

# Econ 103 – Statistics for Economists

## Chapter 6 and 7: Confidence Intervals

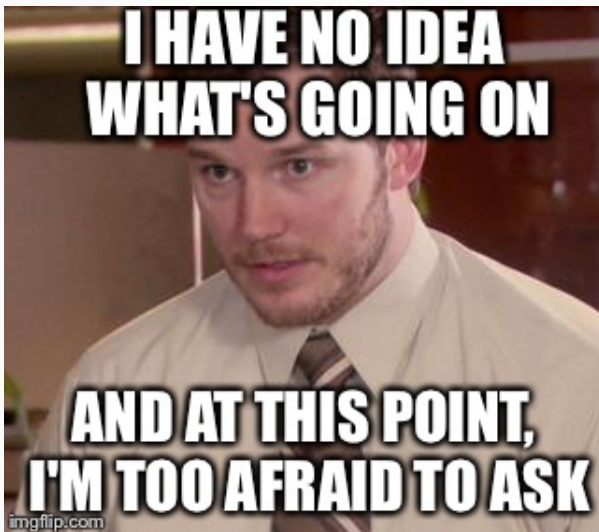
---

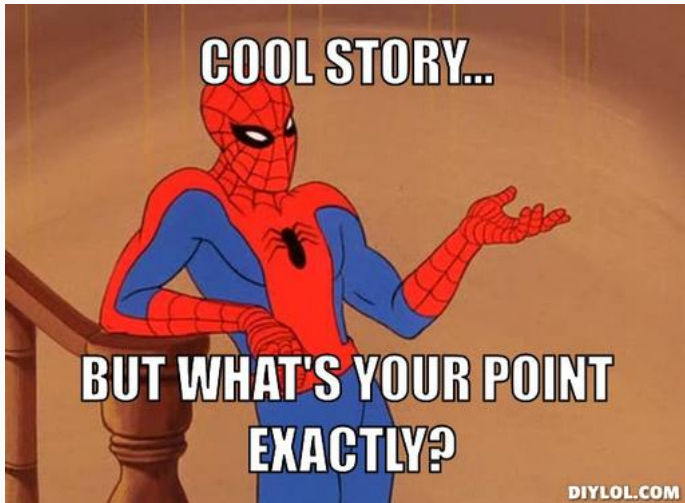
Mallick Hossain

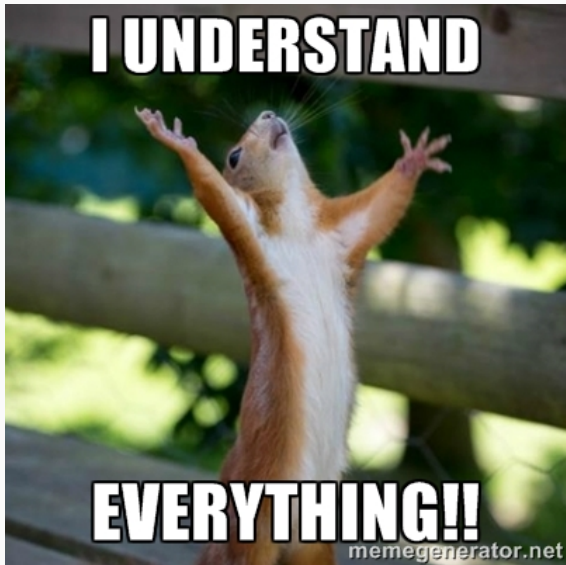
University of Pennsylvania

3 Students

---







# What's the Point?

The goal is to get you closer to the squirrel (or at least Spider-Man)

## Recap and Motivation

---

## What We've Done So Far (Theory Side)

- We spent the past few weeks covering discrete and continuous random variables



## What We've Done So Far (Theory Side)

- We spent the past few weeks covering discrete and continuous random variables
  - You should be very comfortable with each random variable and their associated properties (see random variable handout for a nice (not necessarily exhaustive) summary)

## What We've Done So Far (Theory Side)

- We spent the past few weeks covering discrete and continuous random variables
  - You should be very comfortable with each random variable and their associated properties (see random variable handout for a nice (not necessarily exhaustive) summary)
- We dug into the normal distribution and all of its nice properties

## What We've Done So Far (Theory Side)

- We spent the past few weeks covering discrete and continuous random variables
  - You should be very comfortable with each random variable and their associated properties (see random variable handout for a nice (not necessarily exhaustive) summary)
- We dug into the normal distribution and all of its nice properties
  - The more intuitive the normal RV feels, the easier the rest of the semester will be

## What We've Done So Far (Theory Side)

- We spent the past few weeks covering discrete and continuous random variables
  - You should be very comfortable with each random variable and their associated properties (see random variable handout for a nice (not necessarily exhaustive) summary)
- We dug into the normal distribution and all of its nice properties
  - The more intuitive the normal RV feels, the easier the rest of the semester will be
- Briefly introduced chi-squared, t-, and F-distributions

# What We've Done So Far (Theory Side)

- We spent the past few weeks covering discrete and continuous random variables
  - You should be very comfortable with each random variable and their associated properties (see random variable handout for a nice (not necessarily exhaustive) summary)
- We dug into the normal distribution and all of its nice properties
  - The more intuitive the normal RV feels, the easier the rest of the semester will be
- Briefly introduced chi-squared, t-, and F-distributions
  - You'll see why they are so important today! The wait is over!

## What We've Done So Far (Practical Side)

- Random Sampling:  $X_1, \dots, X_n \sim \text{iid}$
- Use estimator  $\hat{\theta}$  to learn about population parameter  $\theta_0$
- Estimator  $\hat{\theta}$  is a random variable:
  - Distribution of  $\hat{\theta}$  is called *sampling distribution*
  - Bias of an estimator
  - Variance of an estimator
  - Mean-squared Error (MSE) of an estimator
  - Consistency of an Estimator

## Confidence Intervals

What values of  $\theta_0$  are consistent with the data we observed?

## Hypothesis Testing

I think that  $\theta_0 = 0$ . Should I change my mind based on the data?

- Do we expect point estimates to be exactly right?



- Do we expect point estimates to be exactly right?
  - No! As we saw last lecture, our estimate is basically a draw from the distribution of a random variable

# Motivation

- Do we expect point estimates to be exactly right?
  - No! As we saw last lecture, our estimate is basically a draw from the distribution of a random variable
- If we predicted that the S&P 500 would close at \$2150.00 on Monday and it closed at \$2150.88, my point estimate was wrong. Does that mean it's worthless though?

# Motivation

- Do we expect point estimates to be exactly right?
  - No! As we saw last lecture, our estimate is basically a draw from the distribution of a random variable
- If we predicted that the S&P 500 would close at \$2150.00 on Monday and it closed at \$2150.88, my point estimate was wrong. Does that mean it's worthless though?
  - No! It was “close” which can be very informative!

# Motivation

- Do we expect point estimates to be exactly right?
  - No! As we saw last lecture, our estimate is basically a draw from the distribution of a random variable
- If we predicted that the S&P 500 would close at \$2150.00 on Monday and it closed at \$2150.88, my point estimate was wrong. Does that mean it's worthless though?
  - No! It was “close” which can be very informative!
  - Confidence intervals are instrumental in giving us a better idea of what counts as “close.”

## Example

---

## (Above?) Average Joe

Joe is 73 inches tall. Based on a sample of US males aged 20 and over, the Centers for Disease Control (CDC) reported a mean height of about 69 inches in a recent report.

Clearly Joe is taller than the average American male!

Do you agree or disagree?

- (a) Agree
- (b) Disagree
- (c) Not Sure

## Remember: The Sample Mean is Random!

Just because the sample mean is 69 inches it doesn't follow that the population mean is 69 inches!

What Else Should We Consider?

## Remember: The Sample Mean is Random!

Just because the sample mean is 69 inches it doesn't follow that the population mean is 69 inches!

### What Else Should We Consider?

- How big was the sample?



## Remember: The Sample Mean is Random!

Just because the sample mean is 69 inches it doesn't follow that the population mean is 69 inches!

### What Else Should We Consider?

- How big was the sample?
  - If the sample was very small there's a higher chance that it won't be representative of the population as a whole

## Remember: The Sample Mean is Random!

Just because the sample mean is 69 inches it doesn't follow that the population mean is 69 inches!

### What Else Should We Consider?

- How big was the sample?
  - If the sample was very small there's a higher chance that it won't be representative of the population as a whole
  - Why? The variance of the sample mean is *decreasing with sample size* so bigger samples are less noisy.

## Remember: The Sample Mean is Random!

Just because the sample mean is 69 inches it doesn't follow that the population mean is 69 inches!

### What Else Should We Consider?

- How big was the sample?
  - If the sample was very small there's a higher chance that it won't be representative of the population as a whole
  - Why? The variance of the sample mean is *decreasing with sample size* so bigger samples are less noisy.
- How much variability is there in height in the population?

## Remember: The Sample Mean is Random!

Just because the sample mean is 69 inches it doesn't follow that the population mean is 69 inches!

### What Else Should We Consider?

- How big was the sample?
  - If the sample was very small there's a higher chance that it won't be representative of the population as a whole
  - Why? The variance of the sample mean is *decreasing with sample size* so bigger samples are less noisy.
- How much variability is there in height in the population?
  - If everyone is very similar in height, any sample we take will be representative of the population.

## Remember: The Sample Mean is Random!

Just because the sample mean is 69 inches it doesn't follow that the population mean is 69 inches!

### What Else Should We Consider?

- How big was the sample?
  - If the sample was very small there's a higher chance that it won't be representative of the population as a whole
  - Why? The variance of the sample mean is *decreasing with sample size* so bigger samples are less noisy.
- How much variability is there in height in the population?
  - If everyone is very similar in height, any sample we take will be representative of the population.
  - Remember: the variance of the sample mean is *increasing* with the population standard deviation.

# Am I Taller Than The Average American Male?

**Table 1:** Height in inches for Males aged 20 and over (approximate)

Sample Mean	69 inches
Sample Std. Dev.	6 inches
Sample Size	5647
Joe's Height	73 inches

We'll return to this example later.

# Theoretical Example

---

## For Now – Single Population, Normally Distributed

$$X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$$

Later we'll look at more than one population and talk about what happens if Normality doesn't hold.



Suppose  $X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$ . What is the sampling distribution of  $\sqrt{n}(\bar{X}_n - \mu)/\sigma$ ?

- (a)  $N(\mu, \sigma^2)$
- (b)  $N(0, 1)$
- (c)  $N(0, \sigma)$
- (d)  $N(\mu, 1)$
- (e) Not enough information to determine.

## Z-score!

Suppose  $X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$ . From above,

$$\begin{aligned} E[\bar{X}_n] &= \mu \\ \text{Var}(\bar{X}_n) &= \sigma^2/n \\ \Rightarrow SD(\bar{X}_n) &= \sigma/\sqrt{n} \end{aligned}$$

Thus,

$$\sqrt{n}(\bar{X}_n - \mu)/\sigma = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} = \frac{\bar{X}_n - E[\bar{X}_n]}{SD(\bar{X}_n)} \sim N(0, 1)$$

Remember that we call the standard deviation of a sampling distribution the **standard error**, written  $SE$ , so

$$\frac{\bar{X}_n - \mu}{SE(\bar{X}_n)} \sim N(0, 1)$$

# Standard Error vs Standard Deviation

- **Standard Deviation**

- The square root of the variance
- Measures the deviation from the mean

- **Standard Error**

- A specific kind of standard deviation
- This is the standard deviation of the estimator
- For example, if we are estimating the population mean, the standard error tells us how far our estimate is from the actual population mean.

