## Topic: Review (ch 1-3)

$X$ is RV w/ sample space $\mathcal{X}$

Random variables (RVs), $X$ say, are defined by their distribution function (CDF).

$$F_X(x) = F(x) = Pr(X \leq x)$$

Mixture

**Ex)** Discrete

$$X = \begin{cases} -1, & \text{w.p. } 0.2 \\ 0, & \text{w.p. } 0.3 \\ 1, & \text{w.p. } 0.5 \end{cases}$$

$$S = \{-1, 0, 1\}$$

Continuous

$$X \sim Exp\left(\tfrac{1}{5}\right)$$

$$S = (0, \infty)$$

Insurance policy reimburse up to some benefit level, $C$, with some deductible, $d$.

$X$ = policy holders $\sim Exp\left(\tfrac{1}{5}\right)$ loss

$Y$ = payout from Insurance co.

$$= \begin{cases} 0, & x < d \\ x-d, & d \leq x < C+d \\ C, & x \geq C+d \end{cases}$$

$$S_y = \{0, C\} \cup (0, c) = [0, C]$$

$Pr(Y=a) = 0$ if $a \in [d\ c+d)$ so $\longrightarrow$ $f_Y(y) = \tfrac{1}{5}e^{-\tfrac{(y+d)}{5}}I\{0 \leq y \leq C\}$

The CDF is a <u>probability measure/law</u> so,

**(Def)**

1. $Pr(\mathcal{X}) = 1$

2. If $A \subset \mathcal{X}$ then $Pr(A) \geq 0$

3. If mutually disjoint $A_1, A_2, \ldots$

$$Pr\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} Pr(A_i)$$

**Def:** Conditional probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

## Law of Total Probability

For events $B_1, \ldots, B_n$ s.t. $\bigcup_{i=1}^{n} B_i = \mathcal{X}$

If $\Bigg\{$    and $\quad B_i \cap B_j = \emptyset$ , for $i \neq j$

     and $\quad Pr(B_i) > 0$ for all $i$

then For any event $A \in \mathcal{X}$

$$Pr(A) = \sum_{i=1}^{n} Pr(A|B_i) \cdot P(B_i)$$

All RVs have a CDF. Many RVs have a density function as well.

$$\boxed{f_X(x) = f(x) = F_X'(x) = \frac{d}{dx} F_X(x)}$$

Def: Likelihood is the density function but viewed as a function of the parameters.

different from textbook

$$* \quad \boxed{f(x; \theta) = f(\theta|x) = L(\theta|x)}$$

is a function of $\theta$.

Read as "the likelihood for $\theta$, given $X = x$."

Applying the Law of Total Prob. to jointly distributed RVs, $(X, Y)$ yields:

$$f_Y(y) = \int_{\mathcal{X}} f_{Y|X=x}(y|X=x) \cdot f_X(x) \, dx$$

(In the case where both $X, Y$ are continuous)

9/2/22 (week 1)

## Bayes' Rule/Law — combines Law of Tot. Prob w/ def. of conditional prob.

For $A, B_1, \ldots, B_n$ where $B_i$ are disjoint w/ all $B_j$, $i \neq j$, $\bigcup_{i=1}^{n} B_i = \mathcal{X}$ and $P(B_i) > 0$ for all $i$,

we have that

$$Pr(B_j | A) = \frac{Pr(A|B_j) Pr(B_j)}{\sum_{i=1}^{n} Pr(A|B_i) Pr(B_i)}$$

def. of cond.

$$\frac{Pr(B_j \cap A)}{Pr(A)} \xrightarrow[\text{law of tot.}]{} \frac{Pr(A \cap B_j)}{\sum_{i=1}^{n} Pr(A|B_i) Pr(B_i)}$$

def. of cond.

"(Reverse) Conditioning"

## Jointly Distributed RVs

Ex)  $(X, Y)$

$(X_1, X_2)$

$(X_1, X_2, \ldots, X_n)$

Q: What is the sample space for jointly distrbted RVs, say, $X$ w/ sample space $\mathcal{X}$ and $Y$ w/ samp. space $\mathcal{Y}$?

# Distribution Notation :

|  | Both discrete | Both continuous |
|---|---|---|
| joint | $P_{XY}(x,y) = Pr(X=x, Y=y)$ <br> $F(x,y) = Pr(X \leq x, Y \leq y)$ | $Pr\left(\binom{X}{Y} \in A\right) = \iint_A f(x,y)\,dy\,dx$ <br> $F(x,y) = Pr(X \leq x, Y \leq y)$ |
| Marginal | $P_X(x) = \sum_{y \in \mathcal{Y}} P(x,y)$ | $f_X(x) = F_X'(x) = \int_y f(x,y)\,dy$ <br> where $F_X(x) = Pr(X \leq x)$ <br> $= \lim_{y \to \infty} F(x,y)$ <br> $= \int_{-\infty}^{x} \int_{-\infty}^{\infty} f(u,y)\,dy\,du$ |
| conditional | $P_{X|Y}(x|y) = \dfrac{P_{XY}(x,y)}{P_Y(y)}$ | $f_{X|Y}(x|y) = \begin{cases} \dfrac{f_{XY}(x,y)}{f_Y(y)}, & \text{if } 0 < f_Y(y) < \infty \\ 0, & \text{otherwise} \end{cases}$ |

## Special EX

### X discrete, Y continuous

|  |  |  |
|---|---|---|
| marginal | $Pr(X=x) = P_X(x)$ | $F_Y(y) = \sum_{x \in \mathcal{X}} Pr(Y \leq y, X \leq x)$ <br> $f_Y(y) = F_Y'(y)$ |
| conditional | $Pr(X=x \mid Y=y) = \dfrac{f_{Y|X}(y \mid x)\, Pr(X=x)}{f_Y(y)}$ | |

In general, knowing the marginal distrb'n of X and of Y is $\underline{NOT}$ enough information for us to determine the joint distrb'n of $(X, Y)$....

unless....

$\qquad$ X and Y are independent $\qquad$ $X \perp\!\!\!\perp Y$ (abbreviation)

## Def: Independent RVs

For RVs $(X_1, ..., X_n)$ w/ joint distrb'n fnctn
$$F(x_1, ..., x_n),$$
we say $(X_1, ..., X_n)$ are independent RVs if
$$F(x_1, ..., x_n) = F_{X_1}(x_1) \cdot F_{X_2}(x_2) \cdot ... \cdot F_{X_n}(x_n).$$

(It can be shown that this is equivalent to saying that the joint pmf or joint density factors.)

## Indicator fnctn

$$\mathbb{I}\{0 < x < 1\} = \begin{cases} 1, & \text{if } 0 < x < 1 \\ 0, & \text{o/w} \end{cases}$$

$$E(\mathbb{I}\{0 < X < 1\}) = P_r(0 < X < 1) = 1 \cdot P_r(0 < X < 1) + 0 \cdot P_r(X \notin (0,1))$$

# Next week...

- expectation, variance, covariance

- conditional expectation & variance

- moment generating functions
- methods of estimation

# Topic : Review (Ch 4)

**Def:** Moment Generating Function (MGF)

of a discrete RV X is:

$$M(t) = \sum_{x \in \chi} e^{tx} P_X(x)$$

of a continuous RV X is:

$$M(t) = \int_\chi e^{tx} f_X(x) dx$$

The MGF of RV does not always exist (ex. Cauchy)
but when it does, it uniquely determines the RV.
(The Characteristic function, like the CDF, always
exists but is a complex function.)

**Def:** The <u>moments</u> of a RV X are

$$E(X^r) \quad \text{for} \quad r = 1, 2, \dots$$

The $r^{th}$ derivative of MGF, $M(t)$, evaluated
at $t=0$, is the $r^{th}$ moment of X;

Ie. $M^{(r)}(0) = E(X^r)$,

provided $M(t)$ exists in an open interval
containing zero.

The first and second moments of a RV
determine its expectation & variance.

| Expected Values | Discrete | Continuous |
|---|---|---|
| $E(X) =$ | $\sum_{x \in \mathcal{X}} x P_X(x)$ | $\int_{\mathcal{X}} x f_X(x) dx$ |
| $E[g(X)] =$ | $\sum_{x \in \mathcal{X}} g(x) P_X(x)$ | $\int_{\mathcal{X}} g(x) f_X(x) dx$ |
| $E[Y \mid X = x] =$ | $\sum_{y \in \mathcal{Y}} y P_{Y \mid X}(y \mid x)$ | $\int_{\mathcal{Y}} y f_{Y \mid X}(y \mid x) dy$ |
| $E[g(Y) \mid X = x] =$ | $\sum_{y \in \mathcal{Y}} g(y) P_{Y \mid X}(y \mid x)$ | $\int_{\mathcal{Y}} g(y) f_{Y \mid X}(y \mid x) dy$ |

"The expected value is the sum of the possibilities of a RV, times their probabilities."

Note: $E[g(X)] \neq g[E(X)]$

Ex) $X = \begin{cases} 1, & \text{w.p. } \frac{1}{2} \\ 2, & \text{w.p. } \frac{1}{2} \end{cases}$ ; $g(X) = \frac{1}{X}$

$E(X) = 1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{2} = \frac{3}{2}$

$E[g(X)] = \frac{1}{1} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} = \frac{3}{4}$

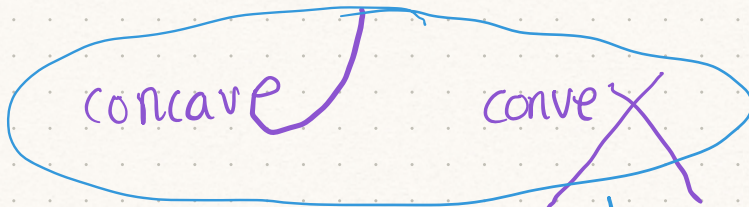$g(E(X)) = g(\frac{3}{2}) = \frac{2}{3}$

However, we do have the following result.

