

STA261 - Module 2

Point Estimation

Rob Zimmerman

University of Toronto

July 12-14, 2022

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 θ or some function of the parameter $\tau(\theta)$ as best we can (whatever that means)
- Example 2.1: X_1, X_2, \dots, X_n = heights of UofT students $\stackrel{iid}{\sim} N(\mu, 3)$.
If we want to estimate μ , maybe take $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n$. Seems reasonable!
 - Event of voting for Candidate X in an election $\sim \text{Bernoulli}(p)$.
Maybe want to estimate $\psi(p) = \log(\frac{p}{1-p})$ "log-odds of p"

Point Estimation

- How do we estimate θ from the observed data \mathbf{x} ?
- Ideally, we want some statistic $T(\mathbf{X})$ such that $T(\mathbf{x})$ will be close to θ
- **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 θ .
- How do we find good point estimators?

Poll Time!

Not a point estimator: $\frac{1}{2}(x_1 + \mu)$

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 θ 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:
 $\mathbb{E}_\theta [\hat{\theta}(\mathbf{X})] \approx \theta$ 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 as $n \rightarrow \infty$.

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:
 - $X \sim N(\mu, \sigma^2) \Rightarrow \mathbb{E}(X) = \mu, \quad \mathbb{E}(X^2) = \mu^2 + \sigma^2$
 - $X \sim \text{Bin}(k, p) \Rightarrow \mathbb{E}[X] = kp, \quad \mathbb{E}[X^2] = kp(1-p) + k^2p^2$
 - $X \sim \text{Poisson}(\lambda) \Rightarrow \mathbb{E}[X] = \lambda, \quad \mathbb{E}[X^2] = \lambda + \lambda^2$
 - $X \sim \text{Exp}(\lambda) \Rightarrow \mathbb{E}[X] = \lambda, \quad \mathbb{E}[X^n] = \frac{n!}{\lambda^n}$ (Exercise!)

$$\text{Var}(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2$$

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 μ
- 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 $\mathbb{E}_{\mathbf{P}}[\hat{\mu}_1(\vec{X})] = \mathbb{E}_{\mathbf{P}}[X_1] = \mu$.
- Can we do better? $\hat{\mu}_n(\vec{X}) := \bar{X}_n \Rightarrow \mathbb{E}_{\mathbf{P}}[\hat{\mu}_n(\vec{X})] = \mu$
and $\text{Var}_{\mathbf{P}}(\hat{\mu}_n(\vec{X})) = \frac{\sigma^2}{n} < \sigma^2 = \text{Var}_{\mathbf{P}}(\hat{\mu}_1(\vec{X}))$
- Now suppose we want to estimate both μ and σ^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 $m_2(\mathbf{X}) \xrightarrow{d} \mu^2 + \sigma^2 = \mathbb{E}[X^2]$ by LN!
- Therefore $m_2(\mathbf{X}) - m_1(\mathbf{X})^2 \xrightarrow{d} \sigma^2$ by the continuous mapping theorem
(STA261 & Module 5)

So take $\hat{\mu}(\vec{X}) = m_1(\vec{X}) = \bar{X}_n$ and $\hat{\sigma}^2(\vec{X}) = m_2(\vec{X}) - m_1(\vec{X})^2$
 $= \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2$

The Method of Moments (MOM)

- 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 θ_j is defined by choosing

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

↳ k^{th} sample moment
 $= \frac{1}{n} \sum_{i=1}^n X_i^k$

Basic principle: LLN + CMT

↑
Although continuity of ψ_j 's is 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 λ .

$$\lambda = [E[X_i]]$$

→ The MOM estimator is $\hat{\lambda}_{\text{mom}}(\vec{x}) = \bar{X}_n$.

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 θ is known. Find the MOM estimator for k .

$$\mathbb{E}_\theta[X_i] = k\theta \Rightarrow k = \frac{\mathbb{E}_\theta[X_i]}{\theta} =: \Psi_i(\mathbb{E}[X])$$

$$\Rightarrow \text{Our MOM estimator is } \hat{k}_{\text{mom}}(\vec{x}) = \Psi_i(\bar{x}_n) = \frac{1}{\theta} \cdot \bar{x}_n.$$

- Could this be a problem?

Yes! There's no reason for $\hat{k}(\vec{x})$ to be a natural number, even though $\mathbb{N} = \mathbb{N}$.

Poll Time!

X_1, \dots, X_n iid $\text{Bin}(k, \theta)$, k known, $\theta \in (0, 1)$.

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

$$\Rightarrow \hat{\theta}_{\text{mean}}(\vec{X}) = \frac{1}{n} \sum_{i=1}^n X_i .$$

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 α .

