

# Lec 5 MATH 341/641

19

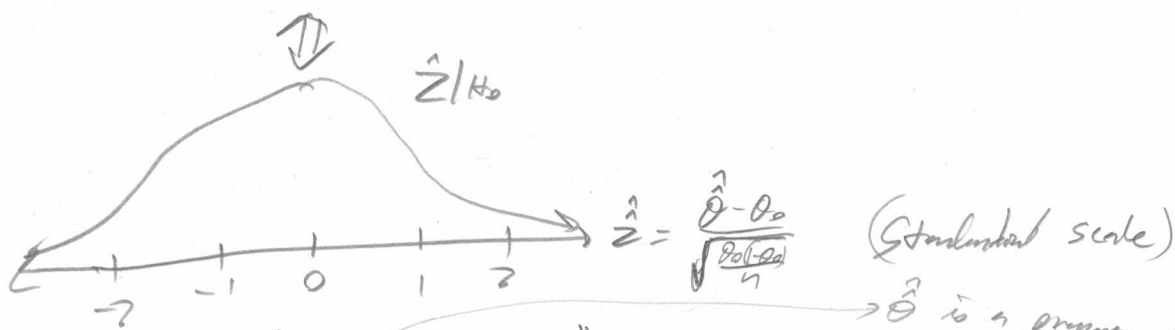
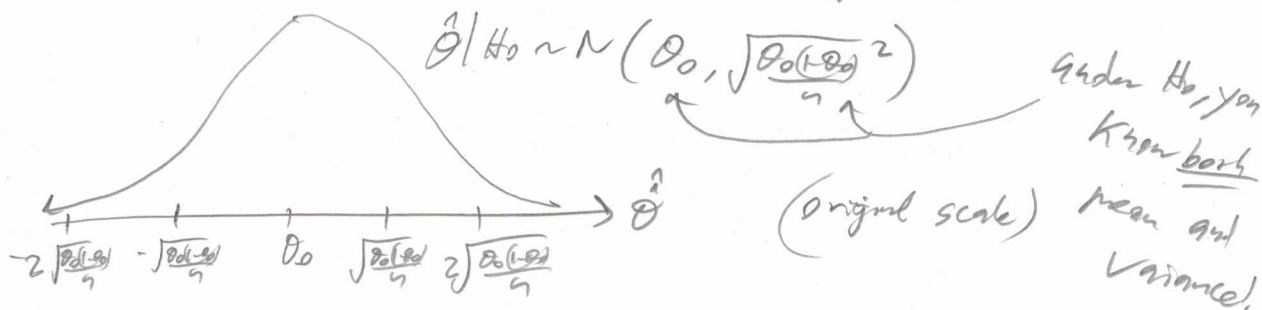
Before we do max. lik. est., let's employ the CLT to see why it's important for statistics. Def: "asymptotically normal estimator" means:

$$\frac{\hat{\theta} - E(\hat{\theta})}{SE[\hat{\theta}]} \xrightarrow{d} N(0,1)$$

Consider OGP:  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(\theta)$  and  $\hat{\theta} = \bar{X}$ . By the CLT,  $\frac{\bar{X} - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} \xrightarrow{d} N(0,1)$

and thus  $\bar{X}$  is asymptotically normal. Once an estimator is asymptotically normal, you know with "large  $n$ "  $\frac{\bar{X} - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} \sim N(0,1)$ . You can use this fact to construct an "approximate test".  $\Rightarrow \bar{X} \sim N(\theta, \frac{\theta(1-\theta)}{n})$

Previously  $\hat{\theta} | H_0 \sim \text{Binom}(n, \theta_0)$  exactly. Now  $\rightarrow$  i.e. the test is approximately  $\alpha$



This test is called the "1-prop z-test". It is approximate.  $\rightarrow \hat{\theta}$  is a proportion of 1's out of  $n$ .

