# Confidence Intervals for Univariate Quantitative Data

# Reminder

The process of statistical analysis:

1. Identify population and parameter you are interested in.

   - Question: What is the average age at which BYU students find out Santa Claus isn't real? Specifically, is the average age at which BYU students find out Santa isn't real older than 8?

   - Parameter: The mean age at which all BYU students find out Santa Claus isn't real. We'll use the Greek letter $\mu$ to denote this value.

2. Collect data

   - A convenience sample of 1575 BYU students who are taking this course and completed the student survey.

3. Posit a statistical model based on information in the sample

   - Explore the data.

   - Posited a normal population model.

4. Draw inference about the population using your model.

# Types of Statistical Inference

3 ways of using sample to make inference about the population:

1. Point Estimation (last lecture notes)

2. Hypothesis Testing (last lecture notes)

3. Confidence Intervals

# An Issue with Hypothesis Testing

A student claims that the average age BYU students learn about Santa Claus is 8. I hypothesize that its older than that. Perform a hypothesis test for these claims.

Step 3 - Draw a conclusion

- Because the $p$-value is small. We say that our data is NOT consistent with the null hypothesis and that the mean is greater than 8.

- "Conclusions" from hypothesis tests are painfully vague!

- On the one hand, if we reject $H_0$, we still don't have a firm conclusion on what the value of the parameter is.

- On the other hand, if we don't reject $H_0$, we can't say $H_0$ is true because we assumed it was true.

# Confidence Intervals

**Goal:** Provide a range of reasonable values that the parameter could be.
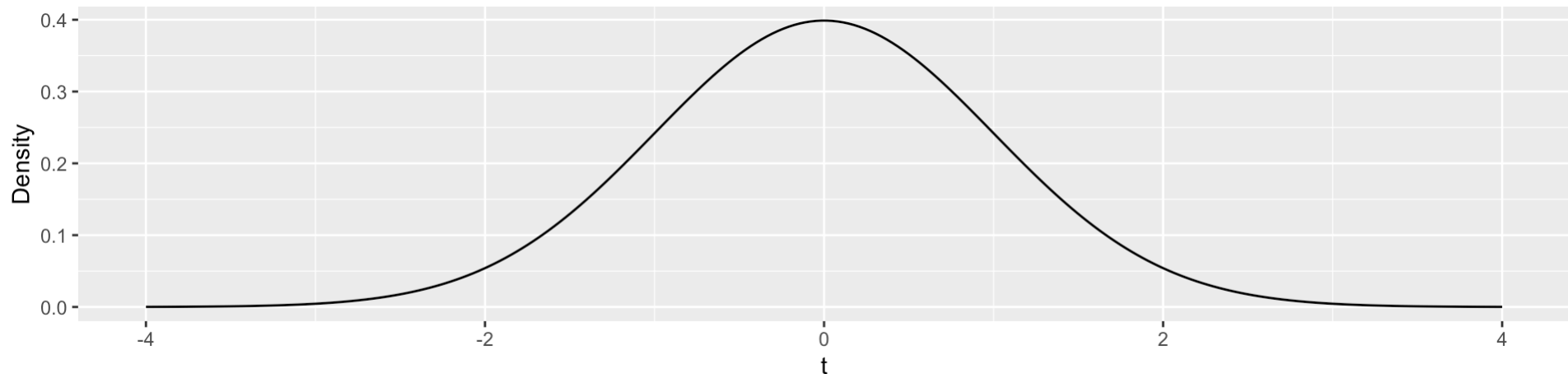
**Tool:** The sampling distribution of $t$.

