

## Lecture 7. Maximum Likelihood Estimation.

## Re-cap of last lecture

$x \sim f_n(x|\theta)$   
plug in random variable

$$L(\theta|x) = \log f_x(x|\theta) \rightarrow \max \theta$$

To find max likelihood, take max derivative

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

when multivar

have matrix

1st info eg.  $E[S(\theta)] = 0$

I diff. at diff. values of  $\theta$   $I(\theta) = \text{Var}(S(\theta)) \rightarrow E[S(\theta)] = 0$

$$= E[S^2(\theta)] - E[S(\theta)]^2$$

2nd info eq:

$$I(\theta) = -\epsilon \left[ \frac{\partial^2 l}{\partial \theta \partial \theta'} \right]$$

## Rao-Cramer bound

$\hat{\theta} = w(x)$  you calculated from  $\hat{\theta}$  data

### Theorem 1 (Rao-Cramer bound).

unbiased  $E\hat{\theta} = \theta$

Let  $X$  be a random data with distribution  $f(x|\theta)$  and information  $I(\theta)$ . Let  $W(X)$  be an estimator of  $\theta$  such that

$$(1) \frac{d}{d\theta} E_\theta[W(X)] = \int W(x) \frac{\partial f(x|\theta)}{\partial \theta} dx$$

$$(2) \text{Var}(W) < \infty.$$

Then

$$\text{Var}(W) \geq \left( \frac{d}{d\theta} E_\theta[W(X)] \right)^2 \frac{1}{I(\theta)}. \Rightarrow \text{Var}(\hat{\theta}) \geq \frac{1}{I(\theta)}$$

In particular, if  $X = (X_1, \dots, X_n)$  is an i.i.d. random sample and  $W$  is unbiased for  $\theta$ , then

$$\text{Var}(W) \geq \frac{1}{I(\theta)} = \frac{1}{n I_1(\theta)}$$

$$I_{\text{total}} = n I_1$$

$$\text{Var}(\hat{\theta}) \geq \frac{1}{n I_1(\theta)}$$



# ML Estimator

Parametric family:

- $X \sim f(x|\theta)$  with  $\theta \in \Theta$ ,  $f(\cdot|\theta)$  is pdf or pmf
- $\ell(\theta|X) = \log f(X|\theta)$  is log-likelihood function

$$\hat{\theta}_{ML} = \arg \max_{\theta \in \Theta} \ell(\theta|x)$$

$$\square \text{ F.O.C } s(\hat{\theta}_{ML}) = \frac{\partial \ell(\hat{\theta}_{ML}|X)}{\partial \theta} = 0$$

$x = (x_1, \dots, x_n)$      $x_i \sim \text{iii}$   $f_i(\cdot|\theta)$  || individual pdf.

pdf.  $f(x|\theta) = \prod_i f_i(x_i|\theta)$

$$\ell(\theta) = \sum_{i=1}^n \ell_i(\theta|x_i) = \sum_{i=1}^n \log f_i(x_i|\theta)$$

Then take

$$\text{FOC} \quad s(\hat{\theta}) = 0$$

$$\frac{\partial \ln}{\partial \theta} = \sum_{i=1}^n s_i(\theta|x_i) = \sum_{i=1}^n s_i(\theta|x_i)$$

## ML Estimator

- If  $X = (X_1, \dots, X_n)$  i.i.d  $f_1(\cdot | \theta_0)$  with  $\theta_0 \in \Theta$
- Then the joint pdf is  $f(x|\theta) = \prod_{i=1}^n f_1(x_i|\theta)$
- The log-likelihood is  $\ell(\theta|x) = \sum_{i=1}^n \log f_1(x_i|\theta)$
- F.O.C.  $\frac{1}{n} \sum_{i=1}^n \frac{\partial \ell_1(\hat{\theta}_{ML}|x_i)}{\partial \theta} = 0$     FOC     $\frac{1}{n} \sum_{i=1}^n (\hat{\theta} | x_i) = 0$   
 $E \hat{\ell}_1(\theta_0)^*$  = 0     $\leftarrow$  population  
eg (1st info. equality)  
\* is a method of moment.

## Theorem 2 (MLE consistency).

In the setting above, assume that

diff value of  $\theta$  give diff. distn  
so diff pdfs.