Suppose  $X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$ . What is the approximate value of the following?

$$P\left(-2 \leq \frac{\bar{X}_n - \mu}{SE(\bar{X}_n)} \leq 2\right)$$

Suppose  $X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$ . What is the approximate value of the following?

$$P\left(-2 \leq \frac{\bar{X}_n - \mu}{SE(\bar{X}_n)} \leq 2\right) \approx 0.95$$

## What happens if I rearrange?

$$P\left(-2 \leq \frac{\bar{X}_n - \mu}{SE(\bar{X}_n)} \leq 2\right) = 0.95$$

## What happens if I rearrange?

$$P\left(-2 \leq \frac{\bar{X}_n - \mu}{SE(\bar{X}_n)} \leq 2\right) = 0.95$$

$$P(-2 \cdot SE \leq \bar{X}_n - \mu \leq 2 \cdot SE) = 0.95$$

## What happens if I rearrange?

$$P\left(-2 \leq \frac{\bar{X}_n - \mu}{SE(\bar{X}_n)} \leq 2\right) = 0.95$$

$$P(-2 \cdot SE \leq \bar{X}_n - \mu \leq 2 \cdot SE) = 0.95$$

$$P(-2 \cdot SE - \bar{X}_n \leq -\mu \leq 2 \cdot SE - \bar{X}_n) = 0.95$$



## What happens if I rearrange?

$$P\left(-2 \leq \frac{\bar{X}_n - \mu}{SE(\bar{X}_n)} \leq 2\right) = 0.95$$

$$P(-2 \cdot SE \leq \bar{X}_n - \mu \leq 2 \cdot SE) = 0.95$$

$$P(-2 \cdot SE - \bar{X}_n \leq -\mu \leq 2 \cdot SE - \bar{X}_n) = 0.95$$

$$P(\bar{X}_n - 2 \cdot SE \leq \mu \leq \bar{X}_n + 2 \cdot SE) = 0.95$$

# Confidence Intervals

---

# Confidence Intervals

## Confidence Interval (CI)

A confidence interval is a range  $(A, B)$  constructed from the **sample data** that has a specified probability of containing a **population parameter**:

$$P(A \leq \theta_0 \leq B) = 1 - \alpha$$

# Confidence Intervals

## Confidence Interval (CI)

A confidence interval is a range  $(A, B)$  constructed from the **sample data** that has a specified probability of containing a **population parameter**:

$$P(A \leq \theta_0 \leq B) = 1 - \alpha$$

## Confidence Level

The **specified probability**, typically denoted  $1 - \alpha$ , is called the confidence level. For example, if  $\alpha = 0.05$  then the confidence level is 0.95 or 95%.

# Confidence Interval for Mean of Normal Population

## Confidence Interval for Mean of Normal Population

The interval  $\bar{X}_n \pm 2\sigma/\sqrt{n}$  has approximately 95% probability of containing the population mean  $\mu$ , provided that:

$$X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$$

# Confidence Interval for Mean of Normal Population

## Confidence Interval for Mean of Normal Population

The interval  $\bar{X}_n \pm 2\sigma/\sqrt{n}$  has approximately 95% probability of containing the population mean  $\mu$ , provided that:

$$X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$$

But What Does This Mean?

## Which quantities are random?

Suppose  $X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$ . Which quantities are random variables?

- (a)  $\mu$  only
- (b)  $\sigma$  and  $\mu$
- (c)  $\sigma$  only
- (d)  $\sigma, \mu$  and  $\bar{X}_n$
- (e)  $\bar{X}_n$  only

## Which quantities are random?

Suppose  $X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$ . Which quantities are random variables?

- (a)  $\mu$  only
- (b)  $\sigma$  and  $\mu$
- (c)  $\sigma$  only
- (d)  $\sigma, \mu$  and  $\bar{X}_n$
- (e)  $\bar{X}_n$  only

What does this mean for our confidence intervals?



## Confidence Interval is a Random Variable!

1.  $X_1, \dots, X_n$  are RVs  $\Rightarrow \bar{X}_n$  is a RV (repeated sampling)

## Confidence Interval is a Random Variable!

1.  $X_1, \dots, X_n$  are RVs  $\Rightarrow \bar{X}_n$  is a RV (repeated sampling)
2.  $\mu, \sigma$  and  $n$  are constants

# Confidence Interval is a Random Variable!

1.  $X_1, \dots, X_n$  are RVs  $\Rightarrow \bar{X}_n$  is a RV (repeated sampling)
2.  $\mu, \sigma$  and  $n$  are constants
3. Confidence Interval  $\bar{X}_n \pm 2\sigma/\sqrt{n}$  is also a RV!

# Meaning of Confidence Interval

## Meaning of Confidence Interval

If we sampled many times we'd get many different sample means, each leading to a **different** confidence interval. Approximately 95% of these intervals will contain  $\mu$ .

## Rough Intuition

What values of  $\mu$  are consistent with the data?

## CI for Population Mean: Repeated Sampling

$$X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$$

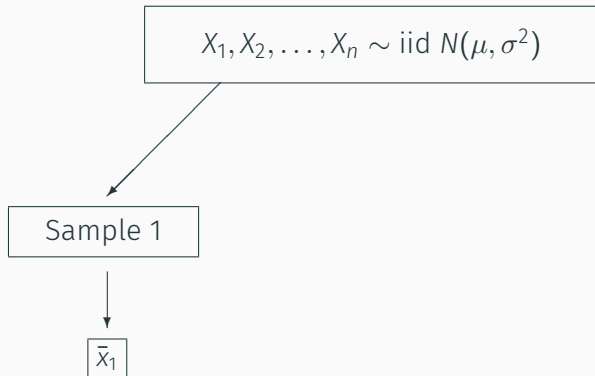
## CI for Population Mean: Repeated Sampling

$$X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$$

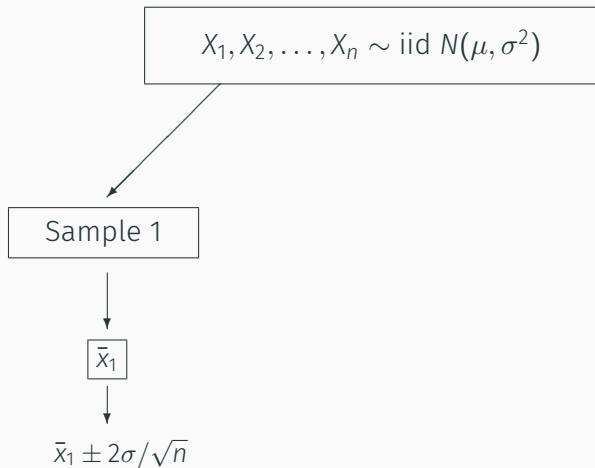


Sample 1

## CI for Population Mean: Repeated Sampling

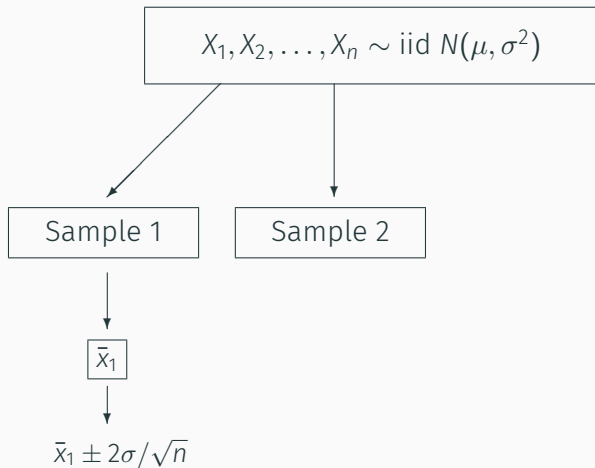


## CI for Population Mean: Repeated Sampling

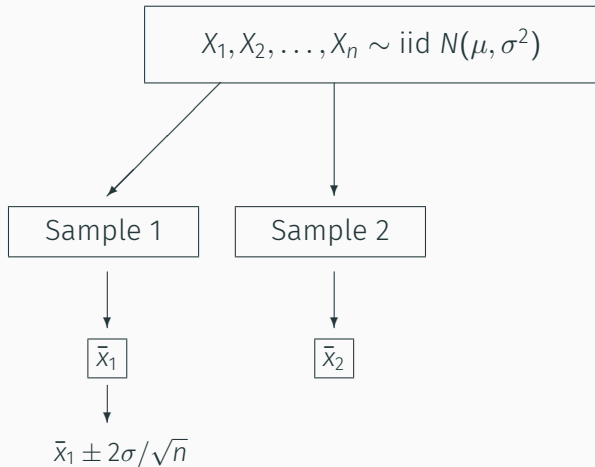




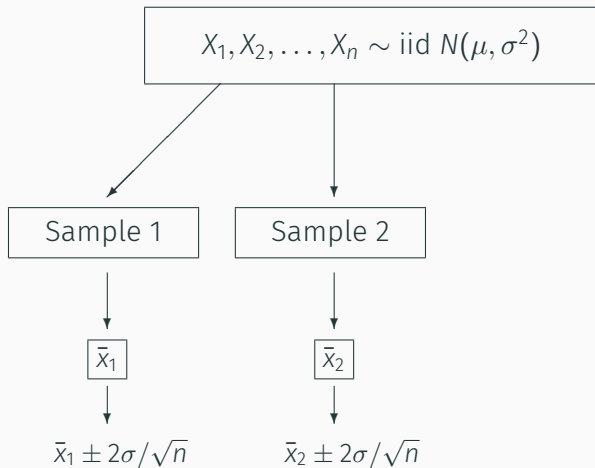
## CI for Population Mean: Repeated Sampling



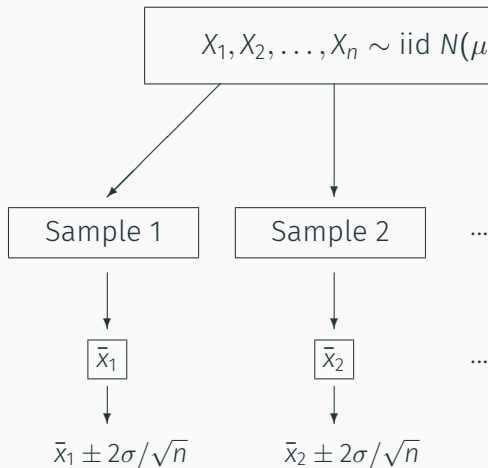
## CI for Population Mean: Repeated Sampling



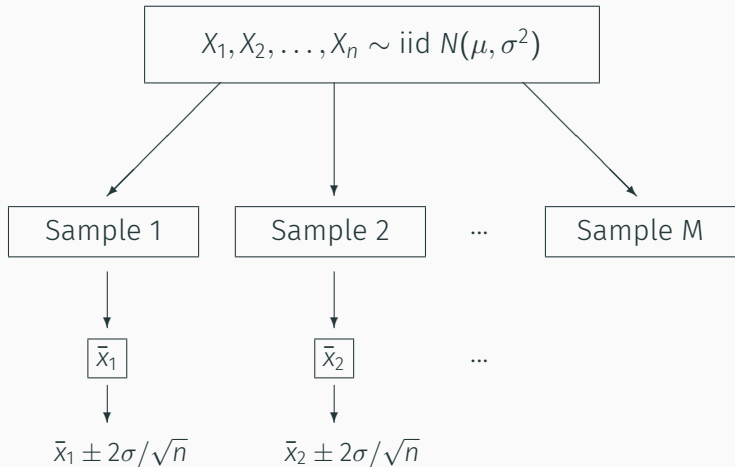
## CI for Population Mean: Repeated Sampling



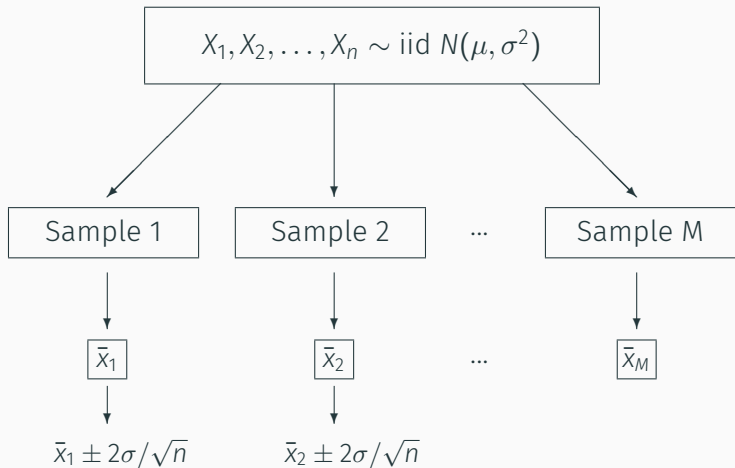
## CI for Population Mean: Repeated Sampling



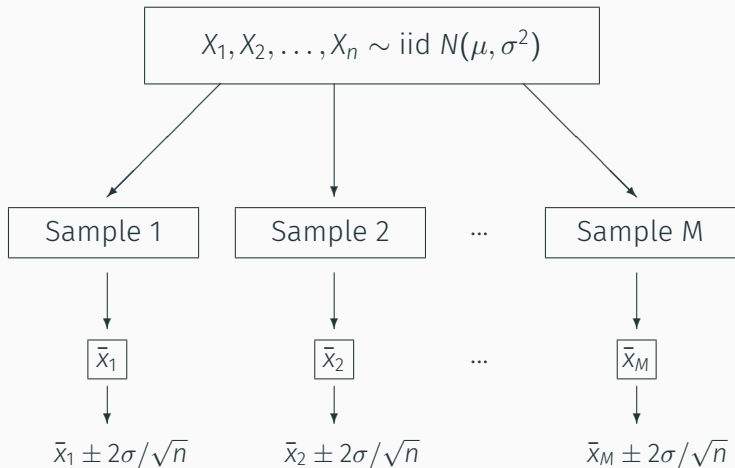
# CI for Population Mean: Repeated Sampling



