_2 | \mathbf{x})$  says that the chance of observing  $\mathbf{X} = \mathbf{x}$  is more likely under  $f_{\theta_1}$  than under  $f_{\theta_2}$  (ie, the likelihood function ranks the elements of  $\Theta$  given the observed data  $\mathbf{x}$ )
- It could be that the likelihood is very small for all  $\theta \in \Theta$ , so knowing  $L(\theta | \mathbf{x})$  for just a single  $\theta$  is useless
- Instead, we want to know how  $L(\theta | \mathbf{x})$  compares to the other  $L(\theta' | \mathbf{x})$ 's

# The Likelihood Principle

- Much of modern statistics revolves around the likelihood function; it will be with us in some form or another for the rest of our course
- The **likelihood principle** states that if two model and data combinations  $L_1(\theta | \mathbf{x})$  and  $L_2(\theta | \mathbf{y})$  are such that  $L_1(\theta | \mathbf{x}) = c(\mathbf{x}, \mathbf{y}) \cdot L_2(\theta | \mathbf{y})$ , then the conclusions about  $\theta$  drawn from  $\mathbf{x}$  and  $\mathbf{y}$  should be identical
- In other words, the likelihood principle says that anything we want to say about  $\theta$  should be based solely on  $L(\cdot | \mathbf{x})$ , regardless of how  $\mathbf{x}$  was actually obtained
- Is this requirement too strong?

Experiment 1: toss a coin ( $P(H) = \theta$ ) 10 times; let  $X = \# \& H \sim \text{Bin}(10, \theta)$ .

- Example 2.7: Observe  $X = 4$ . Then  $L_1(\theta | 4) = \binom{10}{4} \theta^4 (1-\theta)^6$

Experiment 2: toss the same coin until we observe 4 H's; let  $Y = \# \& T$  until then  $\sim \text{NegBin}(4, \theta)$ .  
Observe  $Y = 6$ . Then  $L_2(\theta | 6) = \binom{9}{6} \theta^4 (1-\theta)^6$

Then  $L_1(\theta | 4) \propto L_2(\theta | 6)$ . Likelihood principle says we should be indifferent to Experiment 1 or Experiment 2.

# Maximizing the Likelihood

- Suppose there were some  $\hat{\theta} \in \Theta$  which makes  $L(\hat{\theta} | \mathbf{x})$  the highest; would it be sensible to use that  $\hat{\theta}$  as an estimator?
- If we can maximize  $L(\theta | \mathbf{x})$  with respect to  $\theta$ , the resulting maximizer  $\hat{\theta}$  will be a function of the sample  $\mathbf{x}$
- Example 2.8: Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$ , where  $\theta \in (0, 1)$ . Maximize the likelihood with respect to  $\theta$ .

$$L(\theta | \vec{x}) = f_{\theta}(\vec{x}) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} = \theta^{\sum x_i} (1-\theta)^{n - \sum x_i}$$

We'll soon see that the maximum occurs at  $\hat{\theta} = \bar{x}$ .

So with this idea, a point estimator could be  $\hat{\theta}(\vec{x}) = \bar{x}_n$ .

# Maximum Likelihood Estimation

- **Definition 2.4:** Let  $\mathbf{X} = (X_1, \dots, X_n) \sim f_\theta$ . Let  $L(\theta | \mathbf{x})$  be the likelihood function based on observing  $\mathbf{X} = \mathbf{x}$ . The **maximum likelihood estimate** of  $\theta$  is given by

$$\hat{\theta}(\mathbf{x}) = \operatorname{argmax}_{\theta \in \Theta} L(\theta | \mathbf{x}),$$

ie, the  $\theta \in \Theta$  that maximizes  $L(\theta | \mathbf{x})$

and the **maximum likelihood estimator (MLE)** for  $\theta$  is the point estimator given by  $\hat{\theta}_{\text{MLE}} = \hat{\theta}(\mathbf{x})$ .

equivalently:

A statistic! Has a distribution, an expectation, variance, etc.

$\hat{\theta}(\mathbf{x})$  is such that  $L(\hat{\theta} | \mathbf{x}) \geq L(\theta | \mathbf{x}) \quad \forall \theta \in \Theta$

# Maximum Likelihood: Examples

- Nothing says the distribution needs to have a “nice” functional form
- Example 2.9: Suppose  $\mathcal{X} = \{1, 2, 3\}$  and  $\Theta = \{a, b\}$ , and a parametric family is given by the following table:

	$x = 1$	$x = 2$	$x = 3$
$f_a(x)$	0.3	0.4	0.3
$f_b(x)$	0.1	0.7	0.2

Suppose we observe  $X \sim f_\theta$ . Find the MLE of  $\theta$ .

$$X=1 \Rightarrow f_a(1) > f_b(1) \Rightarrow \hat{\theta}(1) = a$$

$$X=2 \Rightarrow f_a(2) < f_b(2) \Rightarrow \hat{\theta}(2) = b$$

$$X=3 \Rightarrow f_a(3) > f_b(3) \Rightarrow \hat{\theta}(3) = a$$

$$\text{So } \hat{\theta}(x) = \begin{cases} a, & x \in \{1, 3\} \\ b, & x=2 \end{cases} = a \cdot \mathbf{1}_{x \in \{1, 3\}} + b \cdot \mathbf{1}_{x=2}$$

The MLE is therefore  $\hat{\theta}(x) = a \cdot \mathbf{1}_{x \in \{1, 3\}} + b \cdot \mathbf{1}_{x=2}$

# Maximum Likelihood: Examples

- But when the  $f_\theta$  does have a nice form and is continuously differentiable for  $\theta \in \Theta$ , we can use calculus to find the MLE
- Example 2.10: Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$ , where  $\theta \in (0, 1)$ . Find the MLE of  $\theta$ .

$$\begin{aligned} L(\theta | \bar{x}) &= \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} \\ \Rightarrow \frac{dL}{d\theta} &= (\sum x_i) \theta^{\sum x_i - 1} (1-\theta)^{n-\sum x_i} - \theta^{\sum x_i} (n-\sum x_i) (1-\theta)^{n-\sum x_i - 1} \stackrel{\text{set } 0}{=} 0 \\ \Rightarrow (\sum x_i) \theta^{-1} - (n-\sum x_i) (1-\theta)^{-1} &= 0 \\ \Rightarrow \frac{\bar{x}}{1-\bar{x}} = \frac{\theta}{1-\theta} &\Rightarrow \hat{\theta} = \bar{x} \end{aligned}$$

Is this a local max? Need to find  $\frac{d^2L}{d\theta^2}$ , plug in  $\theta = \bar{x}$ , and check if it's  $< 0$ .  
(You can verify). So  $\hat{\theta}(\bar{x}) = \bar{x}_n$  is the MLE for  $\theta$ .

# Maximum Likelihood: Examples

- Suppose that  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$ , where  $\mu \in \mathbb{R}$  and  $\sigma^2$  is known
- What happens if we try to find the MLE of  $\mu$  in the same fashion?

$$\begin{aligned} L(\mu | \vec{x}) &= \prod_{i=1}^n f_\mu(x_i) \\ &= (2\pi\sigma^2)^{-n/2} \cdot \exp\left(-\frac{\sum(x_i - \mu)^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-n/2} \cdot \exp\left(-\frac{1}{2\sigma^2}(\sum x_i^2 - 2\mu\sum x_i + n\mu^2)\right) \\ \frac{dL}{d\mu} &= \underbrace{(2\pi\sigma^2)^{-n/2}}_{\text{never } 0} \cdot \underbrace{\left(\frac{\sum x_i - n\mu}{2\sigma^2}\right)}_{\text{must be } 0} \cdot \exp\left(-\frac{\sum(x_i - \mu)^2}{2\sigma^2}\right) \stackrel{\text{set } 0}{=} 0 \\ \Rightarrow \frac{\sum x_i - n\mu}{2\sigma^2} &= 0 \Rightarrow \hat{\mu} = \frac{1}{n} \sum x_i = \bar{x} \end{aligned}$$

But differentiating  $\frac{dL}{d\mu}$  again would be a nightmare!  
Is there an easier way?

