

# Lecture 21: Chapter 9, Section 1

## Inference for Categorical Variable: Confidence Intervals

---

- 3 Forms of Inference
- Probability vs. Confidence
- Constructing Confidence Interval
- Sample Size; Level of Confidence



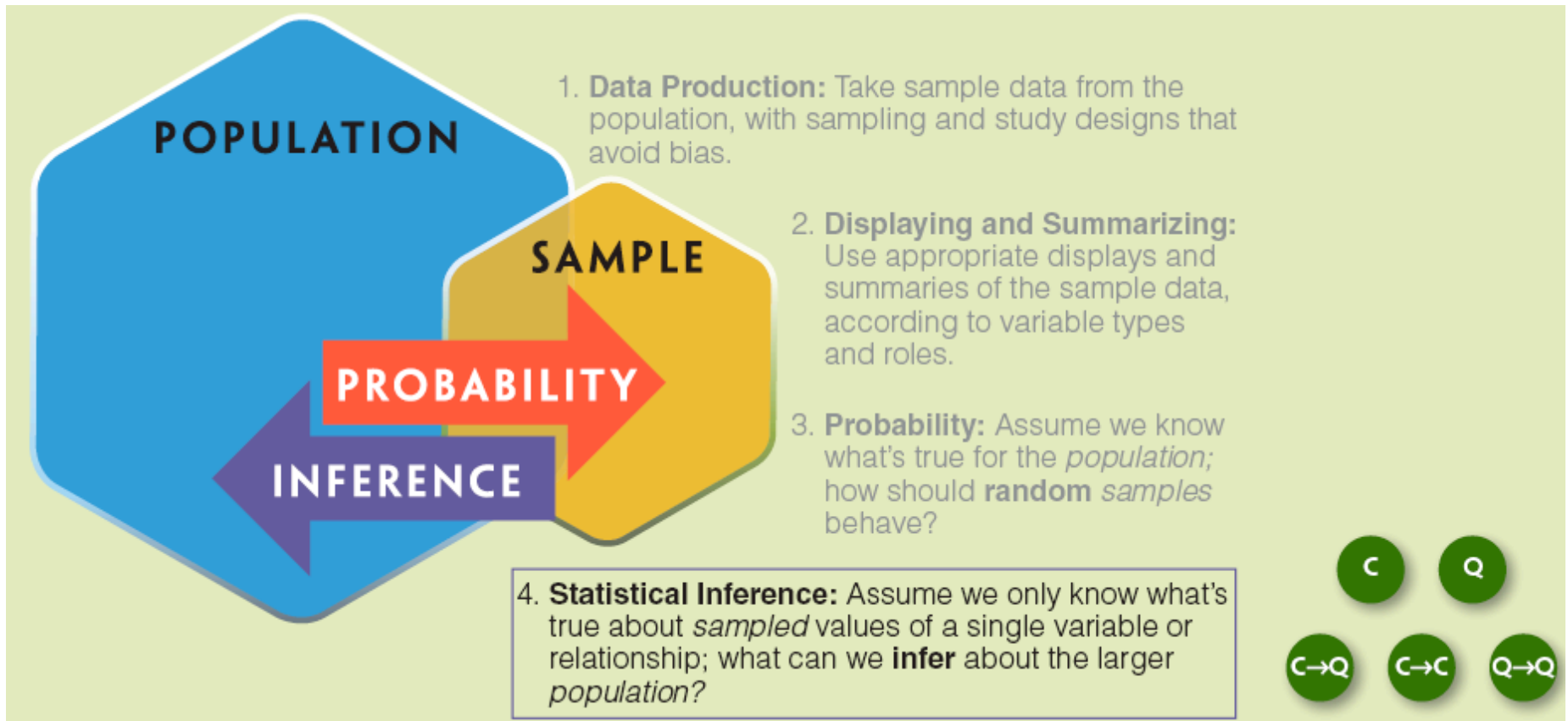
# Looking Back: *Review*

---

## □ 4 Stages of Statistics

- Data Production (discussed in Lectures 1-4)
- Displaying and Summarizing (Lectures 5-12)
- Probability (discussed in Lectures 13-20)
- Statistical Inference
  - 1 categorical
  - 1 quantitative
  - categorical and quantitative
  - 2 categorical
  - 2 quantitative

# Four Processes of Statistics



## Summarizing Categorical Sample Data (*Review*)

---

What proportion of **sampled** students ate breakfast the day of the survey?  $\hat{p} = \frac{X}{n} = \frac{246}{446} = 0.55$

**Looking Back:** In Part 2, we summarized **sample** data for single variables or relationships.

**Looking Ahead:** In Part 4, our goal is to go beyond sample data and draw conclusions about the larger **population** from which the sample was obtained.