In my opinion, there is no reason to use it! The exact test is better! It performs poorly if  $\theta_0$  is close to 0 or 1. Otherwise it performs okay. E.g.

$$H_0: \theta = .529 \quad RET_\alpha = (-2, 2)$$

$$H_1: \theta \neq .529 \Rightarrow SE[\hat{\theta}] = \sqrt{\frac{.529(1-.529)}{21}} = .499$$

$$\alpha = 5\% \quad n = 21$$

$$\hat{\theta} = \frac{10}{21} = .476$$

$$\hat{z} = \frac{.476 - .529}{.499} = -.096$$

Type II error?  
Underpower?  
Fail to Reject

# Procedures for finding estimators

DGP:  $X_1, \dots, X_n \stackrel{iid}{\sim} p(x; \theta_1, \dots, \theta_k)$  for discrete DGP  
or  $f(x; \theta_1, \dots, \theta_k)$  for cont. DGP

DGP  $\xrightarrow{\text{procedure}} \hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$ , a means of estimation  
of all parameters  
of interest

e.g.  $\hat{\theta}_j^{MM} = \beta_j(\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_k)$

this is justified since  $\hat{\mu}_2 \approx E[X^2]$

and all parameters are functions of the moments

Although taken just once in 1800s...

Populated by RA Fisher between 1912-1922 [2]

Another method for finding estimates is the "maximum likelihood method". Here's how this goes for iid data:

$$X_1, \dots, X_n \stackrel{iid}{\sim} DGP(\theta_1, \dots, \theta_K) \stackrel{(PMF)}{=} f(x; \theta_1, \dots, \theta_K) \text{ if discrete} \\ \stackrel{(PDF)}{=} f(x; \theta_1, \dots, \theta_K) \text{ if continuous}$$

Due to independence and identical distributionness,

I'll use this notation only arbitrarily

$$f(x_1, \dots, x_n; \theta_1, \dots, \theta_K) \stackrel{\downarrow}{=} \prod_{i=1}^n f(x_i; \theta_1, \dots, \theta_K)$$

joint density function

The  $\theta_1, \dots, \theta_K$  are the parameters you would need to calculate the density for the data at any point in the support.

But we don't know  $\theta_1, \dots, \theta_K$  and we are trying to estimate it. So, we do the following conceptual inversion:

$$\begin{array}{c} L(\theta_1, \dots, \theta_K) \stackrel{\uparrow}{=} L(\theta_1, \dots, \theta_K; x_1, \dots, x_n) = f(x_1, \dots, x_n; \theta_1, \dots, \theta_K) \stackrel{\downarrow}{=} \prod_{i=1}^n f(x_i; \theta_1, \dots, \theta_K) \\ \begin{array}{ccccccc} \uparrow & & \uparrow & & \uparrow & & \uparrow \\ \text{if iid DGP} & \text{settings you wish to compare} & \text{values} & \text{inputs/settings you wish to compare} & \text{values} & & \text{if iid DGP} \end{array} \end{array}$$

Note:  $L \geq 0$  for continuous probability,  $L \in \{0,1\}$  for discrete

Now we vary  $\theta_1, \dots, \theta_K$  and see which gives us the most likely values. The value is called the "maximum likelihood estimate" (MLE).

$\hat{\theta}_1, \dots, \hat{\theta}_K := \underset{\theta_1, \dots, \theta_K}{\operatorname{argmax}} (L) = \underset{\theta_1, \dots, \theta_K}{\operatorname{argmax}} \left\{ \prod_{i=1}^n f(x_i; \theta_1, \dots, \theta_K) \right\}$   
 value at which the max is achieved  
 What is argmax? go back to precalc.  
 $-x^2 + 4x - 4$   
 $-(x^2 - 4x + 4)$

let  $f(x) = -x^2 + 4x + 1 = -(x-2)^2 + 5$

What is  $\max\{f(x)\} = 5$

$\underset{x}{\operatorname{argmax}} \{f(x)\} = 2 \Rightarrow \underset{x}{\operatorname{argmax}} \{f(x)\} := \{x : f(x) = \max\{f(x)\}\}$

How to find argmax? Set  $f'(x) = 0$  and solve

$f'(x) = -2x + 4 \stackrel{\text{set}}{=} 0 \Rightarrow x_* = 2$ . If you're responsible, you should check  $f''(x_*)$  to ensure it's  $< 0$  so curve is concave (and not convex).  
 $f''(x) = -2 \Rightarrow f''(2) = -2 < 0$  ✓

Note: the argmax is unaffected by a strictly increasing function

Proof: let  $g(x)$  be a function whose domain is always positive  $g'(x) > 0 \forall x$

Find  $\underset{x}{\operatorname{argmax}} \{g(f(x))\} = \underset{x}{\operatorname{argmax}} \left[ \frac{d}{dx} (g(f(x))) \right] = \underset{x}{\operatorname{argmax}} \left[ \underbrace{g'(f(x))}_{\text{some pos}} f'(x) \stackrel{\text{set}}{=} 0 \right]$

$\Rightarrow f'(x) \stackrel{\text{set}}{=} 0$  solution is  $x_* = \underset{x}{\operatorname{argmax}} \{f(x)\}$

Second derivative test...