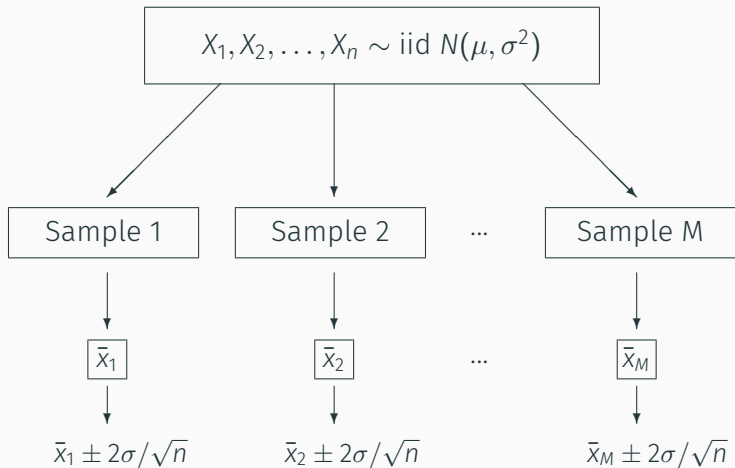
# CI for Population Mean: Repeated Sampling



# CI for Population Mean: Repeated Sampling



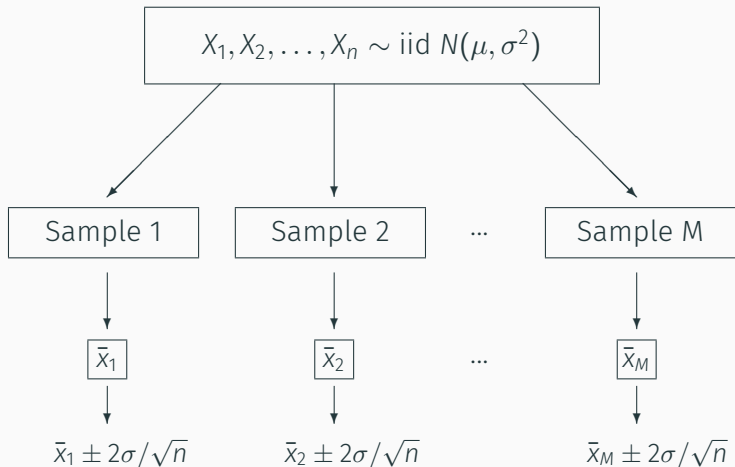
# CI for Population Mean: Repeated Sampling



Repeat  $M$  times  $\rightarrow$  get  $M$  different intervals



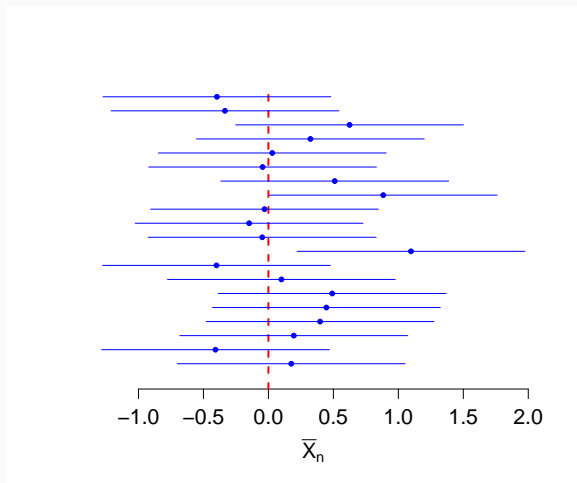
# CI for Population Mean: Repeated Sampling



Repeat  $M$  times  $\rightarrow$  get  $M$  different intervals

Large  $M \Rightarrow$  Approx. 95% of these Intervals Contain  $\mu$

## Simulation Example: $X_1, \dots, X_5 \sim \text{iid } N(0, 1), M = 20$



**Figure 1:** Twenty confidence intervals of the form  $\bar{X}_n \pm 2\sigma/\sqrt{n}$  where  $n = 5$ ,  $\sigma^2 = 1$  and the true population mean is 0.

## Meaning of Confidence Interval for $\theta_0$

$$P(A \leq \theta_0 \leq B) = 1 - \alpha$$

Each time we sample we'll get a different confidence interval, corresponding to different realizations of the random variables  $A$  and  $B$ . If we sample many times, approximately  $100 \times (1 - \alpha)\%$  of these intervals will contain the population parameter  $\theta_0$ .

## True or False?

Suppose

$$X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$$

Then the population mean  $\mu$  has approximately a 95% chance of falling in the interval  $\bar{X}_n \pm 2\sigma/\sqrt{n}$ .

- (a) True
- (b) False

## True or False?

Suppose

$$X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$$

Then the population mean  $\mu$  has approximately a 95% chance of falling in the interval  $\bar{X}_n \pm 2\sigma/\sqrt{n}$ .

- (a) True
- (b) False

FALSE! –  $\mu$  is a constant!

## Margin of Error

When a CI takes the form  $\hat{\theta} \pm ME$ ,  $ME$  is the Margin of Error.

# Confidence Intervals: Some Terminology

## Margin of Error

When a CI takes the form  $\hat{\theta} \pm ME$ ,  $ME$  is the Margin of Error.

## Lower and Upper Confidence Limits

The lower endpoint of a CI is the lower confidence limit (LCL), while the upper endpoint is the upper confidence limit (UCL).

# Confidence Intervals: Some Terminology

## Margin of Error

When a CI takes the form  $\hat{\theta} \pm ME$ ,  $ME$  is the Margin of Error.

## Lower and Upper Confidence Limits

The lower endpoint of a CI is the lower confidence limit (LCL), while the upper endpoint is the upper confidence limit (UCL).

## Width of a Confidence Interval

The distance  $|UCL - LCL|$  is called the width of a CI. This means exactly what it says.



# Margin of Error

---

# What is the Margin of Error

In the preceding example of a 95% confidence interval for the mean of a normal population when the population variance is known, which of these is the margin of error?

(a)  $\sigma/\sqrt{n}$

(b)  $\bar{X}_n$

(c)  $\sigma$

(d)  $2\sigma/\sqrt{n}$

(e)  $1/\sqrt{n}$

# What is the Margin of Error

In the preceding example of a 95% confidence interval for the mean of a normal population when the population variance is known, which of these is the margin of error?

(a)  $\sigma/\sqrt{n}$

(b)  $\bar{X}_n$

(c)  $\sigma$

(d)  $2\sigma/\sqrt{n}$

(e)  $1/\sqrt{n}$

$2\sigma/\sqrt{n}$ , since the CI is  $\bar{X}_n \pm 2\sigma/\sqrt{n}$

## What is the Width?

In the preceding example of a 95% confidence interval for the mean of a normal population when the population variance is known, which of these is the width of the interval?

(a)  $\sigma/\sqrt{n}$

(b)  $2\sigma/\sqrt{n}$

(c)  $3\sigma/\sqrt{n}$

(d)  $4\sigma/\sqrt{n}$

(e)  $5\sigma/\sqrt{n}$

# What is the Width?

In the preceding example of a 95% confidence interval for the mean of a normal population when the population variance is known, which of these is the width of the interval?

- (a)  $\sigma/\sqrt{n}$
- (b)  $2\sigma/\sqrt{n}$
- (c)  $3\sigma/\sqrt{n}$
- (d)  $4\sigma/\sqrt{n}$
- (e)  $5\sigma/\sqrt{n}$

$4\sigma/\sqrt{n}$ , since the CI is  $\bar{X}_n \pm 2\sigma/\sqrt{n}$

## Example: Calculate the Margin of Error

$X_1, \dots, X_{100} \sim \text{iid } N(\mu, 1)$  but we don't know  $\mu$ .  
Want to create a 95% confidence interval for  $\mu$ .

What is the margin of error?

## Example: Calculate the Margin of Error

$X_1, \dots, X_{100} \sim \text{iid } N(\mu, 1)$  but we don't know  $\mu$ .  
Want to create a 95% confidence interval for  $\mu$ .

What is the margin of error?

The confidence interval is  $\bar{X}_n \pm 2\sigma/\sqrt{n}$  so

$$ME = 2\sigma/\sqrt{n} = 2 \cdot 1/\sqrt{100} = 2/10 = 0.2$$

## Example: Calculate the Lower Confidence Limit

$X_1, \dots, X_{100} \sim N(\mu, 1)$  but we don't know  $\mu$ . Want to create a 95% confidence interval for  $\mu$ .

We found that  $ME = 0.2$ . The sample mean  $\bar{x} = 4.9$ . What is the lower confidence limit?



## Example: Calculate the Lower Confidence Limit

$X_1, \dots, X_{100} \sim N(\mu, 1)$  but we don't know  $\mu$ . Want to create a 95% confidence interval for  $\mu$ .

We found that  $ME = 0.2$ . The sample mean  $\bar{x} = 4.9$ . What is the lower confidence limit?

$$LCL = \bar{x} - ME = 4.9 - 0.2 = 4.7$$

## Example: Similarly for the Upper Confidence Limit...

$X_1, \dots, X_{100} \sim N(\mu, 1)$  but we don't know  $\mu$ . Want to create a 95% confidence interval for  $\mu$ .

We found that  $ME = 0.2$ . The sample mean  $\bar{x} = 4.9$ . What is the upper confidence limit?

## Example: Similarly for the Upper Confidence Limit...

$X_1, \dots, X_{100} \sim N(\mu, 1)$  but we don't know  $\mu$ . Want to create a 95% confidence interval for  $\mu$ .

We found that  $ME = 0.2$ . The sample mean  $\bar{x} = 4.9$ . What is the upper confidence limit?

$$UCL = \bar{x} + ME = 4.9 + 0.2 = 5.1$$

## Example: 95% CI for Normal Mean, Popn. Var. Known

$X_1, \dots, X_{100} \sim N(\mu, 1)$  but we don't know  $\mu$ .

95% CI for  $\mu = [4.7, 5.1]$

What values of  $\mu$  are plausible?

## Example: 95% CI for Normal Mean, Popn. Var. Known

$X_1, \dots, X_{100} \sim N(\mu, 1)$  but we don't know  $\mu$ .

95% CI for  $\mu = [4.7, 5.1]$

What values of  $\mu$  are plausible?

The data actually came from a  $N(5, 1)$  Distribution.

Want to be more certain? Use higher confidence level.

What value of  $c$  should we use to get a  $100 \times (1 - \alpha)\%$  CI for  $\mu$ ?

$$P\left(-c \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq c\right) = 1 - \alpha$$

Want to be more certain? Use higher confidence level.

What value of  $c$  should we use to get a  $100 \times (1 - \alpha)\%$  CI for  $\mu$ ?

$$P\left(-c \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq c\right) = 1 - \alpha$$

$$P\left(\bar{X}_n - c\sigma/\sqrt{n} \leq \mu \leq \bar{X}_n + c\sigma/\sqrt{n}\right) = 1 - \alpha$$

Want to be more certain? Use higher confidence level.

What value of  $c$  should we use to get a  $100 \times (1 - \alpha)\%$  CI for  $\mu$ ?

$$P\left(-c \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq c\right) = 1 - \alpha$$

$$P\left(\bar{X}_n - c\sigma/\sqrt{n} \leq \mu \leq \bar{X}_n + c\sigma/\sqrt{n}\right) = 1 - \alpha$$

Take  $c = \text{qnorm}(1 - \alpha/2)$



Want to be more certain? Use higher confidence level.