## Three Types of Inference Problem

---

*In a sample of 446 students, 0.55 ate breakfast.*

1. What is our **best guess** for the population proportion of students who eat breakfast?

**Point Estimate**

2. What **interval** should contain the population proportion of students who eat breakfast?

**Confidence Interval**

3. Is the population proportion of students who eat breakfast **more than half (50%)**?

**Hypothesis Test**

## Behavior of Sample Proportion (*Review*)

---

For **random** sample of size  $n$  from population with  $p$  in category of interest, sample proportion  $\hat{p} = \frac{X}{n}$  has

- **mean**  $p$

→  $\hat{p}$  is *unbiased estimator* of  $p$

(sample must be **random**)



## Example: *Checking if Estimator is Unbiased*

---

- **Background:** Survey produced sample proportion of intro stat students (various ages and times of day) at a university who'd eaten breakfast.
- **Questions:**
  - Is the sample representative of *all* college students? All students at that university?
  - Were the values of the variable (breakfast or not) recorded without bias?
- **Responses:**
  - Differences among college cafeterias, etc. →

---

  - Question not sensitive → \_\_\_\_\_

## Example: *Point Estimate for $p$*

---

- **Background:** In a representative sample of students, 0.55 ate breakfast.
- **Question:** What is our best guess for the proportion of all students at that university who eat breakfast?
- **Response:**  $\hat{p}$  unbiased estimator for  $p \rightarrow$   
\_\_\_\_\_ is best guess for  $p$





## Example: *Point Estimate Inadequate*

---

- **Background:** Our best guess for  $p$ , population proportion eating breakfast, is sample proportion 0.55.
- **Questions:**
  - Are we pretty sure the population proportion is 0.55?
  - By approximately what amount is our guess “off”?
  - Are we pretty sure population proportion is  $> 0.50$ ?
- **Responses:**
  - \_\_\_\_\_
  - \_\_\_\_\_
  - \_\_\_\_\_



# Beyond a Point Estimate

---

Sample proportion from unbiased sample is best estimate for population proportion.

*Looking Ahead: For point estimate we don't need **sample size** or info about **spread**. These are required for **confidence intervals** and hypothesis tests, to quantify how good our point estimate is.*



# Probability vs. Confidence

---

- **Probability:** given population proportion, how does sample proportion behave?
- **Confidence:** given sample proportion, what is a range of plausible values for population proportion?

## Example: *Probability Statement*

- **Background:** If students pick numbers from 1 to 20 at random,  $p=0.05$  should pick #7. For  $n=400$ ,  $\hat{p}$  has
  - mean 0.05
  - standard deviation  $\sqrt{\frac{0.05(1-0.05)}{400}} = 0.01$
  - shape approximately normal
- **Question:** What does the “95” part of the 68-95-99.7 Rule tell us about  $\hat{p}$ ?
- **Response:** Probability is approximately 0.95 that  $\hat{p}$  falls within \_\_\_\_\_ of \_\_\_\_\_.

*Looking Ahead: This statement about **sample proportion** is correct but not very useful for practical purposes. In most real-life problems, we want to draw conclusions about an **unknown population proportion**.*



## **Example:** *How Far is One from the Other?*

---

- **Background:** An instructor can say about his/her position in the classroom:  
“I’m within 10 feet of this particular student.”
- **Question:** What can be said about where that student is in relation to the instructor?
- **Response:**



# Definitions

---

**Margin of Error:** *Distance* around a sample statistic, within which we have reason to believe the corresponding parameter falls.

A common margin of error is 2 s.d.s.

**Confidence Interval** for parameter: *Interval* within which we have reason to believe the parameter falls = **range of plausible values**

A common confidence interval is sample statistic plus or minus 2 s.d.s.