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

$$\Rightarrow \alpha = 3 \cdot E[X]$$

$$\Rightarrow \hat{\alpha}_{\text{mom}}(\bar{x}) = 3 \cdot \bar{X}_n.$$

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 α and β .

$$\Psi_1 = \mathbb{E}(X) = \frac{\alpha}{\beta} \quad \textcircled{1}$$

$$\Psi_2 = \mathbb{E}(X^2) = \frac{\alpha}{\beta^2} + \frac{\alpha^2}{\beta^2} = \frac{\alpha + \alpha^2}{\beta^2}. \quad \textcircled{2}$$

$$\textcircled{1} \Rightarrow \alpha = \Psi_1 \cdot \beta$$

$$\Rightarrow \Psi_2 = \frac{\Psi_1 \beta + \Psi_1^2 \beta^2}{\beta^2} = \frac{\Psi_1}{\beta} + \Psi_1^2$$

$$\Rightarrow \beta = \frac{\Psi_1}{\Psi_2 - \Psi_1^2}$$

$$\Rightarrow \alpha = \frac{\Psi_1^2}{\Psi_2 - \Psi_1^2}$$

Our MOM estimators are

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

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

Method of Moments: Advantages and Disadvantages

GOOD: Very simple and easy to calculate; they exist whenever the moments do

BAD: They may not be in the correct parameter space Θ

BAD: If moments don't exist (ie, Cauchy), it won't work

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

BAD: May not have lowest variance possible among all point estimators of θ
(as will see...)

The Likelihood Function

- **Definition 2.3:** Let $\mathbf{X} \sim f_\theta$, where f_θ is a pdf or pmf in a parametric family. Given the observation $\mathbf{X} = \mathbf{x}$, the **likelihood function for θ** is the function $L(\cdot | \mathbf{x}) : \Theta \rightarrow [0, \infty)$ given by $L(\theta | \mathbf{x}) = f_\theta(\mathbf{x})$. ($L(\theta | \vec{x})$ is a random function of θ)
If $\theta_1, \theta_2 \in \mathbb{R}$ then $L(\theta_1 | \vec{x}) = f_{\theta_1}(\vec{x})$ and $L(\theta_2 | \vec{x}) = f_{\theta_2}(\vec{x})$.
- Interpret this as the “probability” of observing the sample \mathbf{x} , given that the sample came from f_θ NOT $P(\theta = \theta_1 | \vec{X} = \vec{x})$!!! $L(\theta | \vec{x}) \notin [0, 1]$ in general!
- 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_{θ_1} than under f_{θ_2} So $L(\cdot | \vec{x})$ ranks the elements of \mathbb{R} given the observed data.
- It could be that the likelihood is very small for all $\theta \in \Theta$, so knowing $L(\theta | \mathbf{x})$ for just a single θ 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 $L_1(\theta | x)$ and $L_2(\theta | y)$ are such that $L_1(\theta | x) = c(x, y) \cdot L_2(\theta | y)$, then the conclusions about θ drawn from x and y should be identical
 - ie, if $\frac{L_1(\theta | x)}{L_2(\theta | y)}$ is fixed & θ
- In other words, the likelihood principle says that anything we want to say about θ should be based solely on $L(\cdot | x)$, regardless of how x was actually obtained
- Is this requirement too strong?
 - Experiment 1: toss a coin w/ $P(H) = \theta$ 10 times; let $X = \# \& H \sim \text{Bin}(10, \theta)$
 - Say we observe $X=4$. $L_1(\theta | 4) = \binom{10}{4} \theta^4 (1-\theta)^6$
- Example 2.7:
 - Experiment 2: toss the same coin until we observe 4 H's; let $Y = \# \& T$ until that happens.
 - Then $Y \sim \text{NegBin}(4, \theta)$. Say we observe $Y=6$. Then $L_2(\theta | 6) = \binom{9}{6} \theta^4 (1-\theta)^6$
 - Then $L_1(\theta | x=4) \propto L_2(\theta | y=6)$. The Likelihood Principle says we should be indifferent to Exp 1 or Exp 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 θ , 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 θ .

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

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

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 θ is given by

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

and the **maximum likelihood estimator (MLE)** for θ is the point estimator given by $\hat{\theta}_{\text{MLE}} = \hat{\theta}(\mathbf{X})$. ← This is a statistic!

Equivalently, $\hat{\theta}(\mathbf{x})$ is s.t. $L(\hat{\theta}(\mathbf{x}) | \mathbf{x}) \geq L(\theta | \mathbf{x}) \quad \forall \theta \in \Theta$.

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 θ .

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

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

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

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

So $\hat{\theta}_{MLE}(x) = a \cdot \mathbb{1}_{x \in \{1, 3\}} + b \cdot \mathbb{1}_{x=2}$.