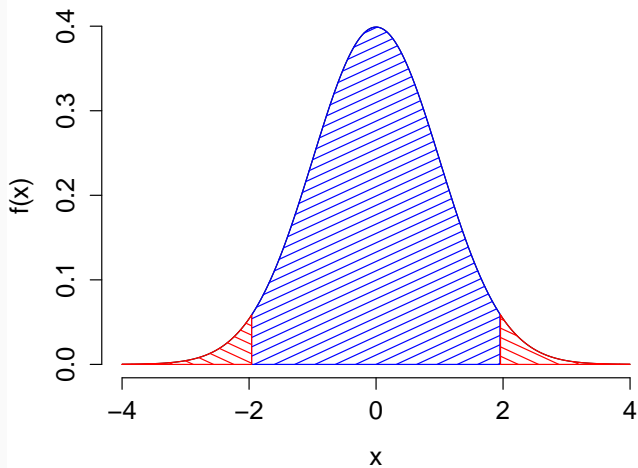
What value of  $c$  should we use to get a  $100 \times (1 - \alpha)\%$  CI for  $\mu$ ?

$$P\left(-c \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq c\right) = 1 - \alpha$$

$$P\left(\bar{X}_n - c\sigma/\sqrt{n} \leq \mu \leq \bar{X}_n + c\sigma/\sqrt{n}\right) = 1 - \alpha$$

Take  $c = \mathbf{qnorm}(1 - \alpha/2)$

$$\bar{X}_n \pm \mathbf{qnorm}(1 - \alpha/2) \times \sigma/\sqrt{n}$$



## Confidence Interval for a Normal Mean, $\sigma$ Known

$$\bar{X}_n \pm \text{qnorm}(1 - \alpha/2) \times \sigma / \sqrt{n}$$

# What Affects the Margin of Error?

$$\bar{X}_n \pm \text{qnorm}(1 - \alpha/2) \times \sigma / \sqrt{n}$$

## Sample Size $n$

ME decreases with  $n$ : bigger sample  $\implies$  tighter interval

## Population Std. Dev. $\sigma$

ME increases with  $\sigma$ : more variable population  $\implies$  wider interval

## Confidence Level $1 - \alpha$

ME increases with  $1 - \alpha$ : higher conf. level  $\implies$  wider interval

# What Affects the Margin of Error?

$$\bar{X}_n \pm \text{qnorm}(1 - \alpha/2) \times \sigma / \sqrt{n}$$

## Sample Size $n$

ME decreases with  $n$ : bigger sample  $\implies$  tighter interval

## Population Std. Dev. $\sigma$

ME increases with  $\sigma$ : more variable population  $\implies$  wider interval

## Confidence Level $1 - \alpha$

ME increases with  $1 - \alpha$ : higher conf. level  $\implies$  wider interval

Conf. Level	90%	95%	99%
$\alpha$	0.1	0.05	0.01
$\text{qnorm}(1 - \alpha/2)$	1.64	1.96	2.56

## But What if $\sigma$ is Unknown?

- What we've done so far assumed that  $\sigma$  was known.
- In real applications this is typically not the case.

## But What if $\sigma$ is Unknown?

- What we've done so far assumed that  $\sigma$  was known.
- In real applications this is typically not the case.
- So what do we do now?

# The Suspense!





## We Don't know $\sigma$ . What to use instead?

$$\bar{X}_n \pm \text{qnorm}(1 - \alpha/2) \times \sigma / \sqrt{n}$$

What about Sample Standard Deviation  $S$ ?

$$P\left(-2 \leq \frac{\bar{X}_n - \mu}{S/\sqrt{n}} \leq 2\right) = 0.95 \text{ ???}$$

Not Quite!

Although  $(\bar{X}_n - \mu)/(\sigma/\sqrt{n}) \sim N(0, 1)$ ,  $S \neq \sigma$ . In fact,  $S$  is an **estimator** of  $\sigma$  so it is a **random variable**!

# What is the sampling distribution?

Suppose  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$

$$\boxed{\frac{\bar{X}_n - \mu}{S/\sqrt{n}} \sim ???}$$

## First Step

What is the sampling distribution of  $S$ ?

## What is the Distribution?

Suppose  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ . What is the distribution of this sum?

$$\sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2$$

- (a)  $\chi^2(n)$
- (b)  $N(\mu, \sigma^2)$
- (c)  $N(0, 1)$
- (d)  $N(\mu, \sigma^2/n)$
- (e)  $\chi^2(1)$

## Towards the Sampling Dist. of $S$

If  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ , then

$$\sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 \sim \chi^2(n)$$

Now:

$$\sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 =$$

## Towards the Sampling Dist. of $S$

If  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ , then

$$\sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 \sim \chi^2(n)$$

Now:

$$\sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 = \left( \frac{n-1}{\sigma^2} \right) \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu)^2 \right]$$

## Towards the Sampling Dist. of $S$

If  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ , then

$$\sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 \sim \chi^2(n)$$

Now:

$$\sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 = \left( \frac{n-1}{\sigma^2} \right) \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu)^2 \right] \sim \chi^2(n)$$

Anything look familiar?

# Sampling Distribution of Sample Variance

Suppose  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ . Then whereas

$$\left( \frac{n-1}{\sigma^2} \right) \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu)^2 \right] \sim \chi^2(n)$$

# Sampling Distribution of Sample Variance

Suppose  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ . Then whereas

$$\left( \frac{n-1}{\sigma^2} \right) \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu)^2 \right] \sim \chi^2(n)$$

Replacing  $\mu$  with  $\bar{X}$  “loses” a degree of freedom

$$\left( \frac{n-1}{\sigma^2} \right) \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \right] =$$



# Sampling Distribution of Sample Variance

Suppose  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ . Then whereas

$$\left(\frac{n-1}{\sigma^2}\right) \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu)^2 \right] \sim \chi^2(n)$$

Replacing  $\mu$  with  $\bar{X}$  “loses” a degree of freedom

$$\left(\frac{n-1}{\sigma^2}\right) \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \right] = \left(\frac{n-1}{\sigma^2}\right) S^2$$

# Sampling Distribution of Sample Variance

Suppose  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ . Then whereas

$$\left(\frac{n-1}{\sigma^2}\right) \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu)^2 \right] \sim \chi^2(n)$$

Replacing  $\mu$  with  $\bar{X}$  “loses” a degree of freedom

$$\left(\frac{n-1}{\sigma^2}\right) \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \right] = \left(\frac{n-1}{\sigma^2}\right) S^2 \sim \chi^2(n-1)$$

Ultimately, we will use this fact to work out the sampling distribution of  $\sqrt{n}(\bar{X}_n - \mu)/S$ , but for now let's take a detour...

Detour

---

## 95% CI for Variance of Normal Population

We know that:

$$\left( \frac{n-1}{\sigma^2} \right) S^2 \sim \chi^2(n-1)$$

## 95% CI for Variance of Normal Population

We know that:

$$\left( \frac{n-1}{\sigma^2} \right) S^2 \sim \chi^2(n-1)$$

First Step: find  $a, b$  such that

$$P \left[ a \leq \left( \frac{n-1}{\sigma^2} \right) S^2 \leq b \right] = 0.95$$

## 95% CI for Variance of Normal Population

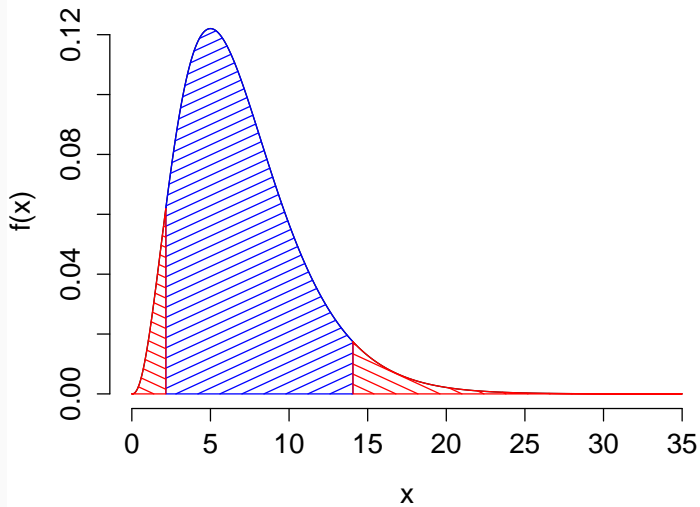
We know that:

$$\left( \frac{n-1}{\sigma^2} \right) S^2 \sim \chi^2(n-1)$$

First Step: find  $a, b$  such that

$$P \left[ a \leq \left( \frac{n-1}{\sigma^2} \right) S^2 \leq b \right] = 0.95$$

Although there are many choices for  $a, b$  that would work, a sensible idea is to put 2.5% in each tail...



What R command should I use to calculate  $a$ ?

$$P\left[a \leq \left(\frac{n-1}{\sigma^2}\right) S^2 \leq b\right] = 0.95$$

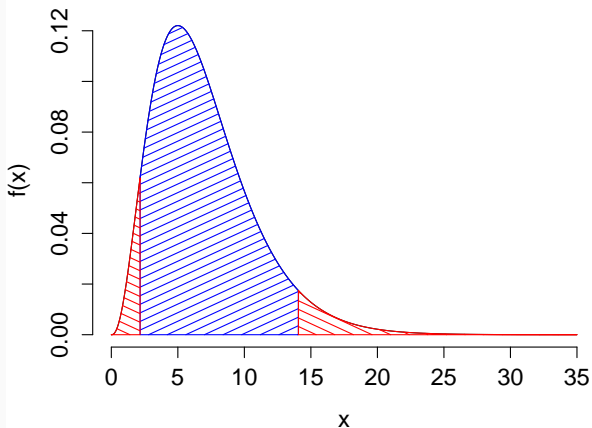
- (a) `qchisq(0.95, df = n - 1)`
- (b) `qchisq(0.025, df = n)`
- (c) `qchisq(0.975, df = n - 1)`
- (d) `qchisq(0.025, df = n - 1)`
- (e) `qchisq(0.975, df = n)`



What R command should I use to calculate  $b$ ?

$$P\left[a \leq \left(\frac{n-1}{\sigma^2}\right) S^2 \leq b\right] = 0.95$$

- (a) `qchisq(0.95, df = n - 1)`
- (b) `qchisq(0.025, df = n)`
- (c) `qchisq(0.975, df = n - 1)`
- (d) `qchisq(0.025, df = n - 1)`
- (e) `qchisq(0.975, df = n)`



```
a = qchisq(0.025, df = n - 1)
b = qchisq(0.975, df = n - 1)
```

## Step 2: After Finding $a, b$ Rearrange

$$P \left[ a \leq \left( \frac{n-1}{\sigma^2} \right) S^2 \leq b \right] = 0.95$$

## Step 2: After Finding $a, b$ Rearrange

$$P \left[ a \leq \left( \frac{n-1}{\sigma^2} \right) S^2 \leq b \right] = 0.95$$

$$P \left[ \frac{a}{(n-1)S^2} \leq \frac{1}{\sigma^2} \leq \frac{b}{(n-1)S^2} \right] = 0.95$$

## Step 2: After Finding $a, b$ Rearrange

$$P \left[ a \leq \left( \frac{n-1}{\sigma^2} \right) S^2 \leq b \right] = 0.95$$

$$P \left[ \frac{a}{(n-1)S^2} \leq \frac{1}{\sigma^2} \leq \frac{b}{(n-1)S^2} \right] = 0.95$$

$$P \left[ \frac{(n-1)S^2}{b} \leq \sigma^2 \leq \frac{(n-1)S^2}{a} \right] = 0.95$$

## Step 2: After Finding $a, b$ Rearrange

$$P \left[ a \leq \left( \frac{n-1}{\sigma^2} \right) S^2 \leq b \right] = 0.95$$

$$P \left[ \frac{a}{(n-1)S^2} \leq \frac{1}{\sigma^2} \leq \frac{b}{(n-1)S^2} \right] = 0.95$$

$$P \left[ \frac{(n-1)S^2}{b} \leq \sigma^2 \leq \frac{(n-1)S^2}{a} \right] = 0.95$$

This CI is *not* symmetric: it *doesn't* take the form  $\hat{\theta} \pm ME$ !

