

# Consistency and Asymptotic Distribution of MLE

Preliminary fact:

(from derivation of alternate  
formula for  $I(\theta)$ )

$$\int_{-\infty}^{\infty} \left[ \frac{\partial}{\partial \theta} (\log f(x|\theta)) \right] f(x|\theta) dx = 0$$

equivalently

$$E \left[ \frac{\partial}{\partial \theta} \log f(x|\theta) \right] = 0$$

when  $X$  has density  $f(x|\theta)$ .

## Heuristic argument for consistency of MLE

Notation:

$X_1, X_2, \dots, X_n$  IID  $f(x|\theta_0)$



true value.

$f(x|\theta)$  is regular one-parameter family of pdf's (or pmf's).

$\hat{\theta}$  denotes MLE of  $\theta$ .

Consistent:  $\hat{\theta} \approx \theta_0$  for large  $n$

$(\hat{\theta} \rightarrow \theta_0 \text{ as } n \rightarrow \infty)$

$\hat{\theta}$  is the value  $\theta$  maximizing

$$\text{log-likelihood } l(\theta) = \sum_{i=1}^n \log f(X_i|\theta)$$

or equivalently

$$\frac{1}{n} l(\theta) = \frac{1}{n} \sum_{i=1}^n \log f(X_i|\theta)$$

This is a sample average (of something).

For large  $n$ , will be close to population average. (by LLN)

LLN: If  $V_1, V_2, \dots, V_n$  IID,  
then  $\bar{V} \approx EV_i$  for large  $n$   
 $(\bar{V} \rightarrow EV_i \text{ as } n \rightarrow \infty)$

Take  $V_i = \log f(x_i | \theta)$  in LLN. (fix  $\theta$ )  
arbitrary

For large  $n$  expect

$$\begin{aligned}\frac{1}{n} l(\theta) &\approx E \log f(x_i | \theta) \\ &= \int (\log f(x | \theta)) f(x | \theta_0) dx \\ &\quad \text{(assume density case.)} \\ &= \Psi(\theta), \text{ (definition)}\end{aligned}$$

If  $\frac{1}{n} l(\theta) \approx \Psi(\theta)$  for "all"  $\theta$ ,

then  $\hat{\theta}$  (the value maximizing  $\frac{1}{n} l(\theta)$ )  
should be close to the value of  $\theta$  which  
maximizes  $\Psi(\theta)$ .

But  $\Psi(\theta)$  is maximized when  $\theta = \theta_0$ :

$$\frac{\partial \Psi(\theta)}{\partial \theta} = \int \frac{\partial}{\partial \theta} (\log f(x | \theta)) f(x | \theta_0) dx = 0 \text{ when } \theta = \theta_0.$$