Jensen's Inequality

For any convex function, $g$, and any RV $X$,

$$g(E[X]) \leq E[g(X)]$$

provided $E[g(X)]$ and $g(E[X])$ exist and are finite.

concave     convex

Incorrect! These a swapped!

Expectation is a __linear__ operator:

$$E\left[\sum_{i=1}^{n}(a_i + b_i X_i)\right] = E\left[(a_1 + b_1 X_1) + (a_2 + b_2 X_2) + \dots + (a_n + b_n X_n)\right]$$

$$= E[a_1 + b_1 X_1] + E[a_2 + b_2 X_2] + \dots + E[a_n + b_n X_n]$$

$$\vdots$$

$$= \sum_{i=1}^{n} a_i + \sum_{i=1}^{n} b_i E[X_i]$$

## Markov's Inequality

If $X$ is a positive RV for which $E(X)$ exists,

$$Pr(X \geq t) \leq \frac{E(X)}{t}, \qquad \text{for any } t \in \mathbb{R}.$$

"mean slasher"

Ex) (of Markov's Inequality)    "variance difference"

## Chebyshev's Inequality

If $X$ is a RV whose first and second moments exist, then f any $t > 0$:

$$Pr\left(|X - E(X)| > t\right) = Pr\left((X - E(X))^2 > t\right)$$

$$\leq \frac{E[(X - E(X))^2]}{t^2}$$

$$= Var(X) / t^2$$

<u>Law of Iterated (Total) Expectation</u>

For RVs X and Y,

$$E[Y|X] \text{ is a RV}$$

because it is a function of X, which is not fixed. It always holds that

$$E[E(Y|X)] = E[Y].$$

(For a proof, see pg. 149.)

Note: $E[Y|X=x]$ is a function of $x$ and is thus <u>NOT</u> a RV since $X=x$ is fixed.

## Variance

If RV X has $E(X) < \infty$ then

$$Var(X) = E\left\{[X - E(X)]^2\right\}$$

$$\vdots$$

$$= E(X^2) - [E(X)]^2$$

Variance is a non-linear operator:

$$Var\left(\sum_{i=1}^{n} a_i + b_i X_i\right) = Var\left((a_1 + b_1 X_1) + (a_2 + b_2 X_2) + \ldots + (a_n + b_n X_n)\right)$$

$$= Var\left(\sum_{i=1}^{n} a_i + \sum_{i=1}^{n} b_i X_i\right)$$

$$= Var\left(\sum_{i=1}^{n} b_i X_i\right)$$

$$= \ldots \, ? \longrightarrow \text{see next} \atop \text{pages}$$

If we are only interested in one RV then:

$$Var(a + bX) = b^2 Var(X)$$

"Variance is the _average_ (squared) _distance_ between the possibilities, of a RV and its _expectation_."

## E-V-E Formula (Iterated Variance)

For RVs $X$ and $Y$, we have that

$$Var(Y) = E[Var(Y|X)] + Var[E(Y|X)]$$

(Proof on pg 151)

## Covariance

If $X, Y$ are jointly distb'td RVs whose expectations exist,

$$\text{Cov}(X, Y) = E\left[(X - E(X))(Y - E(Y))\right]$$
$$\vdots$$
$$= E(XY) - E(X)E(Y)$$

Furthermore,

$$\text{Cov}(X, X) = \text{Var}(X)$$

## Properties of variance + covariance:

Let $U = a + \sum_{i=1}^{n} b_i X_i$, $V = c + \sum_{j=1}^{m} d_j Y_j$

for RVs $X_1, X_2, \ldots, X_n, Y_1, Y_2, \ldots, Y_m$

$$\text{Cov}(U, V) = \sum_{i=1}^{n} \sum_{j=1}^{m} b_i d_j \text{Cov}(X_i, Y_j)$$

In particular,

$$\text{Var}(U) = \text{Var}\left(a + \sum_{i=1}^{n} b_i X_i\right)$$

$$= \text{Var}\left(\sum_{i=1}^{n} b_i X_i\right)$$

$$= \text{Cov}\left(\sum_{i=1}^{n} b_i X_i, \sum_{i=1}^{n} b_i X_i\right)$$

$$= \sum_{i=1}^{n} \sum_{j=1}^{n} b_i b_j \text{Cov}(X_i, X_j)$$

Ex) If $X_1, \ldots, X_n$ are independent (and identically distributed,

(IID)

what is $\text{Var}\left(\sum_{j=1}^{n} X_j\right)$?

$$\text{Var}\left(\sum_{i=1}^{n} X_i\right) = \sum_{i=1}^{n} \text{Var}(X_i)$$

# Topic: Estimation Part I (ch 4+8)

Setting: $X_1, \ldots, X_n$ IID $f(x; \theta)$ is marginal density
$$f_X(x)$$

Q) What's the difference btwn a statistic and an estimate?

(more general)

Both however are fnctns of the random sample.

(targets a particular parameter)

## Deriving an estimator ("Recipes")

## Method 1: Method of moments

Consider the first few moments of the population distb'n

$\mu_1 = E[X^1]$

$\mu_2 = E[X^2]$

$\mu_3 = E[X^3]$

$\vdots$

Create a system of equations that can be solved for the parameter(s) $\theta$

Then, substitute the sample estimates of these moments into solution for $\theta$ above.

So
$$\begin{cases} \hat{\mu}_1 = \frac{1}{n} \sum_{i=1}^{n} x_i^1 \\ \hat{\mu}_2 = \frac{1}{n} \sum_{i=1}^{n} x_i^2 \\ \hat{\mu}_3 = \frac{1}{n} \sum_{i=1}^{n} x_i^3 \\ \vdots \end{cases}$$
take the place of $\mu_1, \mu_2, \ldots$ and the result is the estimator!

# Method 2: Maximize the Likelihood

**Q)** If $(X_1, X_2, \ldots, X_n)$ are an (IID) sample from a population w/ distb'n $F_X(x)$ and density $f_X(x)$, then

what is the joint density of $(X_1, X_2, \ldots, X_n)$?

$$(X_1, X_2, \ldots, X_n) \sim \prod_{i=1}^{n} f_{X_i}(x) = f_X^n(x) = f^n(x; \theta)$$

If we think of this joint distribution as a function of the <u>parameter(s)</u> for fixed (observed) data $(X_1, X_2, \ldots, X_n) = (x_1, x_2, \ldots, x_n)$ then we are referring to the ==likelihood== of the ==parameter(s)==, ==given the data.==

likelihood: $\text{lik}(\theta) = f(x; \theta)$

log-likelihood: $\ell(\theta) = \log\left[\text{lik}(\theta)\right]$

both can be vectors

Once we have a likelihood for $\theta$, often we can maximize this function (w.r.t. $\theta$). The maximum (global) is often a useful estimate for $\theta$.

~ Paradigm Shift! ~

## Method 3: Use Bayes' Theorem

Treat the parameter, $\theta$, as a RV. Come up w/ an initial guess for the distribution of $\theta \sim f_\theta(\theta)$.

Typically a prior is denoted as
$$\theta \sim p(\theta) \quad \text{or} \quad \theta \sim \pi(\theta)$$

Given a likelihood function for $\theta$, conditioned upon the observed data, $\underset{\sim}{x} = (x_1, x_2, \dots, x_n)$, use Bayes' theorem to find the conditional distribution $f(\theta | \underset{\sim}{x})$.

Typically, this posterior density is denoted as
$$\theta | \underset{\sim}{x} \sim \pi(\theta | \underset{\sim}{x})$$

Altogether we have

likelihood : $f(\underset{\sim}{x}; \theta)$

prior : $\pi(\theta)$

posterior : $\pi(\theta|\underset{\sim}{x}) = \dfrac{\pi(\theta) f(\underset{\sim}{x}; \theta)}{\displaystyle\int_{\Theta} \pi(\theta) f(\underset{\sim}{x}; \theta) d\theta}$

Prior distrib'n for $\theta$

likelihood for $\theta$ given data $\underset{\sim}{x} = (x_1, x_2, \ldots, x_n)$

"normalizing" constant

Q) What is $\Theta$?

$\Theta$ is the ==parameter space==

Often, we can ignore the "normalizing" constant and specify the posterior up to proportionality:

$$\pi(\theta|\underset{\sim}{x}) \propto \pi(\theta) f(\underset{\sim}{x}; \theta)$$

"is proportional to"

Note: The entire distb'n of the posterior is a distribution function estimate for $\theta$!

We can derive point estimates for $\theta$ by considering different qualities of the posterior. For example:   posterior mean
                  posterior mode

Q) Are these the only ways to derive an estimator?
no! there are infinite numb of ways
to derive an estimator

Q) How do we know if an estimator is useful?
This is what we'll discuss next!

Setting:

Given an sample $(x_1, x_2, \dots, x_n)$ of RVs that follow a distribution depending on unknown parameter $\theta$, denote

$$\hat{\theta}_n = \hat{\theta}_n(x) \text{ as an estimator for } \theta$$

Note: $\hat{\theta}_n(X)$ is a RV ; $\hat{\theta}_n(x)$ is a fixed constant.