***A Closer Look:*** *A parameter is **not** a R.V. It does **not** obey the laws of probability so we must use the word “confidence”.*

## Example: *Confidence Interval for $p$*

---

- **Background:**  $30/400=0.075$  students picked #7 “at random” from 1 to 20. Let’s assume sample proportion for  $n=400$  has s.d.  $0.01$ .
- **Question:** What can we claim about population proportion  $p$  picking #7?
- **Response:** We’re pretty sure  $p$  is

*Looking Back:* In Part I, we learned about biased samples. The data suggest  $p > 0.05$ : students were apparently biased in favor of #7. Their selections were **haphazard**, not random. If sampling individuals or assigning them to experimental treatments is **not randomized**, then we produce a confidence interval that is **not centered** at  $p$ .



## Level of Confidence Corresponds to Multiplier

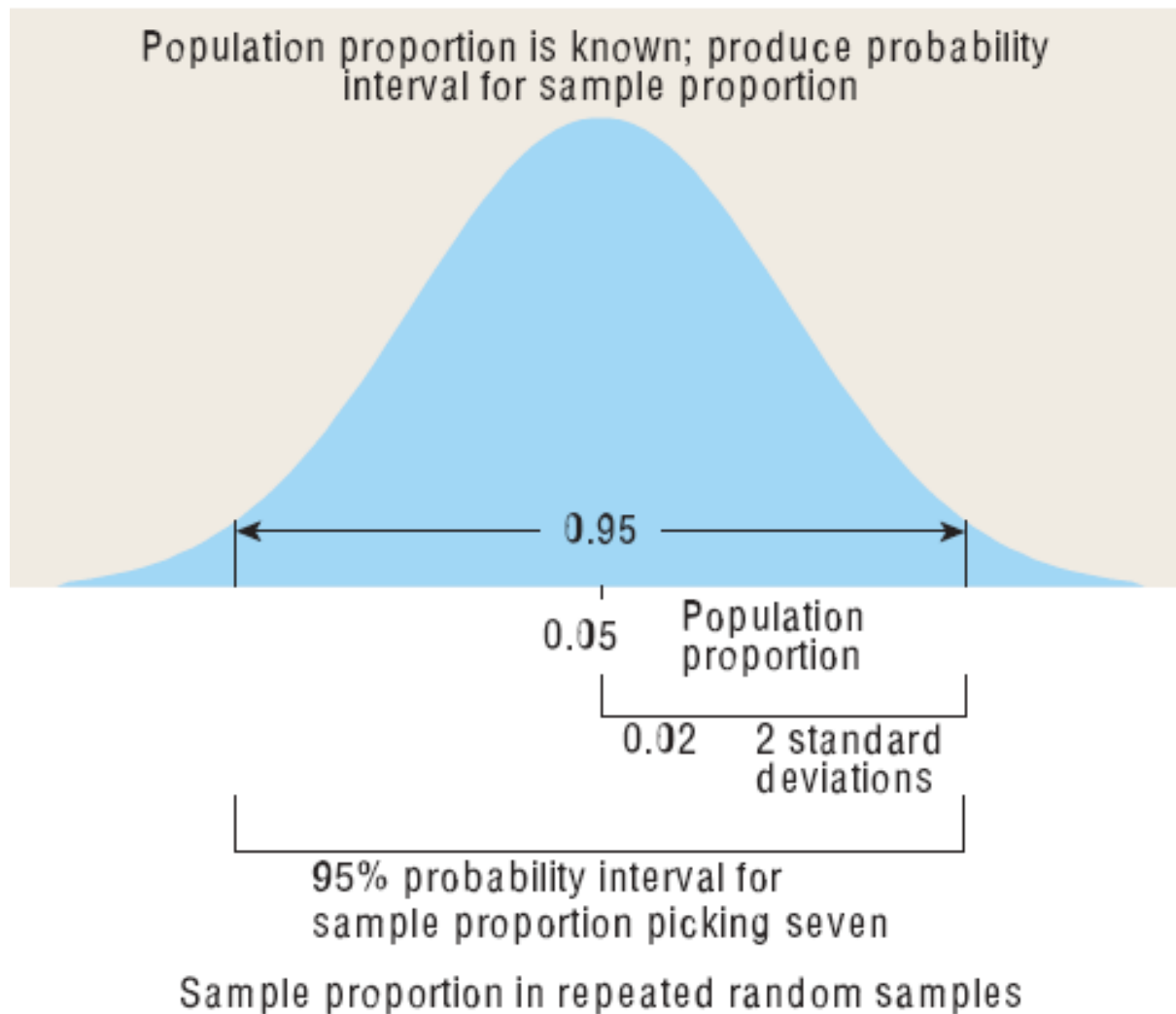
---

By “pretty sure”, we mean “95% confident”, because 95% is the probability of sample proportion within 2 s.d.s of  $p$  (for large enough  $n$ ).

*Looking Back: Our probability statement claimed sample proportion should fall within 2 s.d.s of population proportion. Now, the inference statement claims population proportion should be within 2 s.d.s of sample proportion.*



# Probability Interval for $\hat{p}$ Picking #7



# Confidence Interval for $p$ Picking #7

*We do not sketch a curve showing probabilities for population proportion because it is **not a random variable**.*

Measure sample proportion; produce confidence interval for unknown population proportion

Sample proportion in one random sample

0.02 = margin of error

95% confidence interval for population proportion picking seven  
Unknown population proportion

***A Closer Look:***  
*How do we know the margin of error?*

## Behavior of Sample Proportion (*Review*)

---

For random sample of size  $n$  from population with  $p$  in category of interest, sample proportion  $\hat{p}$  has

- mean  $p$

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

We do inference because  $p$  is unknown; how can we know the standard deviation, which involves  $p$ ?

# Definition

---

**Standard error:** estimated standard deviation of a sampling distribution.

---

We estimate standard deviation of  $\hat{p}$  with standard error  $\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$ .