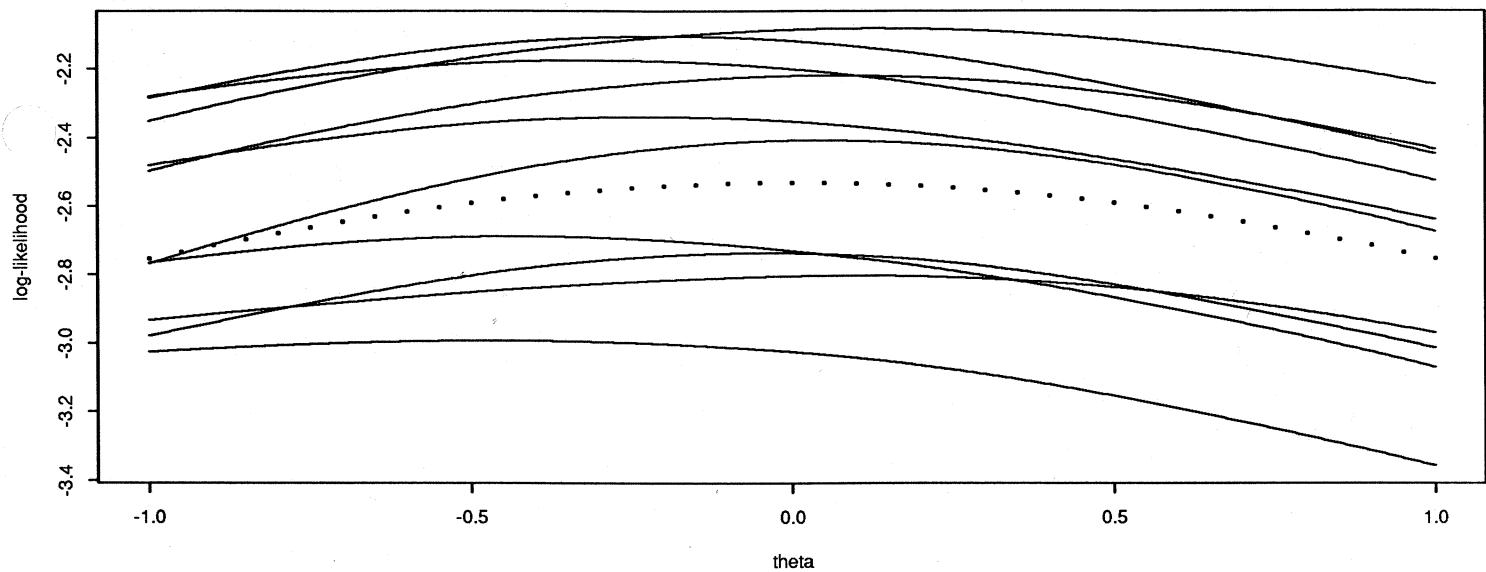
(at least  $\theta_0$  is a stationary point.)

Thus expect

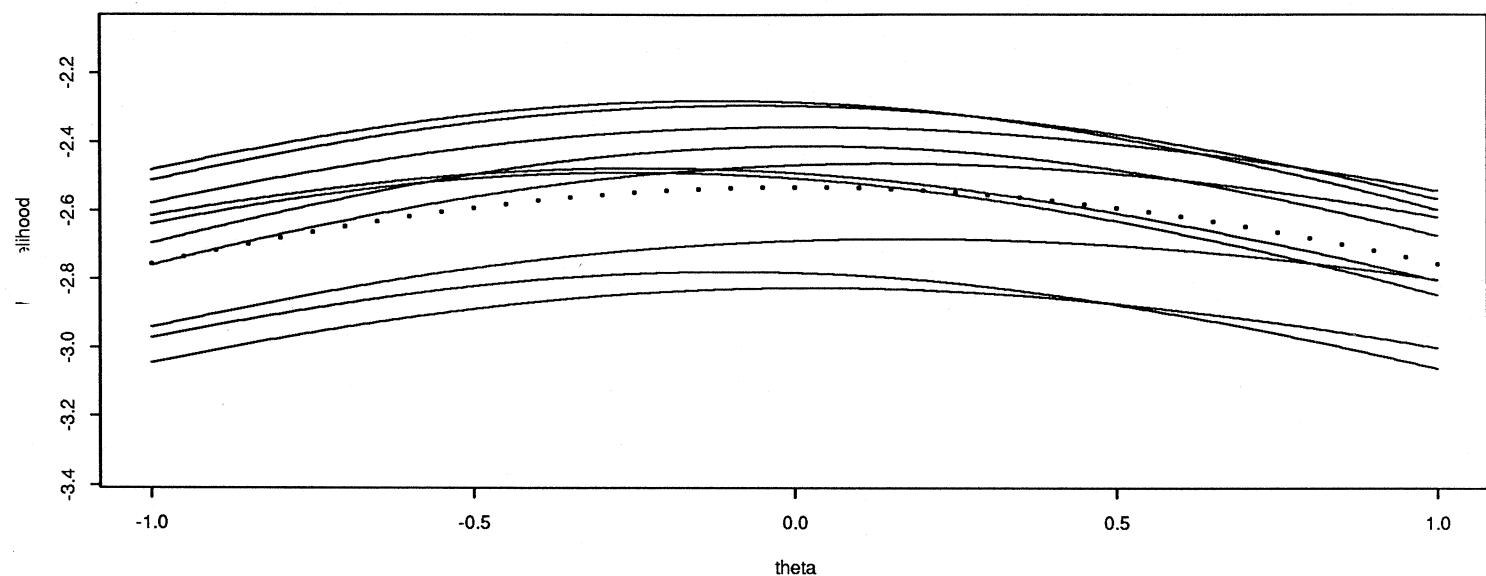
$$\hat{\theta} \approx \theta_0 \text{ for large } n. \quad \text{QED}$$

$\hat{\theta}$  is consistent estimator.