Maximum Likelihood: Examples

- But when the f_θ 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 θ .

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

Is this a local max? Need to find $\frac{d^2 L}{d\theta^2}$, plug in $\hat{\theta} = \bar{x}_n$, and check that $\frac{d^2 L}{d\theta^2} < 0$ there.
(You can verify). So $\hat{\theta}_{MLE}(\vec{x}) = \bar{x}_n$.

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 σ^2 is known
- What happens if we try to find the MLE of μ in the same fashion?

$$\begin{aligned} L(\mu | \bar{x}) &= \prod_{i=1}^n f_\mu(x_i) \\ &= (2\pi\sigma^2)^{-n/2} \cdot \exp\left(-\frac{\sum(x_i - \mu)^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-n/2} \cdot \exp\left(-\frac{1}{2\sigma^2} \left\{ \sum x_i^2 - 2\mu \sum x_i + n\mu^2 \right\}\right) \\ \frac{dL}{d\mu} &= (2\pi\sigma^2)^{-n/2} \cdot \left(\underbrace{\frac{\sum x_i - n\mu}{2\sigma^2}}_{\neq 0} \right) \cdot \exp\left(-\frac{\sum(x_i - \mu)^2}{2\sigma^2}\right) \stackrel{\text{set } 0}{=} 0 \\ &\Rightarrow \hat{\mu} = \frac{1}{n} \sum x_i = \bar{x}. \end{aligned}$$

But differentiating $\frac{dL}{d\mu}$ w.r.t. μ would be a nightmare!
Is there an easier way?!

...yes.

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 it's easier to differentiate a sum than a product (♥ linearify!)

If the data are iid, then $\ell(\theta | \mathbf{x}) = \log(L(\theta | \mathbf{x}))$

$$\begin{aligned}&= \log\left(\prod_{i=1}^n f_\theta(x_i)\right) \\&= \sum_{i=1}^n \log(f_\theta(x_i))\end{aligned}$$

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 σ^2 known. Find the MLE of μ .

$$L(\mu | \vec{x}) = \left(\frac{1}{2\pi\sigma^2} \right)^{-n/2} \cdot \exp\left(-\frac{\sum(x_i - \mu)^2}{2\sigma^2}\right)$$

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

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

$$\Rightarrow \hat{\mu} = \bar{x}_n.$$

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

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

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

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 λ .

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

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

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

$$\Rightarrow \hat{\lambda} = \bar{y}_x$$

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

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

So the MLE of λ 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 α .

$$L(\alpha | \vec{x}) = \prod_{i=1}^n \frac{2^\alpha}{\Gamma(\alpha)} x_i^{\alpha-1} e^{-2x_i} = \frac{2^{na}}{\Gamma(\alpha)^n} \cdot \left(\prod_i x_i \right)^{\alpha-1} \cdot e^{-2 \sum x_i}$$

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

$$\Rightarrow S(\alpha | \vec{x}) = n \cdot \log(2) - ?? + \log\left(\prod_i x_i\right)$$



No closed form! We can't differentiate this!

Maximum Likelihood: More Examples

- What about when θ 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 $\theta = (\mu, \sigma^2)$.

$$\begin{aligned} L(\mu, \sigma^2 | \vec{x}) &= (2\pi\sigma^2)^{-n/2} \cdot \exp\left(-\frac{\sum(x_i - \mu)^2}{2\sigma^2}\right) \\ \Rightarrow l(\mu, \sigma^2 | \vec{x}) &= c - \frac{n}{2} \log(\sigma^2) - \frac{\sum(x_i - \mu)^2}{2\sigma^2} \\ \Rightarrow S(\mu, \sigma^2 | \vec{x}) &= \nabla l = \left(\frac{\partial S}{\partial \mu}, \frac{\partial S}{\partial \sigma^2} \right) = \left(\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu), -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 \right) \stackrel{\text{set}}{=} \vec{0} = (0, 0) \\ \text{solve tediously} \Rightarrow (\hat{\mu}, \hat{\sigma}^2) &= \left(\bar{x}, \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right) \end{aligned}$$

Second derivative test:

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

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

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

The determinant of the Hessian is

$$\begin{vmatrix} \frac{\partial^2 S}{\partial \mu^2} & \frac{\partial^2 S}{\partial \mu \partial \sigma^2} \\ \frac{\partial^2 S}{\partial \sigma^2 \partial \mu} & \frac{\partial^2 S}{\partial (\sigma^2)^2} \end{vmatrix} = \begin{vmatrix} -\frac{n}{\sigma^2} & -\frac{1}{\sigma^4} \sum_{i=1}^n (x_i - \mu) \\ -\frac{1}{\sigma^4} \sum_{i=1}^n (x_i - \mu) & \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \sum_{i=1}^n (x_i - \mu)^2 \end{vmatrix}_{\substack{\mu = \bar{x} \\ \sigma^2 = \hat{\sigma}^2}} \cdots = \frac{1}{\hat{\sigma}^6} \cdot \frac{n^2}{2} > 0. \text{ So } (\hat{\mu}, \hat{\sigma}^2) \text{ is the MLE.}$$