***Looking Ahead:** In many situations throughout inference, when needed information about the **population** is unknown, we substitute known information about the **sample**.*

# Definition

---

**95% confidence interval for  $p$ : (approx.)**

$$\hat{p} \pm 2\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

estimate = sample proportion

standard error

margin of error =  
2 standard errors

**95% confidence interval  
for population proportion**

# Confidence Interval Formula: Conditions

---

**95% confidence interval for  $p$ :** (approx.)

$$\hat{p} \pm 2\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

- Sample must be unbiased  
(otherwise interval is not really centered at  $\hat{p}$  )
- $n$  must be large enough so  $\hat{p}$  is approx. normal  
(otherwise multiplier 2 from 68-95-99.7 Rule is incorrect)
- Population size must be at least  $10n$   
(otherwise formula for s.d., which requires independence, is incorrect)

## Conditions for Normality in Confidence Interval

---

Multiplier 2 from normal dist. approximately correct if  $np$  and  $n(1-p)$  both at least 10.

But  $p$  is unknown so substitute  $\hat{p}$ :

Require

$$n\hat{p} = nX/n = X \geq 10$$

$$n(1 - \hat{p}) = n - nX/n = n - X \geq 10$$

Sample count in ( $X$ ) and out ( $n-X$ ) of category of interest should both be at least 10.



## Example: *Checking Sample Size*

---

- **Background:**  $30/400=0.075$  students picked #7 “at random” from 1 to 20.
- **Question:** Do the data satisfy requirement for approximate normality of sample proportion?
- **Response:**





## Example: *Checking Population Size*

---

- **Background:** To draw conclusions about criminal histories of a city's 750 bus drivers, a random sample of 100 drivers was used.
- **Question:** Is there approximate independence in spite of sampling without replacement, so formula for standard error is accurate?
- **Response:**

## Example: *Revisiting Original Question*

---

- **Background:** In sample of 446 college students, 246 (proportion 0.55) ate breakfast.
- **Question:** Assuming sample is representative, what interval should contain proportion of all students at that university who eat breakfast?
- **Response:** Approx. 95% confidence interval for  $p$  is

