- Homework assignment going up today
- Next in-class lab will be next Wednesday
- I'm trying to get caught up on grading
- I'll be sending out more information this weekend about the final projects
- Starting Ch 4 today. We'll double back to the last bits of 3 later if we have time
- Polling: `rembold-class.ddns.net`

## Friday the 13th

Between 1990 - 1992 researchers in the UK collected data on traffic flow, accidents, and hospital admissions on Friday 13th and the previous Friday, Friday 6th. Below is an excerpt from this data set on traffic flow. We can assume that traffic flow on given day at locations 1 and 2 are independent.

|    | type    | date            | 6th    | 13th   | diff | location |
|----|---------|-----------------|--------|--------|------|----------|
| 1  | traffic | 1990, July      | 139246 | 138548 | 698  | loc 1    |
| 2  | traffic | 1990, July      | 134012 | 132908 | 1104 | loc 2    |
| 3  | traffic | 1991, September | 137055 | 136018 | 1037 | loc 1    |
| 4  | traffic | 1991, September | 133732 | 131843 | 1889 | loc 2    |
| 5  | traffic | 1991, December  | 123552 | 121641 | 1911 | loc 1    |
| 6  | traffic | 1991, December  | 121139 | 118723 | 2416 | loc 2    |
| 7  | traffic | 1992, March     | 128293 | 125532 | 2761 | loc 1    |
| 8  | traffic | 1992, March     | 124631 | 120249 | 4382 | loc 2    |
| 9  | traffic | 1992, November  | 124609 | 122770 | 1839 | loc 1    |
| 10 | traffic | 1992, November  | 117584 | 117263 | 321  | loc 2    |

**Friday the 13th**

- We want to investigate if people's behavior is different on Friday $13^{th}$ compared to Friday $6^{th}$.
- One approach is to compare the traffic flow on these two days.
- $H_0$ : Average traffic flow on Friday $6^{th}$ and $13^{th}$ are equal.
  $H_A$ : Average traffic flow on Friday $6^{th}$ and $13^{th}$ are different.

- We want to investigate if people's behavior is different on Friday 13th compared to Friday 6th.
- One approach is to compare the traffic flow on these two days.
- $H_0$ : Average traffic flow on Friday 6th and 13th are equal.
  $H_A$ : Average traffic flow on Friday 6th and 13th are different.

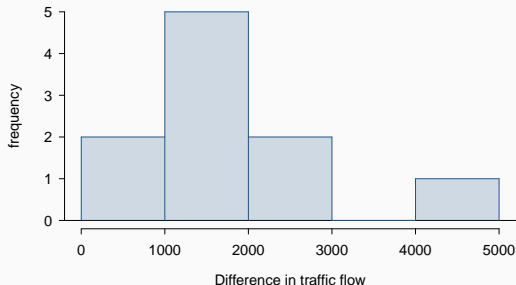Each case in the data set represents traffic flow recorded at the same location in the same month of the same year: one count from Friday 6th and the other Friday 13th. Are these two counts independent?
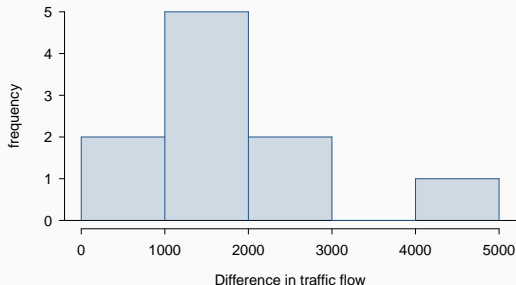
## Conditions

- **_Independence:_** We are told to assume that cases (rows) are independent.

- **_Sample size / skew:_** The sample distribution does not appear to be extremely skewed, but it's very difficult to assess with such a small sample size.

- We do not know $\sigma$ and *n* is too small to assume *s* is a reliable estimate for $\sigma$.

## Conditions

- **_Independence:_** We are told to assume that cases (rows) are independent.

- **_Sample size / skew:_** The sample distribution does not appear to be extremely skewed, but it's very difficult to assess with such a small sample size.



- We do not know $\sigma$ and $n$ is too small to assume $s$ is a reliable estimate for $\sigma$.

So what do we do when the sample size is small?

4

As long as observations are independent, and the population distribution is not extremely skewed, a large sample would ensure that...