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

$$L(\theta | \vec{x}) = \prod_{i=1}^n f_\theta(x_i) = \theta^{-n} \cdot \mathbb{1}_{0 \leq x_1 < x_m \leq \theta} = \mathbb{1}_{0 \leq x_m} \cdot \theta^{-n} \cdot \mathbb{1}_{x_m \leq \theta}$$

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

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

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

Hence $\hat{\theta}_{me}(\vec{x}) = x_m$. But we couldn't use calculus to find it,
because $L(\theta | \vec{x})$ is not differentiable at $\theta = x_m$.

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 β .

all known!

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

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

$$\Rightarrow S(\beta | \vec{y}) = \frac{\sum x_i(y_i - \beta x_i)}{\sigma^2} \stackrel{!}{=} 0$$

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

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

Second derivative test:

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

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

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

Reparameterization

- Instead of θ 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 θ^2 .

$$\text{let } \tau = \theta^2.$$

$$\text{Then } L(\tau | \vec{x}) = \sqrt{\tau}^{\sum x_i} (1 - \sqrt{\tau})^{n - \sum x_i}$$

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

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

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

$$\Rightarrow \hat{\tau} = (\bar{x})^2$$

$$\Rightarrow \hat{\tau}(\vec{x}) = (\bar{x}_n)^2 = (\Theta_{\text{MLE}}(\vec{x}))^2$$

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 Θ , then the MLE of $\tau(\theta)$ is given by $\tau(\hat{\theta}(\mathbf{X}))$.

Proof. Let $\Psi = \tau(\theta)$ so that $\theta = \tau^{-1}(\Psi)$, and let $\hat{\Psi} := \tau^{-1}(\hat{\theta})$.
ie, $\tau(\theta)_{MLE}(\vec{x}) = \tau(\hat{\theta}(\vec{x}))$
"plug-in estimator"

Let the likelihood under θ be $L(\theta | \vec{x})$, and the likelihood under Ψ be $L^*(\Psi | \vec{x})$.

Then for any $\Psi = \tau(\theta) \in \tau(\Theta)$,

$$\begin{aligned} L^*(\hat{\Psi} | \vec{x}) &= f_{\tau^{-1}(\hat{\Psi})}(\vec{x}) \\ &= L(\tau^{-1}(\hat{\Psi}) | \vec{x}) \\ &= L(\hat{\theta} | \vec{x}) \\ &\geq L(\theta | \vec{x}) \\ &= L(\tau^{-1}(\Psi) | \vec{x}) \\ &= f_{\tau^{-1}(\Psi)}(\vec{x}) \\ &= L^*(\Psi | \vec{x}). \end{aligned}$$

Hence $\hat{\Psi}$ maximizes L^* . \square

Discussion prompt: what if τ is not one-to-one?

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 $\log(\frac{x}{1-x})$ is one-to-one on $(0, 1)$, the MLE of $\tau(p)$ is $\log\left(\frac{\bar{X}_n}{1-\bar{X}_n}\right)$

by the invariance property.

Poll Time!

Invariance principle! $\hat{\lambda}_{MLE} = \bar{X}_n$
 $\Rightarrow \underbrace{(-\log(\lambda))}_{MLE} = -\log(\bar{X}_n).$

Maximum Likelihood Estimation

- Maximum likelihood is *by far* the most common method that statisticians use to find point estimates¹
- Maximum likelihood estimators tend to have quite good properties (especially for large sample sizes):

Good: $\hat{\theta}_{MLE}$ is always in Θ

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

Good: (Relatively) easy to easy to implement (usually general-purpose optimization software does it)

"BAD": $\hat{\theta}_{MLE}$ may not have the "right" expectation (e.g., $\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ but $E[\hat{\sigma}_{MLE}^2] \neq \sigma^2$)
↑ in the $N(\mu, \sigma^2)$ model

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

¹ 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 θ
- 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

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 $\tau(\theta)$.
 $\Rightarrow \mathbb{E}_\theta[T(\vec{X})] = \tau(\theta)$

- Example 2.18:

$X_1, \dots, X_n \sim N(\mu, \sigma^2)$, $\mu \in \mathbb{R}$, $\sigma^2 > 0$. Then $T(\vec{X}) = \bar{X}_n$ is unbiased for μ .