# Constructing a Confidence Intervals

> **Theorem: Sampling distribution of t**
>
> If the normal population model is appropriate and the null hypothesis $H_0 : \mu = \mu_0$ is true, then
>
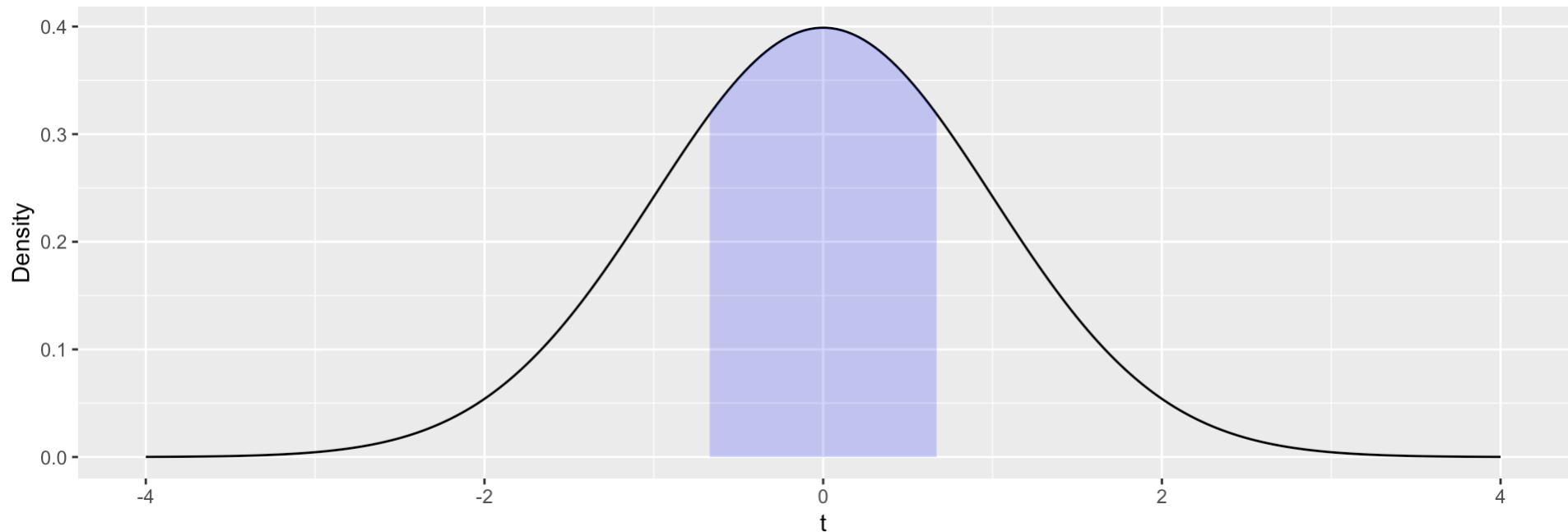> $$t = \frac{\bar{y} - \mu}{s/\sqrt{n}}$$
>
> is a standardized statistic and its sampling distribution is a t-distribution with center $0$, spread $1$ and degrees of freedom $n - 1$ where $n$ is the size of the sample.