## Desirable characteristics for estimators:

- consistency

$\hat{\theta}_n$ is <u>consistent</u> for $\theta$ if, for all $\varepsilon > 0$

$$\lim_{n \to \infty} Pr\left(|\hat{\theta}_n - \theta| > \varepsilon\right) = 0.$$

Q) What type of convergence is this?
this is an example of 'limit' in probability

Ex) (weak) LLN: sample moments $\xrightarrow{P}$ pop. moments

Note: continuous functions preserve consistency

- unbiased

$\hat{\theta}_n$ is <u>unbiased</u> if $E[\hat{\theta}_n] = \theta$,
ie. the center of its sampling distb'n is $\theta$.

# Evaluating an estimator

**Def**: Mean Square Error

If we are targeting parameter $\theta$ w/ an estimator $\hat{\theta}_n$, then

$$MSE(\hat{\theta}_n) = E[(\hat{\theta}_n - \theta)^2]$$

trick is to $\pm E(\hat{\theta}_n)$

$$\vdots$$

$$= (E(\hat{\theta}_n) - \theta)^2 + Var(\hat{\theta}_n)$$

bias$^2$          variance

## Strategies to show consistency:

If $\hat{\theta}_n$ is unbiased — subsititute $E(\hat{\theta}_n)$ in for $\theta$
then apply a limiting $\leq$

If $\hat{\theta}_n$ (potentially) biased — then we have to usually work w/ the CDF of $\hat{\theta}_n$

$$Pr(|\hat{\theta}_n - \theta| > \varepsilon) = Pr(\hat{\theta}_n > \theta + \varepsilon) + Pr(\hat{\theta}_n < \theta - \varepsilon)$$

evaluate seperately

But,
Some times you have to get more creative!
Eg. Jensen's $\leq$ could be used to prove brasedness.
(the strict version)

## Topic - Detour for errata

**Recall**

### Jensen's Inequality

For any ~~convex~~ concave up function, $g$, and any RV $X$,

$$g(E[X]) \leq E[g(X)]$$

provided $E[g(X)]$ and $g(E[X])$ exist and are finite.

Correction: My heuristic for remembering concave/convex doesn't work!

**Instead**

looks like a cave!

convex = concave up
$f'' > 0$

concave down
$f'' < 0$

Q) When is the inequality **strict**?

When the concavity is strict
(no plateaus)

## Now back to properties of estimators...

Ex). Suppose $X_1, \ldots, X_n$ are IID from $U(0, \theta)$. Consider the following estimates and determine if they are consistent.

| Consistent | Unbiased | Estimate |
|---|---|---|
| yes | yes | $\hat{\theta}_1 = 2\bar{X}$ |
| No | yes | $\hat{\theta}_2 = 2X_1$ |
| yes | no | $\hat{\theta}_3 = X_{(n)}$ |
| no | no | $\hat{\theta}_4 = \frac{1}{2}X_1^2$ |

$n^{th}$ order statistic, i.e. largest observation

only use the first observation! (not necessarily the minimum!)

Note:

$X_1$ has density

$$f_{X_1}(x) = \frac{1}{\theta} \, \mathbb{I}\{0 < x \leq \theta\}$$

and CDF

$$F_{X_1}(x) = \Pr(X_1 \leq x) = \frac{x}{\theta} \, \mathbb{I}\{0 < x \leq \theta\}$$

$$\boxed{\hat{\theta}_1 = 2\bar{x}}$$

## Unbiased?

$$E[2\bar{x}] = 2E\left[\frac{1}{n}\sum_{i=1}^{n} X_i\right] = 2 \times \left(\frac{1}{n}\right) \times E[X_1 + X_2 + \dots + X_n]$$

$$= \frac{2}{n} \times \left(E(X_1) + E(X_2) + \dots + E(X_n)\right)$$

$$= \frac{2}{n} \cdot n\, E(X_1) = 2\left(\frac{\theta}{2}\right) = \theta \quad ✓$$

is unbiased

## Consistent?

$\bar{x}$ is consistent for $E[X_1] = \frac{\theta}{2}$. **Why?**

A: b/c sample moments are consistent for population moments (LLN)

$g(\bar{x}) = 2x$ is a continuous function

$\therefore \quad g(\bar{x}) = 2\bar{x}$ is consistent for $2E[X_1] = \theta$

✓ is consistent

$\boxed{\hat{\Theta}_2 = 2X_1}$

## Unbiased?

$$E(\hat{\Theta}_2) = E(2X_1) = 2E(X_1) = 2 \cdot \frac{\Theta}{2} = \Theta \quad \checkmark \quad \text{is unbiased}$$

## Consistent?

$$Pr(|\hat{\Theta}_2 - \Theta| > \varepsilon) = Pr(|\hat{\Theta}_2 - E(\hat{\Theta}_2)| > \varepsilon)$$
$$= Pr(|2X_1 - E(2X_1)| > \varepsilon)$$
$$\le \frac{Var(2X_1)}{\varepsilon^2}$$

> Chebyshev:
> $$P(|X - E(X)| > \varepsilon) \le \frac{Var(X)}{\varepsilon^2}$$

$$= \frac{4 Var(X_1)}{\varepsilon^2} \quad \leftarrow \text{not a fnctn of } n \text{ so not helpful.}$$

Let's try the CDF approach...

$$Pr(|\hat{\Theta}_2 - \Theta| > \varepsilon) = Pr(2X_1 > \Theta + \varepsilon) + Pr(2X_1 < \Theta - \varepsilon)$$
$$= Pr\left(X_1 > \frac{\Theta + \varepsilon}{2}\right) + Pr\left(X_1 < \frac{\Theta - \varepsilon}{2}\right)$$
$$= 1 - Pr\left(X_1 \le \frac{\Theta + \varepsilon}{2}\right) + \left(\frac{\Theta - \varepsilon}{2}\right)\Big/ \Theta$$
$$= 1 - \frac{\left(\frac{\Theta + \varepsilon}{2}\right)}{\Theta} + \frac{\left(\frac{\Theta - \varepsilon}{2}\right)}{\Theta}$$
$$= 1 - \frac{\Theta + \varepsilon}{2\Theta} + \frac{\Theta - \varepsilon}{2\Theta}$$
$$= \frac{2\Theta - \Theta + \varepsilon + \Theta - \varepsilon}{2\Theta}$$
$$= 1 \quad \bigotimes \text{ not consistent}$$

Note: Estimator is function of $X_1$ only... so we really didn't need to do all that work!

$$\boxed{\hat{\Theta}_3 = X_{(n)}}$$

CDF for $\hat{\Theta}_3$: $\quad Pr(\hat{\Theta}_3 \leq x) = Pr(X_{(n)} \leq x)$

$\quad\quad$ by def$^n$ of $\quad = Pr(X_1 \leq x, X_2 \leq x, \ldots, X_n \leq x)$
$\quad\quad X_{(n)}$

$\quad\quad$ by indep $\quad = Pr(X_1 \leq x) \cdot Pr(X_2 \leq x) \cdots Pr(X_n \leq x)$

$\quad\quad$ by identical $\quad = [Pr(X_1 \leq x)]^n$

$\quad\quad\quad\quad = \left(\frac{x}{\theta}\right)^n \mathbb{I}\{0 < x \leq \theta\}$

density for $\hat{\Theta}_3$:

$$f_{\hat{\Theta}_3}(x) = n \cdot x^{n-1} \cdot \frac{1}{\theta} \mathbb{I}\{0 < x \leq \theta\}$$

---

## Unbiased?

$$E[\hat{\Theta}_3] = E[X_{(n)}] = \int_0^\theta \frac{n}{\theta} \cdot x \cdot x^{n-1} \, dx = \frac{n}{\theta} \int_0^\theta x^n \, dx$$

$$= \frac{n}{\theta} \left( \frac{x^{n+1}}{n+1} \Big|_{x=0}^\theta \right)$$

$$= \frac{n}{\theta} \left[ \frac{\theta^{n+1}}{n+1} - 0 \right] = \frac{n\theta^{n+1}}{(n+1)\theta} \quad \bigotimes$$

NOT unbiased

## Consistent?

$$Pr(|\hat{\Theta}_3 - \theta| > \varepsilon) = Pr(X_{(n)} > \theta + \varepsilon) + Pr(X_{(n)} < \theta - \varepsilon)$$

$$= \int_{\theta+\varepsilon}^\theta f_{X_{(n)}}(x) \, dx + \int_0^{\theta-\varepsilon} f_{X_{(n)}}(x) \, dx$$

$(\varepsilon > 0)$
$$= 0 + \int_0^{\theta-\varepsilon} \frac{n}{\theta} x^{n-1} \, dx$$

$$= \frac{n}{\theta} \int_0^{\theta-\varepsilon} x^{n-1} \, dx \quad\quad \left[ \begin{array}{c} \text{can assume} \\ \varepsilon < \theta \end{array} \right]$$

$$= \frac{n}{\theta} \left[ \frac{x^n}{n} \Big|_{x=0}^{\theta-\varepsilon} \right]$$

$$= \frac{1}{\theta} \left[ (\theta-\varepsilon)^n - 0 \right]$$

$$\lim_{n \to \infty} \frac{(\theta-\varepsilon)^n}{\theta} = 0 \quad \text{since } \theta \in (0,1) \quad \checkmark \quad \text{is consistent}$$

$$\boxed{\tilde{\theta}_4 = 1/X_1^2}$$

## Unbiased?

$$E\left[\frac{1}{X_1^2}\right] = \int_0^\theta \frac{1}{X_1^2} f_{X_1}(x_1)\, dx_1 \;=\; \int_0^\theta \frac{1}{X_1^2} \cdot \frac{1}{\theta}\, dx_1 \;=\; \frac{1}{\theta}\int_0^\theta \frac{1}{X_1^2}\, dx_1$$

$$= \frac{1}{\theta}\left[\left.\frac{-1}{X_1}\right|_{X_1=0}^{\theta}\right] \qquad \text{undefined} \; \otimes$$

not unbiased b/c
expectation doesn't exist!