# The Log-Likelihood

- **Definition 2.5:** Given data  $\mathbf{x}$  and a parametric model with likelihood function  $L(\theta | \mathbf{x})$ , the **log-likelihood function** is defined as by

$$\ell(\theta | \mathbf{x}) = \log(L(\theta | \mathbf{x})).$$

- Maximizing the log-likelihood is equivalent to maximizing the likelihood because it's a monotonically increasing function of  $L(\theta | \mathbf{x})$ .
- ...but usually way easier  
... because sums are usually way easier to differentiate than products!  
(we  $\heartsuit$  linearity!)

- If the data are iid, then  $\ell(\theta | \mathbf{x}) = \log(L(\theta | \mathbf{x}))$   
 $= \log\left(\prod_{i=1}^n f_\theta(x_i)\right)$  $= \sum_{i=1}^n \log(f_\theta(x_i))$

# The Score Function

- **Definition 2.6:** Given data  $\mathbf{x}$  and a parametric model with log-likelihood function  $\ell(\theta | \mathbf{x})$ , the **score function** is defined as

$$S(\theta | \mathbf{x}) = \frac{\partial}{\partial \theta} \ell(\theta | \mathbf{x}),$$

when it exists.

- When  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)$  is a vector, this is interpreted as the gradient

$$S(\boldsymbol{\theta} | \mathbf{x}) = \nabla \ell(\boldsymbol{\theta} | \mathbf{x}) = \left( \frac{\partial}{\partial \theta_1} \ell(\boldsymbol{\theta} | \mathbf{x}), \dots, \frac{\partial}{\partial \theta_k} \ell(\boldsymbol{\theta} | \mathbf{x}) \right)$$

- If the likelihood function is nice enough, then any extremum  $\hat{\theta}$  will satisfy the *score equation*  $S(\hat{\theta} | \mathbf{x}) = 0$
- So finding the MLE amounts to finding  $\hat{\theta}$  such that  $S(\hat{\theta} | \mathbf{x}) = 0$  and then checking that  $\hat{\theta}$  is a global maximum

# Maximum Likelihood: More Examples

- **Example 2.11:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$  with  $\mu \in \mathbb{R}$  and  $\sigma^2$  known. Find the MLE of  $\mu$ .

$$L(\mu | \vec{x}) = (2\pi\sigma^2)^{-n/2} \cdot \exp\left(-\frac{\sum(x_i - \mu)^2}{2\sigma^2}\right)$$

$$\Rightarrow l(\mu | \vec{x}) = C - \frac{\sum(x_i - \mu)^2}{2\sigma^2} \quad \text{where } C \text{ is free of } \mu$$

$$\begin{aligned} \Rightarrow S(\mu | \vec{x}) &= \frac{\partial}{\partial \mu} \left( -\frac{\sum x_i^2 + 2\mu \sum x_i - n\mu^2}{2\sigma^2} \right) \\ &= \frac{\sum x_i - n\mu}{\sigma^2} \stackrel{\text{set}}{=} 0 \Rightarrow \hat{\mu} = \bar{x} \end{aligned}$$

Second derivative test:

$$\frac{\partial^2}{\partial \mu^2} S(\mu | \vec{x}) = -\frac{n}{\sigma^2} \Rightarrow \frac{\partial^2}{\partial \mu^2} S(\mu | \vec{x}) \Big|_{\hat{\mu} = \bar{x}} < 0.$$

So  $\hat{\mu}(\vec{x}) = \bar{x}_n$  is the MLE for  $\mu$ .

# Maximum Likelihood: More Examples

- **Example 2.12:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Exp}(\lambda)$  with  $\lambda > 0$ . Find the MLE of  $\lambda$ .

$$L(\lambda | \vec{x}) = \prod_{i=1}^n \lambda e^{-\lambda x_i} = \lambda^n e^{-\lambda \sum x_i} = \lambda^n e^{-\lambda n \bar{x}}.$$

$$\Rightarrow \ell(\lambda | \vec{x}) = n \cdot \log(\lambda) - \lambda n \bar{x}$$

$$\Rightarrow S(\lambda | \vec{x}) = \frac{n}{\lambda} - n \bar{x} \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \hat{\lambda} = \frac{1}{\bar{x}}$$

Second derivative test:  $\frac{\partial}{\partial \lambda} S(\lambda | \vec{x}) = -\frac{n}{\lambda^2}$

$$\Rightarrow \left. \frac{\partial}{\partial \lambda} S(\lambda | \vec{x}) \right|_{\lambda=\frac{1}{\bar{x}}} = -\frac{n}{\left(\frac{1}{\bar{x}}\right)^2} < 0.$$

So the MLE of  $\lambda$  is  $\hat{\lambda}(\vec{x}) = \frac{1}{\bar{x}_n}$

# Maximum Likelihood: More Examples

- Even if the likelihood is smooth and well-behaved, this method doesn't always work
- Example 2.13:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \Gamma(\alpha, 2)$  with  $\alpha > 0$ . Try to find the MLE of  $\alpha$ .

$$\begin{aligned} L(\alpha | \vec{x}) &= \prod_{i=1}^n \frac{2^\alpha}{\Gamma(\alpha)} x_i^{\alpha-1} e^{-2x_i} \\ &= \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left( \prod_{i=1}^n x_i \right)^{\alpha-1} e^{-2\sum x_i} \end{aligned}$$

$$\Rightarrow l(\alpha | \vec{x}) = n\alpha \cdot \log(2) - n \cdot \log(\Gamma(\alpha)) + (\alpha-1) \cdot \log \left( \prod_{i=1}^n x_i \right) + c \quad \text{where } c \text{ is free of } \alpha$$

$$\Rightarrow S(\alpha | \vec{x}) = n \cdot \log(2) - \text{???} + \log \left( \prod_{i=1}^n x_i \right)$$

↑ No closed form! Can't differentiate this!

# Maximum Likelihood: More Examples

- What about when  $\theta$  is multidimensional? We need to bring out our multivariate calculus

(Completed after lecture)

- Example 2.14:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$  with  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ . Find the MLE of  $\theta = (\mu, \sigma^2)$ .

$$L(\mu, \sigma^2 | \vec{x}) = (2\pi\sigma^2)^{-n/2} \cdot \exp\left(-\frac{\sum(x_i - \mu)^2}{2\sigma^2}\right)$$

$$\Rightarrow l(\mu, \sigma^2 | \vec{x}) = c - \frac{n}{2} \log(\sigma^2) - \frac{\sum(x_i - \mu)^2}{2\sigma^2}$$

$$\Rightarrow S(\mu, \sigma^2 | \vec{x}) = \nabla l = \left( \frac{\partial S}{\partial \mu}, \frac{\partial S}{\partial \sigma^2} \right) = \left( \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu), \frac{-n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 \right) \stackrel{\text{set}}{=} \vec{0} = (0, 0)$$

solve tediously  
 $\Rightarrow (\hat{\mu}, \hat{\sigma}^2) = \left( \bar{x}, \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right)$

Second derivative test:

$$\frac{\partial^2 S}{\partial \mu^2} = -\frac{n}{\sigma^2} < 0$$

$$\frac{\partial^2 S}{\partial (\sigma^2)^2} = \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \sum_{i=1}^n (x_i - \mu)^2$$

$$\frac{\partial^2 S}{\partial \mu \partial \sigma^2} = -\frac{1}{\sigma^4} \sum_{i=1}^n (x_i - \mu)$$

The determinant of the Hessian is

$$\begin{vmatrix} \frac{\partial^2 S}{\partial \mu^2} & \frac{\partial^2 S}{\partial \mu \partial \sigma^2} \\ \frac{\partial^2 S}{\partial \mu \partial \sigma^2} & \frac{\partial^2 S}{\partial (\sigma^2)^2} \end{vmatrix} = \begin{vmatrix} -\frac{n}{\sigma^2} & -\frac{1}{\sigma^4} \sum_{i=1}^n (x_i - \mu) \\ -\frac{1}{\sigma^4} \sum_{i=1}^n (x_i - \mu) & \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \sum_{i=1}^n (x_i - \mu)^2 \end{vmatrix}_{\substack{\mu = \hat{\mu} \\ \sigma^2 = \hat{\sigma}^2}}$$

$$\dots = \frac{1}{\hat{\sigma}^6} \cdot \frac{n^2}{2} > 0. \text{ So } (\hat{\mu}, \hat{\sigma}^2) \text{ is the MLE.}$$

Note: we could've also done this by maximizing  $l(\mu, \sigma^2 | \vec{x})$  for fixed  $\sigma^2$  to find  $\hat{\mu}$ , and then maximizing  $l(\hat{\mu}, \sigma^2 | \vec{x})$  in  $\sigma^2$ .

