

this more about maximum likelihood  
time: & Bayesian inference

AMS  
132

next yet more about Bayes  
time:

7 Feb  
17

Step 0 You start  
with a sampling model

for your data set

$Y = (Y_1, \dots, Y_n)$   
~ observed data

Fisher's recipe ①  
for maximum likelihood  
(large-sample) inference in  
"nice" problems (regular)

$\underline{I} = (I_1, \dots, I_n)$   
~ random variables  
representing the  
sampling process  
before it's occurred

Ex. Kaiser Ichi Case Study

marginal sampling dist.  
of obs.  $(I_i | \theta, \mathcal{B}) \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\theta)$   
( $i=1, \dots, n$ )

$\mathcal{B}$  = background  
information,  
assumptions  
& judgments

i.e., for  $Y_i = \{0, 1\}$   
 $P(I_i = Y_i | \theta, \mathcal{B}) = \theta^{Y_i} (1-\theta)^{1-Y_i}$

Step 1 write out the joint

sampling distribution of  $\underline{I} = (I_1, \dots, I_n)$

Ex. (kaiser) By independence ( $\textcircled{I}$  ID),  $\textcircled{2}$

joint dist. of  $(Y_1, \dots, Y_n)$  is product of marginals:

$$\begin{aligned} P(Y_1 = y_1, \dots, Y_n = y_n \mid \theta \mathcal{B}) & \stackrel{\textcircled{I} \text{ ID}}{=} \prod_{i=1}^n P(Y_i = y_i \mid \theta \mathcal{B}) \\ & = \prod_{i=1}^n \theta^{y_i} (1-\theta)^{1-y_i} = \theta^{y_1 + \dots + y_n} (1-\theta)^{n - (y_1 + \dots + y_n)} \\ & = \theta^s (1-\theta)^{n-s}, \text{ where } s = \sum_{i=1}^n y_i. \end{aligned}$$

Before data arrives, this is a function of  $\underline{y} = (y_1, \dots, y_n)$  for fixed  $\theta$ .

Step 2] After data arrives, think of it (the joint sampling dist.) as a function of  $\theta$  for fixed  $\underline{y}$ , and multiply it

by an arbitrary positive constant  $c$ :

$$\ell(\theta \mid \underline{y} \mathcal{B}) = c P(\underline{Y} = \underline{y} \mid \theta \mathcal{B})$$

likelihood function joint ~ sampling dist.

Ex. (Kaiser)  $L(\theta | y, B) = c \theta^S (1-\theta)^{n-S}$  (3)

Important point,  
which I forgot  
to make earlier

Notice that the  
likelihood function  
depends on the data  
vector  $y = (y_1, \dots, y_n)$

only through the sum  $S = \sum_{i=1}^n y_i$ .

Definition  
(Fisher)

If  $L(\theta | y)$  depends on  $y$  only through  
some function  $s(y)$ ,  $S$  is said to be

(a) sufficient statistic for  $\theta$  in the  
sampling model that gave rise to  $L(\theta | y)$ ;  
short hand:  $S$  is sufficient for  $\theta$

Ex. (Kaiser) | Some sufficient statistics  
in the Bernoulli sampling model:

①  $s_1(y) = y$ ; i.e., the entire data vector is sufficient (Fisher was trying to achieve dimensionality reduction, from  $n$  observations down to a smaller number of "pieces of information"; for this purpose  $s_1(y) = y$  is useless) (4)

---

② Suppose that  $n$  is even, and define  $s_2(y) = \left( \sum_{i=1}^{n/2} y_i, \sum_{i=\frac{n}{2}+1}^n y_i \right)$ ; this is also sufficient & is a big dimensionality-reduction improvement on  $s_1(y) = y$ , from  $n$  to 2.

---

③  $s_3(y) = \sum_{i=1}^n y_i$  is even better ( $n$  to 1) <sup>down</sup>

---

If (informal Definition) no sufficient statistic

exists of dimension lower than the  $\Theta$   
one ( $S^*$ , say) you're currently considering,  
 $S^*$  is said to be minimal sufficient.

