

Statistical Inference gives us the ability to answer specific questions about a population.

A statistic is a summary of a sample:  $\bar{X}, \hat{p}, s$   
 Statistics represent parameters which are aspects of the population:  $\mu, p, \sigma$ .

Population (for some variable)

Some mean  $\mu$   
 Some std dev  $\sigma$  } Presumed unknown

Sample

$X_1, X_2, \dots, X_n$

Some mean  $\bar{X}$  (or  $\hat{p}$ )  
 Some std dev  $s$  (or  $\sqrt{\hat{p}(1-\hat{p})}$ )

→ We know the dist of  $\bar{X}$ , sort of

↑ estimation of  $s$

Central Limit Theorem (CLT)

Barring some conditions

Ex:  $n$  should be large ( $n \geq 50$ )

Need a pop without very extreme values

$\bar{X}$  is approximately distributed as

$\bar{X} \sim \text{Normal}(\text{mean} = \mu, \text{std dev} = \frac{\sigma}{\sqrt{n}})$

↑  
pop mean

↑  
pop std dev

Confidence Interval: What is the value of the population mean?

For UAW: sent a survey asking members

↑ Union that represents our TAs, grader and academic student employees

"Are you rent burdened?"

$N = 1,300$   $\hat{p} = 80\% = .8$

→ Spending > 30% of income on rent

3 distributions at play  
Population Distribution

→ dist of entire pop

Sample distribution

→ dist of  $X_1, X_2, \dots, X_n$

→ Dist of  $\hat{p}$

Sampling Dist

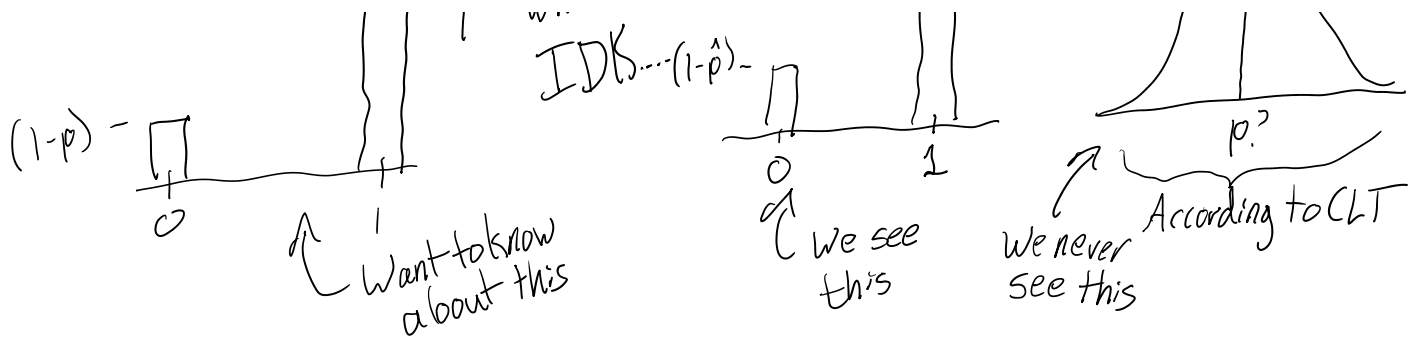
$\sqrt{\frac{p(1-p)}{n}}$

$\Pi - p$

what is  $p$ ?  
 IDs...  $(1-\hat{p}) - \Pi$

$\Pi - \hat{p} = .8$





For a proportion the CLT requires

$n$  large ( $n \geq 50$ ) ✓

$$n \times p \quad n \times (1-p) \geq 10$$

estimate

$$n \times \hat{p} \quad n \times (1-\hat{p}) \geq 10 \quad \checkmark$$

$$1300 \times .8, 1300 \times .2$$

$\hat{p}$  follows a dist  $N(p, \sqrt{\frac{p(1-p)}{n}})$

Estimate  $p$  we use  $\hat{p}$ , but  $\hat{p}$  is wrong

We know  $\hat{p}$  should be close

Thanks to CLT, we know how it should deviate from  $p$ .