- (1)  $\theta_0$  is **identifiable**, i.e. for any  $\theta \neq \theta_0$ , there exists  $x$  such that  $f(x|\theta) \neq f(x|\theta_0)$ ,
- (2) the support of  $f(\cdot|\theta)$  does not depend on  $\theta$ ,
- (3)  $\theta_0$  is an interior point of parameter space  $\Theta$ .

Then  $\hat{\theta}_{ML} \rightarrow_p \theta_0$ . *not on the boundary*

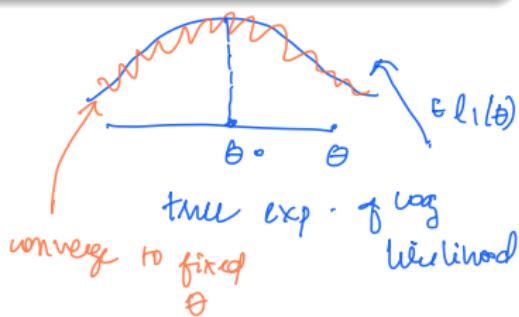
\* Outline of proof:

i)  $\frac{1}{n} \sum_i l_1(\theta|x_i)$   $\xrightarrow{\text{converge}}$   $E l_1(\theta)$   
average do converge to

ii)  $(\text{true}) \quad \theta_0 = \arg \max_{\theta} E l_1(\theta)$

Joining this two,

iii)  $\hat{\theta}_{ML} = \arg \max_{\theta} \frac{1}{n} \sum_i l_1(\theta|x_i) \xrightarrow{P} \theta_0$



e.g.: consistency & common proved using this method

$$\sqrt{n}(\hat{\theta}_{ML} - \theta_0) \Rightarrow N(0, I_1^{-1})$$

### Theorem 3 (MLE asymptotic normality).

In the setting above, assume that conditions (1)-(3) in the MLE consistency theorem hold. In addition, assume that

- (4)  $f_1(x_i|\theta)$  is thrice differentiable with respect to  $\theta$  and we can interchange integration with respect to  $x$  and differentiation with respect to  $\theta$ ,  
\*  $s$  should have small tails
- (5)  $|\partial^3 \log f_1(x_i|\theta)/\partial\theta^3| \leq M(x)$  and  $E[M(X_i)] < \infty$ .

Then

$$\sqrt{n}(\hat{\theta}_{ML} - \theta_0) \Rightarrow N(0, I_1^{-1}(\theta_0))$$

Sketch of Proof

$$0 = \frac{d \ln L(\hat{\theta}_{ML})}{d\theta} \xrightarrow{\text{full likelihood}} \text{FOC we have}$$
$$= \frac{dL(\theta_0)}{d\theta} + \frac{d^2L(\theta^*)}{d\theta^2} \cdot (\hat{\theta}_{ML} - \theta_0) \quad \begin{matrix} \leftarrow \text{Taylor Exp} \\ \uparrow \end{matrix}$$

$$\sqrt{n}(\hat{\theta} - \theta_0) = -\frac{\frac{\partial L(\theta_0)}{\partial\theta} \frac{1}{\sqrt{n}}}{\frac{\partial^2 L(\theta^*)}{\partial\theta^2} \frac{1}{n}} \quad \begin{matrix} \text{some intermediate point between } \theta_0 \text{ and } \hat{\theta}_{ML} \\ \text{always on numerator to normalize} \\ \text{always.} \end{matrix}$$

Now want normalize:

solving once  $\theta_0$  &  $x_i$  are known to  $\theta_0$

$$\textcircled{1} \frac{1}{n} \frac{\partial^2 \ln(\theta^*)}{\partial \theta \partial \theta} = \text{can write as} = \frac{1}{n} \sum_{i=1}^n \frac{\partial^2 \partial l_1(\theta^*|x_i)}{\partial \theta \partial \theta} \text{ random but close to } \theta_0$$

so denominator  $\textcircled{1}$  converge to constant.