Ex. (Kaiser)  $S = \sum_{i=1}^n Y_i$  is a minimal suff. stat.  
in the Bernoulli  
sampling model

Technical note: strictly

speaking ~~of~~ minimal sufficient stat.

in this model is based on  $(\sum_{i=1}^n Y_i, n)$ ,  
but people usually omit mention of  $n$

Note) Minimal sufficiency is not unique;

$S = \sum_{i=1}^n Y_i$  is minimal sufficient in

Kaiser case study, but so is  $S+1$

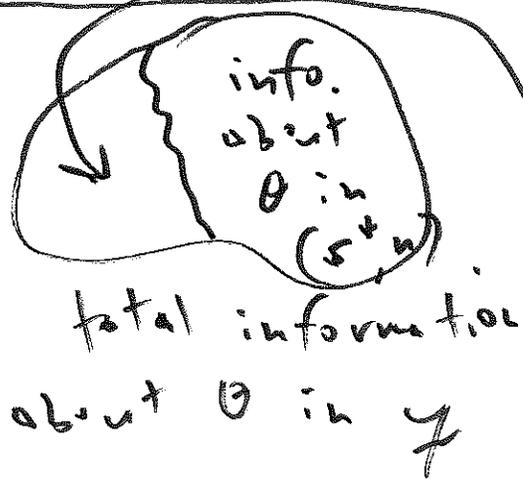
Generally speaking (but not always),

if the unknown quantity  $\theta$  in (b) the sampling model is of dimension  $k$  & a (minimal) suff. stat. exists, it will also be of dimension  $k$  (weird exceptions exist)

Important limitation of the idea of sufficiency

you can throw  $Y$  away & just go forward with  $(S^*, n)$ . This is false:

info useful to criticize your sampling dist. ( &  $\therefore$  your likelihood  $f^n$ ).



ex.  $\vec{y}$  pretend  $(Z_i | \theta) \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta)$  (7)   
 $(i=1, \dots, n)$

$n = 10$

$s=5$	1 0 1 1 0 0 1 0 1 0	IID plausible
$s=5$	0 0 0 0 0 1 1 1 1 1	IID implausible
$s=5$	0 1 0 1 0 1 0 1 0 1	

return to Fisher's recipe

recall: step 2: get lik. fun. from joint. dist.

$l(\theta | \vec{y}) = l(\theta | s, n)$   
 $= \theta^s (1-\theta)^{n-s}$

step 2.5) Also compute

log likelihood function:  $ll(\theta | \vec{y}) = ll(\theta | s, n)$   
 $= \log l(\theta | \vec{y})$

$$\text{here } \ell(\theta | z) = \log \theta^s (1-\theta)^{n-s} \quad (8)$$

$$= (\log \theta^s) + (\log (1-\theta)^{n-s})$$

$$= s \log \theta + (n-s) \log (1-\theta)$$

Find  $\theta$  to

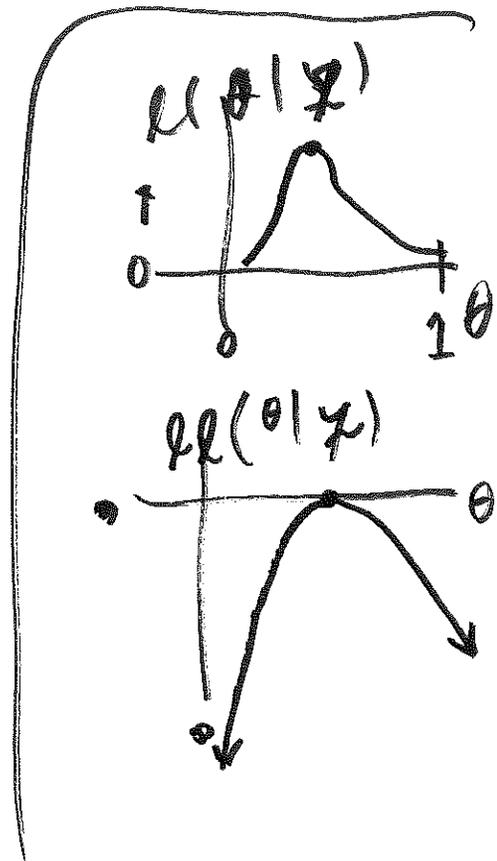
Step 3 } Maximize  $\ell(\theta | s, n)$  or

equivalently to maximize  $\ell(\theta | s, n)$ .

This is usually done  
by differentiating once  
with respect to  $\theta$ ,

setting first partial  
derivative (w.r.t.  $\theta$ ) to

0 \* solving



Note: Always easier to diff.  $ll(\theta|y)$  because log converts products into sums

$\mu$ . (Kaiser)

lik. fn is always a product

$$ll(\theta|s_n) = s \log \theta + (n-s) \log (1-\theta)$$

$$\frac{d}{d\theta} ll(\theta|s_n) = \frac{s}{\theta} + (n-s) \left( \frac{1}{1-\theta} \right) (-1)$$

$$= \frac{s}{\theta} - \frac{n-s}{1-\theta} = \frac{s(1-\theta) - (n-s)\theta}{\theta(1-\theta)}$$

$$= \frac{s - s\theta - n\theta + s\theta}{\theta(1-\theta)} = 0 \quad \text{iff } s = n\theta$$

$$\hat{\theta}_{MLE} = \frac{s}{n} = \bar{y}$$

populat. mean  $\theta = ?$

sample  $\begin{bmatrix} 15 \\ 2 \\ 05 \end{bmatrix}$   $\Rightarrow$   $\begin{bmatrix} 15 \\ 205 \end{bmatrix}$   $n = 112$

mean  $\bar{y} = \frac{s}{n} = \frac{4}{112} = 0.036$

in a fitfully sample mean should be good est. of mean

here  $\hat{\theta}_{MLE} = \frac{4}{112} = 0.036$

Step 4 (10)