## Consistent?

Again, estimator is a function of $X_1$ only.
So what happens as $n \to \infty$?

Nothing. The estimator doesn't change w/
the sample size.

$\otimes$ not consistent

9-16-22

Ex) Stakeholder analysis of using a consistent estimator

$\theta_1$ = dosage that max benefit/min harm

$\theta_2$ = change in B cell counts after using medication

} Possible parameters

$\hat{\theta}_1$ = 10 mg/kg

$\hat{\theta}_2$ = "X" change in fluoresence intensity

} possible estimates

Suppose $\hat{\theta}_n$ = 10mg/kg is a consistent estimator for $\theta$ = dosage that max benefit & min harm

**Choice/Decision**: Decide whether or not to use a drug to treat Systemic Lupus Erythematous within the first few years of diagnosis. Here is an example of a pilot study currently ongoing.

| Stakeholder | Potential results | |
| --- | --- | --- |
| | **Harm** | **Benefit** |
| Medical practitioners<br><br>perscribe ^On dose to patient | possibly not all patients are represented in the population for which we have a sample | for the majority of the population this estimated dosage will be the best dosage |
| Medication users<br><br>take ^On dosage | | |
| | | |

- Example harms: cost of money, time, effort; negative impact to reputations; can be tangible or intangible with immediate or delayed effects
- Example benefits: earning or gaining money; removal of a harm; saved time or effort; improved reputation; demonstration of expertise.

*Source*: Tractenberg, R. E. (2019). Teaching and Learning about ethical practice: The case analysis. https://doi.org/10.31235/OSF.IO/58UMW

# Topic: Estimation Part II (Ch.8)

## Large Sample Theory for MLEs

Setting: $X_1, \ldots, X_n$ IID $f(x; \theta)$

$$\text{lik}(\theta) = \prod_{i=1}^{n} f(x_i; \theta)$$

$$\ell(\theta) = \sum_{i=1}^{n} \ln(f(x_i; \theta))$$

MLE specifically

$\hat{\theta}_n$ = value of $\theta$ that maximizes $\text{lik}(\theta)$

$\theta_0$ = true, unknown value of $\theta$

$n \to \infty$

---

**Def**: The **score** is the gradient (first derivative) of the likelihood fnctn.

$$S(\theta) = \frac{\partial}{\partial \theta} \ell(\theta)$$

rate of change in (log) likelihood

**Note**: $\hat{\theta}_n$ (the MLE for $\theta$, given $\underset{\sim}{X}_{obs}$) is a "zero" of $S(\theta)$
ie. $S(\hat{\theta}_n) = 0$

**Thm**: If $f(x; \theta)$ is "smooth enough", then the MLE is consistent.

**Note**: The expected value of $S(\theta)$ is $\emptyset$ at $\theta = \theta_0$.
b/c...

mistake starts here!

$$E[S(\theta)] = E\left[\frac{\partial}{\partial \theta} \ell(\theta)\right] = \int \left[\frac{\partial}{\partial \theta} \ell(\theta)\right] f(x; \theta) dx$$

$$= \int \left[\frac{1}{f(x; \theta)} \frac{\partial f(x; \theta)}{\partial \theta}\right] f(x; \theta) dx$$

$$= \int \frac{\partial}{\partial \theta} f(x; \theta) dx \underset{*}{=} \frac{\partial}{\partial \theta} \int f(x; \theta) dx = \frac{\partial}{\partial \theta} 1 = 0$$

Fixed version of $E[s(\theta)] = 0$ :

$$E[s(\theta)] = E\left[\frac{\partial}{\partial\theta}\ell(\theta)\right] = \int \cdots \int \left[\frac{\partial}{\partial\theta}\ell(\theta)\right] f(x_1, \ldots, x_n; \theta_0) \, dx_1 \ldots dx_n$$

$$\overbrace{\phantom{= \int \cdots \int}}^{n \text{ times}}$$

at $\theta = \theta_0$.

$$= \int \cdots \int \frac{\frac{\partial}{\partial\theta} f(x_1, \ldots, x_n; \theta)}{f(x_1, \ldots, x_n; \theta)} f(x_1, \ldots, x_n; \theta_0) \, dx_1 \ldots dx_n$$

$$= \int \cdots \int \frac{\partial}{\partial\theta} f(x_1, \ldots, x_n; \theta_0) \, dx_1 \ldots dx_n$$

* If we can
interchange
$\frac{\partial}{\partial\theta}$ & $\int$

$$\hookrightarrow \quad = \frac{\partial}{\partial\theta} \int \cdots \int f(x_1, \ldots, x_n; \theta_0) \, dx_1, \ldots, dx_n$$

$$= \frac{\partial}{\partial\theta}(1)$$

$$= 0 \qquad \qquad \boxed{\phantom{x}}$$

Related to

    * <u>Calc Thm:</u> Leibniz Integral Rule (special case)

$$\frac{d}{dx}\left(\int_a^b f(x, u)\, du\right) = \int_a^b \left[\frac{\partial}{\partial x} f(x, u)\right] du$$

Def: The ==Fisher Information== is the variance of the score.

$$I_n(\theta) = E\left\{\left[\frac{\partial}{\partial\theta}\ell(\theta)\right]^2\right\} \longleftarrow$$

2nd moment of score fnctn

Thm: <u>Information Identity</u>

If $f(x;\theta)$ is "smooth enough", then

$$I_n(\theta) = E\left\{\left[\frac{\partial}{\partial\theta}\ell(\theta)\right]^2\right\} = -E\left[\frac{\partial^2}{\partial\theta^2}\ell(\theta)\right]$$

Thm: <u>Asymptotic Normality of MLEs</u>

If $f(x;\theta)$ is "smooth enough", then

$$\sqrt{n\,I_n(\theta_0)}\,(\hat{\theta}_n - \theta_0) \xrightarrow{d} N(0,1).$$

Q) What does it mean for an estimate to be "optimal"?

**Def:** Suppose $\hat{\theta}_1$ and $\hat{\theta}_2$ are estimators of $\theta$ that have the same bias. Ie. $E[\hat{\theta}_1] - \theta = E[\hat{\theta}_2] - \theta$.
The ==efficiency== of $\hat{\theta}_1$ relative to $\tilde{\theta}_2$ is

$$eff(\hat{\theta}_1, \hat{\theta}_2) = Var(\hat{\theta}_2) / Var(\hat{\theta}_1).$$

**Note:** If we are comparing asy. variance of an estimator, we call this the ==asymptotic relative efficiency.==

**Thm:** <u>Cramér-Rao Inequality</u>
Suppose $X_1, \ldots, X_n$ are IID $f(x;\theta)$, where $f(x;\theta)$ is "smooth enough". Let $T = T(\underset{\sim}{X})$ be an <u>unbiased</u> estimate of $\theta$. Then $Var(T) \geq \dfrac{1}{n I_n(\theta)}$.

cramér-Rao Lower Bound

**Note:** An unbiased estimate w/ variance equal to the $\boxed{CR-LB}$ is said to be effecient.

**Note:** As $n \to \infty$, the MLE is asymptotically effecient.

(Q) Is asymptotic unbiasedness the same thing as consistent? Why/why not?

# Notation

9-21-22

$f(X_i; \theta)$ density for $X_i$

If $X_1, \ldots, X_n$ are IID then the likelihood is:

$$lik(\theta) = \prod_{i=1}^{n} f(X_i, \theta)$$

and the log-likelihood is:

$$\ell(\theta) = \sum_{i=1}^{n} \log\left(f(X_i; \theta)\right)$$

The score function is the gradient of the log-likelihood:

$$\frac{\partial}{\partial \theta} \ell(\theta)$$

The score function has mean zero and variance equal to the Fisher Information

$$I_n(\theta) = E\left\{\left[\frac{\partial}{\partial \theta} \ell(\theta)\right]^2\right\}$$

Info about $\theta$ contained in $(X_1, \ldots, X_n)$

Your textbook considers the score for a single RV, $X$:

$$\frac{\partial}{\partial \theta} \log f(X; \theta)$$

where the Fisher Info is thus

$$I(\theta) = E\left\{\left[\frac{\partial}{\partial \theta} \log f(X; \theta)\right]^2\right\}$$

is the info about $\theta$ contained in $X$ alone.