$X_1, \dots, X_n \sim \text{Bernoulli}(p)$, $p \in (0, 1)$. Then $T(\vec{X}) = \bar{X}_n$ is unbiased for p ,

$$\begin{aligned}\text{because } \text{Bias}_p(T(\vec{X})) &= \mathbb{E}_p[T(\vec{X})] - p = \mathbb{E}_p\left[\frac{1}{n} \sum_{i=1}^n X_i\right] - p \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_p[X_i] - p \\ &= \frac{1}{n} \cdot np - p = p - p = 0.\end{aligned}$$

- Example 2.19:

$X_1, \dots, X_n \sim N(\mu, \sigma^2)$. $\tau(\sigma^2) = \sigma^2$:

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

↑ Assignment 0: $E(S^2) = \text{Var}(X)$

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 $T(X)$ is unbiased for $\tau(\theta) = Y_\theta$.

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

But LHS is unbounded as $\theta \rightarrow 0$, but RHS $\rightarrow T(0) \in \mathbb{R}$. So this can't happen!

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. $\text{MSE}_\theta(T(\vec{x})) = \mathbb{E}_\theta[(T(\vec{x}) - \varphi(\theta))^2]$

$$= \text{Var}_\theta(T(\vec{x}) - \varphi(\theta)) + \mathbb{E}_\theta[(T(\vec{x}) - \varphi(\theta))]^2$$

$$= \text{Var}_\theta(T(\vec{x})) + \text{Bias}_\theta(T(\vec{x}))^2. \quad \square$$

$$\text{Var}(Y) = \mathbb{E}[Y^2] - \mathbb{E}[Y]^2$$

Take $Y = T(\vec{x}) - \varphi(\theta)$.

So among all estimators with a fixed MSE, you must choose between more accuracy (\leftarrow bias) and less precision (\rightarrow variance), or vice versa..

Poll Time!

$$MSE_{\theta}(\hat{T}(\vec{x})) = \text{Bias}_{\theta}(\hat{T}(\vec{x}))^2 + \text{Var}_{\theta}(\hat{T}(\vec{x}))$$

If $\hat{T}(\vec{x})$ is unbiased, then $\text{Bias}_{\theta}(\hat{T}(\vec{x})) = 0 \Rightarrow MSE = \text{Var.}$

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 $\tau(\theta)$.

Handwritten notes:

- A bracket under "best unbiased estimator" has the handwritten note "UGER".
- A bracket under "UMVUE" has the handwritten note "lowest variance out of all unbiased estimators".

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 $\mu_x = E(X)$, $\mu_y = E(Y)$. If $Var(Y) = 0$, then $Y = \mu_y$. Then

$$\text{Cov}(X, Y) = E[(X - \mu_x)(Y - \mu_y)] = 0 \leq 0 = \sqrt{\text{Var}(X) \cdot \text{Var}(Y)}.$$

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

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

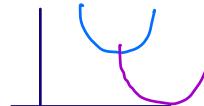
This is quadratic in t and non-negative, so it has at most one real root.

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

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

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

Quadratic formula: $0 = ax^2 + bx + c$
 $\Rightarrow x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$



Equality holds iff $E[(Z - tW)^2] = 0$ for some t

$$\text{iff } Z = tW$$

$$\text{iff } X - \mu_x = t \cdot (Y - \mu_y)$$

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

Moreover, check that $t = \frac{\text{Cov}(X, Y)}{\text{Var}(Y)}$. \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 $W^* = \frac{1}{2}(W + W')$.

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

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

But W^* can't beat a UMVUE, so equality must hold (ie, $\text{Cov}(W, W') = \text{Var}_{\theta}(W)$ &
 $\Rightarrow W' = a(\theta) \cdot W + b(\theta)$. What are $a(\theta)$ and $b(\theta)$?

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

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

Finally, $\tau(\theta) = \mathbb{E}_{\theta}[W]$

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

Hence $W = W'$. \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 θ . 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,
 $\text{Var}_\theta(W_T(\mathbf{X})) \leq \text{Var}_\theta(W(\mathbf{X})).$ ie, conditioning on sufficient statistics never hurts!

Proof. Unbiasedness: $\mathbb{E}_\theta[W_T(\bar{x})] = \mathbb{E}_\theta[\mathbb{E}_\theta[W(\bar{x}) | T(\bar{x})]] \stackrel{\text{"tower rule"}}{=} \mathbb{E}_\theta[W(\bar{x})] = \tau(\theta)$ since W is unbiased.

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

What about sufficiency? If $T(\bar{x})$ weren't sufficient, then $\mathbb{E}_\theta[W(\bar{x}) | T(\bar{x})]$ wouldn't be free of θ
 \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: $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$, $\lambda > 0$.