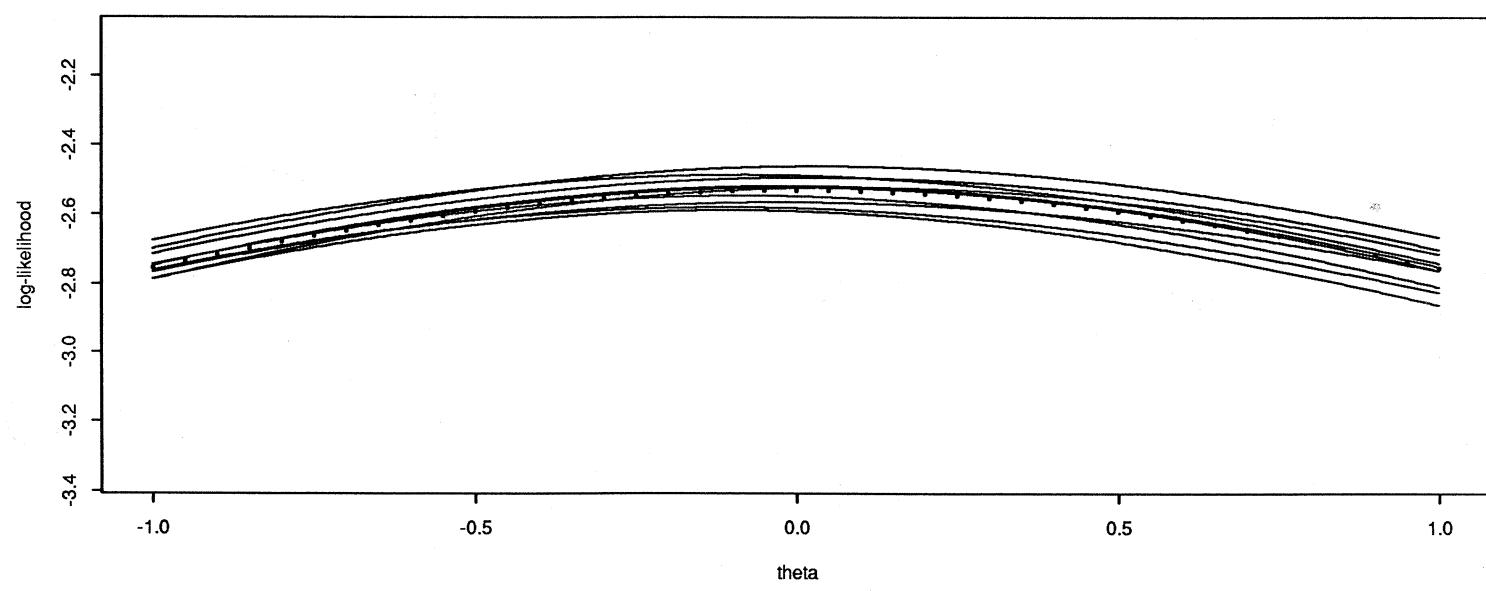
Log-Likelihood Functions for 10 Samples of Size 20



Log-Likelihood Functions for 10 Samples of Size 100



Log-Likelihood Functions for 10 Samples of Size 500



## Asymptotic distn. of MLE

follows (by an elaborate argument)

from the CLT applied to

$$\underbrace{\frac{1}{n} \dot{l}(\theta_0)}_{\text{This is } \bar{V}.} = \frac{1}{n} \sum_{i=1}^n \left. \frac{\partial}{\partial \theta} \log f(x_i | \theta) \right|_{\theta=\theta_0}$$

