$\hat{\theta}_{MM}$ and $\hat{\theta}_{MLE}$ have different MSE's.

Is there a fundamental

limit of estimation?

A minimum MSE?

In general, no; the class of possible estimators is too large!

However, if you limit

it to only the set of unbiased

estimators, there is a

bound called the

Cramer-Rao Lower Bound

(CRLB) discovered in 1945-46.

Unbiased

Estimators that achieve the CRLB are called a

"uniformly minimum variance unbiased estimator" (UMVUE).

In order to prove the CRLB, we begin with the Covariance Inequality: For all rv's $A, B$,

$$\left( \text{Cov}(A,B) \right) = \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad$$

$$\implies \text{Var}(A) \geq \frac{\text{Cov}(A,B)^2}{\text{Var}(B)} = \frac{\left( E[AB] - E[A]E[B] \right)^2}{E[B^2] - E[B]}$$

(from MATH 390)

i.e. a lower bound on the variance of any rv.

_____

Consider DGP: $X_1, \dots X_n \overset{iid}{\sim} p(x;\theta)$ or $f(x;\theta)$

and any unbiased estimator $\hat{\theta}$.

Let $A = \hat{\theta}$, $B = S$, the "score function" for parameter $\theta$.

Let $\alpha \neq \hat{\theta}$. Add Define the "score function"  $S := \frac{\partial}{\partial \theta}\left[\ln f(x_1 \ldots x_i ; \theta)\right]$  (def 1)

$\Rightarrow E[\theta] \neq \theta$

assume 1 unbiased

A property of the DGP, not any estimator.

by chain rule

$= \frac{\frac{\partial}{\partial \theta}\left[f(x_1 \ldots x_i ; \theta)\right]}{f(x_1 \ldots x_i ; \theta)}$  (def 2)

by iid

$= \frac{\partial}{\partial \theta}\left[\ln\left(\prod_{i=1}^{n} f(x_i ; \theta)\right)\right]$  (def 3)

precalc

Identity of derivative operator

$= \frac{\partial}{\partial \theta}\left[\sum_{i=1}^{n} \ln f(x_i ; \theta)\right]$  (def 4)

From def 1,

Since $\mathcal{L} = f(x_1, \ldots, x_i ; \theta)$, $l := \ln(\mathcal{L})$

$= \sum_{i=1}^{n} \frac{\partial}{\partial \theta}\left[\ln\left(f(x_i ; \theta)\right)\right]$  (def 5)

$= \frac{\partial}{\partial \theta}\left[l(\theta ; x_1, \ldots x_n)\right]$  (def 6)

def of $l'$

$= l'(\theta, x_1, \ldots x_n)$  (def 7)

by def 7 and linearity of deriv. operator

$= \frac{\partial}{\partial \theta}\sum l(\theta, x_i) = \sum_{i=1}^{n} l'(\theta ; x_i)$  (def 8)

All X's capital letters since $S$ is a r.v.

We need to find $E[\hat{\theta}S]$, $E(S^2)$, $E(S)$ to prove the CRLB formula.

Since $Cov(\hat{\theta}, S) := E(\hat{\theta}S) - E(\hat{\theta})E(S)$ and $Var(S) = E(S^2) - E(S)^2$

we start with $E[S]$, then $E(S^2)$ then $E(\hat{\theta}S)$ and the substitute to find:

$$Var(\hat{\theta}) \geq \frac{E(\hat{\theta}S) - E(\hat{\theta})E(S)}{E(S^2) - E(S)^2}$$

$$E[S] = E\left[\frac{\frac{\partial}{\partial\theta}(f(X_1,\ldots,X_n;\theta))}{f(X_1,\ldots,X_n;\theta)}\right] = \int_{S_{X_1}}\cdots\int_{S_{X_n}} \frac{\frac{\partial}{\partial\theta}[f(x_1,\ldots,x_n;\theta)]}{f(x_1,\ldots,x_n;\theta)} f(x_1,\ldots,x_n;\theta)\, dx_1\ldots dx_n$$