- the sampling distribution of the mean is nearly normal
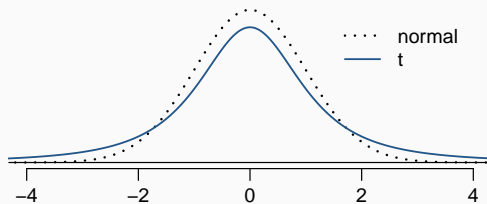- the estimate of the standard error, as $\frac{s}{\sqrt{n}}$, is reliable

## The normality condition

- The CLT, which states that sampling distributions will be nearly normal, holds true for any sample size as long as the population distribution is nearly normal.
- While this is a helpful special case, it's inherently difficult to verify normality in small data sets.
- We should exercise caution when verifying the normality condition for small samples. It is important to not only examine the data but also think about where the data come from.
  - For example, ask: would I expect this distribution to be symmetric, and am I confident that outliers are rare?
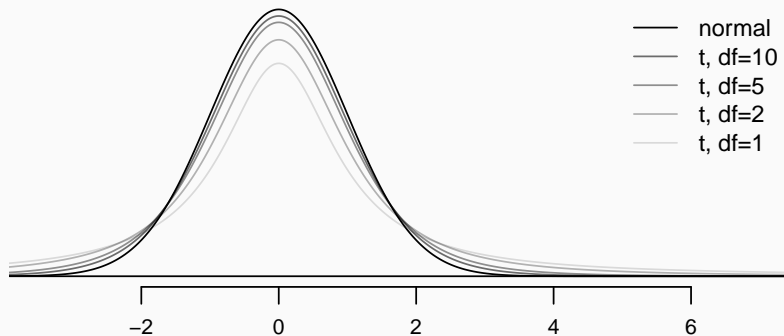
## The $t$ distribution

- When the population standard deviation is unknown (almost always), the uncertainty of the standard error estimate is addressed by using a new distribution: the $t$ **distribution**.
- This distribution also has a bell shape, but its tails are **thicker** than the normal model's.
- Therefore observations are more likely to fall beyond two SDs from the mean than under the normal distribution.
- These extra thick tails are helpful for resolving our problem with a less reliable estimate the standard error (since $n$ is small)
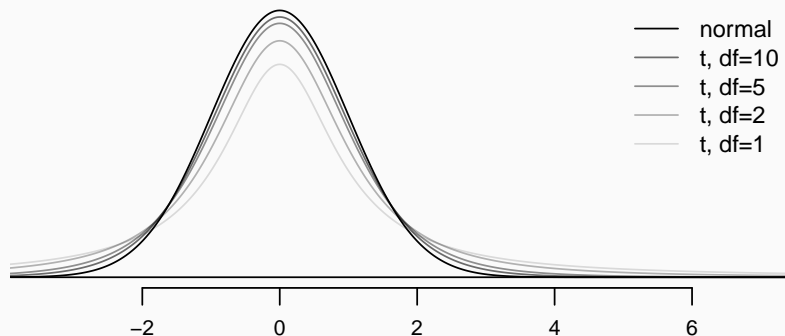
## The *t* distribution (cont.)

- Always centered at zero, like the standard normal (*z*) distribution.
- Has a single parameter: ***degrees of freedom*** (*df*).

- Always centered at zero, like the standard normal (*z*) distribution.
- Has a single parameter: ***degrees of freedom*** (*df*).



What happens to shape of the *t* distribution as *df* increases?

| | type | date | 6th | 13th | diff | location |
|---|---|---|---|---|---|---|
| 1 | traffic | 1990, July | 139246 | 138548 | 698 | loc 1 |
| 2 | traffic | 1990, July | 134012 | 132908 | 1104 | loc 2 |
| 3 | traffic | 1991, September | 137055 | 136018 | 1037 | loc 1 |
| 4 | traffic | 1991, September | 133732 | 131843 | 1889 | loc 2 |
| 5 | traffic | 1991, December | 123552 | 121641 | 1911 | loc 1 |
| 6 | traffic | 1991, December | 121139 | 118723 | 2416 | loc 2 |
| 7 | traffic | 1992, March | 128293 | 125532 | 2761 | loc 1 |
| 8 | traffic | 1992, March | 124631 | 120249 | 4382 | loc 2 |
| 9 | traffic | 1992, November | 124609 | 122770 | 1839 | loc 1 |
| 10 | traffic | 1992, November | 117584 | 117263 | 321 | loc 2 |

$$\downarrow$$

$$\bar{x}_{diff} = 1836$$

$$s_{diff} = 1176$$

9

# Finding the test statistic

## Test statistic for inference on a small sample mean

The test statistic for inference on a small sample ($n < 50$) mean is the $T$ statistic with $df = n - 1$.

$$T_{df} = \frac{\text{point estimate} - \text{null value}}{SE}$$

**Test statistic for inference on a small sample mean**

The test statistic for inference on a small sample ($n < 50$) mean is the $T$ statistic with $df = n - 1$.