$\underbrace{\quad}_{\text{Call this } V_i.}$

CLT says (for large n)

$$\bar{V} \sim \text{approx } N(EV_i, \frac{\text{Var } V_i}{n}).$$

Since

$$EV_i = E \left. \frac{\partial}{\partial \theta} \log f(x_i | \theta) \right|_{\theta=\theta_0} = 0,$$

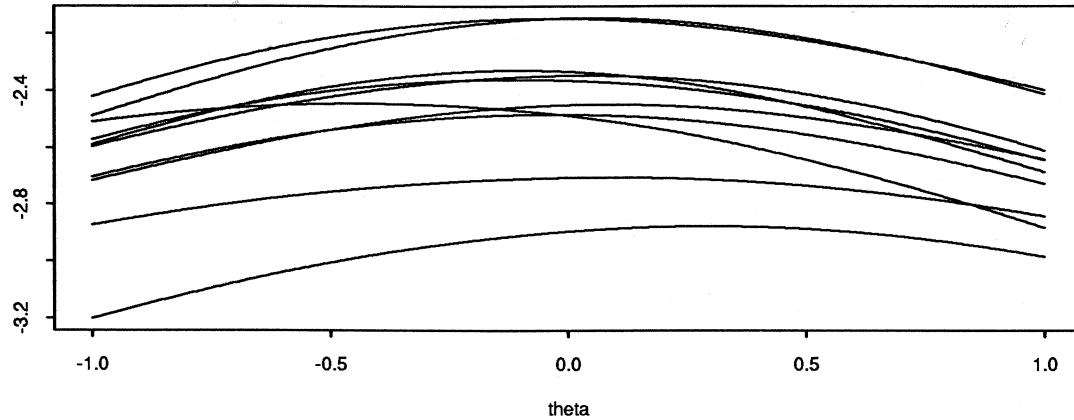
$$\text{and } \text{Var } V_i = \text{Var} \left( \left. \frac{\partial}{\partial \theta} \log f(x_i | \theta) \right|_{\theta=\theta_0} \right) = I(\theta_0),$$

we have  $\frac{1}{n} \dot{l}(\theta_0) \sim \text{approx } N(0, \frac{1}{n} I(\theta_0)).$

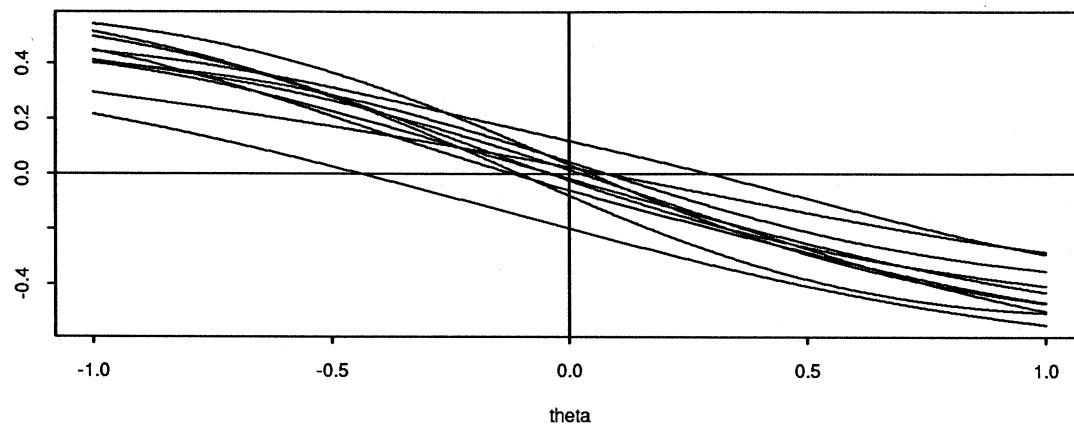
(There is much more.)