Assume derivative and integral can be interchanged (Assum 1).        by definition of JDF/JMF

$$\stackrel{\downarrow}{=} \frac{\partial}{\partial\theta}\left[\int_{S_{X_1}}\cdots\int_{S_{X_n}} f(x_1,\ldots,x_n;\theta)\,dx_1\ldots dx_n\right] = \frac{\partial}{\partial\theta}(1) = 0 \quad \text{(Fact 1a)}$$

$$\boxed{\Rightarrow Var(\hat\theta) \ge \frac{E[\partial S]}{E[S^2]} \quad \text{here on our way!}}$$

def ℓ

by linearity of expectation
by iid

$$\Rightarrow E[S] = E[\ell'(\theta; X_1,\ldots,X_n)] = 0 \quad \text{and} \quad E[S] = E\left[\sum \ell'(\theta;x)\right] = n\, E[\ell'(\theta,x)] = 0$$

$$\Rightarrow E[\ell'(\theta;x)] = 0$$

~~[scribbled out text]~~

def θ

$$(q_1 + \ldots + q_n)^2 = \sum q_i^2 + \sum_{i\ne j} q_i q_j \quad \text{(Fact 1b)}$$

$$Var[S] = E[S^2] - \underbrace{E[S]^2}_{0} = E\left[\left(\sum_{i=1}^n \ell'(\theta;x_i)\right)^2\right] =$$

linearity of expectation

$$E\left[\sum_{i=1}^n \ell'(\theta;x_i)^2 + \sum_{i\ne j}\ell'(\theta;x_i)\ell'(\theta;x_j)\right] = \sum_{i=1}^n E[\ell'(\theta;x_i)^2] + \sum_{i\ne j} E[\ell'(\theta;x_i)\ell'(\theta;x_j)]$$

by iid

Recall $E[UV] = E[U]E[V]$ if $U,V$ indep

$$\mathcal{I}(\theta) = E[S^2]$$

$$= n\, E[\ell'(\theta;x_i)^2] + \left(\sum_{i\ne j}\right) \underbrace{E[\ell'(\theta;x_i)]}_{0} \underbrace{E[\ell'(\theta;x_j)]}_{0} = n\, E[\ell'(\theta;x_i)^2]$$

(by fact 1b)

$$\underbrace{\qquad}_{\mathcal{I}_n(\theta)}$$

$$\Rightarrow Var(\hat\theta) \ge \frac{E[\partial S]}{n\,\mathcal{I}(\theta)}$$

Sometimes Fisher Information is defined this way

(N.O.W. ~~to solve~~ for $E[\hat{\theta}S] = E\left[\hat{\theta} \dfrac{\frac{\partial}{\partial\theta}[f(x_1,\ldots x_n;\theta)]}{f(x_1,\ldots x_n;\theta)}\right]$

def of expectation of vector rv
$\downarrow$

$= \displaystyle\int_{S_X}\cdots\int_{S_X} \hat{\theta}\, \dfrac{\frac{\partial}{\partial\theta}\left[f(x_1,\ldots x_n;\theta)\right]}{f(x_1,\ldots x_n;\theta)}\, f(x_1,\ldots x_n;\theta)\, dx_1\ldots dx_n$

Assum 1

def of expectation vector rv
$\downarrow$

$\overset{\downarrow}{=} \dfrac{\partial}{\partial\theta}\left[\displaystyle\int_{S_X}\cdots\int_{S_X} \hat{\theta}\, f(x_1,\ldots x_n;\theta)\, dx_1\ldots dx_n\right] = \dfrac{\partial}{\partial\theta}\left[E[\hat{\theta}]\right] \overset{\downarrow}{=} \dfrac{\partial}{\partial\theta}(\theta) = 1$

assumpiont unbiasedness of $\hat{\theta}$

we're done

consider starting at some $n^{-1}$!