We have at least two unbiased estimators for λ : \bar{X}_n and S^2 .

But \bar{X}_n is sufficient for λ by Theorem 1.2, so $E[S^2 | \bar{X}_n]$ is better than S^2 itself,

but by Theorem 2.9, it can't be better than $E[\bar{X}_n | \bar{X}_n] = \bar{X}_n$.

Rao-Blackwell: Examples

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

- 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.

Unbiasedness: $\mathbb{E}_\theta[W(\vec{X})] = \mathbb{E}_\theta[\mathbb{1}_{X_1=1}] = P_\theta(X_1=1) = k\theta(1-\theta)^{k-1} = \tau(\theta)$.

Recall that $T(\vec{X}) = \sum_{i=1}^n X_i$ is sufficient for θ .

Let $W_t(\vec{X}) = \mathbb{E}_\theta[W(\vec{X}) | T(\vec{X}) = t]$. How do we use this? Suppose $T(\vec{X}) = t$.

$$\begin{aligned} \text{Then } \mathbb{E}_\theta[W(\vec{X}) | T(\vec{X}) = t] &= P(X_1=1 | \sum_{i=1}^n X_i = t) \\ &= \frac{P(X_1=1 \wedge \sum_{i=1}^n X_i = t)}{P(\sum_{i=1}^n X_i = t)} \\ &= \frac{P(X_1=1 \wedge \sum_{i=2}^n X_i = t-1)}{P(\sum_{i=1}^n X_i = t)} \end{aligned}$$

$$\text{indep } \frac{P(X_1=1) \cdot P(\sum_{i=2}^n X_i = t-1)}{P(\sum_{i=1}^n X_i = t)}$$

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

$$= \frac{k \binom{k(n-1)}{t-1}}{\binom{kn}{t}}$$

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

Suppose $V_T(\vec{X}) = \mathbb{E}_\theta[V(\vec{X}) | T(\vec{X})]$ is also an MLE, where $V(\vec{X})$ is unbiased for $\tau(\theta)$, which we can do by Rao-Blackwell.

The Lehmann-Scheffé Theorem

- 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 $W_T(\mathbf{X}) = \mathbb{E}[W(\mathbf{X}) | T(\mathbf{X})]$ is the unique UMVUE.

Proof. Suppose that $V(\vec{x})$ is a UMVUE of $v(\theta)$. Then $V_T(\vec{x}) := \mathbb{E}[V(\vec{x}) | T(\vec{x})]$ is also unbiased for $v(\theta)$, and $\text{Var}(V_T(\vec{x})) \leq \text{Var}(V(\vec{x}))$ by Rao-Blackwell, so it too is a UMVUE. By Theorem 2.4, $V(\vec{x}) = V_T(\vec{x})$.

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

\Rightarrow By completeness, $P_{\theta}(h(T(\vec{x})) = 0) = 1 \quad \forall \theta \in \Theta$.

$$\Rightarrow W_T(\vec{x}) = V_T(\vec{x})$$

$$= V(\vec{x}).$$

So the UMVUE is $\mathbb{E}[W(\vec{x}) | T(\vec{x})]$. \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
- However, Lehmann-Scheffé assumes that a complete sufficient statistic exists (which isn't always the case, as we know from Module 1), so it doesn't subsume Theorem 2.4
- In fact, there do exist models where UMVUEs exist but complete sufficient statistics don't

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 (μ, σ^2) .

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

Also \bar{X}_n is unbiased for μ , and S^2 is unbiased for σ^2 .

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

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

That's not the MLE of (μ, σ^2) !

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 λ .

We know that \bar{X}_n is unbiased for λ , and it's also a complete sufficient statistic.

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

Poll Time!

What About the Likelihood?

- Rao-Blackwellization and Lehmann-Scheffé tell us how to get the unique UMVUE (if it exists) via complete 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.

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

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

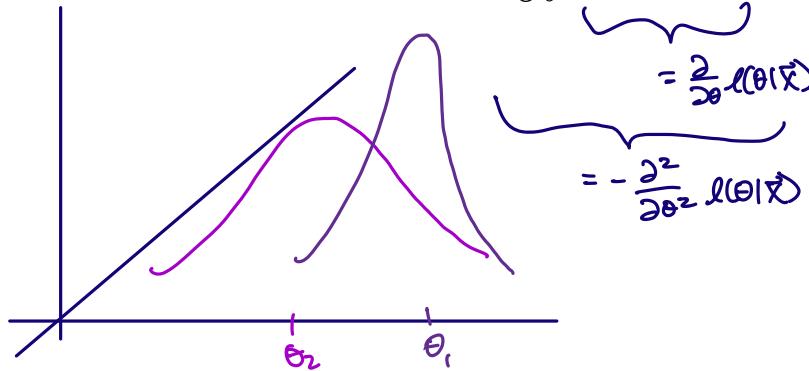
$$I_n(\theta) = \text{Var}_\theta (S(\theta | \mathbf{X})).$$