## 10 Samples of Size 50 from the Cauchy Distrn.

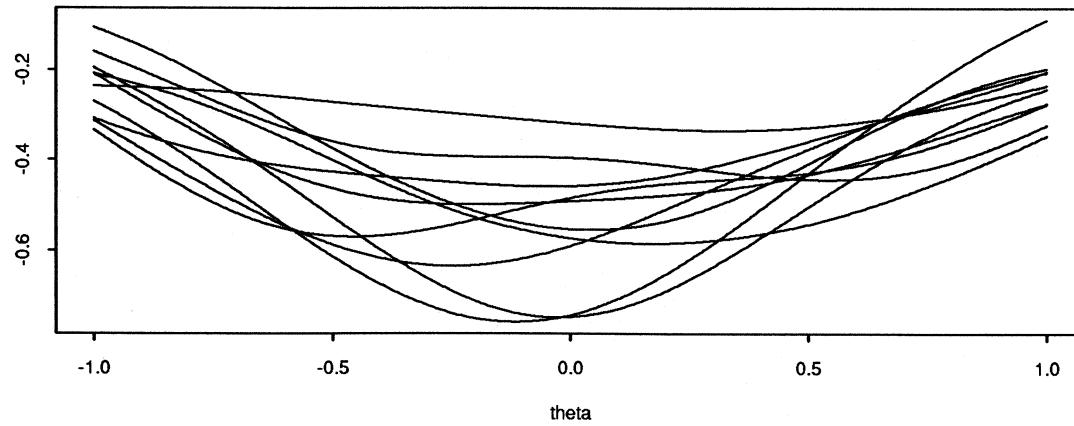
Log-likelihood



First Derivative of Log-likelihood



Second Derivative of Log-likelihood



Heuristic argument  
for asymptotic distn. of MLE

$X_1, X_2, \dots, X_n$  IID  $f(x|\theta_0)$ .

Notation  
↑  
true value

$\hat{\theta}$  denotes MLE.

$l(\theta) = \text{log-likelihood function.}$

MLE  $\hat{\theta}$  is the solution of  $\frac{\partial l}{\partial \theta} = 0$ .

Define  $s(\theta) = \frac{\partial l(\theta)}{\partial \theta}$  ( $= \dot{l}(\theta)$  earlier)

Then  $\hat{\theta}$  is solution of  $s(\theta) = 0$ .

For  $\theta$  near  $\theta_0$

$$s(\theta) \approx s(\theta_0) + s'(\theta_0)(\theta - \theta_0).$$

If  $\hat{\theta}$  is close to  $\theta_0$  (we can show  $\hat{\theta}$  is a consistent estimator of  $\theta$ ),  $\hat{\theta}$  will be  $\approx$  equal to the solution of

$$s(\theta_0) + s'(\theta_0)(\theta - \theta_0) = 0$$

which is

$$\begin{cases} \theta - \theta_0 = -\frac{s(\theta_0)}{s'(\theta_0)} \\ \rightarrow \theta = \theta_0 - \frac{s(\theta_0)}{s'(\theta_0)} \end{cases}$$

$$\text{Thus } \hat{\theta} \approx \theta_0 - \frac{s(\theta_0)}{s'(\theta_0)}.$$

$$\begin{aligned} s(\theta) &= \frac{\partial l(\theta_0)}{\partial \theta} = \frac{\partial}{\partial \theta} \sum_{i=1}^n \log f(x_i | \theta) \Big|_{\theta=\theta_0} \\ &= \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(x_i | \theta) \end{aligned}$$

$$s'(\theta_0) = \frac{\partial s(\theta_0)}{\partial \theta} = \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f(x_i | \theta_0).$$

$$\frac{s(\theta_0)}{s'(\theta_0)} = \frac{\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(x_i | \theta_0)}{\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f(x_i | \theta_0)} \rightarrow \begin{array}{l} \text{Apply CLT} \\ \text{CLT} \end{array} \quad \begin{array}{l} \text{Apply LLN} \\ \text{LLN} \end{array} \quad \begin{array}{l} \text{CLT} \\ \text{LLN} \end{array}$$

Define  $\mathcal{N}$  and  $\mathcal{D}$  to be the numerator and denominator of this fraction.

Suppose  $n$  is large.

By the LLN,

$$\mathcal{D} \equiv \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f(X_i | \theta_0) \approx E \left( \frac{\partial^2}{\partial \theta^2} \log f(X_1 | \theta_0) \right) = -I(\theta_0).$$

Thus

$$\hat{\theta} \approx \theta_0 - \frac{s(\theta_0)}{s'(\theta_0)} = \theta_0 - \frac{\mathcal{N}}{\mathcal{D}} \approx \theta_0 + I(\theta_0)^{-1} \mathcal{N}$$

Note that  $\mathcal{N} = n^{-1}s(\theta_0) = n^{-1}\dot{\ell}(\theta_0)$ .