# Constructing Confidence Intervals
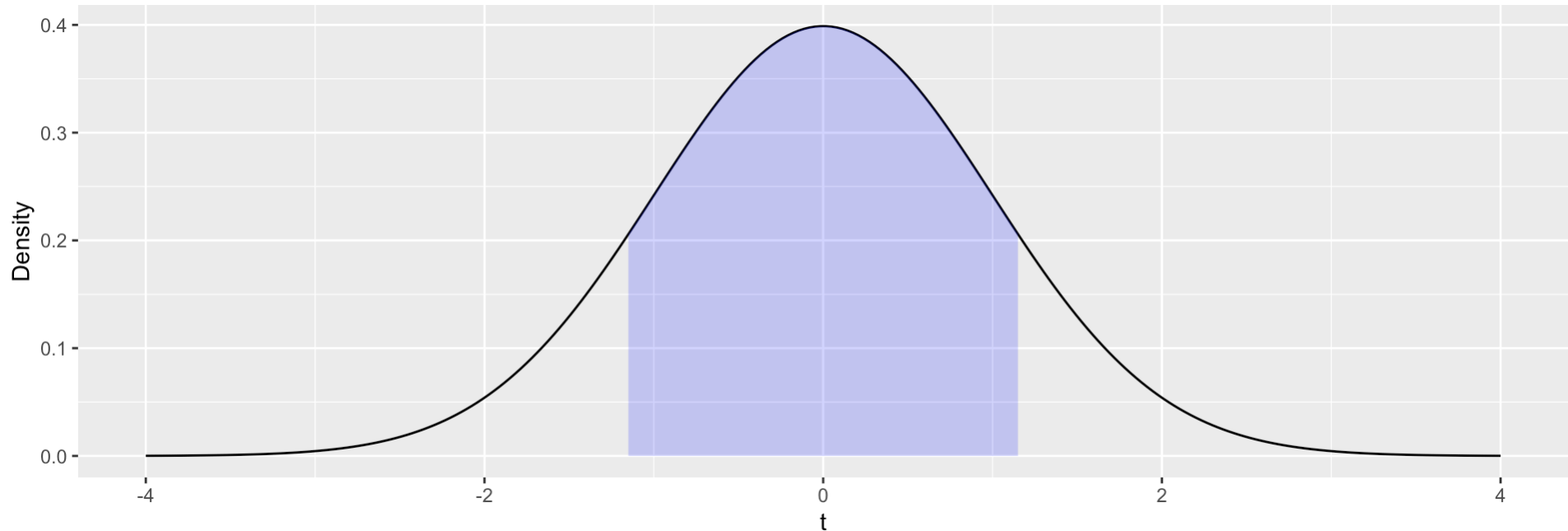
According to the $t$-distribution:

- 50% of possible $t = (\bar{y} - \mu)/(s/\sqrt{(n)})$ values are within 0.67 of 0.

# Constructing Confidence Intervals

According to the $t$-distribution:

- 50% of possible $t = (\bar{y} - \mu)/(s/\sqrt{(n)})$ values are within 0.67 of 0.

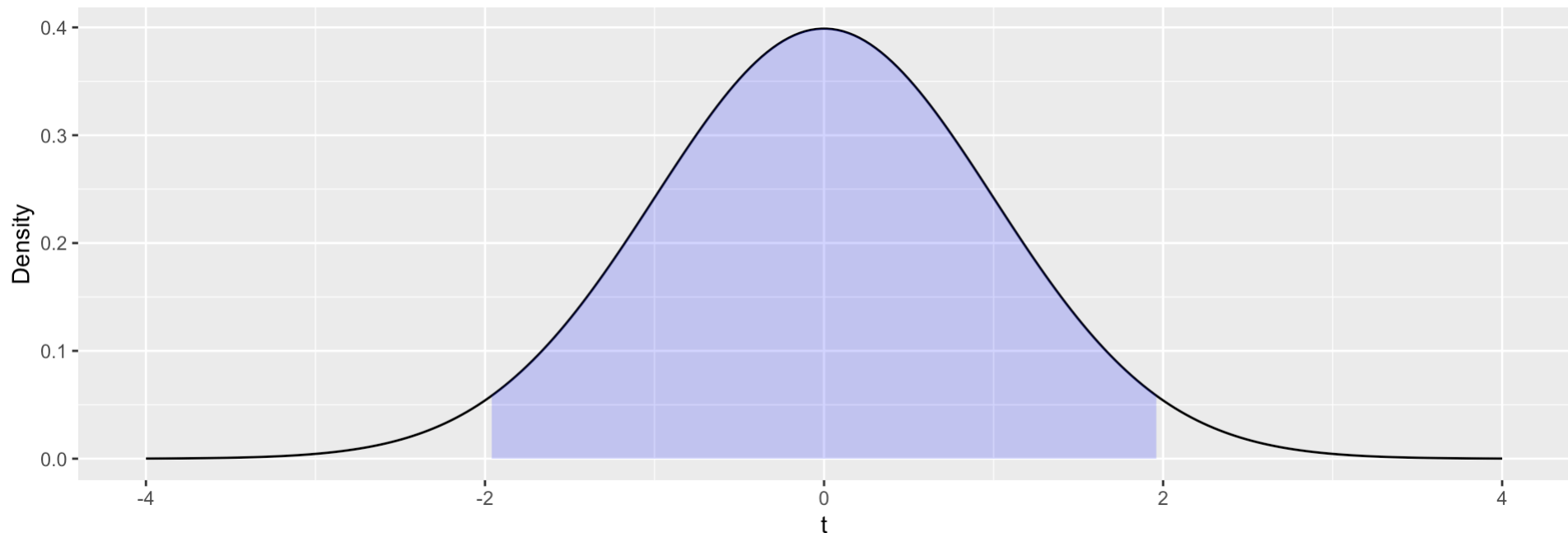- 75% of possible $t = (\bar{y} - \mu)/(s/\sqrt{(n)})$ values are within 1.15 of 0.