## Example: 95% Confidence Interval for Normal Variance

$X_1, \dots, X_{100} \sim N(\mu, \sigma^2)$ . Here  $n - 1 = 99$ , hence

## Example: 95% Confidence Interval for Normal Variance

$X_1, \dots, X_{100} \sim N(\mu, \sigma^2)$ . Here  $n - 1 = 99$ , hence

$$a = \text{qchisq}(0.025, \text{df} = 99) \approx 73$$



## Example: 95% Confidence Interval for Normal Variance

$X_1, \dots, X_{100} \sim N(\mu, \sigma^2)$ . Here  $n - 1 = 99$ , hence

$$a = \text{qchisq}(0.025, \text{df} = 99) \approx 73$$

$$b = \text{qchisq}(0.975, \text{df} = 99) \approx 128$$

## Example: 95% Confidence Interval for Normal Variance

$X_1, \dots, X_{100} \sim N(\mu, \sigma^2)$ . Here  $n - 1 = 99$ , hence

$$a = \text{qchisq}(0.025, \text{df} = 99) \approx 73$$

$$b = \text{qchisq}(0.975, \text{df} = 99) \approx 128$$

From the sample data,  $s^2 = 4.3$

$$LCL = (n - 1)s^2/b =$$

## Example: 95% Confidence Interval for Normal Variance

$X_1, \dots, X_{100} \sim N(\mu, \sigma^2)$ . Here  $n - 1 = 99$ , hence

$$a = \text{qchisq}(0.025, \text{df} = 99) \approx 73$$

$$b = \text{qchisq}(0.975, \text{df} = 99) \approx 128$$

From the sample data,  $s^2 = 4.3$

$$LCL = (n - 1)s^2/b = 99 \times 4.3/128$$

## Example: 95% Confidence Interval for Normal Variance

$X_1, \dots, X_{100} \sim N(\mu, \sigma^2)$ . Here  $n - 1 = 99$ , hence

$$a = \text{qchisq}(0.025, \text{df} = 99) \approx 73$$

$$b = \text{qchisq}(0.975, \text{df} = 99) \approx 128$$

From the sample data,  $s^2 = 4.3$

$$LCL = (n - 1)s^2/b = 99 \times 4.3/128 \approx 3.3$$

## Example: 95% Confidence Interval for Normal Variance

$X_1, \dots, X_{100} \sim N(\mu, \sigma^2)$ . Here  $n - 1 = 99$ , hence

$$a = \text{qchisq}(0.025, \text{df} = 99) \approx 73$$

$$b = \text{qchisq}(0.975, \text{df} = 99) \approx 128$$

From the sample data,  $s^2 = 4.3$

$$LCL = (n - 1)s^2/b = 99 \times 4.3/128 \approx 3.3$$

$$UCL = (n - 1)s^2/a =$$

## Example: 95% Confidence Interval for Normal Variance

$X_1, \dots, X_{100} \sim N(\mu, \sigma^2)$ . Here  $n - 1 = 99$ , hence

$$a = \text{qchisq}(0.025, \text{df} = 99) \approx 73$$

$$b = \text{qchisq}(0.975, \text{df} = 99) \approx 128$$

From the sample data,  $s^2 = 4.3$

$$LCL = (n - 1)s^2/b = 99 \times 4.3/128 \approx 3.3$$

$$UCL = (n - 1)s^2/a = 99 \times 4.3/73$$

## Example: 95% Confidence Interval for Normal Variance

$X_1, \dots, X_{100} \sim N(\mu, \sigma^2)$ . Here  $n - 1 = 99$ , hence

$$a = \text{qchisq}(0.025, \text{df} = 99) \approx 73$$

$$b = \text{qchisq}(0.975, \text{df} = 99) \approx 128$$

From the sample data,  $s^2 = 4.3$

$$LCL = (n - 1)s^2/b = 99 \times 4.3/128 \approx 3.3$$

$$UCL = (n - 1)s^2/a = 99 \times 4.3/73 \approx 5.8$$

## Example: 95% Confidence Interval for Normal Variance

$X_1, \dots, X_{100} \sim N(\mu, \sigma^2)$ . Here  $n - 1 = 99$ , hence

$$a = \text{qchisq}(0.025, \text{df} = 99) \approx 73$$

$$b = \text{qchisq}(0.975, \text{df} = 99) \approx 128$$

From the sample data,  $s^2 = 4.3$

$$LCL = (n - 1)s^2/b = 99 \times 4.3/128 \approx 3.3$$

$$UCL = (n - 1)s^2/a = 99 \times 4.3/73 \approx 5.8$$

95% CI for  $\sigma^2$  is  $[3.3, 5.8]$ . What values are plausible?



## Example: 95% Confidence Interval for Normal Variance

$X_1, \dots, X_{100} \sim N(\mu, \sigma^2)$ . Here  $n - 1 = 99$ , hence

$$a = \text{qchisq}(0.025, \text{df} = 99) \approx 73$$

$$b = \text{qchisq}(0.975, \text{df} = 99) \approx 128$$

From the sample data,  $s^2 = 4.3$

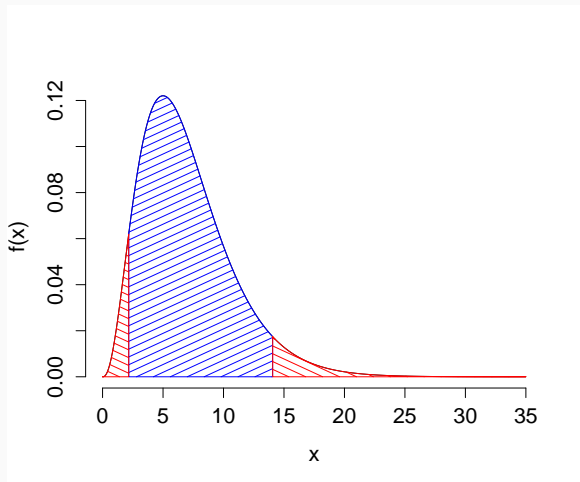
$$LCL = (n - 1)s^2/b = 99 \times 4.3/128 \approx 3.3$$

$$UCL = (n - 1)s^2/a = 99 \times 4.3/73 \approx 5.8$$

95% CI for  $\sigma^2$  is [3.3, 5.8]. What values are plausible?

The actual population variance in this case was 4

## Arbitrary Confidence Level: $(1 - \alpha)$



```
a = qchisq( $\alpha/2$ , df = n - 1)
```

```
b = qchisq( $1 - \alpha/2$ , df = n - 1)
```

## CI for Normal Variance

`a = qchisq( $\alpha/2$ , df = n - 1)`

`b = qchisq( $1 - \alpha/2$ , df = n - 1)`

$$P\left[a \leq \left(\frac{n-1}{\sigma^2}\right) S^2 \leq b\right] = 1 - \alpha$$

$$P\left[\frac{a}{(n-1)S^2} \leq \frac{1}{\sigma^2} \leq \frac{b}{(n-1)S^2}\right] = 1 - \alpha$$

$$P\left[\frac{(n-1)S^2}{b} \leq \sigma^2 \leq \frac{(n-1)S^2}{a}\right] = 1 - \alpha$$

## CI for Normal Variance

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$  and let:

$$a = \text{qchisq}(\alpha/2, \text{df} = n - 1)$$

$$b = \text{qchisq}(1 - \alpha/2, \text{df} = n - 1)$$

Then,

$$\left[ \frac{(n-1)S^2}{b}, \frac{(n-1)S^2}{a} \right]$$

is a  $100 \times (1 - \alpha)\%$  confidence interval for  $\sigma^2$ .

We want to know the Sampling Distribution of  $\sqrt{n}(\bar{X}_n - \mu)/S$  and we just saw that:

$$\left( \frac{n-1}{\sigma^2} \right) S^2 \sim \chi^2(n-1)$$

How can we use this fact to help us?

## Back on Track

---

## What is the Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$ ?

This slide is just algebra:

$$\frac{\bar{X}_n - \mu}{S/\sqrt{n}} =$$

## What is the Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$ ?

This slide is just algebra:

$$\frac{\bar{X}_n - \mu}{S/\sqrt{n}} = \frac{\bar{X}_n - \mu}{S/\sqrt{n}} \cdot \left( \frac{\sigma/\sqrt{n}}{\sigma/\sqrt{n}} \right) =$$



## What is the Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$ ?

This slide is just algebra:

$$\begin{aligned}\frac{\bar{X}_n - \mu}{S/\sqrt{n}} &= \frac{\bar{X}_n - \mu}{S/\sqrt{n}} \cdot \left( \frac{\sigma/\sqrt{n}}{\sigma/\sqrt{n}} \right) = \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \frac{\sigma/\sqrt{n}}{S/\sqrt{n}} \right) \\ &= \end{aligned}$$

## What is the Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$ ?

This slide is just algebra:

$$\begin{aligned}\frac{\bar{X}_n - \mu}{S/\sqrt{n}} &= \frac{\bar{X}_n - \mu}{S/\sqrt{n}} \cdot \left( \frac{\sigma/\sqrt{n}}{\sigma/\sqrt{n}} \right) = \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \frac{\sigma/\sqrt{n}}{S/\sqrt{n}} \right) \\ &= \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \frac{\sigma}{S} \right) =\end{aligned}$$

## What is the Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$ ?

This slide is just algebra:

$$\begin{aligned}\frac{\bar{X}_n - \mu}{S/\sqrt{n}} &= \frac{\bar{X}_n - \mu}{S/\sqrt{n}} \cdot \left( \frac{\sigma/\sqrt{n}}{\sigma/\sqrt{n}} \right) = \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \frac{\sigma/\sqrt{n}}{S/\sqrt{n}} \right) \\ &= \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \frac{\sigma}{S} \right) = \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \sqrt{\frac{n-1}{n-1}} \cdot \sqrt{\frac{\sigma^2}{S^2}} \right) \\ &= \end{aligned}$$

## What is the Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$ ?

This slide is just algebra:

$$\begin{aligned}\frac{\bar{X}_n - \mu}{S/\sqrt{n}} &= \frac{\bar{X}_n - \mu}{S/\sqrt{n}} \cdot \left( \frac{\sigma/\sqrt{n}}{\sigma/\sqrt{n}} \right) = \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \frac{\sigma/\sqrt{n}}{S/\sqrt{n}} \right) \\&= \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \frac{\sigma}{S} \right) = \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \sqrt{\frac{n-1}{n-1}} \cdot \sqrt{\frac{\sigma^2}{S^2}} \right) \\&= \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \sqrt{\frac{(n-1)\sigma^2}{(n-1)S^2}} \right) \\&= \end{aligned}$$

# What is the Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$ ?

This slide is just algebra:

$$\begin{aligned}\frac{\bar{X}_n - \mu}{S/\sqrt{n}} &= \frac{\bar{X}_n - \mu}{S/\sqrt{n}} \cdot \left( \frac{\sigma/\sqrt{n}}{\sigma/\sqrt{n}} \right) = \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \frac{\sigma/\sqrt{n}}{S/\sqrt{n}} \right) \\&= \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \frac{\sigma}{S} \right) = \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \sqrt{\frac{n-1}{n-1}} \cdot \sqrt{\frac{\sigma^2}{S^2}} \right) \\&= \left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right) \left( \sqrt{\frac{(n-1)\sigma^2}{(n-1)S^2}} \right) \\&= \frac{\left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \right)}{\sqrt{\left[ \frac{(n-1)S^2}{\sigma^2} \right] / (n-1)}}\end{aligned}$$

## Distribution of $\sqrt{n}(\bar{X}_n - \mu)/\sigma$

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$  and  $\bar{X}_n$  is the sample mean. Then the sampling distribution of  $\sqrt{n}(\bar{X}_n - \mu)/\sigma$  is

- (a)  $t(n)$
- (b)  $t(n - 1)$
- (c)  $\chi^2(n)$
- (d)  $\chi^2(n - 1)$
- (e)  $N(\mu, \sigma^2)$
- (f)  $N(0, 1)$
- (g)  $N(\mu, \sigma^2/n)$
- (h)  $F(n, n - 1)$

## Distribution of $(n - 1)S^2/\sigma^2$

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$  and  $S^2$  is the sample variance. Then the sampling distribution of  $(n - 1)S^2/\sigma^2$  is

- (a)  $t(n)$
- (b)  $t(n - 1)$
- (c)  $\chi^2(n)$
- (d)  $\chi^2(n - 1)$
- (e)  $N(\mu, \sigma^2)$
- (f)  $N(0, 1)$
- (g)  $N(\mu, \sigma^2/n)$
- (h)  $F(n, n - 1)$

# What is the Sampling Distribution?

Suppose  $Z \sim N(0, 1)$  independent of  $Y \sim \chi^2(n - 1)$ . Then the sampling distribution of  $Z/\sqrt{Y/(n - 1)}$  is

- (a)  $t(n)$
- (b)  $t(n - 1)$
- (c)  $\chi^2(n)$
- (d)  $\chi^2(n - 1)$
- (e)  $N(\mu, \sigma^2)$
- (f)  $N(0, 1)$
- (g)  $N(\mu, \sigma^2/n)$
- (h)  $F(n, n - 1)$



# What is the Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$ ?

From three slides back:

$$\begin{aligned}\frac{\bar{X}_n - \mu}{S/\sqrt{n}} &= \frac{\left(\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}\right)}{\sqrt{\left[\frac{(n-1)S^2}{\sigma^2}\right]/(n-1)}} \\ &= \end{aligned}$$

# What is the Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$ ?

From three slides back:

$$\begin{aligned}\frac{\bar{X}_n - \mu}{S/\sqrt{n}} &= \frac{\left(\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}\right)}{\sqrt{\left[\frac{(n-1)S^2}{\sigma^2}\right]/(n-1)}} \\ &= \frac{N(0,1)}{\sqrt{\frac{\chi^2(n-1)}{n-1}}} \\ &\sim\end{aligned}$$

# What is the Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$ ?

From three slides back:

$$\begin{aligned}\frac{\bar{X}_n - \mu}{S/\sqrt{n}} &= \frac{\left(\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}\right)}{\sqrt{\left[\frac{(n-1)S^2}{\sigma^2}\right]/(n-1)}} \\ &= \frac{N(0,1)}{\sqrt{\frac{\chi^2(n-1)}{n-1}}} \\ &\sim t(n-1)\end{aligned}$$

# What is the Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$ ?

From three slides back:

$$\begin{aligned}\frac{\bar{X}_n - \mu}{S/\sqrt{n}} &= \frac{\left(\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}\right)}{\sqrt{\left[\frac{(n-1)S^2}{\sigma^2}\right]/(n-1)}} \\ &= \frac{N(0,1)}{\sqrt{\frac{\chi^2(n-1)}{n-1}}} \\ &\sim t(n-1)\end{aligned}$$

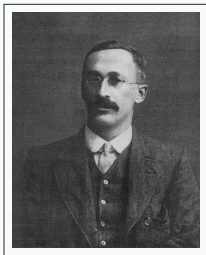
Strictly speaking, need to show that numerator and denominator are independent, but you can take my word for it!