$\Rightarrow Var[\hat{\theta}] \geq \dfrac{1}{n\, I(\theta)} = \dfrac{I(\theta)^{-1}}{n}$   So Fisher Information is super important! We will let it sink in for now by using it and then revisit it conceptually later...

It's a metric about the distribution itself that tells you how ... $\theta$ ... degree $\theta$ is ... able to be estimated.

---

Let's prove some estimators are UMVUE's! First, let's get a more convenient expression for Fisher Information:

Assumption 1 AND lots of
work for HW

$I(\theta) := E\left[\ell'(\theta;x)^2\right] = \cdots = E\left[-\ell''(\theta;x)\right]$

Obp: $X_1,\ldots X_n \overset{iid}{\sim} Bern(\theta)$, $\hat{\theta}=\bar{X}$, $Var[\bar{X}] = \dfrac{\theta(1-\theta)}{n}$. Is this the optimal UMVUE? Let's compute the CRLB.

$f(\theta;x) = p(x;\theta) = \theta^x(1-\theta)^{1-x}$

$\ell(\theta;x) = x\ln(\theta) + (1-x)\ln(1-\theta)$        $-\ell''(\theta;x) = \dfrac{x}{\theta^2} + \dfrac{1-x}{(1-\theta)^2}$

$\ell'(\theta;x) = \dfrac{x}{\theta} + \dfrac{1-x}{1-\theta}$

$\ell''(\theta;x) = -\dfrac{x}{\theta^2} - \dfrac{1-x}{(1-\theta)^2}$    $I(\theta) = E[-\ell''] = \dfrac{E(x)}{\theta^2} + \dfrac{E(1-x)}{(1-\theta)^2} = \dfrac{\theta}{\theta^2} + \dfrac{1-\theta}{(1-\theta)^2} = \dfrac{1}{\theta} + \dfrac{1}{1-\theta}$

$= \dfrac{1-\theta}{\theta(1-\theta)} + \dfrac{\theta}{\theta(1-\theta)} = \dfrac{1}{\theta(1-\theta)} \Rightarrow I(\theta)^{-1} = \theta(1-\theta)$

$$CRLB = \frac{I(\theta)^{-1}}{n} = \underline{\qquad} \quad \frac{\theta(1-\theta)}{n} \implies \hat{\theta} = \overline{X} \text{ is the CMVUE!}$$

Hooray!

Remember, there's no way to find the $\hat{\theta}$ with variance $=$ CRLB. But if we find $\hat{\theta}$ with variance $=$ CRLB, we find the best one!!

$X_1, \ldots, X_n \overset{iid}{\sim} N(\theta_1, \theta_2)$, $\hat{\theta}_1 = \overline{X}$. $Var[\overline{X}] = \frac{\theta_2}{n}$. Is it optimal?

$$\mathcal{L}(\theta; X) = f(X; \theta) = \frac{1}{\sqrt{2\pi \theta_2}} e^{-\frac{1}{2\theta_2}(X-\theta_1)^2} \quad -\frac{X^2}{2\theta_2} + \frac{X\theta_1}{\theta_2} - \frac{\theta_1^2}{2\theta_2}$$

$$\ell(\theta; X) = -\frac{1}{2}\ln(2\pi) - \frac{1}{2}\ln(\theta_2) - \frac{1}{2\theta_2}(X-\theta_1)^2$$

$$\ell'(\theta; X) = \frac{X}{\theta_2} - \frac{\theta_1}{\theta_2}$$

$$\ell''(\theta; X) = -\frac{1}{\theta_2}$$

$$-\ell''(\theta; X) = \frac{1}{\theta_2}$$

no $X$!

$$I(\theta) = E[\quad] = \frac{1}{\theta_2} \implies I(\theta)^{-1} = \theta_2$$

$$CRLB = \frac{I(\theta)^{-1}}{n} = \underline{\qquad} \quad \frac{\theta_2}{n} \implies \hat{\theta}_1 = \overline{X} \text{ is a CMVUE!}$$

Hooray!