As shown earlier using the CLT,

$$\mathcal{N} \sim \text{approx } N(0, n^{-1}I(\theta_0)).$$

Thus we conclude

$$\hat{\theta} \sim \text{approx } N\left(\theta_0, \frac{1}{nI(\theta_0)}\right).$$

Here we have used the fact:

$$X \sim N(\mu, \sigma^2) \Rightarrow aX + b \sim N(a\mu + b, a^2\sigma^2).$$

Fisher Information, CR-bound,  
 Asymp. distn. of MLE's in the  
 Multi-parameter case

Notation:  $\tilde{x} \sim f(\tilde{x}|\theta)$

$$\theta = (\theta_1, \theta_2, \dots, \theta_p)$$

$$\frac{\partial}{\partial \theta} = \begin{pmatrix} \frac{\partial}{\partial \theta_1}, \\ \vdots \\ \frac{\partial}{\partial \theta_p} \end{pmatrix}$$

$$S = \text{vector of scores} = \frac{\partial}{\partial \theta} \log f(\tilde{x}|\theta) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} \log f(\tilde{x}|\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_p} \log f(\tilde{x}|\theta) \end{pmatrix}$$

$p \times 1$  vector

Definition       $\xrightarrow{\text{computed at } \theta} \xrightarrow{\text{evaluated at } \theta} \left. \right\} \text{same } \theta!$

$$I_{\tilde{x}}(\theta) = E(S S^t)$$

$$\underbrace{p \times p}_{\text{matrix}} \quad \underbrace{p \times 1}_{\text{ }} \quad \underbrace{1 \times p}_{\text{ }} \quad \underbrace{p \times p}_{\text{ }}$$

For a vector or matrix, we define expected value in this way:

$$E \begin{pmatrix} Y \\ Z \end{pmatrix} = \begin{pmatrix} EY \\ EZ \end{pmatrix}$$

$$E \begin{pmatrix} W & X \\ Y & Z \end{pmatrix} = \begin{pmatrix} EW & EX \\ EY & EZ \end{pmatrix}$$

Properties:

$$\rightarrow E_{\theta} S = 0$$

$p \times 1 \quad p \times 1$

$$\rightarrow I_{\tilde{X}}(\theta) = \text{Cov}(S) \quad (\text{The variance-covariance matrix of } S)$$

$\rightarrow$  If  $\tilde{X} = (X_1, X_2, \dots, X_n)$  has independent components, then

$$I_{\tilde{X}}(\theta) = I_{X_1}(\theta) + I_{X_2}(\theta) + \dots + I_{X_n}(\theta)$$

$p \times p \quad p \times p \quad p \times p$

→ If  $X_1, X_2, \dots, X_n$  are iid,  
then  $I_{\tilde{X}_n}(\theta) = n I_{X_1}(\theta)$   
 $p \times p$

$$\rightarrow I_{\tilde{X}}(\theta) = E\left(-\frac{\partial^2}{\partial \theta^2} \log f(\tilde{X}|\theta)\right)$$

where we define

$$\frac{\partial^2}{\partial \theta^2} \log f(\tilde{X}|\theta) = \left( \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(\tilde{X}|\theta) \right)$$

= the  $p \times p$  matrix whose  $(i,j)$  entry  
is  $\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(\tilde{X}|\theta)$ .

## Asymptotic distn. of MLE (of $\theta$ )

If  $\hat{\theta}_n = \hat{\theta}_n(x_1, x_2, \dots, x_n)$  is the sequence of MLE's (based on progressively larger samples), then

$$\hat{\theta}_n \sim AN(\theta, (I_{\tilde{X}}(\theta))^{-1})$$

asymptotically  $\searrow$  multivariate normal

which means roughly that

$$\hat{\theta}_n \sim \text{approx } N(\theta, (I_{\tilde{X}}(\theta))^{-1})$$

for large  $n$ .