Consider the log density $\log(f(X;\theta))$:

Q) What is the $1^{st}$ (population) moment?

$$E\left[\underbrace{\log(f(X;\theta))}_{g(X)}\right] = \int \underbrace{[\log(f(x;\theta))]}_{g(X)} \underset{\underset{\text{density for } X!}{\uparrow}}{f(x;\theta_0)}\, dx$$

Q) What is the $1^{st}$ sample moment?

$$\frac{1}{n}\sum_{i=1}^{n} \log(f(x_i;\theta)) = \frac{1}{n}\ell(\theta)$$

Now consider the gradient of the log density $\frac{\partial}{\partial\theta}\log(f(X;\theta))$:

Q) What is the $1^{st}$ (population) moment?

$$E\left[\frac{\partial}{\partial\theta}\log f(X;\theta)\right] = \int \frac{\partial}{\partial\theta}\log f(x;\theta)\, f(x;\theta_0)\, dx$$

$$= \int \frac{\frac{\partial}{\partial\theta}f(x;\theta)}{f(x;\theta)} f(x;\theta_0)\, dx \;\underset{\textcolor{blue}{\bigstar}}{=}\; \int \frac{\partial}{\partial\theta}f(x;\theta_0)\, dx = \frac{\partial}{\partial\theta}\int f(x;\theta_0)\, dx = \frac{\partial}{\partial\theta}(1) = 0$$

Q) What is the (population) variance?

$$\mathrm{Var}\left[\frac{\partial}{\partial\theta}\log f(X;\theta)\right] = E\left\{\left[\frac{\partial}{\partial\theta}\log f(X;\theta)\right]^2\right\} - \left\{E\left[\frac{\partial}{\partial\theta}\log f(X;\theta)\right]\right\}^2$$

$$= I(\theta)$$

$\textcolor{blue}{\bigstar}$ $\textcolor{blue}{\text{at } \theta = \theta_0}$

# Warm-up group work: 9-23-22

## 5 mins
### Identify strategies, stuck points, approachs you tried to solve assigned HW 8 problem

|          | # 1                        | # 2                        | # 3                        |
|----------|----------------------------|----------------------------|----------------------------|
| Sec 1:   | Seth<br>Sherry<br>Miles    | Koji<br>Annie<br>Amy       | Brian<br>Patty<br>Guy<br>Tillie |
| Sec 2:   | Mwangangi<br>Tinashe<br>Zack | Ben H<br>Atesh<br>Jonathan | Joey<br>Jason<br>Rodas     |
|          | Alex<br>Jorge<br>Ben C     | Sarah<br>Hellman<br>Ian    | Nancy<br>Noha<br>Gertrud   |

# Review & consider:
## What strategies/approaches were most useful?

# Sufficiency

Setting: $X_1, \ldots, X_n \overset{IID}{\sim} f(x; \theta)$

$$lik(\theta) = \prod_{i=1}^{n} f(x_i; \theta)$$

$$\ell(\theta) = \sum_{i=1}^{n} \ln(f(x_i; \theta))$$

$\hat{\theta}_n = \hat{\theta}(x_1, \ldots, x_n)$

is an estimator for $\theta$.

Q) Is there an estimator that contains as much information about $\theta$ as the entire sample, $x_1, \ldots, x_n$?

Def: A statistic $T = T(X_1, \ldots, X_n)$ is **sufficient** for parameter $\theta$ if

$$(X_1, \ldots, X_n) | T = t$$

follows a distribution that does not depend on $\theta$.

Thm: Factorization Theorem

Statistic $T(X_1, \ldots, X_n)$ is sufficient for $\theta$

iff

$$f(x_1, \ldots, x_n; \theta) = g[T(x_1, \ldots, x_n); \theta] \cdot h(x_1, \ldots, x_n)$$

likelihood

must involve <u>all</u> of the observed data!

# Exponential Family

The family of probability distb'n functions that have sufficient statistics of the same dimension as the parameter space is called the ==exponential family==.

1-Parameter Exponential family:

$$f(x; \theta) = \exp\left\{ c(\theta) T(x) + d(\theta) + S(x) \right\}$$

for all $x \in A$ where $A \perp\!\!\!\perp \theta$

$K$-parameter Exponential family:

$$f(x; \theta) = \exp\left\{ \sum_{j=1}^{K} c_j(\theta) T_j(x) + d(\theta) + S(x) \right\}$$

for all $x \in A$ where $A \perp\!\!\!\perp \theta$

Note: If $T$ is sufficient for $\theta$, then the MLE is a function of $T$.

We can see this is the case b/c...

$T(x_1, \ldots, x_n)$ sufficient means

$$\text{lik}(\theta) = \underbrace{f(x_1, \ldots, x_n; \theta)}_{\substack{\text{maximize} \\ \text{wrt } \theta}} = \underbrace{g[T(x_1, \ldots, x_n), \theta]}_{\substack{\text{maximize} \\ \text{wrt } \theta}} \cdot \underbrace{h(x_1, \ldots, x_n)}_{\perp\!\!\!\perp \theta}$$

## Thm: Rao-Blackwell Theorem

Let $\hat{\theta}$ be an estimator for $\theta$ s.t. $E(\hat{\theta}^2) < \infty$. If $T$ is sufficient for $\theta$ and if $\tilde{\theta} = E[\hat{\theta} | T]$, then, for all $\theta$,

$$MSE \longrightarrow E\left[(\tilde{\theta} - \theta)^2\right] \leq E\left[(\hat{\theta} - \theta)^2\right].$$

Furthermore, the inequality is strict unless $\hat{\theta} = \tilde{\theta}$.

Note: If an estimator is not a function of a sufficient statistic, and if a sufficient statistic exists, then the estimator can be improved!

# Group Work:
## Dissecting Proofs

## Example : Information Identity

Define $I(\theta) = E\left\{\left[\frac{\partial}{\partial\theta}\log f(X;\theta)\right]^2\right\}$

If $f(\cdot)$ is "smooth enough", the we have

$$E\left\{\left[\frac{\partial}{\partial\theta}\log f(X;\theta)\right]^2\right\} = -E\left[\frac{\partial^2}{\partial\theta^2}\log f(X;\theta)\right].$$

1. Confusing steps?

combining identities in a useful way

how does $\frac{\partial}{\partial\theta}\int\left[\frac{\partial}{\partial\theta}\log f(x;\theta)\right]f(x;\theta)dx = \int\left[\frac{\partial^2}{\partial\theta^2}\log f(x;\theta)\right]f(x;\theta)dx$
$+ \int\left[\frac{\partial}{\partial\theta}\log f(x;\theta)\right]^2 f(x;\theta)dx$

2. Useful techniques?

the fact that $\int f(x;\theta)dx = 1$ ; swapping $\frac{\partial}{\partial\theta}$ and $\int\cdot dx$
and
rearrange $\frac{\partial}{\partial\theta}\log f(x;\theta) = \frac{\partial/\partial\theta\, f(x;\theta)}{f(x;\theta)}$

3. Narrative?

use property of         + take 2nd
density fnctns             deriv                    + rearrange
       +                                                      identities
swap diff. & integ.    applying calc. rules    = result

# Example: Working Thru Steps of Cramér-Rao

(For me, these were the most confusing steps in this proof.)

pg 30

$$E[ZT] = E\left\{ \left[ \sum_{i=1}^{n} \frac{\partial}{\partial\theta} \log(f(X_i;\theta)) \right] T(X_1,...,X_n) \right\} =$$

$$\int \cdots \int t(x_1,...,x_n) \left[ \sum_{i=1}^{n} \frac{\partial}{\partial\theta} \log\left(f(x_i;\theta)\right) \right] f(x_1,...,x_n;\theta) \, dx_1 \cdots dx_n$$

$\underbrace{\qquad}$ $n$ times

$$=$$

$$\int \cdots \int t(x_1,...,x_n) \left[ \sum_{i=1}^{n} \frac{\partial}{\partial\theta} \log\left(f(x_i;\theta)\right) \right] \left[ \prod_{j=1}^{n} f(x_j;\theta) \, dx_j \right]$$

and note $\displaystyle\sum_{i=1}^{n} \frac{\frac{\partial}{\partial\theta} f(x_i;\theta)}{f(x_i;\theta)} \prod_{j=1}^{n} f(x_j;\theta) =$

$$\frac{\frac{\partial}{\partial\theta} f(x_1;\theta)}{f(x_1;\theta)} \Big( f(x_1;\theta) f(x_2;\theta) \cdots f(x_n;\theta) \Big)$$

$$+ \frac{\frac{\partial}{\partial\theta} f(x_2;\theta)}{f(x_2;\theta)} \Big( f(x_1;\theta) \cdots f(x_n;\theta) \Big)$$

$$+ \cdots + \frac{\frac{\partial}{\partial\theta} f(x_n;\theta)}{f(x_n;\theta)} \Big( f(x_1;\theta) \cdots f(x_n;\theta) \Big)$$

$$E[Z] = E\left[ \sum_{i=1}^{n} \frac{d}{d\theta} \log\left( f(X_i; \theta) \right) \right]$$

$$= E\left[ \sum_{i=1}^{n} \frac{\frac{d}{d\theta} f(X_i; \theta)}{f(X_i, \theta)} \right]$$

$$= \sum_{i=1}^{n} E\left[ \frac{\frac{d}{d\theta} f(X_i; \theta)}{f(X_i; \theta)} \right]$$

$$= \sum_{i=1}^{n} \left\{ \int \left[ \frac{\frac{d}{d\theta} f(x_i; \theta)}{f(x_i; \theta)} \right] f(x_i; \theta_0) \, dx_i \right\}$$

$$\overset{at\ \theta=\theta_0}{=} \sum_{i=1}^{n} \left\{ \int \frac{d}{d\theta} f(x_i; \theta_0) \, dx_i \right\}$$

$$= \sum_{i=1}^{n} \left\{ \frac{d}{d\theta} \int f(x_i; \theta_0) \, dx_i \right\}$$