If we get  $\rightarrow$

$$\frac{\partial^2 l_1}{\partial \theta \partial \theta} (\theta_0|x_i) = -I_1$$

future info

$$\textcircled{2} \frac{1}{\sqrt{n}} \frac{\partial l}{\partial \theta} (\theta_0) = \frac{1}{\sqrt{n}} \sum_i \frac{\partial l_1}{\partial \theta} (\theta_0|x_i)$$

$\sim$  iid random variable

$$= \frac{1}{\sqrt{n}} \sum_i \underbrace{s_1(\theta_0|x_i)}_{Y_i}$$

$$\frac{1}{\sqrt{n}} \sum_i Y_i \stackrel{D}{\rightarrow} \sqrt{n} \left( \frac{1}{n} \sum_i Y_i - 0 \right) \Rightarrow N(0, I_1)$$

② numerator converge to Gaussian

$$\begin{aligned} EY_i &= E s_1 = 0 \\ \text{Var} Y_i &= \text{Var}(s_1) = I_1 \end{aligned}$$

$$\text{So, } \sqrt{n} (\hat{\theta} - \theta_0) \Rightarrow -\frac{1}{I_1} \cdot N(0, I_1) = N(0, \frac{I_1}{I_1}) = N(0, \frac{1}{I_1})$$

## Examples

"kinda":  $\hat{\theta}_{ML} - \theta_0 \stackrel{\text{?}}{\sim} N(0, \frac{1}{I_n}) = N(0, \frac{1}{I_n})$

$$\text{Var}(\hat{\theta}_{ML} - \theta_0) \propto \frac{1}{I_n}$$

MLE is asymptotically efficient.

Example 1)  $x_1, \dots, x_n \sim_{iid} \exp(\lambda) - \text{expo}$

$$x_i \sim \exp(\lambda) \quad f_1(x_i | \lambda) = \lambda \exp(-\lambda x)$$

$$\begin{cases} n \geq 0 \\ \lambda \geq 0 \end{cases}$$

$$L(\lambda) = \log \lambda - \lambda x_i \rightarrow \text{one obs log likelihood}$$

$$\ln(\lambda) = n \log \lambda - \lambda \sum x_i \rightarrow \text{for all } \lambda \rightarrow \max_{\lambda}$$

Taking FOC:  $\frac{n}{\lambda} - \sum x_i \text{ & so } \hat{\lambda} = \frac{n}{\sum x_i} = \frac{1}{\bar{x}}$

b) Exam  
what is the  
distribution of  $\hat{\lambda}$

$$E X_i = \frac{1}{\lambda} = \int_0^\infty x \lambda e^{-\lambda x} dx$$

\* One way of getting asymptotics:

$$\bar{x} \xrightarrow{P} E X_i = \frac{1}{\lambda}, \text{ by CLT} \quad \hat{\lambda}_{m_2} = \frac{1}{\bar{x}} \xrightarrow{P} \frac{1}{E X_i} = \frac{1}{1/\lambda} = \lambda$$

$$\hat{\lambda}_{m_2} \rightarrow \lambda_0 \text{ (consistent)}$$

By Delta method,  $g(\bar{x}) = \frac{1}{\bar{x}}$

$$\sqrt{n} \left( g(\bar{x}) - g\left(\frac{1}{\lambda}\right) \right) = \sqrt{n} \left( \hat{\lambda} - \lambda_0 \right) \xrightarrow{\text{calc. variance.}} N(0, ?) \quad I_1 = \frac{1}{\lambda^2}$$

\* Second way:

$$\sqrt{n} (\hat{\lambda}_{m_2} - \lambda) \xrightarrow{D} N(0, \frac{1}{I_1(\lambda)})$$

$$\text{so } \sqrt{n} (\hat{\lambda}_{m_2} - \lambda) \xrightarrow{D} N(0, \lambda^2)$$

$$\text{Var}(X_i) > \frac{1}{\lambda^2}$$

Can calc by int. of part twice.

$$\frac{\partial l_1}{\partial \theta} = \frac{1}{\lambda} - x_i \quad \frac{\partial^2 l_1}{\partial \theta \partial \theta} = -\frac{1}{\lambda^2}$$