This works since  $\mu$  and  $\sigma^2$  aren't functions of each other.

# Maximum Likelihood: More Examples

- The likelihood may not be differentiable, but that doesn't mean it can't be maximized
- Example 2.15:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(0, \theta)$  with  $\theta > 0$ . Find the MLE of  $\theta$ .

$$\begin{aligned} L(\theta | \vec{x}) &= \prod_{i=1}^n f_\theta(x_i) = \theta^{-n} \cdot \mathbb{1}_{0 \leq x_{(1)} < x_{(n)} \leq \theta} \\ &= \mathbb{1}_{0 \leq x_{(1)}} \cdot \theta^{-n} \cdot \mathbb{1}_{x_{(n)} \leq \theta} \end{aligned}$$

If  $\theta = x_{(n)}$ , then  $L(x_{(n)} | \vec{x}) = \mathbb{1}_{0 \leq x_{(1)} < x_{(n)}} \cdot x_{(n)}^{-n}$

If  $\theta > x_{(n)}$ , then  $L(\theta | \vec{x}) = \mathbb{1}_{0 \leq x_{(1)}} \cdot \theta^{-n} < \mathbb{1}_{0 \leq x_{(1)}} \cdot x_{(n)}^{-n} = L(x_{(n)} | \vec{x})$

If  $\theta < x_{(n)}$ , then  $L(\theta | \vec{x}) = \mathbb{1}_{0 \leq x_{(1)}} \cdot \theta^{-n} \cdot 0 = 0 < L(x_{(n)} | \vec{x})$

Hence  $\hat{\theta}(\vec{x}) = x_{(n)}$  is the MLE.

But we couldn't use calculus to find it, because  $L(\theta | \vec{x})$  is not differentiable at  $\theta = x_{(n)}$

(E/R Figure 6.2.2)

# Regression Through the Origin

- **Example 2.16:** Let  $Y_1, Y_2, \dots, Y_n$  be independent where  $Y_i \sim \mathcal{N}(\beta x_i, \sigma^2)$  with  $\beta \in \mathbb{R}$ ,  $x_i \in \mathbb{R}$ , and  $\sigma^2 > 0$ . Find the MLE of  $\beta$ .

all known!

$$L(\beta | \vec{y}) = \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \cdot \exp\left(-\frac{(y_i - \beta x_i)^2}{2\sigma^2}\right)$$
$$= (2\pi\sigma^2)^{-n/2} \cdot \exp\left(-\frac{\sum(y_i - \beta x_i)^2}{2\sigma^2}\right)$$

$$\Rightarrow l(\beta | \vec{y}) = c - \frac{\sum(y_i - \beta x_i)^2}{2\sigma^2} \quad \text{where } c \text{ is free of } \beta$$

Second derivative test:

$$\frac{\partial^2}{\partial \beta^2} S(\beta | \vec{y}) = -\frac{\sum x_i^2}{\sigma^2} < 0 \quad \forall \beta \in \mathbb{R}$$

$$\text{Hence } \hat{\beta}(\vec{y}) = \frac{\sum x_i y_i}{\sum x_i^2}.$$

$$\Rightarrow \sum x_i(y_i - \beta x_i) = 0$$

$$\Rightarrow \hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

- This is a particular case of **linear regression**; see Assignment 2 for more

# Reparameterization

- Instead of  $\theta$  itself, what if we want to find the MLE of some one-to-one function of the parameter  $\tau(\theta)$ ?
- Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$ , where  $\theta \in (0, 1)$ . Find the MLE of  $\theta^2$ .

Let  $\tau = \theta^2$ .

Then  $L(\tau | \vec{x}) = \tau^{n\bar{x}} (1-\tau)^{n-n\bar{x}}$

$$\Rightarrow l(\tau | \vec{x}) = n\bar{x} \cdot \log(\tau) + (n-n\bar{x}) \cdot \log(1-\tau)$$

$$\Rightarrow S(\tau | \vec{x}) = \frac{n\bar{x}}{2\tau} + \frac{n-n\bar{x}}{2(1-\tau)} \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \sqrt{\tau} = \bar{x}$$

$$\Rightarrow \tau = (\bar{x})^2 \quad \psi(\bar{x}) = (\bar{x}_n)^2 = (\Theta(\vec{x}))^2$$

Exercise: second derivative test.

# Reparameterization

- That wasn't a coincidence
- Theorem 2.1 (**Invariance Property**): If  $\hat{\theta}(\mathbf{X})$  is an MLE of  $\theta \in \Theta$  and  $\tau(\cdot)$  is one-to-one on  $\Theta$ , then the MLE of  $\tau(\theta)$  is given by  $\tau(\hat{\theta}(\mathbf{X}))$ .  
$$\hat{\tau}(\theta) = \tau(\hat{\theta}) \quad \text{"plug-in estimator"}$$

Proof. Let the likelihood under  $\theta$  be  $L(\theta | \vec{x}) = f_\theta(\vec{x})$ .

Let the likelihood under  $\tau = \tau(\theta)$  be  $L^*(\tau | \vec{x}) = g_\tau(\vec{x})$ .

If  $\hat{\theta}$  is the MLE of  $\theta$ , then

$$\begin{aligned} L^*(\tau(\hat{\theta}) | \vec{x}) &= g_{\tau(\hat{\theta})}(\vec{x}) = f_{\hat{\theta}}(\vec{x}) = L(\hat{\theta} | \vec{x}) \\ &\stackrel{\substack{\curvearrowleft \\ \hat{\theta} \text{ is MLE}}}{\geq} L(\theta | \vec{x}) \\ &= L^*(\tau(\theta) | \vec{x}) \quad \forall \theta \in \Theta. \end{aligned}$$

Hence  $\tau(\hat{\theta})$  is the MLE under the  $\tau$  parameterization.  $\square$

# Reparameterization

- **Example 2.17:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$  where  $p \in (0, 1)$ . Find the MLE of  $\tau(p) = \log\left(\frac{p}{1-p}\right)$ .

From Ex 2.11,  $\hat{p}_{\text{MLE}}(\vec{x}) = \bar{X}_n$ .

Since  $\log\left(\frac{p}{1-p}\right)$  is one-to-one, the MLE of  $\tau(p)$  is  $\log\left(\frac{\bar{X}_n}{1-\bar{X}_n}\right)$  by the invariance property.

# Poll Time!

# Maximum Likelihood Estimation

- Maximum likelihood is *by far* the most common method that statisticians use to find point estimates<sup>1</sup>
- Maximum likelihood estimators tend to have quite good properties (especially for large sample sizes):

GOOD:  $\hat{\theta}$  is always in  $\Theta$  (by definition)

GOOD: Widely applicable; requires very few assumptions (don't need moments, etc.)

GOOD: (Relatively) easy to implement numerically (usually just general-purpose optimization software)

"BAD":  $\hat{\theta}(\bar{x})$  may not have the "right" expectation, especially when  $n$  is small

( $\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$  does not have expectation  $\sigma^2$ )

BAD: numerical optimizers have pitfalls (e.g., starting values, poor in high dimensions, etc.)