## Punchline: Sampling Distribution of $\sqrt{n}(\bar{X}_n - \mu)/S$

If  $X_1, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$ , then

$$\frac{\bar{X}_n - \mu}{S/\sqrt{n}} \sim t(n-1)$$

# Who was “Student?”



*“Student” is the pseudonym used in 19 of 21 published articles by William Sealy Gosset, who was a chemist, brewer, inventor, and self-trained statistician, agronomer, and designer of experiments ... [Gosset] worked his entire adult life ... as an experimental brewer for one employer: Arthur Guinness, Son & Company, Ltd., Dublin, St. James’s Gate. Gosset was a master brewer and rose in fact to the top of the top of the brewing industry: Head Brewer of Guinness. [Source](#)*

# Three Key Sampling Distributions

Suppose that  $X_1, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$ . Then:

$$\left( \frac{n-1}{\sigma^2} \right) S^2 \sim \chi^2(n-1)$$

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

$$\frac{\bar{X}_n - \mu}{S/\sqrt{n}} \sim t(n-1)$$

## CI for Mean of Normal Distribution, Popn. Var. Unknown

Same argument as we used when the variance was known, except with  $t(n - 1)$  rather than standard normal distribution:

$$P\left(-c \leq \frac{\bar{X}_n - \mu}{S/\sqrt{n}} \leq c\right) = 1 - \alpha$$

$$P\left(\bar{X}_n - c \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X}_n + c \frac{S}{\sqrt{n}}\right) = 1 - \alpha$$

$$c = \text{qt}(1 - \alpha/2, \text{df} = n - 1)$$

$$\bar{X}_n \pm \text{qt}(1 - \alpha/2, \text{df} = n - 1) \frac{S}{\sqrt{n}}$$



# Comparison of CIs for Mean of Normal Distribution

$$X_1, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$$

Known Population Std. Dev. ( $\sigma$ )

$$\bar{X}_n \pm \text{qnorm}(1 - \alpha/2) \frac{\sigma}{\sqrt{n}}$$

Unknown Population Std. Dev. ( $\sigma$ )

$$\bar{X}_n \pm \text{qt}(1 - \alpha/2, \text{df} = n - 1) \frac{S}{\sqrt{n}}$$

# Standard Error vs. Estimator of Standard Error

## Standard Error

Recall that the standard deviation of the sampling distribution of an estimator is called the *standard error* (SE) of that estimator.

## Example: Standard Error of the Mean

$$SE(\bar{X}_n) = \sqrt{\text{Var}(\bar{X}_n)} = \sigma/\sqrt{n}$$

## Estimator of Standard Error of the Mean

Whereas  $\sigma/\sqrt{n}$  *is* the standard error of the mean,  $S/\sqrt{n}$  is an *estimator* of this quantity:  $\widehat{SE}(\bar{X}_n) = S/\sqrt{n}$

## Writing the CIs in terms of Actual and Estimated SE

$$X_1, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$$

Known Population Std. Dev. ( $\sigma$ )

$$\bar{X}_n \pm \text{qnorm}(1 - \alpha/2) \text{SE}(\bar{X}_n)$$

Unknown Population Std. Dev. ( $\sigma$ )

$$\bar{X}_n \pm \text{qt}(1 - \alpha/2, \text{df} = n - 1) \widehat{\text{SE}}(\bar{X}_n)$$

## Comparison of Normal and $t$ CIs

**Table 2:** Values of  $qt(1 - \alpha/2, df = n - 1)$  for various choices of  $n$  and  $\alpha$ .

$n$	1	5	10	30	100	$\infty$
$\alpha = 0.10$	6.31	2.02	1.81	1.70	1.66	1.64
$\alpha = 0.05$	12.71	2.57	2.23	2.04	1.98	1.96
$\alpha = 0.01$	63.66	4.03	3.17	2.75	2.63	2.58

Recall that as  $n \rightarrow \infty$ ,  $t(n - 1) \rightarrow N(0, 1)$

In a sense, using the  $t$ -distribution involves making a “small-sample correction.” In other words, it is only when  $n$  is fairly small that this makes a practical difference for our confidence intervals.

# Is Joe Taller Than The Average American Male?

Assuming the population is normal,

$$\bar{X}_n \pm \text{qt}(1 - \alpha/2, \text{df} = n - 1) \widehat{SE}(\bar{X}_n)$$

What is the approximate value of  
 $\text{qt}(1 - 0.05/2, \text{df} = 5646)$ ?

Sample Mean	69 inches
Sample Std. Dev.	6 inches
Sample Size	5647
Joe's Height	73 inches

$$\begin{aligned}\widehat{SE}(\bar{X}_n) &= s/\sqrt{n} \\ &= 6/\sqrt{5647} \\ &\approx 0.08\end{aligned}$$

# Is Joe Taller Than The Average American Male?

Sample Mean	69 inches
Sample Std. Dev.	6 inches
Sample Size	5647
Joe's Height	73 inches

$$\begin{aligned}\widehat{SE}(\bar{X}_n) &= s/\sqrt{n} \\ &= 6/\sqrt{5647} \\ &\approx 0.08\end{aligned}$$

Assuming the population is normal,

$$\bar{X}_n \pm qt(1 - \alpha/2, df = n - 1) \widehat{SE}(\bar{X}_n)$$

What is the approximate value of  $qt(1 - 0.05/2, df = 5646)$ ?

For large  $n$ ,  $t(n - 1) \approx N(0, 1)$ , so the answer is approximately 2

# Is Joe Taller Than The Average American Male?

Sample Mean	69 inches
Sample Std. Dev.	6 inches
Sample Size	5647
Joe's Height	73 inches

$$\begin{aligned}\widehat{SE}(\bar{X}_n) &= s/\sqrt{n} \\ &= 6/\sqrt{5647} \\ &\approx 0.08\end{aligned}$$

Assuming the population is normal,

$$\bar{X}_n \pm qt(1 - \alpha/2, df = n - 1) \widehat{SE}(\bar{X}_n)$$

What is the approximate value of  $qt(1 - 0.05/2, df = 5646)$ ?

For large  $n$ ,  $t(n - 1) \approx N(0, 1)$ , so the answer is approximately 2

What is the ME for the 95% CI?

# Is Joe Taller Than The Average American Male?

Sample Mean	69 inches
Sample Std. Dev.	6 inches
Sample Size	5647
Joe's Height	73 inches

$$\begin{aligned}\widehat{SE}(\bar{X}_n) &= s/\sqrt{n} \\ &= 6/\sqrt{5647} \\ &\approx 0.08\end{aligned}$$

Assuming the population is normal,

$$\bar{X}_n \pm qt(1 - \alpha/2, df = n - 1) \widehat{SE}(\bar{X}_n)$$

What is the approximate value of  $qt(1 - 0.05/2, df = 5646)$ ?

For large  $n$ ,  $t(n - 1) \approx N(0, 1)$ , so the answer is approximately 2

What is the ME for the 95% CI?

$$ME \approx 0.16 \implies 69 \pm 0.16$$



Stop Here for Midterm

---

## Two-sample Problem

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu_x, \sigma_x^2)$  independently of  $Y_1, \dots, Y_m \sim \text{iid } N(\mu_y, \sigma_y^2)$ . What is  $E[\bar{X}_n - \bar{Y}_m]$ , the expectation of the sampling distribution of the difference of sample means?

- (a)  $\mu_x$
- (b)  $\mu_x - \mu_y$
- (c)  $\mu_y$
- (d)  $\mu_x + \mu_y$
- (e) 0

## Two-sample Problem

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu_x, \sigma_x^2)$  independently of  $Y_1, \dots, Y_m \sim \text{iid } N(\mu_y, \sigma_y^2)$ . What is  $E[\bar{X}_n - \bar{Y}_m]$ , the expectation of the sampling distribution of the difference of sample means?

- (a)  $\mu_x$
- (b)  $\mu_x - \mu_y$
- (c)  $\mu_y$
- (d)  $\mu_x + \mu_y$
- (e) 0

$$E[\bar{X}_n - \bar{Y}_m] = E[\bar{X}_n] - E[\bar{Y}_m] = \mu_x - \mu_y$$

## Two-sample Problem

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu_x, \sigma_x^2)$  independently of  $Y_1, \dots, Y_m \sim \text{iid } N(\mu_y, \sigma_y^2)$ . What is  $\text{Var}[\bar{X}_n - \bar{Y}_m]$ , the variance of the sampling distribution of the difference of sample means?

- (a)  $\sigma_x^2 - \sigma_y^2$
- (b)  $\sigma_x^2 + \sigma_y^2$
- (c)  $\sigma_x^2/n + \sigma_y^2/m$
- (d)  $\sigma_x^2/n - \sigma_y^2/m$
- (e) 1

## Two-sample Problem

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu_x, \sigma_x^2)$  independently of  $Y_1, \dots, Y_m \sim \text{iid } N(\mu_y, \sigma_y^2)$ . What is  $\text{Var}[\bar{X}_n - \bar{Y}_m]$ , the variance of the sampling distribution of the difference of sample means?

- (a)  $\sigma_x^2 - \sigma_y^2$
- (b)  $\sigma_x^2 + \sigma_y^2$
- (c)  $\sigma_x^2/n + \sigma_y^2/m$
- (d)  $\sigma_x^2/n - \sigma_y^2/m$
- (e) 1

By independence:  $\text{Var}[\bar{X}_n - \bar{Y}_m] = \text{Var}[\bar{X}_n] + \text{Var}[\bar{Y}_m] = \frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}$

## Two-sample Problem

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu_x, \sigma_x^2)$  independently of  $Y_1, \dots, Y_m \sim \text{iid } N(\mu_y, \sigma_y^2)$ . What is the **sampling distribution** of  $\bar{X}_n - \bar{Y}_m$ , the difference of sample means?

- (a)  $\chi^2$
- (b)  $t$
- (c)  $F$
- (d) Normal

## Two-sample Problem

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu_x, \sigma_x^2)$  independently of  $Y_1, \dots, Y_m \sim \text{iid } N(\mu_y, \sigma_y^2)$ . What is the **sampling distribution** of  $\bar{X}_n - \bar{Y}_m$ , the difference of sample means?

- (a)  $\chi^2$
- (b)  $t$
- (c)  $F$
- (d) Normal

**Normal**, by independence and linearity property of normal distributions.

## Sampling Distribution of $\bar{X}_n - \bar{Y}_m$

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu_x, \sigma_x^2)$  independently of  $Y_1, \dots, Y_m \sim \text{iid } N(\mu_y, \sigma_y^2)$ . Then,

$$(\bar{X}_n - \bar{Y}_m) \sim N\left(\mu_x - \mu_y, \frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}\right)$$

$$\frac{(\bar{X}_n - \bar{Y}_m) - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}}} \sim N(0, 1)$$

$$\text{Shorthand: } SE(\bar{X}_n - \bar{Y}_m) = \sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}}$$



## CI for Difference of Population Means, $\sigma_x^2, \sigma_y^2$ Known

$$\frac{(\bar{X}_n - \bar{Y}_m) - (\mu_x - \mu_y)}{SE(\bar{X}_n - \bar{Y}_m)} \sim N(0, 1)$$

Thus, we construct a  $100 \times (1 - \alpha)\%$  CI for  $\mu_x - \mu_y$  as follows:

$$(\bar{X}_n - \bar{Y}_m) \pm \text{qnorm}(1 - \alpha/2) SE(\bar{X}_n - \bar{Y}_m)$$

Where  $SE(\bar{X}_n - \bar{Y}_m) = \sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}}$

## Calculate the ME for the Difference of Means

I generated independent random samples of size 25 from two normal distributions in R. One had a population standard deviation of 4 and the other had a population standard deviation of 3. The sample means were approximately 4.2 and 3.1.