$\frac{d^2}{dx^2} [g(f(x))] = \frac{d}{dx} [g'(f(x)) f'(x)] = g'(f(x)) f''(x) + f'(x)^2 g''(f(x))$   
 $\left|_{x=x_*} \begin{matrix} > 0 & < 0 \end{matrix} \right| = \underbrace{g'(f(x_*))}_{>0} \underbrace{f''(x_*)}_{<0} + \underbrace{f'(x_*)^2}_{=0} \underbrace{g''(f(x_*))}_{?} < 0$



Let us use the  $\ln$  function. The log makes everything easier since

$$l := \ln(L) = \ln\left(\prod_{i=1}^n f(x_i; \theta_1, \dots, \theta_K)\right) = \sum_{i=1}^n \ln(f(x_i; \theta_1, \dots, \theta_K))$$

and sums are easy to take derivatives of

To get the MLE's, we take  $\arg\max_{\theta \in \Theta} \{l\}$ . If local max's exist, we use the following system of eq's. (from space)

$$\sum_{i=1}^n \frac{\partial}{\partial \theta_1} [\ln(f(x_i; \theta_1, \dots, \theta_K))] \stackrel{\text{set}}{=} 0$$

$$\sum \frac{\partial}{\partial \theta_2} [\dots] \stackrel{\text{set}}{=} 0$$

$$\vdots$$

$$\sum \frac{\partial}{\partial \theta_K} [\dots] \stackrel{\text{set}}{=} 0$$

If local max's don't exist, we need to work harder.

e.g.  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bern}(\theta)$  ( $K=1$  dimension). Here:

$$\sum_{i=1}^n \frac{\partial}{\partial \theta} \ln(p(x_i; \theta)) = \sum_{i=1}^n \frac{\partial}{\partial \theta} [\ln(\theta^{x_i} (1-\theta)^{1-x_i})]$$

$$= \sum_{i=1}^n \frac{\partial}{\partial \theta} [x_i \ln(\theta) + (1-x_i) \ln(1-\theta)] = \sum_{i=1}^n \frac{x_i}{\theta} - \frac{1-x_i}{1-\theta}$$

$$= \frac{\sum x_i}{\theta} - \frac{n - \sum x_i}{1-\theta} \stackrel{\text{set}}{=} 0 \Rightarrow (1-\theta)(\sum x_i) = (\theta)(n - \sum x_i)$$

$$\Rightarrow \sum x_i - \theta \sum x_i = \theta n - \theta \sum x_i \Rightarrow \sum x_i = \theta (\sum x_i + n - \sum x_i)$$

$$\Rightarrow \hat{\theta}_{MLE} = \frac{\sum x_i}{n} = \bar{X} \quad \text{which was our estimate and the true estimate.}$$

eg.  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta_1, \theta_2)$ . The MLE's  $\hat{\theta}_1^{MLE}$   $\hat{\theta}_2^{MLE}$  (5)

Let's do  $\theta_1$  first:

$$\sum_{i=1}^n \frac{\partial}{\partial \theta_1} \left[ \ln(f(x_i; \theta_1, \theta_2)) \right] = \sum_{i=1}^n \frac{\partial}{\partial \theta_1} \left[ \ln \left( \frac{1}{\sqrt{2\pi\theta_2}} e^{-\frac{1}{2\theta_2}(x_i - \theta_1)^2} \right) \right]$$

$$= \sum_{i=1}^n \frac{\partial}{\partial \theta_1} \left[ -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\theta_2) - \frac{1}{2\theta_2} (x_i - \theta_1)^2 \right] = -\frac{x_i^2}{2\theta_2} + \frac{x_i\theta_1}{\theta_2} - \frac{\theta_1^2}{2\theta_2}$$

$$= \sum_{i=1}^n \frac{x_i}{\theta_2} - \frac{\theta_1}{\theta_2} = \frac{\sum x_i}{\theta_2} - \frac{n\theta_1}{\theta_2} \stackrel{\text{set}}{=} 0 \Rightarrow \hat{\theta}_1^{MLE} = \bar{X}$$