Compute  $\vec{SE}(\hat{\theta}_{MLE}) = ?$

$$\vec{SE}_{RS}(\hat{\theta}_{MLE}) = \sqrt{\hat{V}_{RS}(\hat{\theta}_{MLE})}$$

↑  
repeated  
sampling

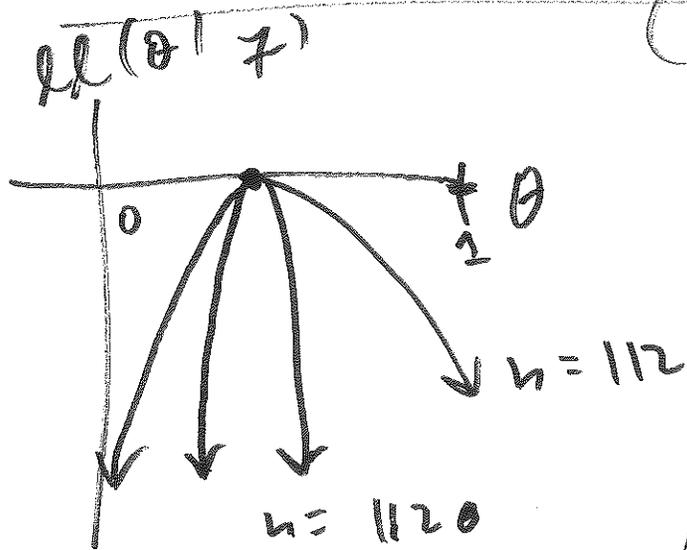
Note: we already know  
in earlier problem that

$$\begin{aligned}\hat{V}_{RS}(\hat{\theta}_{MLE}) &= \hat{V}_{RS}(\bar{Y}) \\ &= \frac{\hat{\theta}(1-\hat{\theta})}{n}\end{aligned}$$

So

$$\vec{SE}_{RS}(\hat{\theta}_{MLE}) = \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}} = 0.017$$

Fisher had a general-purpose method ⑪  
 for computing  $\hat{V}_{RS}(\hat{\theta}_{MLE})$ :



(observed) Fisher information

$$\hat{I}(\hat{\theta}_{MLE}) =$$

$$\left[ - \frac{d^2}{d\theta^2} l(\theta | z) \right]_{\theta = \hat{\theta}_{MLE}}$$

Fisher showed that

$$\hat{V}_{RS}(\hat{\theta}_{MLE}) = \hat{I}^{-1}(\hat{\theta}_{MLE})$$

$$\hat{SE}_{RS}(\hat{\theta}_{MLE}) = \sqrt{\hat{I}^{-1}(\hat{\theta}_{MLE})}$$

$$l(\theta | s, n) = s \log \theta + (n-s) \log (1-\theta)$$

$$\frac{d}{d\theta} \ell(\theta | s, n) = \frac{s}{\theta} - \frac{n-s}{1-\theta} \quad (12)$$

$$\frac{d^2}{d\theta^2} \ell(\theta | s, n) = -\frac{s}{\theta^2} - \frac{n-s}{(1-\theta)^2}$$

simplify

$$\downarrow$$

$$= \frac{n}{\theta(1-\theta)}$$

$\int_0$

$$\hat{I}(\hat{\theta}_{MLE}) = \frac{n}{\hat{\theta}_{MLE} (1-\hat{\theta}_{MLE})}$$

$$\neq \text{SE}_{FS}(\hat{\theta}_{MLE}) = \frac{1}{\sqrt{I(\hat{\theta}_{MLE})}} = \frac{1}{\sqrt{\frac{n}{\hat{\theta}_{MLE} (1-\hat{\theta}_{MLE})}}} = \sqrt{\frac{\hat{\theta}_{MLE} (1-\hat{\theta}_{MLE})}{n}}$$