$$= \sum_{i=1}^{n} \left\{ \frac{d}{d\theta} (1) \right\}$$

$$= \frac{\partial}{\partial \theta} f(x_1; \theta) \left[ f(x_2; \theta) \cdots f(x_n; \theta) \right]$$

$$+ \frac{\partial}{\partial \theta} f(x_2; \theta) \left[ f(x_1, \theta) f(x_3, \theta) \cdots f(x_n, \theta) \right]$$

$$+ \cdots$$

$$+ \frac{\partial}{\partial \theta} f(x_n; \theta) \left[ f(x_1; \theta) d \; f(x_2; \theta) \cdots f(x_{n-1}; \theta) \right)$$

$$= \frac{\partial}{\partial \theta} \prod_{i=1}^{n} f(x_i, \theta) .$$

Hence

$$E[ZT] =$$

$$\int \cdots \int t(x_1, \ldots, x_n) \left[ \sum_{i=1}^{n} \frac{\partial}{\partial \theta} \log \left( f(x_i; \theta) \right) \right] \left[ \prod_{j=1}^{n} f(x_j; \theta) \, dx_j \right]$$

$$= \int \cdots \int t(x_1, \ldots, x_n) \frac{\partial}{\partial \theta} \prod_{i=1}^{n} f(x_i; \theta) \, dx_i$$

$$= \frac{\partial}{\partial \theta} \int \cdots \int t(x_1, \ldots, x_n) \prod_{i=1}^{n} f(x_i; \theta) \, dx_i \cdots dx_n$$

$$= \frac{\partial}{\partial \theta} E[T(X_1, \ldots, X_n)]$$

**Bayesian Estimation/Prediction**

$\theta_o$ = prob. that bball player successfully makes a shot

prior: $\pi(\theta) \sim U[0,1]$

obs. data: 2 successful shots in a row

assume: outcomes (of shots) are independent

(a) What is the posterior density of $\theta$?

(b) What would you estimate is the probability
that this player makes a third shot?

---

What is the (probability) model for the data?

Let $x = \begin{cases} 0, \text{ miss} \\ 1, \text{ score} \end{cases}$  $X \sim Bern(\theta_o)$

$$P_r(X = x) = \theta^x (1-\theta)^{1-x}$$

Now we can evaluate the likelihood for $\theta$
given the observed outcomes (data):

$x_1 = 1, x_2 = 1$

$$P_r(X_1 = 1, X_2 = 1) = P_r(X_1 = 1) \cdot P_r(X_2 = 1)$$

$$= \theta^1 (1-\theta)^{1-1} \cdot \theta^1 (1-\theta)^{1-1}$$

$$= \theta^2$$

What is the prior density on $\theta$?

$$f_\theta(\theta) = 1 \cdot \mathbb{I}\{0 \le \theta \le 1\} = \pi(\theta)$$

Now we can evaluate the posterior, conditioned upon the observed data:

$$\pi(\theta \mid x_1=1, x_2=1) = \frac{\pi(\theta) \cdot f(x_1=1, x_2=1; \theta)}{\int_0^1 \pi(\theta) f(x_1=1, x_2=1; \theta) d\theta}$$

$$= \frac{1 \cdot \mathbb{I}\{0 \leq \theta \leq 1\} \cdot \theta^2}{\int_0^1 1 \cdot \theta^2 d\theta}$$

$$\vdots$$

$$= \frac{\theta^2}{\theta^3/3 \big|_{\theta=0}^1} \quad \mathbb{I}\{0 \leq \theta \leq 1\} = \cdots = \boxed{3\theta^2},$$
for $0 \leq \theta \leq 1$

Finally, we can check our answer by verifying that $\underline{\int \pi(\theta \mid \underline{x}) d\theta = 1}$:

$$\int_0^1 3\theta^2 d\theta = \cdots = 1$$

Part (b) is a question about how to use the posterior to estimate the true value of $\theta$.

$$E(\theta \mid x_1=1, x_2=1) = \int_0^1 \theta \cdot \pi(\theta \mid x_1=1, x_2=1) d\theta$$

$$= \int_0^1 3\theta^3 d\theta$$

# Group Work Results
## for Dissecting Proofs Worksheet

## Cramér-Rao Inequality

Most confusing steps: $E[z] = 0$

$$Cov(z, T) = E[zT]$$

Jee        Example: Working Thru Steps of Cramér-Rao above!

Tricks & techniques:   chain rule

Leibniz rule for diff & int.

properties of Score &

definition of Fisher info.

Story:

## Rao-Blackwell Thm

Most confusing steps: "$\text{Var}(\hat{\theta}|T) = 0$ only if..."

understanding what is meant by $\tilde{\theta}$.

how does comparing MSE's come down to comparing variances?

Note: $E(\hat{\theta}) = E[E[\hat{\theta}|T]]$ by law of iterated expectat.

so $\hat{\theta}$ and $\tilde{\theta} = E[\hat{\theta}|T]$ have the same bias!

Also note: If $\hat{\theta}$ is a function of $T$, then
$\tilde{\theta} = E[\hat{\theta}|T] = E[\hat{\theta}(T)|T]$ is not random!

Tricks & techniques: law of iterated expectation
and E-V-E property of conditional variance

Story:

# Factorization Thm

Most confusing steps: 
$$\frac{Pr(\underset{\sim}{X}=x, T=t)}{Pr(T=t)} = \frac{h(\underset{\sim}{x})}{\sum\limits_{T(\underset{\sim}{x})=t} h(\underset{\sim}{x})}$$

how to get $g(t;\theta) \sum\limits_{T(\underset{\sim}{x})=t} h(\underset{\sim}{x})$ ?

Suppose $X_1, \dots X_n$ are continuous over sample space $\mathcal{X}$. Then

$$Pr(\underset{\sim}{X}=\underset{\sim}{x}, T=t) = Pr(X_1=x_1, X_2=x_2, \dots, X_n=x_n, T(X_1, \dots, X_n)=t)$$

$$= \int \dots \int_A f(x_1, x_2, \dots, x_n; \theta) dx_1 \dots dx_2$$

where $A$ is $\{\underset{\sim}{x} \in \mathcal{X} : T(x_1, \dots, x_n)=t\}$

(by assumption) $= \int \dots \int_A g(T(x_1, \dots, x_n)) \cdot h(x_1, \dots, x_n) dx_1 \dots dx_n$

Tricks & techniques:

expand joint density terms.

manipulate sums.

assume A, deduce B. then assume B, deduce A.

Story:

# Topic: Estimation Part III
### (ch 8)

## Confidence Intervals — quantify the uncertainty inherent to point estimation using properties of random sampling from an assumed model

↳ indirect assessment of uncertainty

For IID data

$$(X_1, \ldots, X_n) \sim \prod_{i=1}^{n} f(x_i; \theta_0)$$

↳ assumed model

parameter fixed, unknown always a constant

Recall

$$\hat{\Theta}_n = \hat{\Theta}(X_1, \ldots, X_n) \quad \text{is a point estimate for } \theta_0$$

↳ is random, has a sampling distb'n

but

$$\hat{\Theta}_n = \hat{\Theta}(x_1, \ldots, x_n) \quad \text{is the point estimate evaluated for observed data.}$$

↳ is fixed, data has been observed

Similarly,
A confidence interval for $\theta_0$ is a <u>random interval</u> ... until the data is observed.

↳ random b/c it is a fnctn of $X_1, \ldots, X_n$

## Process:

Use the sampling dist'bn of $\hat{\Theta}_n$ (in particular the sampling variance of $\hat{\Theta}_n$) to identify a lower bound (LB) and upper bound (UB) on the most plausible values for $\Theta_0$.

## Interpretation:

Although we say we are $(1-\alpha) \times 100\%$ confident that the true value of $\Theta$ (ie. $\Theta_0$) lies w/in $[LB, UB]$, what we mean is something a bit more involved...

Based on the assumed model for the data, the probability that the random interval $[LB(\hat{\Theta}_n), UB(\hat{\Theta}_n)]$ contains the value of $\Theta$ that generated the data, $\Theta_0$, is $(1-\alpha)$.

## Tips & techniques:

Often, it is useful to plot the density (or mass) function for the sampling dist'bn of $\hat{\Theta}_n$ to identify which dist'bn quantiles to use in the CI.

# Example of exact and approximate CIs

## HW 8 #2b

$X_1, \ldots, X_n \overset{\text{IID}}{\sim} \text{Exp}(\tau)$

> Note: This version is consistent w/ the parameterization in your textbook

> (the "rate" parameterization vs. "scale")

$$\text{lik}(\theta) = \prod_{i=1}^{n} f(x_i; \tau) = \prod_{i=1}^{n} \left[ \tau e^{-\tau x_i} \, \mathbb{I}\{x_i \geq 0\} \right] = \tau^n e^{-\tau \sum_{i=1}^{n} x_i} \, \mathbb{I}\{x_{(1)} \geq 0\}$$

$\hat{\Theta}_{MLE} = \bar{X} \quad \longleftarrow$

> To do: use this sampling distrb'n to find a $(1-\alpha)100\%$ CI for $\tau$.