# How to calculate Fisher Information

- If  $I_1(\theta)$  is continuous in  $\theta$  (which is needed for asymptotic results), then  $(I_1(\hat{\theta}_{ML}))^{-1}$  is consistent for  $(I_1(\theta_0))^{-1}$ ; usually,
- $\hat{I} = -\frac{1}{n} \sum_{i=1}^n \frac{\partial^2 \ell_1(\hat{\theta}|X_i)}{\partial \theta^2}$  is consistent for  $I_1(\theta_0)$   $\hat{I} \rightarrow I_1(\theta_0)$
- $\hat{I} = \frac{1}{n} \sum_{i=1}^n \left( \frac{\partial \ell_1(\hat{\theta}|X_i)}{\partial \theta} \right)^2$  is consistent for  $I_1(\theta_0)$

$$I_1(\theta_0) = E \left( \left[ \frac{\partial \ell_1(\theta_0|X_i)}{\partial \theta} \right]^2 \right) = -E \left[ \frac{\partial^2 \ell_1}{\partial \theta^2} (\theta_0|X_i) \right]$$

↑  
Challenge 1

Exp is the integral. Take avg  
 $\theta \approx \frac{1}{n} \sum (1|x_i)$   
white space

Challenge 2 is being at  $\theta_0$   
Doubt - b/c mle is consistent

We know from previous proof:  $\sqrt{n}(\hat{\theta}_{ML} - \theta_0) \xrightarrow{D} N(0, \frac{I(1)}{I^2(2)})$

$$I_{(2)}^{-1} \quad I_{(1)} \quad I_{(2)}^1$$

# How to calculate Fisher Information

Another idea is parametric bootstrap:

- For  $b = 1, \dots, B$  simulate sample  $X_b^* = (X_{1b}^*, \dots, X_{nb}^*)$  as i.i.d. draws from  $f_1(x_i | \hat{\theta}_{ML})$  *draw samples from  $\hat{\theta}_{ML}$*
- Find MLE using sample  $X_b^*$ , denote it  $\hat{\theta}_b^*$
- Calculate the sample variance of  $(\hat{\theta}_1^*, \dots, \hat{\theta}_B^*)$

$$\text{Var}(\hat{\theta}_{ML}) \approx \frac{1}{B-1} \sum_{b=1}^B (\hat{\theta}_b^* - \bar{\theta}^*)^2$$

$\bar{\theta}^* \in \hat{\sigma}_{ML}$

# When MLE asymptotic theory fails us

**Example 1** If support depends on the parameter.

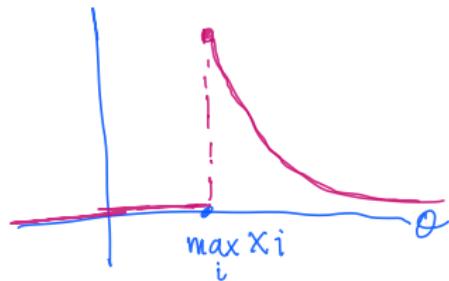
Consider  $X_1, \dots, X_n$  be i.i.d. from  $U[0, \theta]$ .

$$f(\ln|\theta|) = \begin{cases} \frac{1}{\theta} & \text{for } \pi \in [0, \theta] \\ 0 & \text{otherwise} \end{cases}$$

joint pdf

$$f(n) = \frac{1}{\theta^n} \prod_{i=1}^n \max_{\theta} x_i \leq \theta$$

↖ decaying



$$\hat{\theta}_{ML} = \max_i x_i \xrightarrow{n \rightarrow \infty} \theta_0 \text{ (consistent)}$$

-  $\hat{\theta}_{ML} \leq \theta_0$  (always on one side)  
so not be gaussian

super consistent est.  $n(\hat{\theta}_{ML} - \theta_0) \xrightarrow{\text{like exp}} \text{special distribution}$  (not normal)

Violation of mle theory; assumptions violated.

## When MLE asymptotic theory fails us

when parameter on boundary

**Example 2** If the true parameter value  $\theta_0$  is on the boundary of  $\Theta$ .

Consider  $X_1, \dots, X_n$  i.i.d.  $N(\mu, 1)$  with  $\mu \geq 0$ .

$\mu \geq 0$  is like a constraint

$X_1, \dots, X_n$  iid  $N(\mu, 1)$ )

$$\ln(\mu) = \frac{n}{2} \log(2\pi) - \frac{1}{2} \sum (x_i - \mu)^2 \rightarrow \max_{\mu \geq 0}$$

over parameter space

FOC:  $\hat{\mu} = \bar{x}$

true restricted optimum  $\hat{\mu}_M = \begin{cases} \bar{x} & ; \text{if } \bar{x} \geq 0 \\ 0 & ; \text{otherwise} \end{cases}$

if  $\mu_0 = 0$ ,  $\sqrt{n} (\hat{\mu}_M - \mu_0) = \begin{cases} \bar{x} & ; \text{if } \bar{x} \geq 0 \\ 0 & ; \text{otherwise} \end{cases}$



with  $p = 1/2$   $\bar{x} \leq 0$  and  $\hat{\mu}_M = 0$

"parameter on boundary problem"

"moment inequality problem"

still active area of research

# When MLE asymptotic theory fails us

problem of FE.