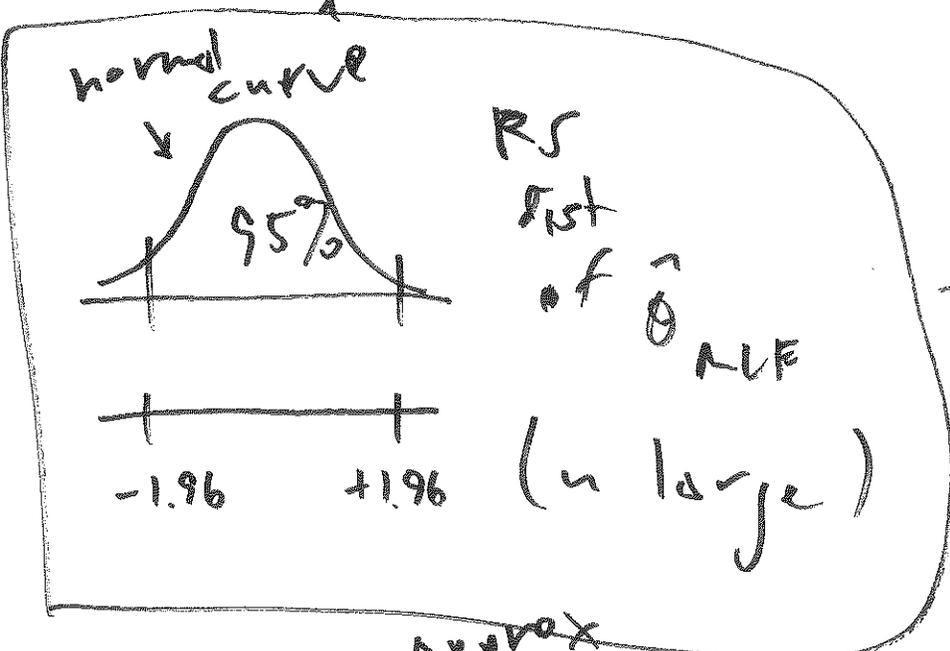
$$= \sqrt{\frac{\hat{\theta}_{MLE} (1-\hat{\theta}_{MLE})}{n}}$$

$\approx 0.017$  ← give or take for  $\hat{\theta}_{MLE}$

as  $n \uparrow$ ,  
 $\hat{I} \uparrow$  at a rate  
 proportional  
 to  $n$

Step 5] approx.  $100(1-\alpha)\%$  confidence (CI) (13)

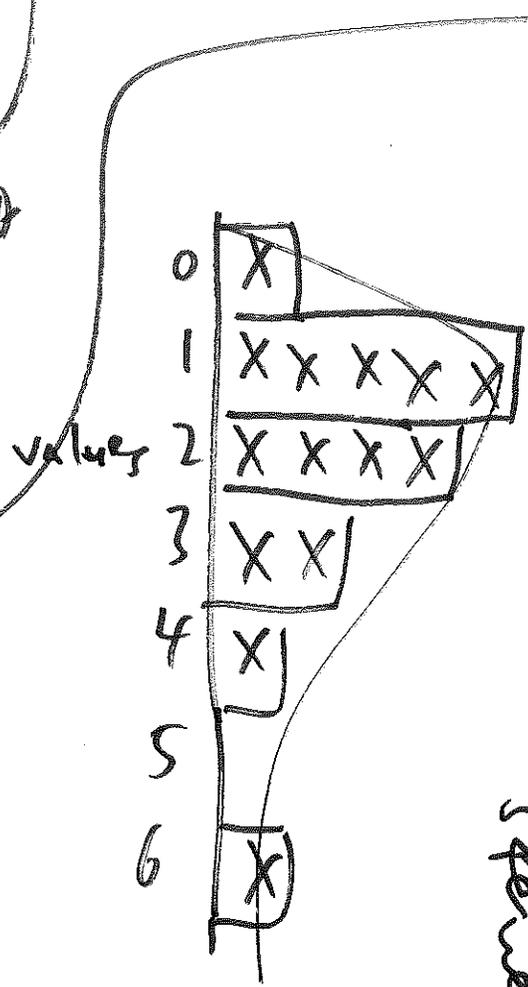
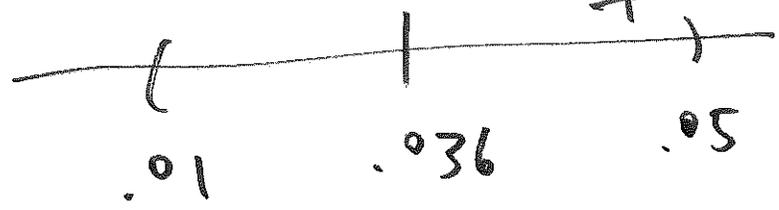
interval for  $\theta$ :  $\alpha = .05 \rightarrow 95\%$  CI



$$\hat{\theta}_{MLE} \pm 1.96 \hat{SE}(\hat{\theta}_{MLE})$$

$$.036 \pm 2(.017)$$

approx.  $95\%$  CI for  $\theta$



1 2 1 4 1 2 2 0  
 3 6 2 1 7

skewed