- When in doubt, it's usually a good idea to use maximum likelihood if you can

---

<sup>1</sup> Assuming those statisticians aren't Bayesians – more on that in Module 6

# Evaluating Estimators

- Back to the idea of what makes a point estimator “good”
  - From now on, we focus on point estimators of  $\tau(\theta)$ , rather than  $\theta$
  - It turns out there's a much more convenient way to assess the quality of a point estimator estimator than our earlier thoughts
  - Consider the *error* (or *absolute deviation*) of an estimator  $|T(\mathbf{X}) - \tau(\theta)|$ , which is of course a random variable
  - It's too much to ask for this to *always* be small; some random sample  $\mathbf{X}_j$  may be an “outlier”, so that  $T(\mathbf{X}_j)$  is far from  $\tau(\theta)$
  - But we can ask for it to be small on average i.e.,  $E_\theta[|T(\vec{\mathbf{x}}) - \tau(\theta)|]$  small.

# Mean-Squared Error

- In other words, it's reasonable to ask for  $\mathbb{E}_\theta [|T(\mathbf{X}) - \tau(\theta)|]$  to be small
- That's fine, but it turns out that for mathematical reasons, it's much more convenient to ask for the *squared error*  $(T(\mathbf{X}) - \tau(\theta))^2$  to be small on average
- **Definition 2.7:** Let  $T(\mathbf{X})$  be an estimator for  $\tau(\theta)$ . The **mean-squared error (MSE)** is defined as

$$\text{MSE}_\theta (T(\mathbf{X})) = \mathbb{E}_\theta [(T(\mathbf{X}) - \tau(\theta))^2].$$

- So why not look for the  $T(\mathbf{X})$  that minimizes the MSE for all  $\theta \in \Theta$ ?
- Because unfortunately, such a  $T(\mathbf{X})$  almost never exists
- Let's try to restrict the class of estimators under consideration to one where minimizers of the MSE are easier to find

# Bias

- Definition 2.8: The **bias** of a point estimator  $T(\mathbf{X})$  is defined as

$$\text{Bias}_\theta(T(\mathbf{X})) = \mathbb{E}_\theta[T(\mathbf{X})] - \tau(\theta).$$

If  $\text{Bias}_\theta(T(\mathbf{X})) = 0$ , then  $T(\mathbf{X})$  is said to be an **unbiased estimator** of  $\tau(\theta)$ .

$$\Rightarrow \mathbb{E}_\theta[T(\bar{\mathbf{x}})] = \tau(\theta)$$

- Example 2.18:

$X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ . Then  $T(\bar{\mathbf{x}}) = \bar{X}_n$  is unbiased for  $\mu$ .

$X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$ . Then  $T(\bar{\mathbf{x}}) = \bar{X}_n$  is unbiased for  $p$ ,

$$\text{because } \text{Bias}_p(\bar{X}_n) = \mathbb{E}_p[\bar{X}_n] - p$$

$$= \frac{1}{n} \cdot np - p$$

$$= 0$$

- Example 2.19:

$X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ .  $\tau(\sigma^2) = \sigma^2$

$$\text{Bias}_{\sigma^2}(\hat{\sigma}^2_{\text{MLE}}) = \text{Bias}_{\sigma^2}\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right) = \left(\frac{n-1}{n}\right) \sigma^2 - \sigma^2 = \frac{\sigma^2}{n} \neq 0. \text{ So the MLE is biased!}$$

# Unbiased Estimators Don't Always Exist

- **Example 2.20:** Let  $X \sim \text{Bernoulli}(\theta)$ , where  $\theta \in (0, 1)$ . There exists no unbiased estimator of  $\tau(\theta) = \frac{1}{\theta}$ .

Suppose  $T(x)$  were unbiased for  $\tau_\theta$ .

$$\begin{aligned}\frac{1}{\theta} &= E_\theta[T(X)] = T(0) \cdot P_\theta(X=0) + T(1) \cdot P_\theta(X=1) \\ &= T(0) \cdot (1-\theta) + T(1) \cdot \theta \quad \forall \theta \in (0,1). \quad \times\end{aligned}$$

Argument 1:  $\frac{1}{\theta} = \frac{1}{1-(1-\theta)} = \sum_{j=0}^{\infty} (1-\theta)^j = 1 + (1-\theta) + (1-\theta)^2 + \dots$

But we can't have a polynomial  $\times$   
equal to an infinite series!

Argument 2:  $\frac{1}{\theta}$  is unbounded for  $\theta \in (0,1)$  but ~~it~~ is clearly bounded.

# The Bias-Variance Tradeoff

- Theorem 2.2 (**Bias-Variance Tradeoff**): If a point estimator  $T(\mathbf{X})$  has a finite second moment, then

$$\text{MSE}_\theta(T(\mathbf{X})) = \text{Bias}_\theta(T(\mathbf{X}))^2 + \text{Var}_\theta(T(\mathbf{X})).$$

Proof.

$$\begin{aligned}\text{MSE}_\theta(T(\vec{\mathbf{x}})) &= \mathbb{E}_\theta\{ (T(\vec{\mathbf{x}}) - \varepsilon(\theta))^2 \} \\ &= \text{Var}_\theta(T(\vec{\mathbf{x}}) - \varepsilon(\theta)) + \mathbb{E}\{ T(\vec{\mathbf{x}}) - \varepsilon(\theta) \}^2 \\ &= \text{Var}_\theta(T(\vec{\mathbf{x}})) + \text{Bias}_\theta(T(\vec{\mathbf{x}}))^2.\end{aligned}$$

□

So among all estimators with a fixed MSE, you can choose between more accuracy (< bias)  
+ less precision ( $\Rightarrow$  variance)  
Or vice versa.

# Poll Time!

# Best Unbiased Estimation

- So let's restrict our attention to the class of unbiased estimators, and *then* choose the one (or ones?) with the lowest MSE
- Equivalently, choose the unbiased estimator (or estimators?) with the lowest variance
- Definition 2.9:** An unbiased estimator  $T^*(\mathbf{X})$  of  $\tau(\theta)$  is a **best unbiased estimator** of  $\tau(\theta)$  if

$$\text{Var}_\theta(T^*(\mathbf{X})) \leq \text{Var}_\theta(T(\mathbf{X})) \quad \text{for all } \theta \in \Theta$$

where  $T(\mathbf{X})$  is any other unbiased estimator of  $\tau(\theta)$ . A best unbiased estimator is also called a **uniform minimum variance unbiased estimator (UMVUE)** of  $\tau(\theta)$ .

*HOE ( $\hat{\theta}$ )*      *lowest variance out of all unbiased estimators*

When we say "best", this is what we mean!

# Questions That We Will Answer

- How do we know whether or not an estimator  $T(\mathbf{X})$  is a UMVUE for  $\tau(\theta)$ ?
- How do we find a UMVUE for  $\tau(\theta)$ ?
- Are UMVUEs unique?

# An Ubiquitous Inequality in Mathematics

- Theorem 2.3 (**Cauchy-Schwarz Inequality**): Let  $X$  and  $Y$  be random variables, each having finite, nonzero variance. Then

$$|\text{Cov}(X, Y)| \leq \sqrt{\text{Var}(X) \text{Var}(Y)}.$$

Furthermore, if  $\text{Var}(Y) > 0$ , then equality is attained if and only if  $X$  and  $Y$  are linearly related. *In particular, if  $X = \frac{\text{Cov}(X, Y)}{\text{Var}(Y)} \cdot Y + a$  for some  $a \in \mathbb{R}$ .*

*Proof.*

Proof: let  $\mu_x = \mathbb{E}[X]$ ,  $\mu_y = \mathbb{E}[Y]$ .

$$\begin{aligned} \text{If } \text{Var}(Y) = 0, \text{ then } Y &= \mu_y. \text{ Then } \text{Cov}(X, Y) = \mathbb{E}[(X - \mu_x)(Y - \mu_y)] \\ &= 0 = \sqrt{\text{Var}(X) \cdot \text{Var}(Y)} \end{aligned}$$

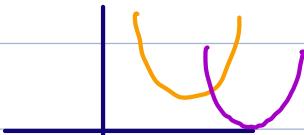
So suppose  $\text{Var}(Y) > 0$ . Let  $Z = X - \mu_x$  and  $W = Y - \mu_y$ .