Ui hard  
to get with  
 $\sigma^2$  can

## Example 3 Incidental parameter problem (number of parameters grow)

Consider

$$X_i = \begin{pmatrix} X_{1i} \\ X_{2i} \end{pmatrix} \stackrel{iid}{\sim} N \left( \begin{pmatrix} \mu_i \\ \mu_i \end{pmatrix}, \begin{pmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{pmatrix} \right)$$

for  $i = 1, \dots, n$ , Two independent scores for two wind-up tests

$$\ln = n \log(2\pi) \rightarrow n \log \sigma^2 - \frac{1}{2} \sigma^2 \in [(x_{1i} - \mu_i)^2 + (x_{2i} - \mu_i)^2]$$

$$\hookrightarrow \max_{\mu_i} y \sigma^2$$

$$\hat{\mu}_i = \frac{x_{1i} + x_{2i}}{2}$$

$$x_{1i} - \hat{\mu}_i = x_{1i} - \frac{x_{1i} + x_{2i}}{2} = \frac{x_{1i} - x_{2i}}{2}$$

$$\frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{\sigma^2} + \frac{1}{2\sigma^4} \mathbb{E} \left[ \underbrace{(x_{1i} - \hat{\mu}_i)^2}_{\frac{(x_{1i} - x_{2i})^2}{4}} + \underbrace{(x_{2i} - \hat{\mu}_i)^2}_{\frac{(x_{1i} - x_{2i})^2}{4}} \right]$$

$$\hat{\sigma}^2 = \frac{1}{2n} \sum_i \left\{ (x_{1i} - \hat{\mu}_i)^2 + (x_{2i} - \hat{\mu}_i)^2 \right\} = \frac{1}{2n} \sum_i \frac{2}{4} (x_{1i} - x_{2i})^2$$

$$\hat{\sigma}^2 = \frac{1}{4n} \sum (x_{1i} - x_{2i})^2 - \text{now see whether } \hat{\sigma}^2 \text{ is}$$

consistent or not-

$$\begin{aligned} & \frac{1}{4} E (x_{1i} - x_{2i})^2 = \frac{1}{4} \text{var}(x_{1i} - x_{2i}) \\ &= \frac{1}{4} (\sigma^2 + \sigma^2) = \frac{\sigma^2}{2} \end{aligned}$$

$\therefore \hat{\sigma}^2 \rightarrow \frac{\sigma^2}{2}$  not to  $\sigma^2$  so the est. is not consistent.

\* when more parameters for est, not be consistent  
\* degrees of freedom also increase.

## Quasi-MLE

$$x_1, \dots, x_n \sim N(\mu, \sigma^2)$$
$$\hat{\sigma}^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

- $X = (X_1, \dots, X_n)$  i.i.d from  $g(x_i)$ .
- We assumed  $X_i \sim f_1(x_i | \theta)$ .
- What would happen if we do MLE?
- MLE estimates a “pseudo-true” parameter value  $\theta_0$ :

$$\theta_0 = \arg \max_{\theta} \int \log[f_1(x_i | \theta)] g(x_i) dx_i = \arg \max_{\theta} \mathbb{E}_g \log f_1(X_i | \theta).$$

$$l = \sum_{i=1}^n x_i \log p + (1-x_i) \log(1-p) - \log(\sigma(x_i)) - \frac{1}{2} \log(2\pi) - \left[ \frac{(y_i - u(x_i))^2}{2\sigma^2(x_i)} \right]_{0,1}$$

$$\hat{l}(\sigma(0), \sigma(1), \mu(0), \mu(1)) =$$

$$\sum_{i=1}^n x_i \log p + (1-x_i) \log(1-p) - \log(\sigma(x_i))$$

$$- \frac{1}{2} \log(2\pi) - \frac{(y_i - u(0))}{2\sigma(0)}$$