# Constructing Confidence Intervals

According to the $t$-distribution:

- 50% of possible $t = (\bar{y} - \mu)/(s/\sqrt{(n)})$ values are within 0.67 of 0.

- 75% of possible $t = (\bar{y} - \mu)/(s/\sqrt{(n)})$ values are within 1.15 of 0.

- 95% of possible $t = (\bar{y} - \mu)/(s/\sqrt{(n)})$ values are within 1.96 of 0.

# Constructing Confidence Intervals

Generally, C% of the time,

$$0 - t^{\star} < \underbrace{\frac{\bar{y} - \mu}{s/\sqrt{n}}}_{t} < 0 + t^{\star}$$

- But, we aren't interested in what $t$ is between, we are interested in what $\mu$ is between. So, lets rearrange this inequality using our algebra skills…

# Constructing Confidence Intervals

Generally, C% of the time,

$$0 - t^\star < \underbrace{\frac{\bar{y} - \mu}{s/\sqrt{n}}}_{t} < 0 + t^\star$$

Rearranging this inequality, we get

$$\bar{y} - t^\star \frac{s}{\sqrt{n}} < \mu < \bar{y} + t^\star \frac{s}{\sqrt{n}}$$

so that

$$\bar{y} \pm t^\star \frac{s}{\sqrt{n}}$$

is an interval estimate for $\mu$.

# The $t$-Confidence Interval for $\mu$

> **Theorem: Sampling distribution of t**
>
> *If the normal model is appropriate, a C% confidence interval for $\mu$ is*
>
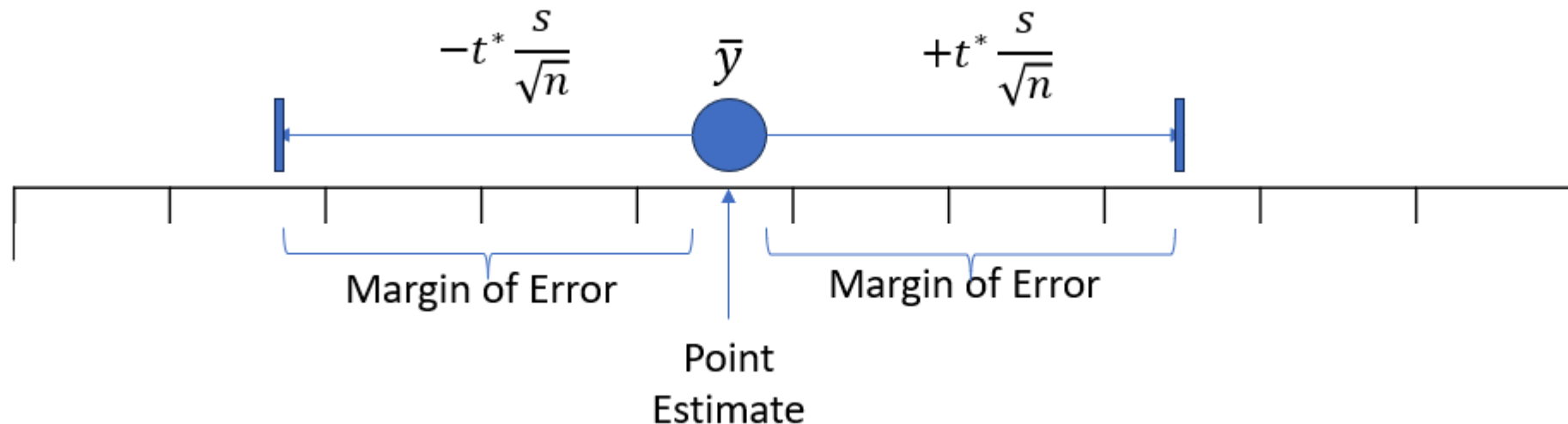> $$\bar{y} \pm t^{\star}\frac{s}{\sqrt{n}}$$

Terminology:

- $t^{\star}$ is a multiplier that corresponds with your chosen percentage $C$.