Recall: In iid case  $I_{\tilde{X}}(\theta) = n I_{X_1}(\theta)$ .

Estimate  $I_{\tilde{X}}(\theta)$  by  $I_{\tilde{X}}(\hat{\theta}_n)$  or

$$\left( -\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(\tilde{X} | \theta) \right) \Big|_{\theta = \hat{\theta}_n}$$

## Multi-parameter CRLB

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^T.$$

- ① If  $E \hat{\boldsymbol{\theta}} = \boldsymbol{\theta}$  (implies  $E \hat{\theta}_i = \theta_i$  for all  $i$ )  
then  $\text{Var}(\hat{\boldsymbol{\theta}}) - (I_{\hat{\boldsymbol{\theta}}}(\boldsymbol{\theta}))^{-1}$  is non-negative definite.

so that  $\text{Var}(\hat{\theta}_i) \geq (I^{-1})_{ii}$

where  $I = I_{\hat{\boldsymbol{\theta}}}(\boldsymbol{\theta})$ .

- ② If  $E \hat{\theta}_i = \theta_i$  for a particular  $i$ ,

then (by fixing  $\theta_j$  for  $j \neq i$  and using  
the CRLB for the one-parameter case)

$$\text{Var}(\hat{\theta}_i) \geq (I_{ii})^{-1} = \frac{1}{E \left( \frac{\partial}{\partial \theta_i} \log f(\mathbf{x} | \boldsymbol{\theta}) \right)^2}$$

- ③ If  $E \hat{\boldsymbol{\theta}} = \boldsymbol{\theta}$ , then

$$\text{Var}(\hat{\theta}_i) \geq (I^{-1})_{ii} \geq (I_{ii})^{-1}$$

Best you can do if ... ↓

All parameters  
are unknown.

↓  
 $\theta_j$  are known  
for  $j \neq i$ .

Example: Normal( $\mu, \sigma^2$ ) dist.

$$f(x|\underbrace{\mu, \sigma^2}_{\Theta}) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\log f = -\frac{1}{2} \log(2\pi\sigma) - \frac{(x-\mu)^2}{2\sigma^2}$$

$$\begin{pmatrix} \frac{\partial}{\partial \mu} \log f \\ \frac{\partial}{\partial \sigma} \log f \end{pmatrix} = \begin{pmatrix} \frac{x-\mu}{\sigma} \\ -\frac{1}{2\sigma} + \frac{(x-\mu)^2}{2\sigma^2} \end{pmatrix}$$

$$\frac{\partial}{\partial \Theta} \log f(x|\Theta)$$

$$\begin{pmatrix} \frac{\partial^2 l}{\partial \mu^2} & \frac{\partial^2 l}{\partial \mu \partial \sigma} \\ \frac{\partial^2 l}{\partial \sigma \partial \mu} & \frac{\partial^2 l}{\partial \sigma^2} \end{pmatrix} = \begin{pmatrix} -\frac{1}{\sigma} & -\frac{(x-\mu)}{\sigma^2} \\ -\frac{(x-\mu)}{\sigma^2} & \frac{1}{2\sigma^2} - \frac{(x-\mu)^2}{2\sigma^3} \end{pmatrix}$$

$$I(\Theta) = -E\left(\frac{\partial}{\partial \Theta} (\frac{\partial}{\partial \Theta})^T l\right) = \begin{pmatrix} \frac{1}{\sigma} & 0 \\ 0 & \frac{1}{2\sigma^2} \end{pmatrix}$$

$$I^{-1} = \begin{pmatrix} 5 & 0 \\ 0 & 25^2 \end{pmatrix} = \begin{pmatrix} \sigma^2 & 0 \\ 0 & 20^4 \end{pmatrix}$$

For an unbiased estimate of  $\mu$

$$(E_{\mu, \sigma^2} W = \mu)$$

$$\text{Var } W \geq \frac{\sigma^2}{n} \quad (\text{achieved by } W = \bar{X})$$