Given: $\sum_{i=1}^{n} X_i \sim \text{Gamma}(n, \tau)$

Derive: $\bar{X} \sim \text{Gamma}(n, n\tau)$

> All changes are highlighted in green.
>
> View the notes for 9-28-22 to see the other version



Gamma$(n, n\tau)$ density

$\frac{\alpha}{2}$

lower $\frac{\alpha}{2}^{th}$ quantile

$\frac{\alpha}{2}$

lower $(1-\frac{\alpha}{2})^{th}$ quantile (ie. upper $(\frac{\alpha}{2})^{th}$ quantile)
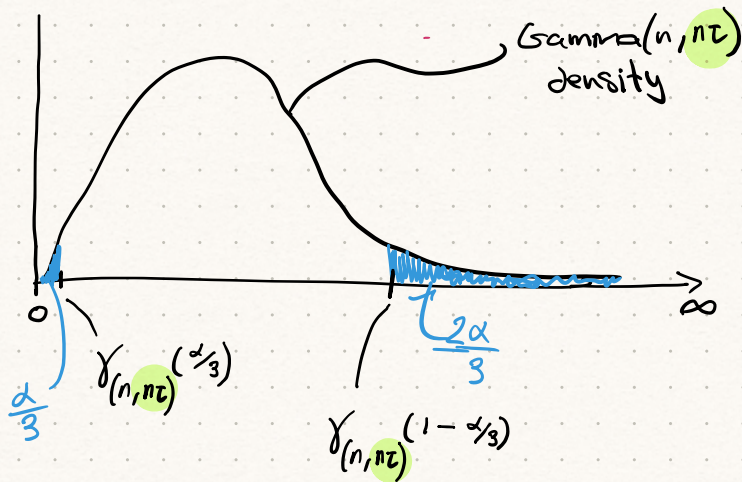
Notation:

$$\gamma_{(n, n\tau)}\left(\frac{\alpha}{2}\right)$$

Notation:

$$\gamma_{(n, n\tau)}\left(1 - \frac{\alpha}{2}\right)$$

Note, we could asymmetrically
choose the quantiles, e.g.



Gamma$(n, n\tau)$
density

$\gamma_{(n, n\tau)}(\alpha/3)$

$\frac{\alpha}{3}$

$\frac{2\alpha}{3}$
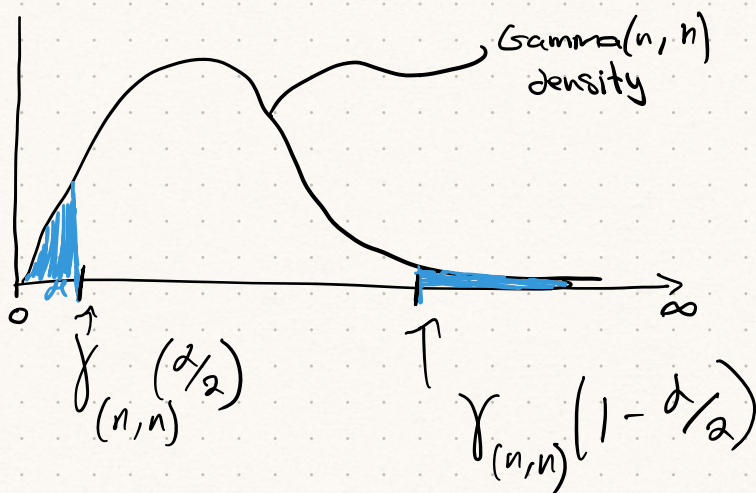
$\gamma_{(n, n\tau)}(1 - \alpha/3)$

But, in either case,
since $\tau$ is unknown, we can't find these
exact quantiles. Instead, we'll try to find
a way to express this idea in term of quantiles
from a distribution w/ no unknown parameters

Using properties of the Gamma distrb'n
we note that

$$\frac{1}{\bar{X}} \sim Gamma(n, n)$$

This is called
a "pivot" b/c
the distrb'n does
NOT depend on
any unknowns.

Hence



Gamma(n, n)
density

$$\gamma_{(n,n)}\left(\frac{\alpha}{2}\right) \qquad \gamma_{(n,n)}\left(1 - \frac{\alpha}{2}\right)$$

and these quantiles don't depend on
any unknowns!

Eg. In R: $\gamma_{(n,n)}\left(\frac{\alpha}{2}\right)$ is found w/ the code

"qgamma$\left(\frac{\alpha}{2}, shape = n, rate = n, lower.tail = T\right)$"

So we have, by definition of quantiles

$$\Pr\left(\gamma_{(n,n)}(\alpha/2) \le \tau\bar{X} \le \gamma_{(n,n)}(1-\alpha/2)\right)$$

$$= \Pr\left(\frac{\gamma_{(n,n)}(\alpha/2)}{\bar{X}} \le \tau \le \frac{\gamma_{(n,n)}(1-\alpha/2)}{\bar{X}}\right)$$

$$= 1-\alpha$$

Hence

$$\left[\frac{\gamma_{(n,n)}(\alpha/2)}{\bar{X}} \;,\; \frac{\gamma_{(n,n)}(1-\alpha/2)}{\bar{X}}\right]$$

Note:
If we invert this, we get the same answer as before (w/ the scale parameterization)

is a $(1-\alpha)100\%$ CI for $\tau$.

$X_1, \ldots, X_n \overset{IID}{\sim} \text{Exp}(\tau)$

$\text{lik}(\theta) = \prod\limits_{i=1}^{n} f(x_i; \tau) = \prod\limits_{i=1}^{n} \left[ \tau e^{-\tau x_i} \mathbb{I}\{x_i \geq 0\} = \tau^n e^{-\tau \sum\limits_{i=1}^{n} x_i} \mathbb{I}\{x_{(1)} \geq 0\}\right.$

$\hat{\Theta}_{MLE} = \overline{X}$ ⟵

To do: Use the CLT to find an approx. $(1-\alpha)100\%$ CI for $\tau$.

**Note: This version is consistent w/ the parameterization in your textbook**

CLT:

$$\frac{\frac{1}{n}\sum\limits_{i=1}^{n} X_i - E[X_1]}{\sqrt{\dfrac{\text{Var}(X_1)}{n}}} \overset{n \to \infty}{\Longrightarrow} N(0,1) \quad \text{for } \overset{IID}{\sim} \text{ sample } X_1, \ldots, X_n$$
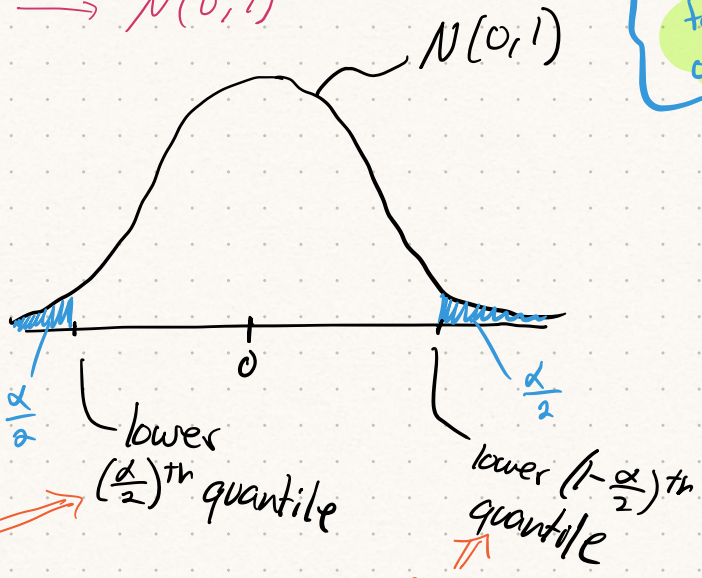
$E[X_1] = \dfrac{1}{\tau}, \quad \text{Var}(X_1) = \dfrac{1}{\tau^2}$

Thus we have

$$\frac{\hat{\tau}_{MLE} - \frac{1}{\tau}}{\left(\frac{1}{\tau^2 n}\right)^{1/2}} \overset{n \to \infty}{\longrightarrow} N(0,1)$$

is a pivot!

**All changes are highlighted in green.**

**View the notes for 9-28-22 to see the other version**

$N(0,1)$

lower $\left(\frac{\alpha}{2}\right)^{th}$ quantile

lower $\left(1-\frac{\alpha}{2}\right)^{th}$ quantile

$\frac{\alpha}{2}$       $\frac{\alpha}{2}$

Notation:

$\mathcal{Z}\left(\frac{\alpha}{2}\right)$        $\mathcal{Z}\left(1-\frac{\alpha}{2}\right)$

By definition of quantile:

$$\Pr\left(\mathcal{Z}\left(\tfrac{\alpha}{2}\right) \leq \frac{\hat{\tau}_{MLE} - \frac{1}{\tau}}{(\tau^2 n)^{-1/2}} \leq \mathcal{Z}\left(1 - \tfrac{\alpha}{2}\right)\right)$$

$$= \Pr\left(\mathcal{Z}\left(\tfrac{\alpha}{2}\right) \leq \tau\sqrt{n}\left(\bar{X} - \tfrac{1}{\tau}\right) \leq \mathcal{Z}\left(1 - \tfrac{\alpha}{2}\right)\right)$$

$$= \Pr\left(\mathcal{Z}\left(\tfrac{\alpha}{2}\right) \leq \tau\sqrt{n}\,\bar{X} - \sqrt{n} \leq \mathcal{Z}\left(1 - \tfrac{\alpha}{2}\right)\right)$$

$$= \Pr\left(\mathcal{Z}\left(\tfrac{\alpha}{2}\right) + \sqrt{n} \leq \sqrt{n}\,\tau\,\bar{X} \leq \mathcal{Z}\left(1 - \tfrac{\alpha}{2}\right) + \sqrt{n}\right)$$

$$= \Pr\left(\frac{\mathcal{Z}\left(\tfrac{\alpha}{2}\right) + \sqrt{n}}{\bar{X}\sqrt{n}} \leq \tau \leq \frac{\mathcal{Z}\left(1 - \tfrac{\alpha}{2}\right) + \sqrt{n}}{\bar{X}\sqrt{n}}\right)$$

$$= \Pr\left(\frac{\frac{\mathcal{Z}\left(\tfrac{\alpha}{2}\right)}{\sqrt{n}} + 1}{\bar{X}} \leq \tau \leq \frac{\frac{\mathcal{Z}\left(1 - \tfrac{\alpha}{2}\right)}{\sqrt{n}} + 1}{\bar{X}}\right)$$

$$= 1 - \alpha$$

Hence
$$\left[\frac{\mathcal{Z}\left(\tfrac{\alpha}{2}\right)/\sqrt{n} + 1}{\bar{X}}, \quad \frac{\mathcal{Z}\left(1 - \tfrac{\alpha}{2}\right)/\sqrt{n} + 1}{\bar{X}}\right]$$

is a $(1-\alpha)100\%$ approx. CI for $\tau$.