Calculate the ME for a 95% confidence interval for the difference of population means.

## Calculate the ME for the Difference of Means

I generated independent random samples of size 25 from two normal distributions in R. One had a population standard deviation of 4 and the other had a population standard deviation of 3. The sample means were approximately 4.2 and 3.1.

Calculate the ME for a 95% confidence interval for the difference of population means.

$$SE = \sqrt{\frac{3^2}{25} + \frac{4^2}{25}} = \frac{\sqrt{9 + 16}}{5} = 1$$

$$ME = \text{qnorm}(1 - 0.05/2) \times SE \approx 2 \times SE = 2$$

## Calculate the LCL for the Difference of Means

I generated independent random samples of size 25 from two normal distributions in R. One had a population standard deviation of 4 and the other had a population standard deviation of 3. The sample means were approximately 4.2 and 3.1.

Calculate the LCL for a 95% confidence interval for the difference of population means.

## Calculate the LCL for the Difference of Means

I generated independent random samples of size 25 from two normal distributions in R. One had a population standard deviation of 4 and the other had a population standard deviation of 3. The sample means were approximately 4.2 and 3.1.

Calculate the LCL for a 95% confidence interval for the difference of population means.

$$LCL = (4.2 - 3.1) - ME = 1.1 - 2 = -0.9$$

## Calculate the UCL for the Difference of Means

I generated independent random samples of size 25 from two normal distributions in R. One had a population standard deviation of 4 and the other had a population standard deviation of 3. The sample means were approximately 4.2 and 3.1.

Calculate the UCL for a 95% confidence interval for the difference of population means.

## Calculate the UCL for the Difference of Means

I generated independent random samples of size 25 from two normal distributions in R. One had a population standard deviation of 4 and the other had a population standard deviation of 3. The sample means were approximately 4.2 and 3.1.

Calculate the UCL for a 95% confidence interval for the difference of population means.

$$UCL = (4.2 - 3.1) + ME = 1.1 + 2 = 3.1$$

## Calculate the UCL for the Difference of Means

I generated independent random samples of size 25 from two normal distributions in R. One had a population standard deviation of 4 and the other had a population standard deviation of 3. The sample means were approximately 4.2 and 3.1.

Calculate the UCL for a 95% confidence interval for the difference of population means.

$$UCL = (4.2 - 3.1) + ME = 1.1 + 2 = 3.1$$

95% Confidence Interval:  $(-0.9, 3.1)$



## Calculate the UCL for the Difference of Means

I generated independent random samples of size 25 from two normal distributions in R. One had a population standard deviation of 4 and the other had a population standard deviation of 3. The sample means were approximately 4.2 and 3.1.

Calculate the UCL for a 95% confidence interval for the difference of population means.

$$UCL = (4.2 - 3.1) + ME = 1.1 + 2 = 3.1$$

95% Confidence Interval:  $(-0.9, 3.1)$

The actual population means were 4 and 3, respectively

## What if $\sigma_x^2, \sigma_y^2$ are Unknown?

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu_x, \sigma_x^2)$  independently of  $Y_1, \dots, Y_m \sim \text{iid } N(\mu_y, \sigma_y^2)$ . Then,

$$\frac{(\bar{X}_n - \bar{Y}_m) - (\mu_x - \mu_y)}{\sqrt{\frac{S_x^2}{n} + \frac{S_y^2}{m}}} \sim t(\nu)$$

Formula for  $\nu$  is Complicated and You Don't Need to Know it

Two possibilities:

1. Have R find the correct value of  $\nu$  for us
2. If  $m, n$  are large enough, approximately standard normal.

## Case of Equal, Unknown Variances

The book considers a case where  $\sigma_x^2 = \sigma_y^2 = \sigma^2$ , that is a common unknown variance. This is a **very dangerous assumption**. It is almost certainly false and can throw off our results in a serious way. You are not responsible for this case.

# Sampling Distributions Under Normality: One-sample

Suppose that  $X_1, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$ . Then:

$$\left( \frac{n-1}{\sigma^2} \right) S^2 \sim \chi^2(n-1)$$

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

$$\frac{\bar{X}_n - \mu}{S/\sqrt{n}} \sim t(n-1)$$

## Sampling Distributions Under Normality: Two-sample

Suppose  $X_1, \dots, X_n \sim \text{iid } N(\mu_x, \sigma_x^2)$  independently of  $Y_1, \dots, Y_m \sim \text{iid } N(\mu_y, \sigma_y^2)$ . Then:

$$\frac{(\bar{X}_n - \bar{Y}_n) - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}}} \sim N(0, 1)$$

$$\frac{(\bar{X}_n - \bar{Y}_m) - (\mu_x - \mu_y)}{\sqrt{\frac{S_x^2}{n} + \frac{S_y^2}{m}}} \sim t(\nu)$$

But what if the population  
isn't Normal?

# The Central Limit Theorem

Suppose that  $X_1, \dots, X_n$  are a random sample from a population with unknown mean  $\mu$ . Then, provided that  $n$  is *sufficiently large*, the sampling distribution of  $\bar{X}_n$  is approximately  $N\left(\mu, \widehat{SE}(\bar{X}_n)^2\right)$ , even if the underlying population is *non-normal*.

In Other Words...

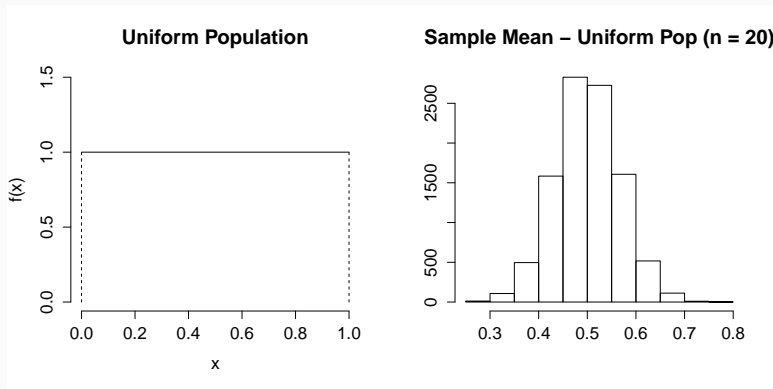
$$\frac{\bar{X}_n - \mu}{\widehat{SE}(\bar{X}_n)} \approx N(0, 1)$$

Use this to create *approximate* CIs for population mean!

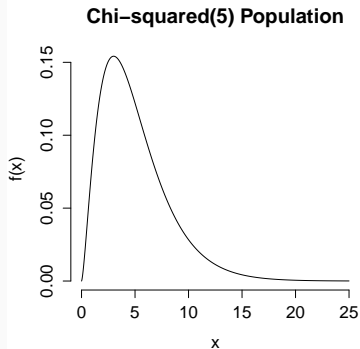
You should be amazed by  
this.



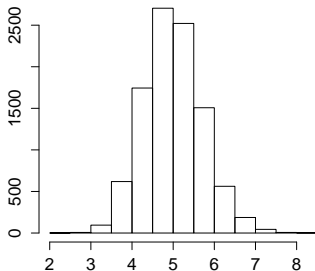
## Example: Uniform(0,1) Population, $n = 20$



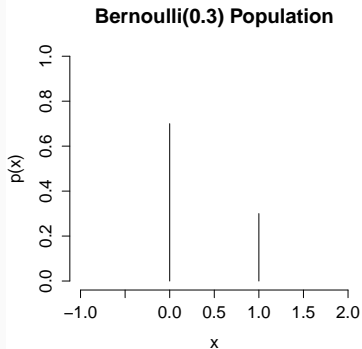
## Example: $\chi^2(5)$ Population, $n = 20$



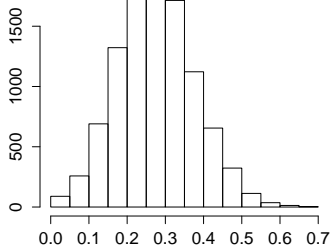
**Sample Mean – Chisq(5) Pop (n=20)**



## Example: Bernoulli(0.3) Population, $n = 20$



**Sample Mean – Ber(0.3) Pop ( $n = 20$ )**



# Who is the Chief Justice of the US Supreme Court?

- (a) Harry Reid
- (b) John Roberts
- (c) William Rehnquist
- (d) Stephen Breyer

# Are US Voters Really That Ignorant?

## The Data

Of 771 registered voters polled, only 39% correctly identified John Roberts as the current chief justice of the US Supreme Court.

## Research Question

Is the majority of voters unaware that John Roberts is the current chief justice, or is this just sampling variation?

Assume Random Sampling...

# Confidence Interval for a Proportion

**What is the appropriate probability model for the sample?**

$X_1, \dots, X_n \sim \text{iid Bernoulli}(p)$ , 1 = Know Roberts is Chief Justice

**What is the parameter of interest?**

$p$  = Proportion of voters *in the population* who know Roberts is Chief Justice.

**What is our estimator?**

Sample Proportion:  $\hat{p} = (\sum_{i=1}^n X_i)/n$

## Sample Proportion *is* the Sample Mean!

$X_1, \dots, X_n \sim \text{iid Bernoulli}(p)$

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n$$

## Sample Proportion *is* the Sample Mean!

$X_1, \dots, X_n \sim \text{iid Bernoulli}(p)$

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n$$

$$E[\hat{p}] = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E[X_i] = \frac{np}{n} = p$$



# Sample Proportion *is* the Sample Mean!

$X_1, \dots, X_n \sim \text{iid Bernoulli}(p)$

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n$$

$$E[\hat{p}] = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E[X_i] = \frac{np}{n} = p$$

$$\text{Var}(\hat{p}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n}$$

# Sample Proportion *is* the Sample Mean!

$X_1, \dots, X_n \sim \text{iid Bernoulli}(p)$

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n$$

$$E[\hat{p}] = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E[X_i] = \frac{np}{n} = p$$

$$\text{Var}(\hat{p}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n}$$

$$SE(\hat{p}) = \sqrt{\text{Var}(\hat{p})} = \sqrt{\frac{p(1-p)}{n}}$$

# Sample Proportion *is* the Sample Mean!

$X_1, \dots, X_n \sim \text{iid Bernoulli}(p)$

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n$$

$$E[\hat{p}] = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E[X_i] = \frac{np}{n} = p$$

$$\text{Var}(\hat{p}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n}$$

$$SE(\hat{p}) = \sqrt{\text{Var}(\hat{p})} = \sqrt{\frac{p(1-p)}{n}}$$

$$\widehat{SE}(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

# Central Limit Theorem Applied to Sample Proportion

## Central Limit Theorem: Intuition

Sample means are approximately normally distributed provided the sample size is large even if the population is non-normal.

### CLT For Sample Mean

$$\frac{\bar{X}_n - \mu}{\widehat{SE}(\bar{X}_n)} \approx N(0, 1)$$

### CLT for Sample Proportion

$$\frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \approx N(0, 1)$$

In this example, the population is Bernoulli( $p$ ) rather than normal. The sample mean is  $\hat{p}$  and the population mean is  $p$ .

## Approximate 95% CI for Population Proportion

$$\frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \approx N(0, 1)$$

$$P\left(-2 \leq \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \leq 2\right) \approx 0.95$$

$$P\left(\hat{p} - 2\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + 2\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right) \approx 0.95$$

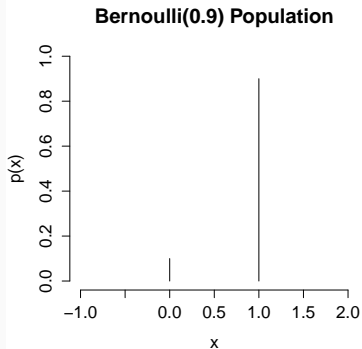
## $100 \times (1 - \alpha)$ CI for Population Proportion ( $p$ )

$X_1, \dots, X_n \sim \text{iid Bernoulli}(p)$

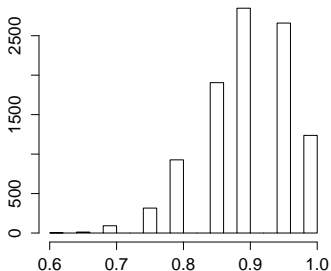
$$\hat{p} \pm \text{qnorm}(1 - \alpha/2) \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

Approximation based on the CLT. Works well provided  $n$  is large and  $p$  isn't too close to zero or one.