Let  $t \in \mathbb{R}$ .

$$\begin{aligned} \text{Then } \mathbb{E}[(Z - tW)^2] &= \mathbb{E}[Z^2] - 2t \cdot \mathbb{E}[ZW] + t^2 \cdot \mathbb{E}[W^2] \\ &= \text{Var}(X) - 2t \cdot \text{Cov}(X, Y) + t^2 \cdot \text{Var}(Y) \end{aligned}$$

This is a quadratic in  $t$  which is non-negative, so it has at most

one real root.



Therefore,

$$4 \cdot \text{Cov}(X, Y)^2 - 4 \cdot \text{Var}(X) \cdot \text{Var}(Y) \leq 0$$

$$\Rightarrow \text{Cov}(X, Y)^2 \leq \text{Var}(X) \cdot \text{Var}(Y)$$

$$\Rightarrow |\text{Cov}(X, Y)| \leq \sqrt{\text{Var}(X) \cdot \text{Var}(Y)}$$

Equality holds iff  $\mathbb{E}[(Z - tW)^2] = 0$  for some  $t$ .

But  $(Z - tW)^2 \geq 0$ , so this holds iff  $Z = tW$

$$\Leftrightarrow X - \mu_x = tY - t\mu_y$$

$$\Leftrightarrow X = tY - t\mu_y + \mu_x \quad (\text{i.e., } X \text{ and } Y \text{ are linearly related}).$$

Moreover,  $\text{Cov}(X, Y) = \mathbb{E}[ZW] = \mathbb{E}[tW^2] = t \cdot \mathbb{E}[W^2] = t \cdot \text{Var}(Y)$

$$\Rightarrow t = \frac{\text{Cov}(X, Y)}{\text{Var}(Y)} \quad \square$$

# UMVUEs Are Unique

- Theorem 2.4: If a UMVUE exists for  $\tau(\theta)$ , then it is unique.

Proof. Let  $W$  and  $W'$  be two UMVUEs of  $\tau(\theta)$ . Let  $W^* = \frac{1}{2}(W + W')$ .

Then  $W^*$  is unbiased for  $\tau(\theta)$ , and moreover

$$\begin{aligned}\text{Var}_\theta(W^*) &= \frac{1}{4} \cdot \text{Var}_\theta(W) + \frac{1}{4} \cdot \text{Var}_\theta(W') + \frac{1}{2} \cdot \text{Cov}_\theta(W, W') \\ &\leq \frac{1}{4} \text{Var}_\theta(W) + \frac{1}{4} \cdot \text{Var}_\theta(W') + \frac{1}{2} \sqrt{\text{Var}_\theta(W) \cdot \text{Var}_\theta(W')} \quad \text{by Cauchy-Schwarz} \\ &= \text{Var}_\theta(W) \quad \text{since all variances are equal (by ass.)} \quad \text{ie, } \text{Cov}_\theta(W, W') = \text{Var}_\theta(W) \quad \textcircled{1}\end{aligned}$$

But  $W^*$  can't be a UMVUE, so equality must hold  $\Rightarrow W' = a(\theta) \cdot W + b(\theta)$ .

What are  $a(\theta)$  and  $b(\theta)$ ?

$$\textcircled{1} \text{ implies } \text{Var}_\theta(W) = \text{Cov}_\theta(W, W')$$

$$\begin{aligned}&= \text{Cov}_\theta(W, a(\theta) \cdot W + b(\theta)) \\ &= \text{Cov}_\theta(W, a(\theta) \cdot W) \\ &= a(\theta) \cdot \text{Var}_\theta(W) \\ \Rightarrow a(\theta) &= 1\end{aligned}$$

$$\begin{aligned}\text{Finally, } \tau(\theta) &= \mathbb{E}[W'] \\ &= \mathbb{E}[1 \cdot W + b(\theta)] \\ &= \tau(\theta) + b(\theta) \\ \Rightarrow b(\theta) &= 0.\end{aligned}$$

$$\text{Hence } W = W'. \quad \square$$

# The Rao-Blackwell Theorem

- It turns out that sufficiency can help us in our search for the UMVUE in powerful ways
- Theorem 2.5 (**Rao-Blackwell**): Let  $W(\mathbf{X})$  be unbiased for  $\tau(\theta)$ , and let  $T(\mathbf{X})$  be sufficient for  $\theta$ . Define  $W_T(\mathbf{X}) = \mathbb{E}_\theta [W(\mathbf{X}) | T(\mathbf{X})]$ . Then  $W_T(\mathbf{X})$  is also an unbiased point estimator of  $\tau(\theta)$ , and moreover,  
 $\text{Var}_\theta (W_T(\mathbf{X})) \leq \text{Var}_\theta (W(\mathbf{X})).$

Proof. Unbiasedness:  $\mathbb{E}_\theta [W_T(\vec{x})] = \mathbb{E}_\theta [\underbrace{\mathbb{E}_\theta [W(\vec{x}) | T(\vec{x})]}_{\text{tower property}}] = \mathbb{E}_\theta [W(\vec{x})] = \tau(\theta).$

Moreover,  $\text{Var}_\theta (W(\vec{x})) = \mathbb{E}_\theta [\underbrace{\text{Var}_\theta (W(\vec{x}) | T(\vec{x}))}_{\geq 0}] + \text{Var}_\theta (\mathbb{E}_\theta [W(\vec{x}) | T(\vec{x})])$   
 $\geq \text{Var}_\theta (\mathbb{E}_\theta [W(\vec{x}) | T(\vec{x})])$   
 $= \text{Var}_\theta (W_T(\vec{x})). \quad \square$

What about sufficiency? If  $T(\vec{x})$  weren't sufficient for  $\theta$ , then

$\mathbb{E}_\theta [W(\vec{x}) | T(\vec{x})]$  wouldn't be free of  $\theta \rightarrow$  not a point estimator!

# Interpreting Rao-Blackwellization

- The process of replacing an estimator with its conditional expectation (with respect to a sufficient statistic) is called **Rao-Blackwellization**
- Theorem 2.5 says that we can always improve on (or at least make no worse) any unbiased estimator  $W(\mathbf{X})$  with a second moment by Rao-Blackwellizing it
- Example 2.21: Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Poisson}(\lambda)$ ,  $\lambda > 0$ .

We have two unbiased estimators for  $\lambda$ :  $\bar{X}_n$  and  $S^2$ .

But  $\bar{X}_n$  is sufficient for  $\lambda$  by Theorem 1.2, so

$E[S^2(\bar{X}_n)]$  is better than  $S^2$  itself, but by Theorem 2.4, it can't be better than  $E[\bar{X}_n | \bar{X}_n] = \bar{X}_n$ .

# Rao-Blackwell: Examples

$$\begin{aligned} \sum_{i=1}^n X_i &\sim \text{Bin}(nK, \theta) \\ \Rightarrow \sum_{i=2}^n X_i &\sim \text{Bin}((n-1)K, \theta) \end{aligned}$$

- **Example 2.22:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bin}(k, \theta)$ , where  $\theta \in (0, 1)$  and  $k$  is known. Let  $\tau(\theta) = k\theta(1-\theta)^{k-1}$ . Show that  $W(\mathbf{X}) = \mathbb{1}_{X_1=1}$  is unbiased for  $\tau(\theta)$ , and then Rao-Blackwellize it.

Unbiased:  $\mathbb{E}_\theta[W(\vec{X})] = P_\theta(X_1=1) = k\theta(1-\theta)^{k-1} = \tau(\theta)$ . Recall  $T(\vec{X}) = \sum_{i=1}^n X_i$  is sufficient for  $\theta$ .

Take  $W_T(\vec{X}) = \mathbb{E}\{W(\vec{X}) | T(\vec{X})\}$ . How do we use this?

Suppose  $T(\vec{X}) = t$ .