---

For an unbiased estimate of  $\sigma^2$

$$\text{Var } W \geq \frac{20^4}{n} \quad (\text{not achieved exactly.})$$

$S^2$  is best unbiased

$$\left( \begin{pmatrix} \bar{X} \\ S^2 \end{pmatrix} \sim AN \left( \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix}, \begin{pmatrix} \frac{\sigma^2}{n} & 0 \\ 0 & \frac{20^4}{n} \end{pmatrix} \right) \right)$$

and

$$S^2 = \frac{\sigma^2}{n-1} \chi^2_{n-1}$$

Limiting distn.  
of MLE

so that

$$\text{Var } S^2 = \frac{20^4}{(n-1)}$$

Note :

$$\text{Var} \left( \frac{1}{n} \sum (X_i - \mu)^2 \right) = \frac{2\sigma^4}{n}$$

$$E \frac{1}{n} \sum (X_i - \mu)^2 = \sigma^2.$$

Achieves CR-bound, but not legitimate estimator if  $\mu$  is unknown.

Example: Gamma( $\alpha, \beta$ )

$$f(x|\alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-x/\beta}$$

$$\log f = -\log \Gamma(\alpha) - \alpha \log \beta + (\alpha-1) \log x - x/\beta$$

$$\begin{pmatrix} \frac{\partial}{\partial \alpha} \log f \\ \frac{\partial}{\partial \beta} \log f \end{pmatrix} = \begin{pmatrix} -\psi(\alpha) - \log \beta + \log x \\ -\frac{\alpha}{\beta} + \frac{x}{\beta^2} \end{pmatrix}$$

$$\begin{pmatrix} \frac{\partial^2}{\partial \alpha^2} & \frac{\partial}{\partial \alpha} \frac{\partial}{\partial \beta} \\ \frac{\partial^2}{\partial \alpha \partial \beta} & \frac{\partial}{\partial \beta} \end{pmatrix} \log f = \begin{pmatrix} -\psi'(\alpha) & -1/\beta \\ -1/\beta & \frac{\alpha}{\beta^2} - \frac{2x}{\beta^3} \end{pmatrix}$$

$$I(\Theta) = -E \left( \frac{\partial}{\partial \theta} \right) = \begin{pmatrix} +\psi'(\alpha) & +1/\beta \\ +1/\beta & 2/\beta^2 \end{pmatrix}$$

[Note:  $E X = \alpha \beta$ ]

$$I^{-1}(\Theta) = \begin{pmatrix} \alpha/\beta^2 & -1/\beta \\ -1/\beta & \psi'(\alpha) \end{pmatrix} \begin{pmatrix} \alpha \psi'(\alpha) - 1 \\ \beta^2 \end{pmatrix}$$

$$I^{-1}(\theta) = \begin{pmatrix} \alpha & -\beta \\ -\beta & \beta^2 \psi'(\alpha) \end{pmatrix} / (\alpha \psi'(\alpha) - 1)$$

CR bound for unbiased estimates of  $\beta$ .

$$\begin{aligned} \text{Var}(\hat{\beta}) &\geq \frac{1}{n} (I^{-1}(\theta))_{22} \geq \frac{1}{n} (I(\theta))_{22}^{-1} \\ &\geq \underbrace{\frac{1}{n} \left( \frac{\beta^2 \psi'(\alpha)}{\alpha \psi'(\alpha) - 1} \right)}_{\frac{\beta^2}{\alpha n} \cdot \left[ \frac{\psi'(\alpha)}{\psi'(\alpha) - \frac{1}{\alpha}} \right]} \geq \underbrace{\frac{1}{n} \frac{1}{\alpha \beta^2}}_{\frac{\beta^2}{\alpha n}} \end{aligned}$$

If  $\alpha$  is known, lower bound is achieved.

$$E\left(\frac{\bar{X}}{\alpha}\right) = \beta$$

$$\text{Var}\left(\frac{\bar{X}}{\alpha}\right) = \frac{1}{\alpha^2} \frac{\text{Var} X}{n} = \frac{\alpha \beta^2}{n \alpha^2} = \frac{\beta^2}{n \alpha}$$

If  $\alpha$  must be estimated, there is a variance penalty which does not vanish asymptotically ( $n \rightarrow \infty$ ).