$$T_{df} = \frac{\text{point estimate} - \text{null value}}{SE}$$

*In context...*

$$\text{point estimate} \;\; = \;\; \bar{x}_{diff} = 1836$$

## Finding the test statistic

### Test statistic for inference on a small sample mean

The test statistic for inference on a small sample ($n < 50$) mean is the $T$ statistic with $df = n - 1$.

$$T_{df} = \frac{\text{point estimate} - \text{null value}}{SE}$$

*In context...*

$$
\begin{aligned}
\text{point estimate} &= \bar{x}_{diff} = 1836 \\
SE &= \frac{s_{diff}}{\sqrt{n}} = \frac{1176}{\sqrt{10}} = 372
\end{aligned}
$$

## Finding the test statistic

### Test statistic for inference on a small sample mean

The test statistic for inference on a small sample ($n < 50$) mean is the $T$ statistic with $df = n - 1$.

$$T_{df} = \frac{\text{point estimate} - \text{null value}}{SE}$$

*In context...*

$$
\begin{aligned}
\text{point estimate} &= \bar{x}_{diff} = 1836 \\
SE &= \frac{s_{diff}}{\sqrt{n}} = \frac{1176}{\sqrt{10}} = 372 \\
T &= \frac{1836 - 0}{372} = 4.94
\end{aligned}
$$

## Finding the test statistic

### Test statistic for inference on a small sample mean

The test statistic for inference on a small sample ($n < 50$) mean is the $T$ statistic with $df = n - 1$.

$$T_{df} = \frac{\text{point estimate} - \text{null value}}{SE}$$

*In context...*

$$
\begin{aligned}
\text{point estimate} &= \bar{x}_{diff} = 1836 \\
SE &= \frac{s_{diff}}{\sqrt{n}} = \frac{1176}{\sqrt{10}} = 372 \\
T &= \frac{1836 - 0}{372} = 4.94 \\
df &= 10 - 1 = 9
\end{aligned}
$$

_____

*Note: Null value is 0 because in the null hypothesis we set $\mu_{diff} = 0$.*

**Finding the p-value**

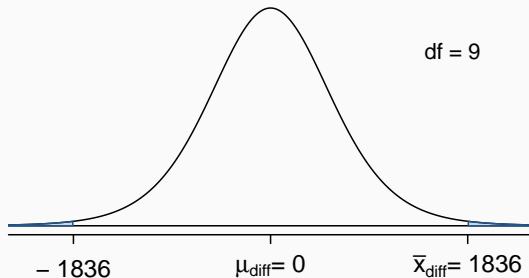- The p-value is, once again, calculated as the tail area under the *t* distribution.
- Using R:

```
> 2 * pt(4.94, df = 9, lower.tail = FALSE)

[1] 0.0008022394
```

- Using a web app:
  https://gallery.shinyapps.io/dist_calc/
- Or when these aren't available, we can use a *t*-table.

## Finding the p-value (cont.)

| one tail | 0.100 | 0.050 | 0.025 | 0.010 | 0.005 |
|---|---|---|---|---|---|
| two tails | 0.200 | 0.100 | 0.050 | 0.020 | 0.010 |
| df   6 | 1.44 | 1.94 | 2.45 | 3.14 | 3.71 |
| 7 | 1.41 | 1.89 | 2.36 | 3.00 | 3.50 |
| 8 | 1.40 | 1.86 | 2.31 | 2.90 | 3.36 |
| 9 | 1.38 | 1.83 | 2.26 | 2.82 | 3.25 |
| 10 | 1.37 | 1.81 | 2.23 | 2.76 | 3.17 |

df = 9

$-1836$        $\mu_{diff} = 0$        $\bar{x}_{diff} = 1836$

# Finding the p-value (cont.)

| one tail | 0.100 | 0.050 | 0.025 | 0.010 | 0.005 |
|---|---|---|---|---|---|
| two tails | 0.200 | 0.100 | 0.050 | 0.020 | 0.010 |
| df   6 | 1.44 | 1.94 | 2.45 | 3.14 | 3.71 |
| 7 | 1.41 | 1.89 | 2.36 | 3.00 | 3.50 |
| 8 | 1.40 | 1.86 | 2.31 | 2.90 | 3.36 |
| 9 | 1.38 | 1.83 | 2.26 | 2.82 | 3.25 |
| 10 | 1.37 | 1.81 | 2.23 | 2.76 | 3.17 |



df = 9

$-1836$      $\mu_{diff} = 0$      $\bar{x}_{diff} = 1836$

# Finding the p-value (cont.)

| one tail | 0.100 | 0.050 | 0.025 | 0.010 | 0.005 | |
|---|---|---|---|---|---|---|
| two tails | 0.200 | 0.100 | 0.050 | 0.020 | 0.010 | → |
| df  6 | 1.44 | 1.94 | 2.45 | 3.14 | 3.71 | |
| 7 | 1.41 | 1.89 | 2.36 | 3.00 | 3.50 | |
| 8 | 1.40 | 1.86 | 2.31 | 2.90 | 3.36 | |
| 9 | 1.38 | 1.83 | 2.26 | 2.82 | 3.25 | → |
| 10 | 1.37 | 1.81 | 2.23 | 2.76 | 3.17 | |