Let's do  $\theta_2$

$$\sum \frac{\partial}{\partial \theta_2} \left[ -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\theta_2) - \frac{1}{2\theta_2} (x_i - \theta_1)^2 \right]$$

$$= \sum -\frac{1}{2\theta_2} + \frac{1}{2\theta_2^2} (x_i - \theta_1)^2 = -\frac{n}{2\theta_2} + \frac{\sum (x_i - \theta_1)^2}{2\theta_2^2} \stackrel{\text{set}}{=} 0$$

system of eq's

$$\Rightarrow \sum (x_i - \theta_1)^2 = n\theta_2 \Rightarrow \theta_2 = \frac{1}{n} \sum (x_i - \theta_1)^2 \Rightarrow \hat{\theta}_2^{MLE} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$$

Since  $w$  is the same for  $\hat{\theta}$  and  $\hat{\theta}^*$ , define the MLE

"Maximum likelihood estimator" as the rv that takes  $\hat{\theta}$ ,  $\hat{\theta}^{MLE}$ .

$\hat{\theta}^{MLE}$  is a value, a pt. estimate and  $\hat{\theta}^{MLE}$  is its sampling distr.

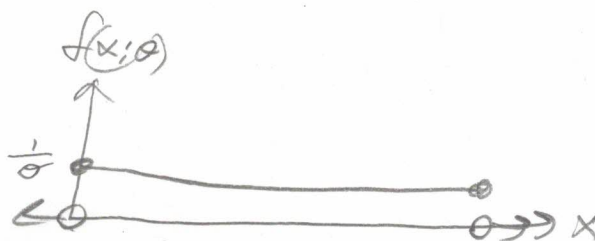
let's do a different example. I'm not expecting you to be able to do this yourself, only follow along. The results of this will (eventually) be a whole new topic. Let  $X_1, \dots, X_n \stackrel{iid}{\sim} U(0, \theta)$ .

We already derived the silly  $\hat{\theta}_{MLE} = 2\bar{X}$ . Can the MLE do "better"?

$$\sum \frac{\partial}{\partial \theta} [\ln(f(x; \theta))] = \sum \frac{\partial}{\partial \theta} [\ln(\frac{1}{\theta})] = -\sum \frac{1}{\theta} = -\frac{n}{\theta} \stackrel{\text{set}}{=} 0$$

NO SOL!!!!

Something is wrong! Yes, this MLE is not a local max. It's a global max.

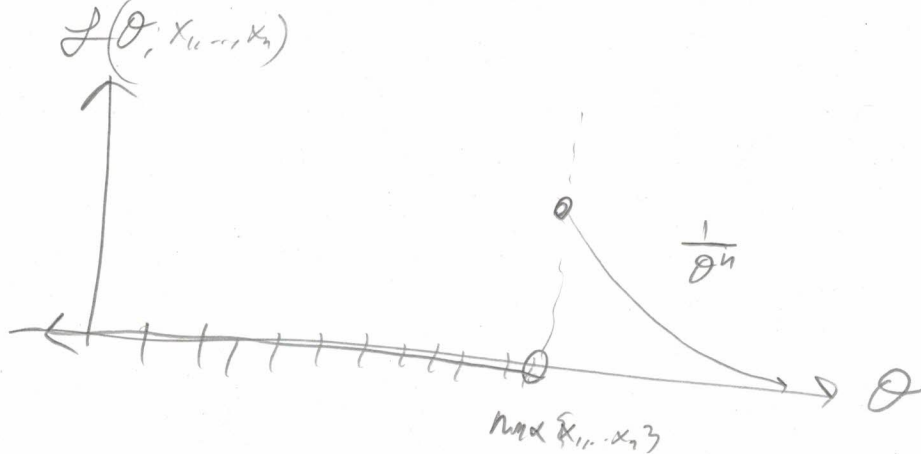


$$f(x; \theta) = \begin{cases} \frac{1}{\theta} & \text{if } 0 \leq x \leq \theta \\ 0 & \text{o/t} \end{cases}$$