- **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

$$J_n(\mathbf{X}) = -\frac{\partial}{\partial \theta} S(\theta | \mathbf{X}_{\bullet}) \Big|_{\theta=\hat{\theta}_{\text{MLE}}} \quad (\text{blue bracket})$$

$$= \frac{\partial^2}{\partial \theta^2} \ell(\theta | \mathbf{x})$$

$$= -\frac{\partial^2}{\partial \theta^2} \ell(\theta | \mathbf{x})$$



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 λ .

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

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

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

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

For $J_n(\vec{x})$, recall $\hat{\lambda}_{\text{MLE}}(\vec{x}) = \bar{X}_n$.

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

$$\begin{aligned} J_n(\vec{x}) &= \left. \frac{\sum x_i}{\lambda^2} \right|_{\lambda=\bar{X}_n} \\ &= \frac{n \bar{X}_n}{(\bar{X}_n)^2} = \frac{n}{\bar{X}_n} = \frac{n^2}{\sum x_i}. \end{aligned}$$

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 σ^2 is known. Calculate the observed and expected Fisher information for μ .

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

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

Recall $\hat{\mu}_{\text{MLE}}(\vec{x}) = \bar{X}_n$. Then

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

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

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

Then

$$\text{Var}_\theta(T(\mathbf{X})) \geq \frac{\left(\frac{d}{d\theta} \mathbb{E}_\theta [T(\mathbf{X})] \right)^2}{I_n(\theta)}.$$

In particular, if $T(\mathbf{X})$ is unbiased for $\tau(\theta)$ and $\tau(\cdot)$ is differentiable on Θ , 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, let $U(\vec{x}) = S(\theta(\vec{x})) = \frac{\partial}{\partial \theta} \ell(\theta(\vec{x}))$.

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

$$= \underbrace{\mathbb{E}_{\theta}[T(\vec{x}) \cdot S(\theta(\vec{x}))]}_{①} - \underbrace{\mathbb{E}_{\theta}[T(\vec{x})] \cdot \mathbb{E}_{\theta}[S(\theta(\vec{x}))]}_{②} = \frac{d}{d\theta} \mathbb{E}_{\theta}[T(\vec{x})] - 0 = \frac{d}{d\theta} \mathbb{E}_{\theta}[T(\vec{x})]. \quad ③$$

$$\begin{aligned} ① &= \int_{\mathbb{R}^n} T(\vec{x}) \cdot S(\theta(\vec{x})) \cdot f_{\theta}(\vec{x}) d\vec{x} \\ &= \int_{\mathbb{R}^n} T(\vec{x}) \left(\frac{\partial}{\partial \theta} \ell(\theta(\vec{x})) \right) f_{\theta}(\vec{x}) d\vec{x} \\ &= \int_{\mathbb{R}^n} T(\vec{x}) \left(\frac{1}{f_{\theta}(\vec{x})} \cdot \frac{\partial}{\partial \theta} f_{\theta}(\vec{x}) \right) f_{\theta}(\vec{x}) d\vec{x} \\ &= \int_{\mathbb{R}^n} T(\vec{x}) \cdot \frac{1}{\partial \theta} f_{\theta}(\vec{x}) d\vec{x} \\ &= \int \frac{1}{\partial \theta} (T(\vec{x}) \cdot f_{\theta}(\vec{x})) d\vec{x} \\ &\stackrel{\text{def}}{=} \frac{d}{d\theta} \int_{\mathbb{R}^n} T(\vec{x}) \cdot f_{\theta}(\vec{x}) d\vec{x} \\ &= \frac{d}{d\theta} \mathbb{E}_{\theta}[T(\vec{x})] \end{aligned}$$

$$\begin{aligned} ② &= \int_{\mathbb{R}^n} \left(\frac{\partial}{\partial \theta} \log(f_{\theta}(\vec{x})) \right) f_{\theta}(\vec{x}) d\vec{x} \\ &= \int_{\mathbb{R}^n} \frac{1}{f_{\theta}(\vec{x})} \left(\frac{\partial}{\partial \theta} f_{\theta}(\vec{x}) \right) f_{\theta}(\vec{x}) d\vec{x} \\ &= \int_{\mathbb{R}^n} \frac{1}{\partial \theta} f_{\theta}(\vec{x}) d\vec{x} \\ &\stackrel{\text{def}}{=} \frac{d}{d\theta} \int_{\mathbb{R}^n} f_{\theta}(\vec{x}) d\vec{x} \\ &= \frac{d}{d\theta} 1 \\ &= 0 \end{aligned}$$