- The "$t^{\star}\frac{s}{\sqrt{n}}$" part is referred to as the margin of error.

- Note: the margin of error is equal to the $t^{\star}$ value times the standard error ($s/\sqrt{n}$).
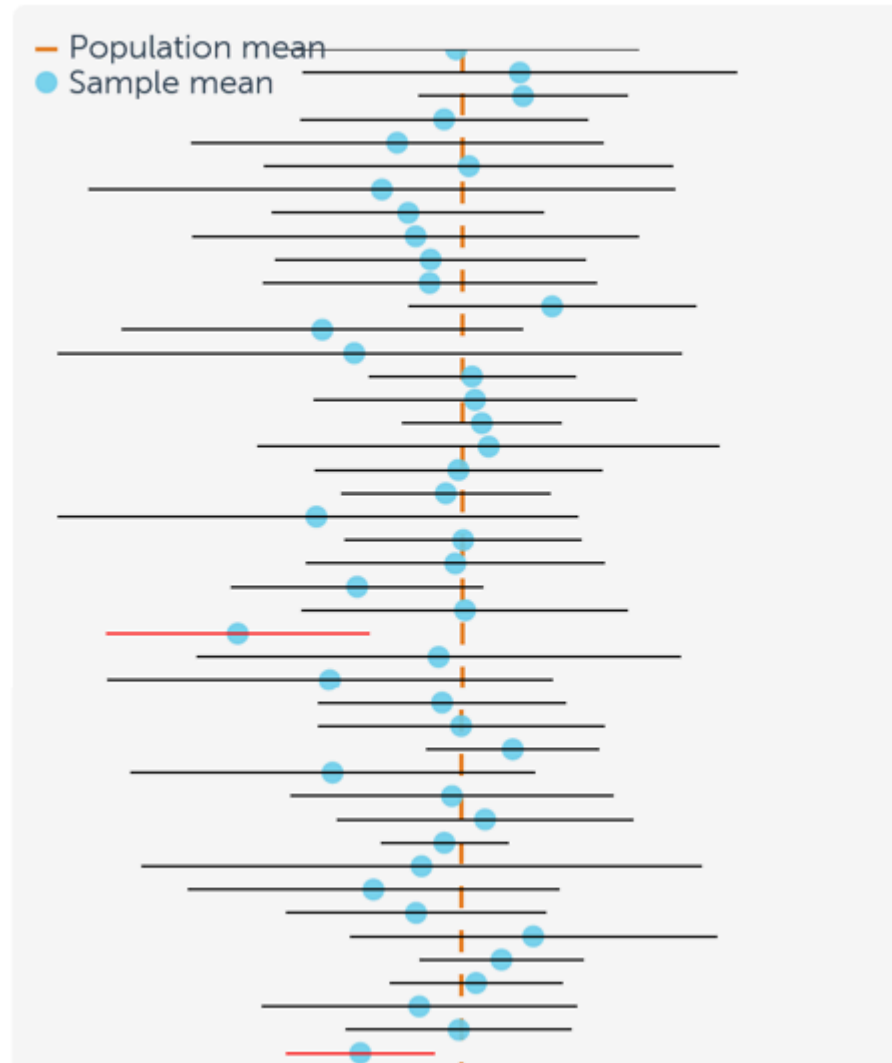
# The $t$-Confidence Interval for $\mu$

Interpreting a confidence interval:

- We are C% confident that $\mu$ is between $\bar{y} - t^\star \frac{s}{\sqrt{n}}$ and $\bar{y} + t^\star \frac{s}{\sqrt{n}}$.

- We have to say "confident" to reflect our belief or uncertainty that $\mu$ is between $\bar{y} - t^\star \frac{s}{\sqrt{n}}$ and $\bar{y} + t^\star \frac{s}{\sqrt{n}}$ (because it might not be).

- When we say C% confident we mean that, of all possible samples we could get from the population, C% of those samples will give an interval that captures $\mu$.