### *Looking Back:*

*Earlier we wondered if a majority of students eat breakfast. The interval suggests this is the case, since it is entirely above 0.50.*

$$\hat{p} \pm 2\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

## Example: *Role of Sample Size*

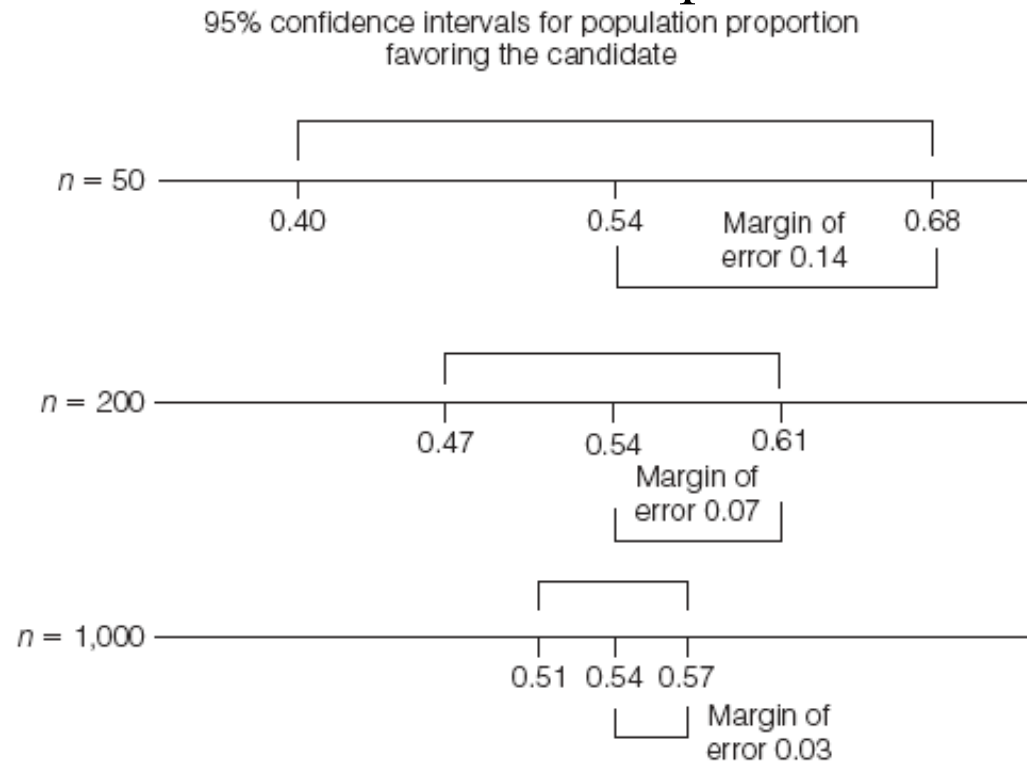
- **Background:** 95% confidence intervals based on sample proportion 0.54 from various sample sizes:

Sample Size $n$	Standard Error of $\hat{p}$	Margin of Error	95% Confidence Interval
50	$\sqrt{\frac{0.54(1 - 0.54)}{50}} = 0.070$	$2(0.070) = 0.14$	(0.40, 0.68)
200	$\sqrt{\frac{0.54(1 - 0.54)}{200}} = 0.035$	$2(0.035) = 0.07$	(0.47, 0.61)
1,000	$\sqrt{\frac{0.54(1 - 0.54)}{1,000}} = 0.016$	$2(0.016) = 0.03$	(0.51, 0.57)

