

# Lecture 19

## Topics:

### I) Estimation continued.

The estimation examples we have been looking at are a bit trivial. Here are a few more interesting ones.

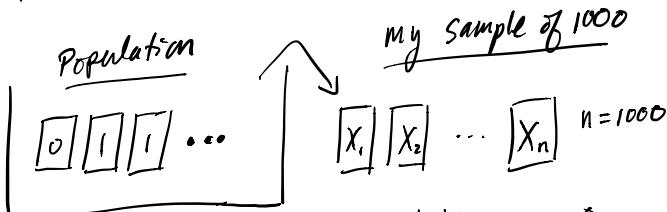
## Estimating a population Proportion

e.g. I just interviewed 1000 UC Davis students & only 34% said they liked statistics.

How reliable is this estimate? Could it plausibly be off by as much as  $\pm 20\%$  (which would mean it could still be true that more than half UC Davis students like statistics : )

Let's analyze the MSE of the estimate.

Here is a box model picture of what we did:



A ticket for each UCD student.

if they like stats

if not

Let  $p$  be the proportion of "likes" in the box

Each  $X_i$  is 1 or 0

Let  $\hat{p}$  be the proportion of likes in the sample

(1)

A couple things to notice.

(2)

1) The 34% I obtained is related to  $\hat{p}$  by

$$0.34 = \hat{p}.$$

2) Since each  $X_i$  is 0 or 1

$$X_1 + \dots + X_n = \underbrace{340}_{\text{the number of likes in my sample.}}$$

$$3) \hat{p} = \frac{X_1 + \dots + X_n}{n} = \bar{X}.$$

Let's find the  $E$  & var each  $X_i$ :

$x$	$P_{X_i}(x)$	$x P_{X_i}(x)$	$x^2 P_{X_i}(x)$
1	$p$	$p$	0
0	$1-p$	0	0

$$\therefore E(X_i) = p$$

$$\text{var}(X_i) = E(X_i^2) - (E(X_i))^2$$

$$= p - p^2 = p(1-p).$$

Let's find the bias, var & MSE of  $\hat{p}$ .

$$\text{bias}(\hat{p}) = E(\hat{p} - p)$$

$$= E(\hat{p}) - p$$

$$= E\left(\frac{X_1 + \dots + X_n}{n}\right) - p$$

$$= \frac{E(X_1) + \dots + E(X_n)}{n} - p$$

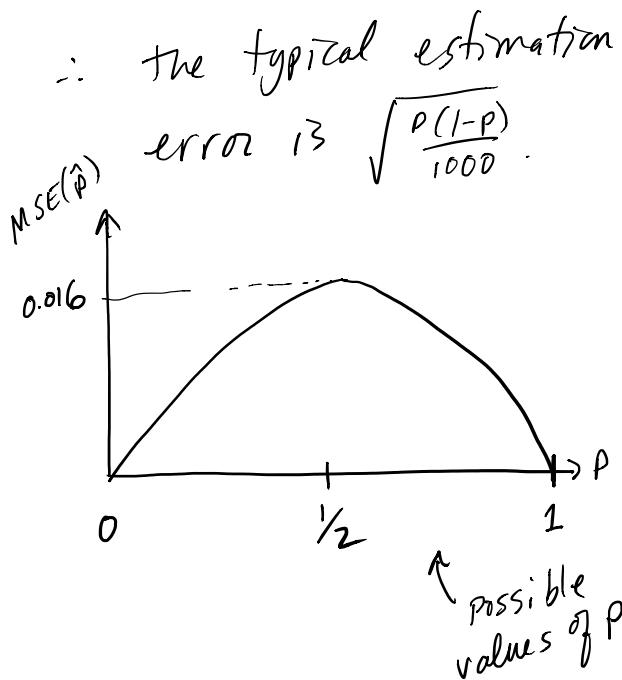
$$= \frac{p + \dots + p}{n} - p$$

$$= \frac{np}{n} - p = 0 \quad \begin{matrix} \leftarrow \\ \text{zero bias so } \hat{p} \end{matrix} \quad \begin{matrix} \text{is accurate "on average."} \\ \text{is accurate "on average."} \end{matrix}$$

$$\begin{aligned}\text{var}(\hat{p}) &= \text{var}\left(\frac{X_1 + \dots + X_n}{n}\right) \\ &= \frac{\text{var}(X_1) + \dots + \text{var}(X_n)}{n^2} \quad \text{since } X_i \text{'s are indep!} \\ &= \frac{p(1-p) + \dots + p(1-p)}{n^2} \\ &= \frac{n p(1-p)}{n^2} = \frac{p(1-p)}{n} = \frac{p(1-p)}{1000}.\end{aligned}$$

since  $n=1000$   
in this case.

$$\begin{aligned}\therefore \text{MSE}(\hat{p}) &= \text{var}(\hat{p}) + (\text{bias}(\hat{p}))^2 \\ &= \frac{p(1-p)}{1000} + 0\end{aligned}$$