Then  $\mathbb{E}\{W(\vec{X}) | T(\vec{X}) = t\}$

$$= P_\theta(X_1=1 | \sum_{i=1}^n X_i = t)$$

$$= \frac{P_\theta(X_1=1 \wedge \sum_{i=1}^n X_i = t)}{P_\theta(\sum_{i=1}^n X_i = t)}$$

$$= \frac{P_\theta(X_1=1 \wedge \sum_{i=2}^n X_i = t-1)}{P_\theta(\sum_{i=1}^n X_i = t)}$$

$$\begin{aligned} &\stackrel{\text{indep}}{=} \frac{P_\theta(X_1=1) \cdot P_\theta(\sum_{i=2}^n X_i = t-1)}{P_\theta(\sum_{i=1}^n X_i = t)} \\ &= \frac{k\theta(1-\theta)^{k-1} \cdot \binom{k(n-1)}{t-1} \cdot \theta^{t-1} (1-\theta)^{k(n-1)-(t-1)}}{\binom{kn}{t} \cdot \theta^t (1-\theta)^{kn-t}} \\ &= \frac{k \binom{k(n-1)}{t-1}}{\binom{kn}{t}}. \end{aligned}$$

So our Rao-Blackwellized estimator is

$$W_T(\vec{X}) = \frac{k \binom{k(n-1)}{\sum_{i=1}^n X_i - 1}}{\binom{kn}{\sum_{i=1}^n X_i}}.$$

# The Lehmann-Scheffé Theorem

"based on  $T(\bar{X})$ "

- Theorem 2.6 (**Lehmann-Scheffé Theorem**): Let  $W(\mathbf{X})$  be unbiased for  $\tau(\theta)$  and let  $T(\mathbf{X})$  be a complete sufficient statistic, for all  $\theta \in \Theta$ . Then  $W_T(\mathbf{X}) = \mathbb{E}_{\bar{X}}[W(\mathbf{X}) | T(\mathbf{X})]$  is the unique UMVUE.

Proof. Suppose  $V_T(\bar{X}) = \mathbb{E}[V(\bar{X}) | T(\bar{X})]$  is also a UMVUE, where  $V(\bar{X})$  is unbiased for  $\tau(\bar{X})$  (which we can do by Rao-Blackwell).

$$\begin{aligned} 0 &= \mathbb{E}_\theta[V_T(\bar{X})] - \mathbb{E}_\theta[W_T(\bar{X})] \\ &= \mathbb{E}_\theta[\mathbb{E}_\theta[V(\bar{X}) | T(\bar{X})] - \mathbb{E}_\theta[W(\bar{X}) | T(\bar{X})]] \\ &= \mathbb{E}_\theta[\underbrace{\mathbb{E}_\theta[V(\bar{X}) - W(\bar{X}) | T(\bar{X})]}_{=: h(T(\bar{X}))}] \\ &= \mathbb{E}_\theta[h(T(\bar{X}))] \quad \forall \theta \in \Theta \end{aligned}$$

$\Rightarrow$  By completeness,  $\mathbb{P}_\theta(h(T(\bar{X})) = 0) = 1 \quad \forall \theta \in \Theta$ .

$$\Rightarrow V_T(\bar{X}) = W_T(\bar{X}) \quad \forall \theta \in \Theta.$$

So the UMVUE is  $\mathbb{E}\{W(\bar{X}) | T(\bar{X})\}$ .  $\square$

# More On Lehmann-Scheffé

- This is a bit startling
- If we take some unbiased estimator and condition it on a complete sufficient statistic, then the resulting estimator is the UMVUE
- As such, if we find an unbiased estimator  $T(\mathbf{X})$  of  $\tau(\theta)$  which is also a complete sufficient statistic, then we're done

Also, if  $V(\tilde{\mathbf{x}})$  and  $W(\tilde{\mathbf{x}})$  are both unbiased for  $\tau(\theta)$  and  $T(\tilde{\mathbf{x}})$  is complete, then  $E\{V(\tilde{\mathbf{x}}) | T(\tilde{\mathbf{x}})\} = E\{W(\tilde{\mathbf{x}}) | T(\tilde{\mathbf{x}})\}$ .

Eg: Poisson:  $E\{\bar{X}_n | \bar{X}_n\} = E\{S^2 | \bar{X}_n\}$

Eg: If  $T(\tilde{\mathbf{x}})$  is complete and  $r(T(\tilde{\mathbf{x}}))$  is unbiased for  $\tau(\theta)$ ,  
then  $E_r[r(T(\tilde{\mathbf{x}})) | T(\tilde{\mathbf{x}})] = r(T(\tilde{\mathbf{x}}))$  is best for  $\tau(\theta)$ .

# Lehmann-Scheffé: Examples

- **Example 2.23:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$  with  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ . Find the UMVUE of  $(\mu, \sigma^2)$ .

We know that  $(\bar{X}_n, S^2)$  is complete (eg, Ex 1.29, Theorem 1.28, A1 Q18).

Also  $\bar{X}_n$  is unbiased for  $\mu$ ,  $S^2$  is unbiased for  $\sigma^2$ .

$\Rightarrow (\bar{X}_n, S^2)$  is unbiased for  $(\mu, \sigma^2)$ .

By Lehmann-Scheffé,  $T(\bar{x}) = (\bar{X}_n, S^2)$  is the UMVUE for  $(\mu, \sigma^2)$ .

(That's not the MLE!)

## Lehmann-Scheffé: Examples

- **Example 2.24:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$ , where  $\lambda > 0$ . Find the UMVUE of  $\lambda$ .

We know that  $\bar{X}_n$  is unbiased for  $\lambda$ , and also a complete sufficient statistic.

By Lehmann-Scheffé, the UMVUE of  $\lambda$  is  $\mathbb{E}[\bar{X}_n | \bar{X}_n] = \bar{X}_n$ .

# Poll Time!

# What About the Likelihood?

- Rao-Blackwellization and Lehmann-Scheffé tell us how to get the unique UMVUE (if it exists) via sufficient statistics
- The likelihood wasn't involved
- It turns out there exists a very helpful tool that helps us with finding the UMVUE (if it exists) by exploiting the likelihood
- It doesn't always work, but when it does, it works like a charm
- But we need several auxiliary results to produce it

# The Covariance Inequality

- Theorem 2.7 (**Covariance Inequality**): Let  $T(\mathbf{X})$  and  $U(\mathbf{X})$  be two statistics such that  $0 < \mathbb{E}_\theta [T(\mathbf{X})^2], \mathbb{E}_\theta [U(\mathbf{X})^2] < \infty$  for all  $\theta \in \Theta$ . Then

$$\text{Var}_\theta (T(\mathbf{X})) \geq \frac{\text{Cov}_\theta (T(\mathbf{X}), U(\mathbf{X}))^2}{\text{Var}_\theta (U(\mathbf{X}))} \quad \text{for all } \theta \in \Theta.$$

Equality holds if and only if

$$T(\mathbf{X}) = \mathbb{E}_\theta [T(\mathbf{X})] + \frac{\text{Cov}_\theta (T(\mathbf{X}), U(\mathbf{X}))}{\text{Var}_\theta (U(\mathbf{X}))} (U(\mathbf{X}) - \mathbb{E}_\theta [U(\mathbf{X})])$$

almost surely. ← Can ignore!