- **Question:** What happens as  $n$  increases?
- **Response:**

# Example: *A Common Margin of Error*

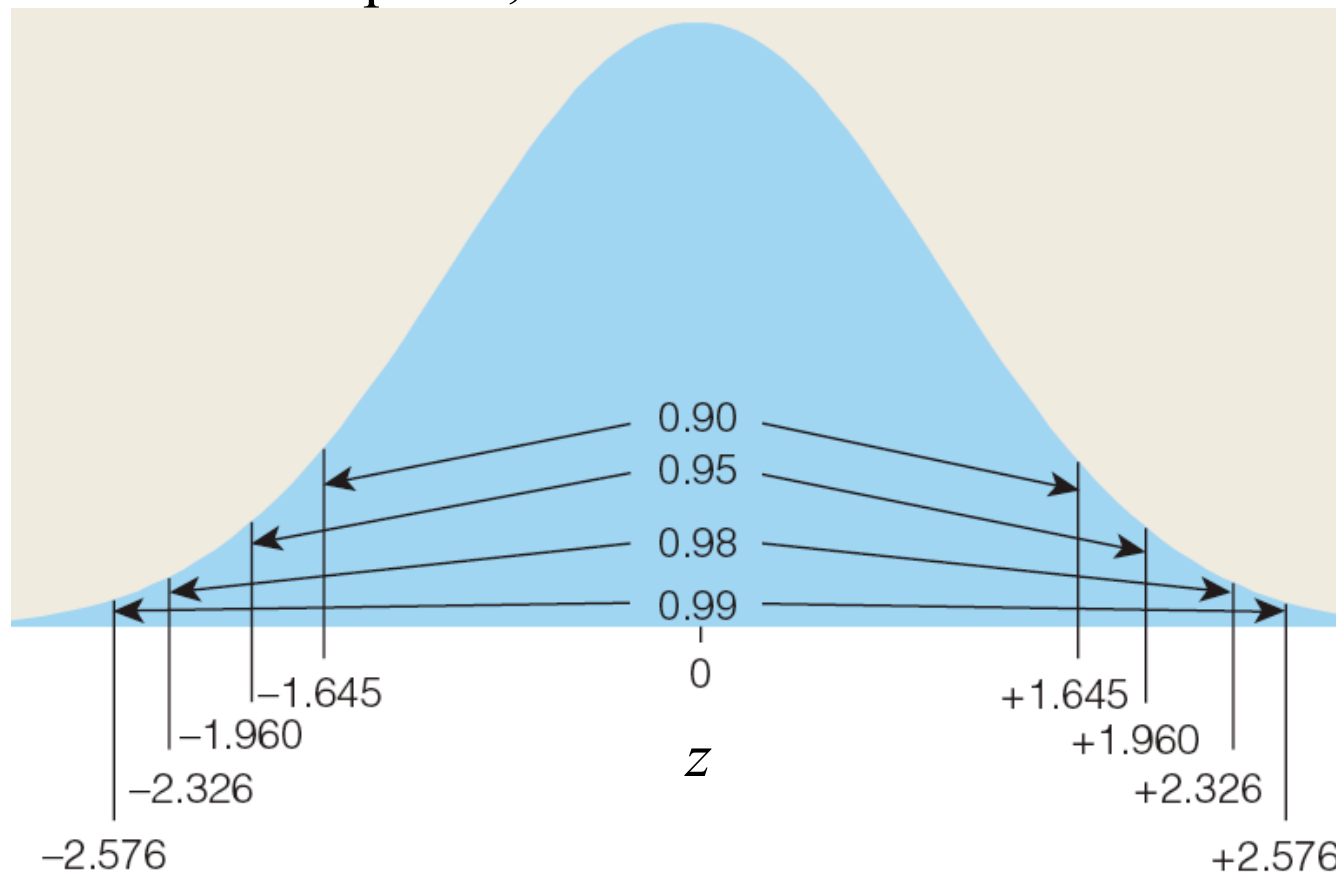
- **Background:** Pollsters most often report a 3% error margin.



- **Question:** What is the most common sample size for polls?
- **Response:** Approximately \_\_\_\_\_.

# Other Levels of Confidence

Confidence level 95% uses multiplier 2. Other levels use other multipliers, based on normal curve:



## Other Levels of Confidence

---

Confidence level 95% uses multiplier 2. Other levels use other multipliers, based on normal curve.

More precise multiplier for 95% is 1.96 instead of 2.