$\therefore$  if 20% listed stats  $\hat{p}$  would be typically wrong by about  $\pm \sqrt{\frac{2(1-2)}{1000}} = \pm 0.12$

Here are some possible values of  $p$  with the corresponding typical  $\hat{p}$  estimation errors

Possible true $p$	typical est error for $\hat{p}$
0.30	0.0145
0.40	0.0155
0.50	0.0158
0.60	0.0155
0.70	0.0145

So if  $p=0.5$  we would expect  $\hat{p}$  to be around

$$0.5 \pm \underbrace{0.0158}_{\text{typical discrepancy}} \text{ in } \hat{p}$$

our  $\hat{p}$  is not like this at all so it is not plausible that  $p=0.5$ .

### Estimating the difference of two population averages

e.g. Do Trump supporters make more money (on average) than Clinton supporters?

Sample 200 Trump supporters

Got \$41671 as the average income, per year, of these 200

Sample 150 Clinton supporters

Got \$44250 as the average income, per year, of these 150.

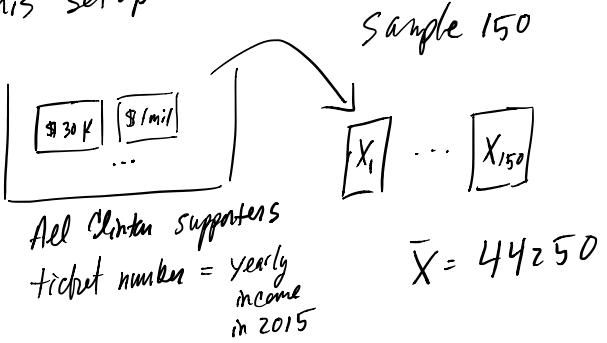
The difference

$$44250 - 41671 = 2579$$

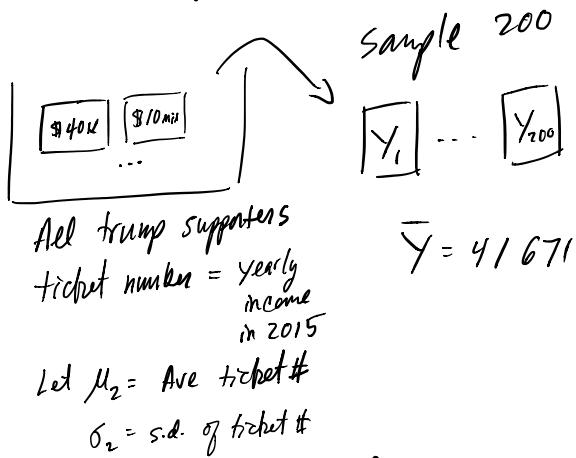
Suggest Clinton supporters make more. Is this real?

Find the MSE of the estimate 2579.

we need two boxes to describe  
this setup. (4)



Let  $\mu_1$  = Ave ticket #  
 $\sigma_1$  = s.d. of ticket #



Our estimate of the difference  $\mu_1 - \mu_2$

$$\text{is } \bar{X} - \bar{Y} = 2579.$$

Let's find  $MSE(\bar{X} - \bar{Y})$  to see how far wrong 2579 could be.

$$\begin{aligned} \text{bias}(\bar{X} - \bar{Y}) &= E\left(\bar{X} - \bar{Y} - (\mu_1 - \mu_2)\right) \\ &\stackrel{\text{est}}{=} E(\bar{X}) - E(\bar{Y}) - (\mu_1 - \mu_2) \\ &= E(\bar{X}) - E(\bar{Y}) - (\mu_1 - \mu_2) \\ &\stackrel{\text{MF1}}{=} \mu_1 - \mu_2 \quad \text{both by MF1} \\ &= 0 \end{aligned}$$

$$\text{var}(\bar{X} - \bar{Y}) = \text{var}(\bar{X}) + \text{var}(\bar{Y}) \quad (5)$$

since  $\bar{X}$  &  $\bar{Y}$  are indep.

$$\text{var}(\bar{X}) = \frac{\sigma_1^2}{150} + \frac{\sigma_2^2}{200}$$

$$MSE(\bar{X} - \bar{Y}) = \frac{\sigma_1^2}{150} + \frac{\sigma_2^2}{200}$$

Suppose the Bureau of labor stats says the overall s.d. of income is about \$30K.

∴ if we assume both  $\sigma_1$  &  $\sigma_2$  are about \$130K then

$$\sqrt{MSE(\bar{X} - \bar{Y})} = \sqrt{\frac{30000^2}{150} + \frac{30000^2}{200}} \\ = 3098$$

∴ we got  $\bar{X} - \bar{Y} = 2579$  but it could be wrong by about  $\pm 3098$  (assuming  $\sigma_1 = \sigma_2 =$  Bureau of labor stats value).

∴ the fact that we observed a difference in ave income b/w the two samples this difference could have been due to sampling fluctuation, also called sampling variability.