Here: $\theta_2, n \overset{iid}{\sim} N(\theta_1, \theta_2)$

Let's do a demo now to try and understand $\mathcal{L}, \ell, \ell', -\ell''$, $I(\theta)$ and how they all related.

Inverse Fisher Information $I(\theta)^{-1}$ of a DGP measures a fundamental limit on the difficulty of estimating $\theta$. If $I(\theta)^{-1}$ is large there's not a lot of "information" in $X$ about $\theta$. $\iff$ If $I(\theta)$ large $\implies$ there's a lot of information in $X$ about $\theta$.

The Last of the 3 goals of statistical inference: confidence sets

Point estimation focused on best guess of $\theta$ e.g. $\hat{\theta} = .176$

Confidence sets focuses on a range of possible $\theta$'s e.g. $[.175, .977]$ or $[.42, .53]$. The confidence sets answers the question "how sure are you of this pt. estimate $\hat{\theta}$? If set has a tight bound $\Rightarrow$ we're sure of $\hat{\theta}$. If wide bands $\Rightarrow$ we're not sure of $\hat{\theta}$.

Define: an "interval estimate" is

$$\left[ W_L(x_1, \ldots x_n), W_U(x_1 \ldots x_n) \right] = [\hat{\theta}_L, \hat{\theta}_U]$$

Where $W_L, W_U$ are two statistical functions s.t. $W_L < W_U$ for all possible datasets.

An "interval estimate" is $\left[ W_L(x_1, \ldots x_n), W_U(x_1, \ldots x_n) \right]$ which is a random interval $= [\hat{\theta}_L, \hat{\theta}_U]$
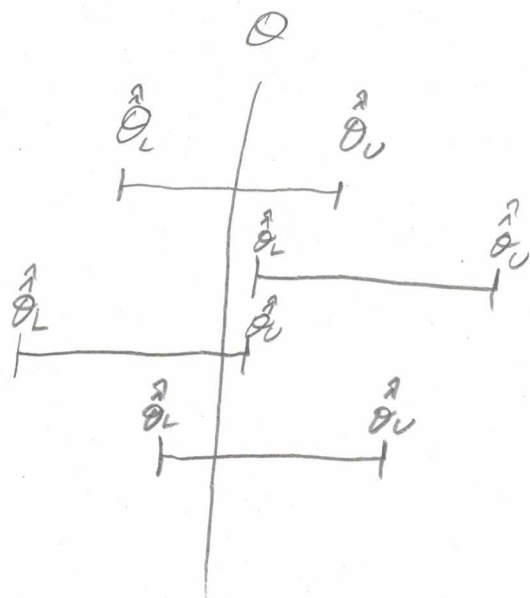
Define: The "coverage probability" of an interval estimate is

$$P\left( \theta \in [\hat{\theta}_L, \hat{\theta}_U] \mid \theta \right)$$

why given $\theta$? Because you need to know the true distr's of $\hat{\theta}_L, \hat{\theta}_U$ to compute coverage explicitly.

Coverage Prob. is best illustrated as follows:



Draws #1

Draws #2

Draws #3

Draws #4

⋮

The coverage prob. is computed over every possible draws.
If this were every draws, cov. prob = 75%.

Def: An confidence band estimator with cov. prob. $1-\alpha$ for param $\theta$

$$\hat{CI}_{\theta, 1-\alpha} := \left[\hat{\theta}_{\frac{\alpha}{2}}, \hat{\theta}_{1-\frac{\alpha}{2}}\right]$$

A two-sided confidence interval estimate ("just 'Confidence Interval'")
corresponding to the above confidence interval estimator is

$$\hat{CI}_{\theta, 1-\alpha} := \left[\hat{\theta}_{\frac{\alpha}{2}}, \hat{\theta}_{1-\frac{\alpha}{2}}\right]$$