## Confidence Interval

Best estimate  $\pm$  Range based on the variability of the sampling dist (via CLT)

$$\hat{p} \pm Z \times \text{Std Error of } \hat{p}$$

Based on standard normal distribution

Std error is like a standard deviation but for a statistic

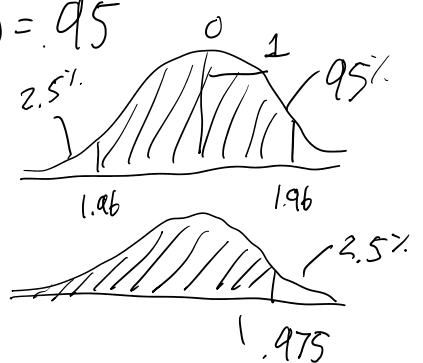
$$\hat{p} \pm Z \cdot \sqrt{\frac{p(1-p)}{n}} \approx \hat{p} \pm Z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

???

$Z$  is based on the level of confidence

I want a 95% CI, so I want a  $Z$  such that  $P(-Z < z < Z) = .95$

$$Z = 1.96$$



$$\hat{p} = .8 \quad n = 1300$$

$$\hat{p} \pm Z_{.975} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$.8 \pm 1.96 \cdot \sqrt{\frac{.8 \cdot (1-.8)}{1300}}$$

std error

$$1.271 \cdot .01271$$

$$\underbrace{1.271^2}_{\text{Margin of error}} = .01271$$

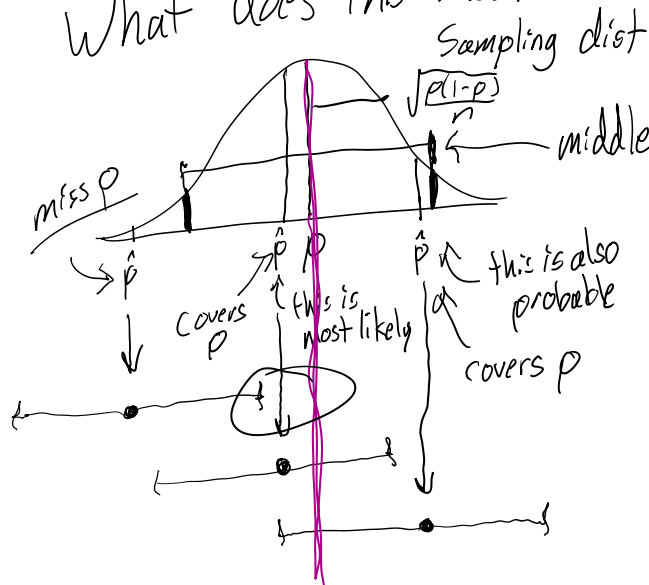
Margin of error  
 $\approx 2.491\%$

1.975

95% Confidence Interval of  $(77.51\%; 82.49\%)$

We are 95% confident that the proportion UAW members at UW facing rent burden is between 77.51% and 82.49%.

What does this mean?



middle 95% is  $z_{.975} \sqrt{\frac{p(1-p)}{n}}$  wide  
 $\approx MoE$

A specific CI may or may not cover the true mean.

We never know which situation were in... why is this useful

If  $\hat{p}$  comes from the middle 95% of the sampling dist } happens 95% of the time  
 $p$  will be in  $\hat{p} \pm z_{.975} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

If  $\hat{p}$  comes from outside the middle 95% } happens 5% of the time  
 $p$  is not in  $\hat{p} \pm z_{.975} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

When we say "We are 95% confident" we're saying

"The process we used to construct this interval should cover true mean 95% of the time".

I don't know  $\mu$  (or  $p$ ) specifically, but the process I used should produce an interval that covers it 95% of the time.

should produce an interval  
time.

Does not mean there a 95% chance  $\mu$  is between (77.5%; 82.5%)