# Bayesian
## Credible Intervals

— quantify our personal feelings of uncertainty about the value of a parameter that generated the observed data based on an assumed model

↱ direct assessment of uncertainty

If $X_1, \ldots, X_n$ are IID

$$(X_1, \ldots, X_n) \sim \prod_{i=1}^{n} f(X_i ; \Theta)$$

$$\Theta \sim \pi(\Theta) \leftarrow$$

both parts form the assumed model

parameter described as a RV

The observed data $(x_1, \ldots, x_n)$ are realized values from the joint distb'n $\prod_{i=1}^{n} f(x_i ; \Theta_0)$.

fixed, unknown value of $\Theta$ that "produced" the observed data

The goal of Bayesian inference is to use the data to describe plausible values for $\Theta_0$ though a posterior distb'n

$$\pi(\Theta | x_1, \ldots, x_n)$$

Data is fixed, NOT random

A credible interval for $\Theta$ is a <u>random interval</u>, <u>always</u>.

↱ random b/c it is a function of a RV w/ density $\pi(\Theta | x_1, \ldots, x_n)$

## Process:

Use the posterior distb'n of $\Theta$ (given the observed data) to identify a lower bound (LB) and upper bound (UB) on the most plausible values for $\Theta_0$.

We choose LB and UB based directly upon quantiles of the posterior.

## Interpretation:

We say a $W\%$ credible interval $[LB, UB]$, contains $\Theta_0$ w/ probability $W$.

Although this is easier to interpret than a confidence interval, what's harder to communicate is the rationale behind the posterior distb'n.

# Example derivation of a Bayesian credible interval

## HW 10 #2

100 items randomly sampled ⎫ Data
3 defects found ⎭

To do: use Beta prior to derive posterior dist'n for $\theta$ and then find a credible interval for $\theta$.

$\theta_o$ = proportion of total defective items in the population

---

$\text{lik}(\theta) = \binom{100}{3} \theta^3 (1-\theta)^{100-3}$  if we let $X = \begin{cases} 0, \text{ not defective} \\ 1, \text{ defective} \end{cases}$

where $X \sim \text{Bern}(\theta)$.

Given
$\pi(\theta) \sim \text{Beta}(a,b)$ means $\pi(\theta) = \dfrac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$

is the probability dist'n we are going to use to express our uncertainty about $\theta_o$.

---

> Detour for Notes on $\Gamma(\cdot)$ function
>
> For positive integer $a$: $\Gamma(a) = (a-1)!$
> $$\Gamma(a+1) = a\,\Gamma(a)$$
>
> For any $a$ besides negative integers or zero: $\Gamma(a) = \dfrac{\Gamma(a+1)}{a}$
>
> In general, $\Gamma(a) = \int_0^\infty t^{a-1} e^{-t}\,dt.$

---

With $\text{lik}(\theta)$ and $\pi(\theta)$ we can now find the posterior density:

$$\pi(\theta \mid x_{obs}) = \frac{\text{lik}(\theta)\,\pi(\theta)}{\int_{(A)} \text{lik}(\theta)\,\pi(\theta)\,d\theta}$$

$y =$ # of successes out of 100 trials

$$\pi(\theta/y=3) = \frac{\binom{100}{3}\theta^3(1-\theta)^{100-3} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}\theta^{a-1}(1-\theta)^{b-1}}{\int_0^1 \binom{100}{3}\theta^3(1-\theta)^{100-3} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}\theta^{a-1}(1-\theta)^{b-1}d\theta}$$
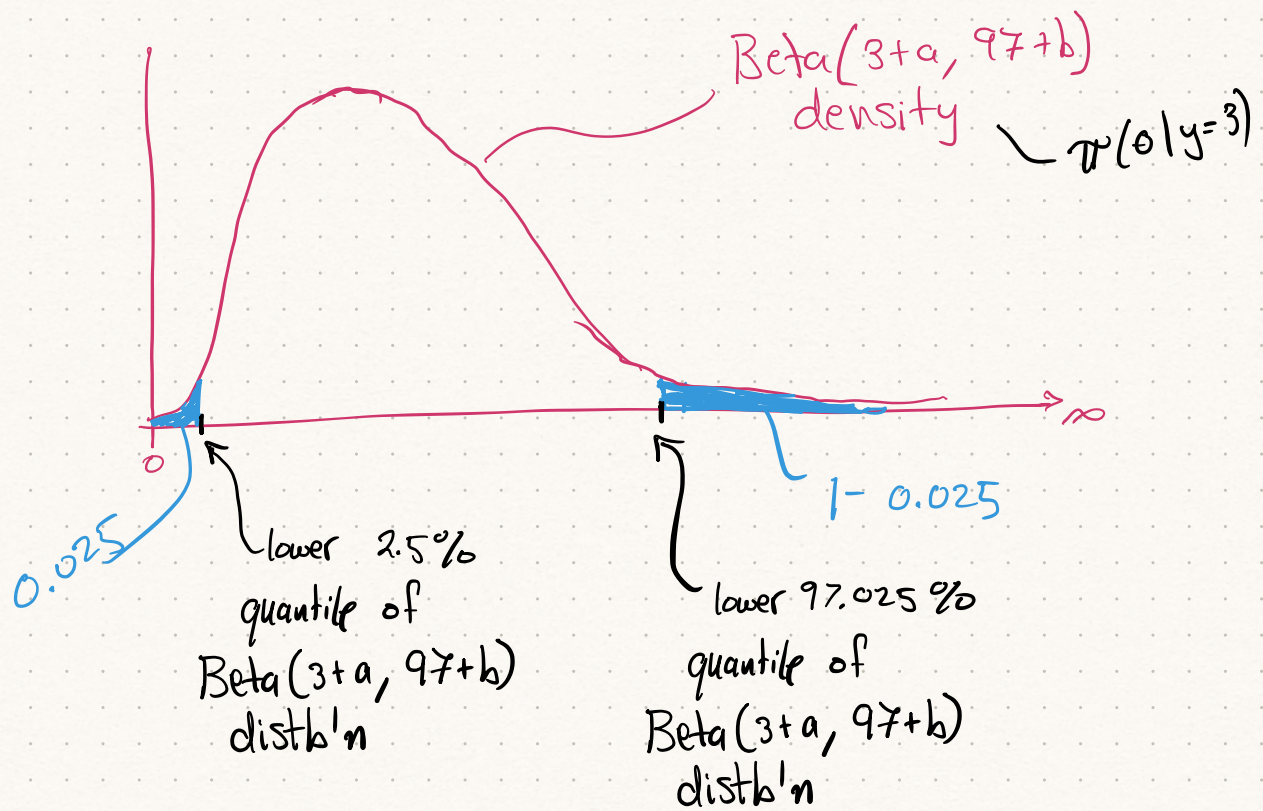
$$= \frac{\theta^{3+a-1}(1-\theta)^{100-3+b-1}}{\int_0^1 \theta^{3+a-1}(1-\theta)^{100-3+b-1}d\theta}$$

looks like $\text{Beta}(3+a, 97+b)$

$$\pi(\theta/y=3) \sim \text{Beta}(3+a, 97+b)$$

So $\theta/x=3 \sim \text{Beta}(3+a, 97+b)$ is the posterior distribution for $\theta$, given the observed data.

For given values of $a$ and $b$, we can find any quantiles we may want!

Beta$(3+a, 97+b)$
density

$\pi(\theta|y=3)$

0.025

lower 2.5%
quantile of
Beta$(3+a, 97+b)$
distb'n

$1 - 0.025$

lower 97.025%
quantile of
Beta$(3+a, 97+b)$
distb'n

In R: qbeta(0.025, shape1 $= 3+a$, shape2 $= 97+b$, lower.tail $= T$)

What we're doing is using the shape of the
posterior density to find an interval that
describes the most typical values for $\theta_o$.
Such credible intervals may also be called
==highest posterior density regions== (hpd for short).

For $W = 95\%$, say,
If $a = b = 1$ then a 95% credible interval for $\theta_o$ is
$$[0.013, 0.842], \text{ but}$$
if $a = 0.5, b = 5$ then a 95% credible interval for $\theta_o$ is
$$[0.0001, 0.4096].$$

# Group Work:

Create a mind-map relating as many theorems from ch. 8 as you can.

- MLE is consistent
- Identity for Fisher Info
- Asymptotic normality of MLE
- Cramér-Rao lower bound
- Factorization thm for sufficient stats
- MLE is a function of sufficient stat
- Rao-Blackwell Theorem for estimation w/ sufficient statistics