Level	Multiplier
90%	1.645
95%	1.960
98%	2.326
99%	2.576

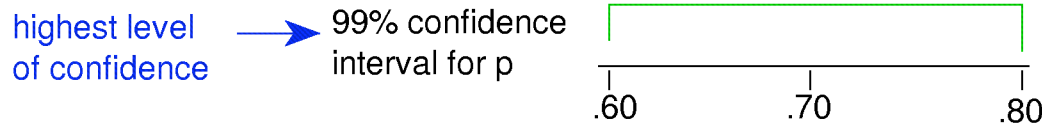
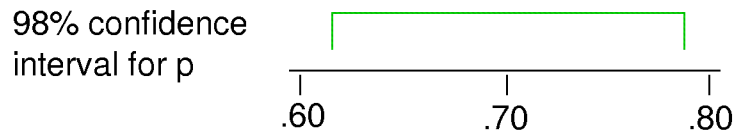
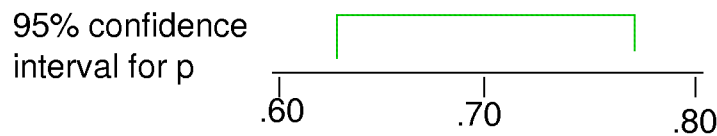
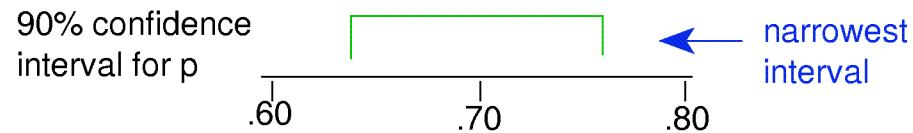
## Example: *Other Levels of Confidence*

---

- **Background:** Of 108 students in committed relationships, 0.70 said they took comfort by sniffing out-of-town partner's clothing. Standard error can be found to be 0.04.
- **Question:** How do 90%, 95%, 98%, 99% confidence intervals compare?
- **Response:**
  - 90% C.I. is \_\_\_\_\_  $= (0.63, 0.77)$
  - 95% C.I. is \_\_\_\_\_  $= (0.62, 0.78)$
  - 98% C.I. is \_\_\_\_\_  $= (0.61, 0.79)$
  - 99% C.I. is \_\_\_\_\_  $= (0.60, 0.80)$

# Example: *Other Levels of Confidence*

Intervals get \_\_\_\_\_ as confidence level increases:







## Confidence Interval and Long-Run Behavior

---

Repeatedly set up 95% confidence interval for proportion of heads, based on 20 coinflips.

In the long run, 95% of the intervals should contain population proportion of heads, 0.5.

# Confidence Interval and Long-Run Behavior

20 coin flips

TTTTTHTHTTHHTHTHTTTHH  
proportion of heads  $9/20 = .45$

HTTHHTHTTTHTHTTTTHHT  
proportion of heads  $8/20 = .40$

TTTHHHHHHTHTHTHTTTT  
proportion of heads  $12/20 = .60$

THTHHHTHTHTHTHTHTHTT  
proportion of heads  $15/20 = .75$

- repeated
- flips of 20
- coins

95% confidence interval

.15 .20 .25 .30 .35 .40 .45 .50 .55 .60 .65 .70 .75 .80 .85

.15 .20 .25 .30 .35 .40 .45 .50 .55 .60 .65 .70 .75 .80 .85

.15 .20 .25 .30 .35 .40 .45 .50 .55 .60 .65 .70 .75 .80 .85

.15 .20 .25 .30 .35 .40 .45 .50 .55 .60 .65 .70 .75 .80 .85

.15 .20 .25 .30 .35 .40 .45 .50 .55 .60 .65 .70 .75 .80 .85

in the long run  
95% of intervals contain  $p = .50$   
5% of intervals do not contain .50

## Example: *Confidence in the Long Run*

---

- **Background:** “President-elect Barack Obama's campaign strategists weren't the only ones vindicated Tuesday. Pollsters came out looking pretty good, too. Of 27 polls of Pennsylvania voters released in the campaign's final two weeks, only seven missed Obama's 10.3-point victory by more than their margins of error. Obama's national victory of about 6 points was within the error margins of 16 of the 21 national polls released in the final week.”
- **Question:** Should pollsters be pleased with success rates of  $20/27=16/21=75\%$  ?
- **Response:**

[pittsburghlive.com/x/pittsburghtrib/news/cityregion/s\\_597288.html](http://pittsburghlive.com/x/pittsburghtrib/news/cityregion/s_597288.html)



# Lecture Summary

## *(Inference for Proportions: Confidence Interval)*

---

- 3 forms of inference; focus on confidence interval
- Probability vs. confidence
- Constructing confidence interval
  - Margin of error based on standard error
  - Conditions
- Role of sample size
- Confidence at other levels
- Confidence interval in the long run