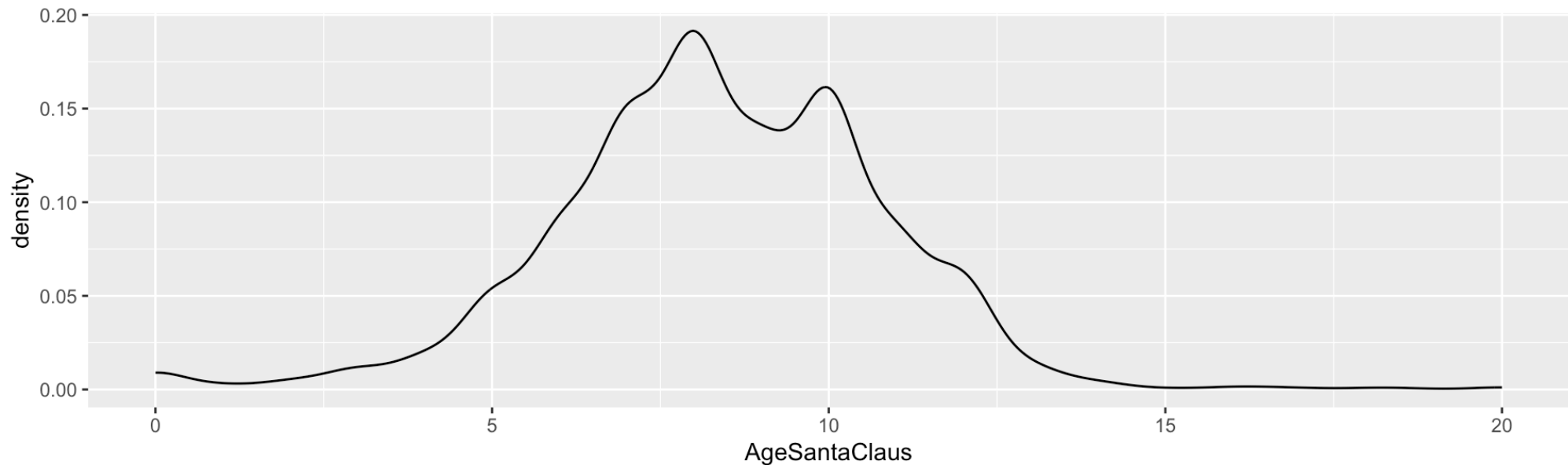# The $t$-Confidence Interval for $\mu$



95% confidence intervals

# Example: Santa Claus

What is the average age at which BYU students find out Santa Claus isn't real? Construct a 99% confidence interval for the average age of *all* BYU students when they found out the truth about Santa (which we'll denote by $\mu$).

- Step 1: Collect data (already done)

- Step 2: Check to make sure I can actually use the $t$ confidence interval (see if the normal model is appropriate).

# Example: Santa Claus

What is the average age at which BYU students find out Santa Claus isn't real? Construct a 99% confidence interval for the average age of *all* BYU students when they found out the truth about Santa (which we'll denote by $\mu$).

- Step 1: Collect data (already done)

- Step 2: Check to make sure I can actually use the $t$ confidence interval (see if the normal model is appropriate).

- Step 3: Have a computer build the confidence interval (I'll show you how to do this in a minute)

$$(8.2, 8.52)$$

- Step 4: Conclude - We are 99% confident that the average age of *all* BYU students when they found out Santa wasn't real is between 8.2 and 8.52.

# Example: Chlorine in Swimming Pools

From the previous chlorine analysis, using a 93% confidence interval, help the pool technician determine the average chlorine content across the whole pool.

Step 0 - Open up the course analysis app

Step 1 - Collect data (done)

Step 2 - Check to see if the $t$-distribution is appropriate.

Step 3 - Calculate the interval.

Step 4 - Draw conclusions

# Using the Tool

1. Make sure you are in the one mean section of the tool

**Stat 121 Analysis Tool**

- Exploratory Data Analysis
- Normal Probability Calculator
- Central Limit Theorem
- Analysis for Mean  ‹
  - » One Mean
  - » Two Means
  - » ANOVA
- Analysis For Proportions  ‹
- Regression  ‹

## One-Sample T Test for Means

### 1) Dataset Selection

**Data Selection**
- ● Use Preexisting Dataset
- ○ Upload Your Own Dataset

**Select dataset:**

| Chlorine ▼ |
|---|

2. Choose the dataset you are working with

Description: Data on the chlorine content (in ppm) in a pool.