$$f(x; \theta) = \frac{1}{\theta} \mathbb{1}_{x \in [0, \theta]}$$

$$\Rightarrow \prod_{i=1}^n f(x_i; \theta) = \left\{ \frac{1}{\theta^n} \mathbb{1}_{x_1 \in [0, \theta]} \mathbb{1}_{x_2 \in [0, \theta]} \dots \mathbb{1}_{x_n \in [0, \theta]} \right.$$

$$\Rightarrow \mathcal{L}(\theta; x_1, \dots, x_n) = \left\{ \frac{1}{\theta^n} \mathbb{1}_{\theta > \max(x_1, \dots, x_n)} \right.$$



What is the maximum likelihood estimate?  $\hat{\theta}_{MLE} = \max \{x_1, \dots, x_n\}$

Var  $\hat{\theta}_{MLE}$  proof in Math 380 class

$$= \theta^2 \frac{n}{(n+1)^2(n+2)}$$

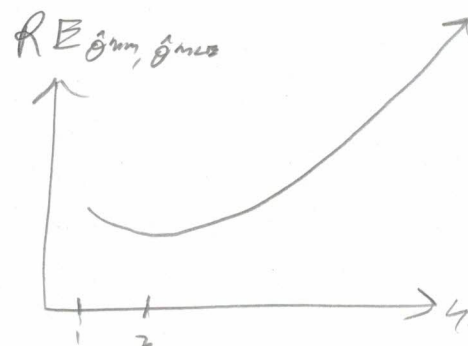
$$X \sim U(a, b) \Rightarrow \text{Var}(X) = \frac{(b-a)^2}{12}$$

$$\text{Var}(\hat{\theta}_{MM}) = \text{Var}(2\bar{X}) = 4 \text{Var}(\bar{X}) = \frac{4}{n} \text{Var}(X) = \frac{4}{n} \frac{\theta^2}{12} = \frac{\theta^2}{3n}$$

$$\text{Hence, } \text{Var}(\hat{\theta}_{MLE}) < \text{Var}(\hat{\theta}_{MM})$$

Relative Efficiency (RE)

$$RE = \frac{\text{Var}(\hat{\theta}_{MM})}{\text{Var}(\hat{\theta}_{MLE})} = \frac{\frac{\theta^2}{3n}}{\theta^2 \frac{n}{(n+1)^2(n+2)}} = \frac{(n+1)^2(n+2)}{3n^2}$$



$\hat{\theta}_{MLE}$  gets better and better!!

Should we use

$\text{MSE}(\hat{\theta}_{MM})$  ? Yes But  $\hat{\theta}_{MM}$  is unbiased.

$\text{MSE}(\hat{\theta}_{MLE})$  we can prove  $\hat{\theta}_{MLE}$  has very small bias

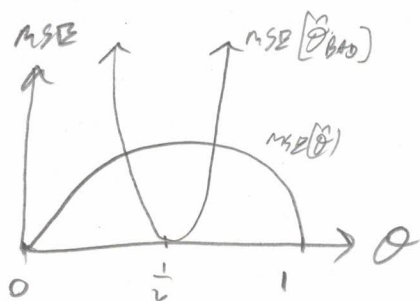
which goes to 0 as  $n \rightarrow \infty$ , so the picture is the same.



There are lots of estimators! Questions:

- ① Is there a theoretical minimum MSE when estimating  $\theta$  for a given DGP?
- ② If so, is there a procedure to find this estimator with theoretically minimum MSE?

The answer to # is NO. Why? Recall  $\hat{\theta}_{\text{BAYO}} = \frac{1}{2}$   
vs  $\hat{\theta} = \bar{X}$  for DGP:  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bern}(\theta)$



The class of estimators is a set which is uncountably infinite and the choice of  $\theta$  matters!

So let's limit the set and ask again.

① Among all unbiased estimators, is there a theoretical minimum var = MSE? Yes! It is called the Cramer-Rao Lower Bound (CRLB) proven in 1945-1946. This estimator is called the "uniformly minimum variance unbiased estimator" (UMVUE)

② Is there a way to derive the UMVUE from a DGP?

No! But if you find a  $\hat{\theta}$  with the CRLB  $\Rightarrow \hat{\theta}$  is the UMVUE!