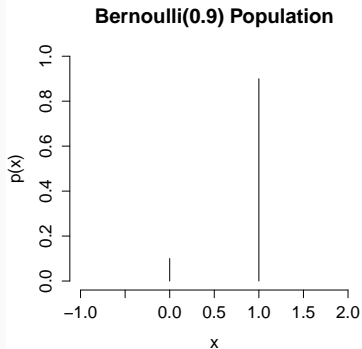
## Example: Bernoulli(0.9) Population, $n = 20$



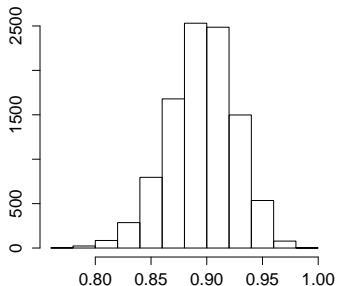
**Sample Mean – Ber(0.9) Pop ( $n = 20$ )**



## Example: Bernoulli(0.9) Population, $n = 100$



**Sample Mean – Ber(0.9) Pop ( $n = 100$ )**





## Approximate 95% CI for Population Proportion

39% of 771 Voters Polled Correctly Identified Chief Justice Roberts

$$\begin{aligned}\widehat{SE}(\widehat{p}) &= \sqrt{\frac{\widehat{p}(1 - \widehat{p})}{n}} = \sqrt{\frac{(0.39)(0.61)}{771}} \\ &\approx 0.018\end{aligned}$$

What is the ME for an approximate 95% confidence interval?

## Approximate 95% CI for Population Proportion

39% of 771 Voters Polled Correctly Identified Chief Justice Roberts

$$\begin{aligned}\widehat{SE}(\hat{p}) &= \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} = \sqrt{\frac{(0.39)(0.61)}{771}} \\ &\approx 0.018\end{aligned}$$

What is the ME for an approximate 95% confidence interval?

$$ME \approx 2 \times \widehat{SE}(\bar{X}_n) \approx 0.04$$

## Approximate 95% CI for Population Proportion

39% of 771 Voters Polled Correctly Identified Chief Justice Roberts

$$\begin{aligned}\widehat{SE}(\hat{p}) &= \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} = \sqrt{\frac{(0.39)(0.61)}{771}} \\ &\approx 0.018\end{aligned}$$

What is the ME for an approximate 95% confidence interval?

$$ME \approx 2 \times \widehat{SE}(\bar{X}_n) \approx 0.04$$

What can we conclude?

Approximate 95% CI: (0.35, 0.43)

# Are Republicans Better Informed Than Democrats?

Of the 239 Republicans surveyed, 47% correctly identified John Roberts as the current chief justice. Only 31% of the 238 Democrats surveyed correctly identified him. Is this difference meaningful or just sampling variation?

Again, assume random sampling.

# Confidence Interval for a Difference of Proportions

**What is the appropriate probability model for the sample?**

$X_1, \dots, X_n \sim \text{iid Bernoulli}(p)$  independently of

$Y_1, \dots, Y_m \sim \text{iid Bernoulli}(q)$

**What is the parameter of interest?**

The difference of population proportions  $p - q$

**What is our estimator?**

The difference of sample proportions:  $\hat{p} - \hat{q}$  where:

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i \quad \hat{q} = \frac{1}{m} \sum_{i=1}^m Y_i$$

# Difference of Sample Proportions $\hat{p} - \hat{q}$ and the CLT

## What We Have

Approx. sampling dist. for *individual* sample proportions from

$$\text{CLT: } \hat{p} \approx N\left(p, \widehat{SE}(\hat{p})^2\right), \quad \hat{q} \approx N\left(q, \widehat{SE}(\hat{q})^2\right)$$

# Difference of Sample Proportions $\hat{p} - \hat{q}$ and the CLT

## What We Have

Approx. sampling dist. for *individual* sample proportions from

$$\text{CLT: } \hat{p} \approx N\left(p, \widehat{SE}(\hat{p})^2\right), \quad \hat{q} \approx N\left(q, \widehat{SE}(\hat{q})^2\right)$$

## What We Want

Sampling Distribution of the *difference*  $\hat{p} - \hat{q}$

# Difference of Sample Proportions $\hat{p} - \hat{q}$ and the CLT

## What We Have

Approx. sampling dist. for *individual* sample proportions from

$$\text{CLT: } \hat{p} \approx N\left(p, \widehat{SE}(\hat{p})^2\right), \quad \hat{q} \approx N\left(q, \widehat{SE}(\hat{q})^2\right)$$

## What We Want

Sampling Distribution of the *difference*  $\hat{p} - \hat{q}$

## Use Independence of the Two Samples

$$\hat{p} - \hat{q} \approx N\left(p - q, \widehat{SE}(\hat{p})^2 + \widehat{SE}(\hat{q})^2\right)$$



# Difference of Sample Proportions $\hat{p} - \hat{q}$ and the CLT

## What We Have

Approx. sampling dist. for *individual* sample proportions from

$$\text{CLT: } \hat{p} \approx N\left(p, \widehat{SE}(\hat{p})^2\right), \quad \hat{q} \approx N\left(q, \widehat{SE}(\hat{q})^2\right)$$

## What We Want

Sampling Distribution of the *difference*  $\hat{p} - \hat{q}$

## Use Independence of the Two Samples

$$\hat{p} - \hat{q} \approx N\left(p - q, \widehat{SE}(\hat{p})^2 + \widehat{SE}(\hat{q})^2\right)$$

$$\implies \widehat{SE}(\hat{p} - \hat{q}) =$$

# Difference of Sample Proportions $\hat{p} - \hat{q}$ and the CLT

## What We Have

Approx. sampling dist. for *individual* sample proportions from

$$\text{CLT: } \hat{p} \approx N\left(p, \widehat{SE}(\hat{p})^2\right), \quad \hat{q} \approx N\left(q, \widehat{SE}(\hat{q})^2\right)$$

## What We Want

Sampling Distribution of the *difference*  $\hat{p} - \hat{q}$

## Use Independence of the Two Samples

$$\hat{p} - \hat{q} \approx N\left(p - q, \widehat{SE}(\hat{p})^2 + \widehat{SE}(\hat{q})^2\right)$$

$$\implies \widehat{SE}(\hat{p} - \hat{q}) = \sqrt{\widehat{SE}(\hat{p})^2 + \widehat{SE}(\hat{q})^2} =$$

# Difference of Sample Proportions $\hat{p} - \hat{q}$ and the CLT

## What We Have

Approx. sampling dist. for *individual* sample proportions from

$$\text{CLT: } \hat{p} \approx N\left(p, \widehat{SE}(\hat{p})^2\right), \quad \hat{q} \approx N\left(q, \widehat{SE}(\hat{q})^2\right)$$

## What We Want

Sampling Distribution of the *difference*  $\hat{p} - \hat{q}$

## Use Independence of the Two Samples

$$\hat{p} - \hat{q} \approx N\left(p - q, \widehat{SE}(\hat{p})^2 + \widehat{SE}(\hat{q})^2\right)$$

$$\implies \widehat{SE}(\hat{p} - \hat{q}) = \sqrt{\widehat{SE}(\hat{p})^2 + \widehat{SE}(\hat{q})^2} = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n} + \frac{\hat{q}(1 - \hat{q})}{m}}$$

## Approx. 95% CI for Difference of Population Proportions

$$\frac{(\hat{p} - \hat{q}) - (p - q)}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n} + \frac{\hat{q}(1-\hat{q})}{m}}} \approx N(0, 1)$$

$$P \left( -2 \leq \frac{(\hat{p} - \hat{q}) - (p - q)}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n} + \frac{\hat{q}(1-\hat{q})}{m}}} \leq 2 \right) \approx 0.95$$

$$(\hat{p} - \hat{q}) \pm \text{qnorm}(1 - \alpha/2) \sqrt{\frac{\hat{p}(1 - \hat{p})}{n} + \frac{\hat{q}(1 - \hat{q})}{m}}$$

## $100 \times (1 - \alpha)$ CI for Diff. of Popn. Proportions ( $p - q$ )

$X_1, \dots, X_n \sim \text{iid Bernoulli}(p)$  indep.  $Y_1, \dots, Y_n \sim \text{iid Bernoulli}(q)$

$$(\hat{p} - \hat{q}) \pm \text{qnorm}(1 - \alpha/2) \sqrt{\frac{\hat{p}(1 - \hat{p})}{n} + \frac{\hat{q}(1 - \hat{q})}{m}}$$

Approximation based on the CLT. Works well provided  $n, m$  large and  $p, q$  aren't too close to zero or one.

# ME for approx. 95% for Difference of Proportions

47% of 239 Republicans vs. 31% of 238 Democrats identified Roberts

## Republicans

$$\hat{p} = 0.47$$

$$n = 239$$

# ME for approx. 95% for Difference of Proportions

47% of 239 Republicans vs. 31% of 238 Democrats identified Roberts

## Republicans

$$\hat{p} = 0.47$$

$$n = 239$$

$$\widehat{SE}(\hat{p}) = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \approx 0.032$$

# ME for approx. 95% for Difference of Proportions

47% of 239 Republicans vs. 31% of 238 Democrats identified Roberts

## Republicans

$$\hat{p} = 0.47$$

$$n = 239$$

$$\widehat{SE}(\hat{p}) = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \approx 0.032$$

## Democrats

$$\hat{q} = 0.31$$

$$m = 238$$



# ME for approx. 95% for Difference of Proportions

47% of 239 Republicans vs. 31% of 238 Democrats identified Roberts

## Republicans

$$\hat{p} = 0.47$$

$$n = 239$$

$$\widehat{SE}(\hat{p}) = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \approx 0.032$$

## Democrats

$$\hat{q} = 0.31$$

$$m = 238$$

$$\widehat{SE}(\hat{q}) = \sqrt{\frac{\hat{q}(1 - \hat{q})}{m}} \approx 0.030$$

# ME for approx. 95% for Difference of Proportions

47% of 239 Republicans vs. 31% of 238 Democrats identified Roberts

## Republicans

$$\hat{p} = 0.47$$

$$n = 239$$

$$\widehat{SE}(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \approx 0.032$$

## Democrats

$$\hat{q} = 0.31$$

$$m = 238$$

$$\widehat{SE}(\hat{q}) = \sqrt{\frac{\hat{q}(1-\hat{q})}{m}} \approx 0.030$$

## Difference: (Republicans - Democrats)

$$\hat{p} - \hat{q} = 0.47 - 0.31 = 0.16$$

# ME for approx. 95% for Difference of Proportions

47% of 239 Republicans vs. 31% of 238 Democrats identified Roberts

## Republicans

$$\hat{p} = 0.47$$

$$n = 239$$

$$\widehat{SE}(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \approx 0.032$$

## Democrats

$$\hat{q} = 0.31$$

$$m = 238$$

$$\widehat{SE}(\hat{q}) = \sqrt{\frac{\hat{q}(1-\hat{q})}{m}} \approx 0.030$$

## Difference: (Republicans - Democrats)

$$\hat{p} - \hat{q} = 0.47 - 0.31 = 0.16$$

$$\widehat{SE}(\hat{p} - \hat{q}) = \sqrt{\widehat{SE}(\hat{p})^2 + \widehat{SE}(\hat{q})^2} \approx 0.044$$

# ME for approx. 95% for Difference of Proportions

47% of 239 Republicans vs. 31% of 238 Democrats identified Roberts

## Republicans

$$\hat{p} = 0.47$$

$$n = 239$$

$$\widehat{SE}(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \approx 0.032$$

## Democrats

$$\hat{q} = 0.31$$

$$m = 238$$

$$\widehat{SE}(\hat{q}) = \sqrt{\frac{\hat{q}(1-\hat{q})}{m}} \approx 0.030$$

## Difference: (Republicans - Democrats)

$$\hat{p} - \hat{q} = 0.47 - 0.31 = 0.16$$

$$\widehat{SE}(\hat{p} - \hat{q}) = \sqrt{\widehat{SE}(\hat{p})^2 + \widehat{SE}(\hat{q})^2} \approx 0.044 \implies ME \approx 0.09$$

# ME for approx. 95% for Difference of Proportions

47% of 239 Republicans vs. 31% of 238 Democrats identified Roberts

## Republicans

$$\hat{p} = 0.47$$

$$n = 239$$

$$\widehat{SE}(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \approx 0.032$$

## Democrats

$$\hat{q} = 0.31$$

$$m = 238$$

$$\widehat{SE}(\hat{q}) = \sqrt{\frac{\hat{q}(1-\hat{q})}{m}} \approx 0.030$$

## Difference: (Republicans - Democrats)

$$\hat{p} - \hat{q} = 0.47 - 0.31 = 0.16$$

$$\widehat{SE}(\hat{p} - \hat{q}) = \sqrt{\widehat{SE}(\hat{p})^2 + \widehat{SE}(\hat{q})^2} \approx 0.044 \implies ME \approx 0.09$$

Approximate 95% CI (0.07, 0.25)

What can we conclude?