Proof. Apply Cauchy-Schwarz to  $X = T(\vec{X})$  and  $Y = U(\vec{X})$ ,  
and square everything. □

# The Fisher Information

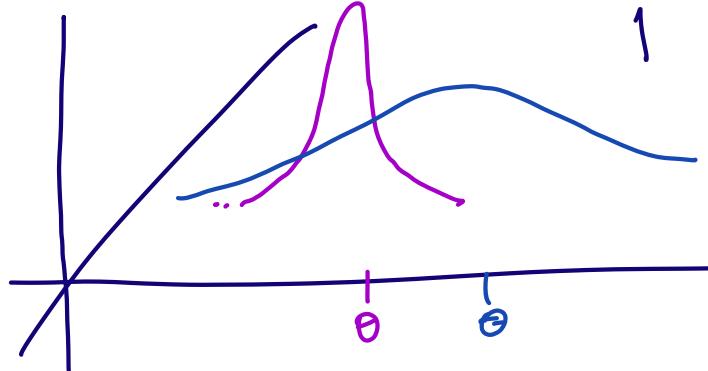
- **Definition 2.10:** Let  $\mathbf{X} = (X_1, \dots, X_n) \sim f_\theta$ , and let  $S(\theta | \mathbf{x})$  be the score function for the parametric model. The **(expected) Fisher information** is the function  $I_n : \Theta \rightarrow [0, \infty)$  defined by

$$I_n(\theta) = \text{Var}_\theta(S(\theta | \mathbf{X})).$$

↑ measures the "curvature" of  
the log-likelihood surface  
at  $\theta$

- **Definition 2.11:** Let  $\mathbf{X} = (X_1, \dots, X_n) \sim f_\theta$ , and let  $S(\theta | \mathbf{x})$  be the score function for the parametric model. The **observed Fisher information** is the function  $J_n : \mathcal{X}^n \rightarrow [0, \infty)$  defined by

$$J_n(\mathbf{X}) = -\frac{\partial}{\partial \theta} S(\theta | \mathbf{X}_{\text{obs}}) \Big|_{\theta=\hat{\theta}_{\text{MLE}}(\mathbf{X})}.$$



# The Fisher Information: Examples

- **Example 2.25:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$ , where  $\lambda > 0$ . Calculate the observed and expected Fisher information for  $\lambda$ .

$$L(\lambda | \vec{x}) = \prod_i \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$$

$$\Rightarrow l(\lambda | \vec{x}) = \sum x_i \log(\lambda) - n\lambda + c, \text{ where } c \text{ is free of } \lambda$$

$$\Rightarrow S(\lambda | \vec{x}) = \frac{\sum x_i}{\lambda} - n$$

$$\begin{aligned} I_n(\lambda) &= \text{Var}_\lambda(S(\lambda | \vec{x})) \\ &= \text{Var}_\lambda\left(\frac{\sum x_i}{\lambda} - n\right) \\ &= \frac{1}{\lambda^2} \cdot \text{Var}_\lambda(\sum x_i) \\ &= \frac{n\lambda}{\lambda^2} = \frac{n}{\lambda}. \end{aligned}$$

Recall  $\hat{\lambda}_{MLE} = \bar{X}_n$ .

$$\text{Then } -\frac{\partial}{\partial \lambda} S(\lambda | \vec{x}) = \frac{\sum x_i}{\lambda^2}$$

$$\Rightarrow J_n(\bar{X}) = \left. \frac{\sum x_i}{\lambda^2} \right|_{\lambda=\bar{X}_n}$$

$$= \frac{n\bar{X}_n}{(\bar{X}_n)^2} = \frac{n}{\bar{X}_n} = \frac{n^2}{\sum x_i}.$$

# The Fisher Information: Examples

- **Example 2.26:** 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. Calculate the observed and expected Fisher information for  $\mu$ .

From Ex 2.12,  $S(\mu | \vec{x}) = \frac{\sum x_i - np}{\sigma^2}$

$$\begin{aligned} I_n(\mu) &= \text{Var}_\mu \left( \frac{\sum x_i - np}{\sigma^2} \right) \\ &= \frac{1}{\sigma^4} \cdot \text{Var} \left( \sum x_i \right) \\ &= \frac{n}{\sigma^2}. \end{aligned}$$

$$\hat{\mu}_{\text{MLE}} = \bar{x}_n. \text{ So}$$

$$\begin{aligned} J_n(\vec{x}) &= -\frac{\partial}{\partial \mu} S(\mu | \vec{x}) \Big|_{\mu=\bar{x}_n} \\ &= \frac{n}{\sigma^2} \Big|_{\mu=\bar{x}_n} \\ &= \frac{n}{\sigma^2}. \end{aligned}$$

They're usually not the same!

# The Cramér-Rao Lower Bound (CRLB)

- Theorem 2.8 (**Cramér-Rao Lower Bound**): Let  $\mathbf{X} = (X_1, \dots, X_n) \sim f_\theta$ , and let  $T(\mathbf{X})$  be any estimator such that  $= \frac{d}{d\theta} \int_{\mathcal{X}} \tau(x) \cdot f_\theta(x) dx$

$$\text{Var}_\theta(T(\mathbf{X})) < \infty \quad \text{and} \quad \underbrace{\frac{d}{d\theta} \mathbb{E}_\theta [T(\mathbf{X})]}_{= \frac{d}{d\theta} \int_{\mathcal{X}} \tau(x) f_\theta(x) dx} = \int_{\mathcal{X}} \frac{\partial}{\partial \theta} [T(\mathbf{x}) f_\theta(\mathbf{x})] d\mathbf{x}.$$

Then

$$\text{Var}_\theta(T(\mathbf{X})) \geq \frac{\left( \frac{d}{d\theta} \mathbb{E}_\theta [T(\mathbf{X})] \right)^2}{I_n(\theta)}. \quad \begin{aligned} \mathbb{E}_\theta [\tau(\bar{x})] \\ = \tau(\theta) \\ \Rightarrow \frac{d}{d\theta} (\mathbb{E}_\theta [\tau(\bar{x})]) = \tau'(\theta) \end{aligned}$$

In particular, if  $T(\mathbf{X})$  is unbiased for  $\tau(\theta)$  and  $\tau(\cdot)$  is differentiable on  $\Theta$ , then

$$\text{Var}_\theta(T(\mathbf{X})) \geq \frac{(\tau'(\theta))^2}{I_n(\theta)}.$$

*Proof.*

# The Cramér-Rao Lower Bound

Proof: In the Covariance Inequality, choose  $U(\vec{x}) = S(\theta|\vec{x}) = \frac{\partial}{\partial \theta} l(\theta|\vec{x})$ .

Then  $\text{Cov}_{\theta}(T(\vec{x}), S(\theta|\vec{x}))$

$$= \underbrace{\mathbb{E}_{\theta}[T(\vec{x}) \cdot S(\theta|\vec{x})]}_{\textcircled{1}} - \mathbb{E}_{\theta}[T(\vec{x})] \cdot \mathbb{E}_{\theta}[S(\theta|\vec{x})] = \frac{1}{d\theta} \mathbb{E}_{\theta}[T(\vec{x})].$$

$$\begin{aligned}\textcircled{1} &= \int_{\vec{x}} T(\vec{x}) \left( \frac{\partial}{\partial \theta} \log(f_{\theta}(\vec{x})) \right) \cdot f_{\theta}(\vec{x}) d\vec{x} \\ &= \int_{\vec{x}} T(\vec{x}) \cdot \left( \frac{1}{f_{\theta}(\vec{x})} \cdot \frac{\partial}{\partial \theta} f_{\theta}(\vec{x}) \right) \cdot f_{\theta}(\vec{x}) d\vec{x} \\ &= \int_{\vec{x}} T(\vec{x}) \cdot \frac{\partial}{\partial \theta} f_{\theta}(\vec{x}) d\vec{x} \\ \text{ass.} &= \frac{1}{d\theta} \int_{\vec{x}} T(\vec{x}) \cdot f_{\theta}(\vec{x}) d\vec{x} \\ &= \frac{1}{d\theta} \mathbb{E}_{\theta}[T(\vec{x})].\end{aligned}$$

$$\begin{aligned}\textcircled{2} &= \int_{\vec{x}} \left( \frac{\partial}{\partial \theta} \log(f_{\theta}(\vec{x})) \right) \cdot f_{\theta}(\vec{x}) d\vec{x} \\ &= \int_{\vec{x}} 1 \cdot \frac{\partial}{\partial \theta} f_{\theta}(\vec{x}) d\vec{x} \\ \text{ass.} &= \frac{1}{d\theta} \int_{\vec{x}} f_{\theta}(\vec{x}) d\vec{x} \\ &= \frac{1}{d\theta} 1 = 0.\end{aligned}$$

Also  $\text{Var}_{\theta}(S(\theta|\vec{x})) \stackrel{\text{def}}{=} I_n(\theta)$ . Plug into the Covariance Inequality to get the result.  $\square$