Sample size: 30

☐ Display Dataset

[ Select This Dataset ]

### 2) Select Variables

**Please select the variable you wish to test (MUST be quantitative):**

| Chlorine ▼ |
|---|

3. Choose the variable that you want to analyze

[ Proceed to EDA ]

# Using the Tool

# Using the Tool



**4) Performing the Test and Calculating the CI**

Null Value:

2

*If only doing a conf int, you can ignore this*

Which sided hypothesis do you...

<

*If only doing a conf int, you can ignore this*

Confidence Level:

0.5                                                                    0.95        0.99

0.5        0.55        0.6        0.65        0.7        0.75        0.8        0.85        0.9        0.95        0.99

*Set the right level*

```
t Test for H0: Mean( Chlorine ) =  2
 Alternative Hypothesis =  less
 y-bar =  1.5392
 t Test statistic = -4.234661
 p-value =  0.0001054179
 95% Conf. Int.:  1.316646 1.761754
```

*The confidence interval*

# Example: Chlorine in Swimming Pools

From the previous chlorine analysis, using a 93% confidence interval, help the pool technician determine the average chlorine content across the whole pool.

Step 0 - Open up the [course analysis app](#)

Step 1 - Collect data (done)

Step 2 - Check to see if the $t$-distribution is appropriate.

- The density plot (or histogram) was normal.

Step 3 - Calculate the interval.

- (1.33448, 1.74392)

Step 4 - Draw conclusions

- We are 93% confident that the average chlorine content across the whole pool is between 1.33 and 1.74.

# Nuances of Confidence Intervals

What do we do if the sampling distribution of $t$ doesn't apply (most likely because the normal population model doesn't apply)?

> **Central Limit Theorem**
>
> *If the normal population model is not appropriate <u>BUT you have a large sample size</u>, the sampling distribution of t is still approximately a t-distribution with center $0$, spread $1$ and degrees of freedom $n-1$.*

Remember: The farther away from a normal model you are, the larger the sample size you will need in order to use the $t$-distribution.

- See the Central Limit Theorem part of the course analysis app

# Nuances of Confidence Intervals

Confidence Level & Margin of Error

- Confidence level = the % confident you want to be
- Margin of Error = $t^\star \frac{s}{\sqrt{n}}$ = amount above and below point estimate we think $\mu$ might be

Important relations:

- As confidence level increases so does margin of error (size of the interval is larger)
- A 100% confidence interval is $(-\infty, \infty)$.
- As sample size goes up, margin of error goes down (good thing)
- Choose confidence level to balance width and confidence in the interval (95% is actually a pretty good balance most of the time)

# Nuances of Confidence Intervals

## Confidence Intervals

- Give a range of reasonable values for the parameter

- Useful if you don't have a hypothesis

- Useful if you are trying to estimate the parameter value

## Hypothesis Tests

- A single conclusion about the validity of a hypothesis.

- Commonly used to assess a "difference" or an "effect".

- Answers yes/no questions about the population

# Nuances of Confidence Intervals

Connection: You can use a CI to perform a 2-sided Hypothesis Test

Example: A 99% confidence interval for the average age at which all BYU students learn the truth about Santa Claus is 8.2 and 8.52.

- Based on this interval, can we say that $\mu$ is different than 9.5? Why or why not?

- Based on this interval, can we say that $\mu$ is different than 8.4? Why or why not?

# Nuances of Confidence Intervals

Connection: You can use a CI to perform a 2-sided Hypothesis Test

Example: A 99% confidence interval for the average age at which all BYU students learn the truth about Santa Claus is 8.2 and 8.52.

- Based on this interval, can we say that $\mu$ is different than 9.5?
  - Yes because 9.5 is not in the interval
- Based on this interval, can we say that $\mu$ is different than 8.4?
  - No because it is in the interval.
- <u>Rules:</u> A C% CI corresponds to a two-sided hypothesis test using $\alpha = (1 - C/100)$ (for example, 95% and 0.05 or 90% and 0.1)

# Key Terminology

- $t$-confidence interval
- confidence level
- margin of error
- confident
- relationship between hypothesis testing and CIs