Also, by def, $\text{Var}_{\theta}(S(\theta(\vec{x}))) = I_n(\theta)$. Plug ③ into the covariance inequality and we're done! \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 $X_1, \dots, X_n \sim N(\mu, 1)$, don't try $T(\bar{X}) = \bar{X}_1 / \bar{X}_n$. 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: $\text{Unif}(0, \theta) \Rightarrow \text{Support } \Theta = (0, \theta) \text{ depends on } \theta$
 $\Rightarrow \frac{\partial}{\partial \theta} \mathbb{E}[T(\bar{X})] \neq \int_0^\theta \left(\frac{\partial^2}{\partial \theta^2} T(\bar{x}) \cdot \frac{1}{\theta} \right) d\bar{x}$ 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. $I_n(\theta) = \text{Var}_{f_\theta}(S(\theta | \mathbf{X}))$ by def

$$\begin{aligned} &= \mathbb{E}[S(\theta | \mathbf{X})^2] - \underbrace{\mathbb{E}[S(\theta | \mathbf{X})]^2}_{=0 \text{ from the prof \& Q&A}} \\ &= \mathbb{E}[S(\theta | \mathbf{X})^2]. \end{aligned}$$

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

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

More Easing

- Theorem 2.11 (**Second Bartlett Identity**): If $X \sim f_\theta$ and f_θ 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_θ is in an exponential family) then

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

Proof.

$$\begin{aligned} \text{RHS} &= -\mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \left(\frac{\partial}{\partial \theta} \log(f_\theta(\theta | X)) \right) \right] \\ &= -\mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \left(\frac{1}{f_\theta(\theta)} \cdot \frac{\partial}{\partial \theta} f_\theta(\theta) \right) \right] \\ &= -\mathbb{E}_\theta [\dots - \dots] \end{aligned}$$

Exercise: you finish off!

Use the assumption somewhere..

Efficiency

- **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 ^{unbiased} efficient estimator exists, then it must be the UMVUE
- But an efficient estimator doesn't always exist, as we'll soon see

Efficiency: Examples

- **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 $T(\mathbf{X}) = \bar{X}_n$ is an efficient estimator for μ .

We need to calculate the CRLB for estimators of μ , and also $\text{Var}_\mu(T(\mathbf{X}))$, and show that they're equal.

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

What about the CRLB? Numerator: $\left(\frac{d}{d\mu} E_\mu[T(\bar{X})]\right)^2 = \left(\frac{d}{d\mu} \mu\right)^2 = 1$.

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

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

So $T(\bar{X}) = \bar{X}_n$ is efficient for μ .

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 holds in the CRB iff

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

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

$$\text{iff } S(\theta|\vec{x}) = \underbrace{\frac{I_\theta(\theta)}{[\tau'(\theta)]^2}}_{=: a(\theta)} (\tau(\vec{x}) - \tau(\theta)). \quad \square$$

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 σ^2 .

If there did exist one, say $T(\vec{x})$, then there would also be some $a(\sigma^2)$ s.t.

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

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

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

But μ is unknown! So $T(\vec{x})$ is not a point estimator.

So no efficient estimator of σ^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.

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

However, consider $W(\vec{x}) = \mathbb{1}_{X_1=0}$, which is unbiased for $\tau(\lambda)$. Also, $T(\vec{x}) = \bar{X}_n$ is a CSS for λ .

By Lehmann-Scheffé, $\mathbb{E}[W(\vec{x}) | T(\vec{x})] = \mathbb{P}(X_1=1 | \bar{X}_n)$ is the UMVUE of $\tau(\lambda)$.

How do we use it? $n\bar{X}_n = \sum_{i=1}^n X_i \sim \text{Poisson}(n\lambda)$ *Exercise: use mfp!* Check that $\vec{X} | \sum X_i = t$ has pmf $\binom{t}{x_1, \dots, x_n} \left(\frac{1}{n}\right)^{x_1} \cdots \left(\frac{1}{n}\right)^{x_n} \sim \text{Multinomial}\left(t; \frac{1}{n}, \dots, \frac{1}{n}\right)$.

$\Rightarrow X_1 | \sum X_i = t \sim \text{Bin}(t, \frac{1}{n})$.

Therefore, $\mathbb{E}[\mathbb{1}_{X_1=0} | \bar{X}_n] = \mathbb{E}[\mathbb{1}_{X_1=0} | \sum X_i] = \mathbb{P}(X_1=0 | \sum X_i) = \left(1 - \frac{1}{n}\right)^{\sum x_i}$.