# The Cramér-Rao Lower Bound Conditions

- Unfortunately, the conditions of the Cramér-Rao Lower Bound don't always hold
- The first says that our estimator must actually have a variance to minimize, which seems reasonable
- Example 2.27: If  $X_1, \dots, X_n \sim N(\theta, \sigma^2)$ , don't try  $T(\vec{x}) = \frac{x_1}{x_n}$ .  
*Won't work!*
- The second says that we need to be able to push a derivative inside an integral, which is more subtle
- When would this condition fail to hold?
- Example 2.28:  $\text{Unif}(0, \theta) \Rightarrow$  support depends on  $\theta$   
 $\Rightarrow \frac{d}{d\theta} E_\theta[T(\vec{x})] \neq \int_0^\theta (\frac{2}{\theta} T(\vec{x}) \cdot \frac{1}{\theta}) d\vec{x}$   
*in general.*

# Easing the Computation

- Theorem 2.9: Under the conditions of Theorem 2.8,

$$I_n(\theta) = \mathbb{E}_\theta [S(\theta | \mathbf{X})^2].$$

Proof.  $I_n(\theta) = \text{Var}_\theta(S(\theta | \mathbf{x}))$

$$\begin{aligned} &= \mathbb{E}_\theta [S(\theta | \mathbf{x})^2] - \underbrace{\mathbb{E}_\theta [S(\theta | \mathbf{x})]}_{=0 \text{ from CLRB proof}}^2 \\ &= \mathbb{E}_\theta [S(\theta | \mathbf{x})^2]. \end{aligned}$$

□

- Theorem 2.10: If  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$  and conditions of Theorem 2.8 hold,

$$I_n(\theta) = n\mathbb{E}_\theta [S(\theta | X)^2]. \quad \text{Exercise!}$$

# More Easing

- Theorem 2.11 (**Second Bartlett Identity**): If  $X \sim f_\theta$  and  $f_\theta$  satisfies

$$\frac{d}{d\theta} \mathbb{E}_\theta [S(\theta | X)] = \int_{\mathcal{X}} \frac{\partial}{\partial \theta} [S(\theta | x) f_\theta(x)] dx,$$

(which is true when  $f_\theta$  is in an exponential family) then

$$\mathbb{E}_\theta [S(\theta | X)^2] = -\mathbb{E}_\theta \left[ \frac{\partial}{\partial \theta} S(\theta | X) \right].$$

Proof. RHS =  $-\mathbb{E}_\theta \left[ \frac{\partial}{\partial \theta} (\frac{1}{f_\theta(x)} \log(f_\theta(x))) \right]$

$$= -\mathbb{E}_\theta \left[ \frac{\partial}{\partial \theta} \left( \frac{1}{f_\theta(x)} \cdot \frac{\partial}{\partial \theta} f_\theta(x) \right) \right]$$

$$= -\mathbb{E}_\theta [\dots - \dots]$$

You finish off! Hint: use the assumption somewhere!

□

# Efficiency

unbiased

- **Definition 2.12:** An estimator  $T(\mathbf{X})$  of  $\tau(\theta)$  that attains the Cramér-Rao Lower Bound is called an **efficient estimator of  $\tau(\theta)$** .
- What's the connection between UMVUEs and efficient estimators?
- If an efficient estimator exists, then it must be the UMVUE
- But an efficient estimator doesn't always exist, as we'll soon see

# Efficiency: Examples

(Completed after lecture)

- **Example 2.29:** 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  $T(\mathbf{X})$  is an efficient estimator for  $\mu$ .

$$\text{I.I. } \bar{X}_n.$$

We need to calculate the CRLB and  $\text{Var}_\mu(\bar{X}_n)$ .

Know  $\text{Var}_\mu(\bar{X}_n) = \sigma^2/n$ .

What about the CRLB?

Numerator:  $\left( \frac{d}{d\mu} \mathbb{E}_\mu[T(\bar{X}_n)] \right)^2 = \left( \frac{d}{d\mu} \mu \right)^2 = 1$ .

Denominator:  $I_\mu(\mu) = n/\sigma^2$  from Example 2.26.

So the CRLB is  $1/n\sigma^2 = \frac{\sigma^2}{n} = \text{Var}_\mu(\bar{X}_n)$ . Hence  $\bar{X}_n$  is efficient for  $\mu$ .

# A Criterion for Efficiency

- Is there a better way to find efficient estimators than simply making an educated guess?
- Theorem 2.12: Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$  satisfy the conditions of Theorem 2.8. An unbiased estimator  $T(\mathbf{X})$  of  $\tau(\theta)$  is efficient if and only if there exists some function  $a : \Theta \rightarrow \mathbb{R}$  such that

$$S(\theta | \mathbf{x}) = a(\theta)[T(\mathbf{x}) - \tau(\theta)].$$

Proof. From the Covariance inequality, equality in the CRLB iff

$$T(\vec{x}) = \mathbb{E}_\theta[T(\vec{x})] + \frac{\text{Cov}(\tau(\vec{x}), S(\theta|\vec{x}))}{\text{Var}_\theta(S(\theta|\vec{x}))} (S(\theta|\vec{x}) - \mathbb{E}_\theta[S(\theta|\vec{x})])$$

$$= \tau(\theta) + \frac{\tau'(\theta)^2}{I_n(\theta)} \cdot S(\theta|\vec{x})$$

if  $S(\theta|\vec{x}) = \underbrace{\frac{I_n(\theta)}{\tau'(\theta)^2}}_{=: a(\theta)} (\tau(\vec{x}) - \tau(\theta)). \quad \square$

# Efficiency: Examples

- **Example 2.30:** 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 there exists no efficient estimator of  $\sigma^2$ .

If there were, then there'd be some  $T(\bar{x})$  a( $\sigma^2$ ) s.t.

$$S(\sigma^2 | \bar{x}) = a(\sigma^2) \cdot (T(\bar{x}) - \sigma^2).$$

$$\text{But } S(\sigma^2 | \bar{x}) = \frac{n}{2\sigma^4} \left( \sum_{i=1}^n \frac{(x_i - \mu)^2}{n} - \sigma^2 \right).$$

By Theorem 2.12, the only candidate is  $T(\bar{x}) = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}$ .

But  $\mu$  is unknown! This isn't a point estimator!

So no efficient estimator of  $\sigma^2$  exists.

# Efficiency: Examples

- If an unbiased point estimator is efficient, then it's the UMVUE – but the converse is not true in general
- Example 2.31: Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$ , where  $\lambda > 0$ . Show that an efficient estimator of  $\tau(\lambda) = \mathbb{P}_\lambda(X = 0)$  does not exist, and find its UMVUE.

$$S(\lambda | \vec{x}) = \frac{\sum x_i - n}{\lambda} = \frac{\sum x_i}{\lambda} - n + e^{-\lambda} - e^{-\lambda}. \text{ Clearly no } T(\vec{x}) \text{ exists for Theorem 2.12, so no efficient estimator exists.}$$

However, consider  $W(\vec{x}) = \mathbb{1}_{X_i=0}$ , which is unbiased for  $\tau(\lambda)$ . Also,  $T(\vec{x}) = \bar{X}_n$  is complete + sufficient. By Lehmann-Scheffé,  $E_\lambda[W(\vec{x}) | T(\vec{x})] = \mathbb{P}_\lambda(X_i=0 | \bar{X}_n)$  is the UMVUE for  $\tau(\lambda)$ .

How do we use it?  $n \bar{X}_n = \sum X_i \sim \text{Pois}(n\lambda)$ . Check that  $\vec{X} | (\sum X_i = \sum x_i)$

has pdf  $\binom{\sum x_i}{x_1, \dots, x_n} \left(\frac{1}{n}\right)^{x_1} \cdots \left(\frac{1}{n}\right)^{x_n} \sim \text{Multinomial}(\sum x_i, \frac{1}{n}, \dots, \frac{1}{n})$

$\Rightarrow X_i | (\sum X_i = \sum x_i) \sim \text{Bin}(\sum x_i, \frac{1}{n})$ .

Therefore,  $E_\lambda[\mathbb{1}_{X_i=0} | \bar{X}_n] = \mathbb{P}_\lambda(X_i=0 | \sum X_i) = \left(1 - \frac{1}{n}\right)^{\sum x_i}$ .