df = 9

$-1836$         $\mu_{diff} = 0$         $\bar{x}_{diff} = 1836$

## Finding the p-value (cont.)

| one tail | 0.100 | 0.050 | 0.025 | 0.010 | 0.005 | |
|---|---|---|---|---|---|---|
| two tails | 0.200 | 0.100 | 0.050 | 0.020 | 0.010 | → |
| df  6 | 1.44 | 1.94 | 2.45 | 3.14 | 3.71 | |
| 7 | 1.41 | 1.89 | 2.36 | 3.00 | 3.50 | |
| 8 | 1.40 | 1.86 | 2.31 | 2.90 | 3.36 | |
| 9 | 1.38 | 1.83 | 2.26 | 2.82 | 3.25 | → |
| 10 | 1.37 | 1.81 | 2.23 | 2.76 | 3.17 | |



df = 9

− 1836          $\mu_{diff}$ = 0          $\bar{x}_{diff}$ = 1836

What is the conclusion of the hypothesis test?

12

## Finding the p-value (cont.)

| one tail | 0.100 | 0.050 | 0.025 | 0.010 | 0.005 | |
|---|---|---|---|---|---|---|
| two tails | 0.200 | 0.100 | 0.050 | 0.020 | 0.010 | → |
| df    6 | 1.44 | 1.94 | 2.45 | 3.14 | 3.71 | |
| 7 | 1.41 | 1.89 | 2.36 | 3.00 | 3.50 | |
| 8 | 1.40 | 1.86 | 2.31 | 2.90 | 3.36 | |
| 9 | 1.38 | 1.83 | 2.26 | 2.82 | 3.25 | → |
| 10 | 1.37 | 1.81 | 2.23 | 2.76 | 3.17 | |



df = 9

What is the conclusion of the hypothesis test?

The data provide convincing evidence of a difference between traffic flow on Friday $6^{th}$ and $13^{th}$.

$- 1836$     $\mu_{diff} = 0$     $\bar{x}_{diff} = 1836$

12

**What is the difference?**

- We concluded that there is a difference in the traffic flow between Friday $6^{th}$ and $13^{th}$.
- But it would be more interesting to find out what exactly this difference is.
- We can use a confidence interval to estimate this difference.

**Confidence interval for a small sample mean**

- Confidence intervals are always of the form

  point estimate ± *ME*

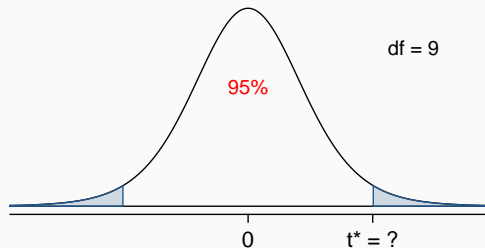**Confidence interval for a small sample mean**

- Confidence intervals are always of the form

$$\text{point estimate} \pm ME$$

- ME is always calculated as the product of a critical value and SE.

**Confidence interval for a small sample mean**

- Confidence intervals are always of the form

  point estimate $\pm$ *ME*

- ME is always calculated as the product of a critical value and SE.
- Since small sample means follow a *t* distribution (and not a *z* distribution), the critical value is a $t^\star$ (as opposed to a $z^\star$).
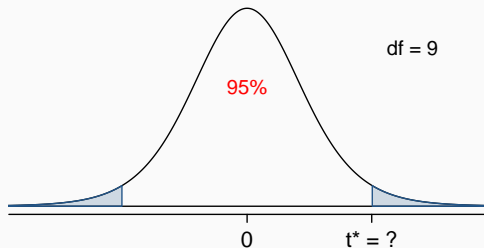
  point estimate $\pm\, t^\star \times SE$

# Finding the critical $t$ ($t^\star$)



df = 9

95%

0      t* = ?

$n = 10$, $df = 10 - 1 = 9$, $t^\star$ is at the intersection of row $df = 9$ and two tail probability 0.05.

| one tail | 0.100 | 0.050 | 0.025 | 0.010 | 0.005 |
|---|---|---|---|---|---|
| two tails | 0.200 | 0.100 | 0.050 | 0.020 | 0.010 |
| df    6 | 1.44 | 1.94 | 2.45 | 3.14 | 3.71 |
| 7 | 1.41 | 1.89 | 2.36 | 3.00 | 3.50 |
| 8 | 1.40 | 1.86 | 2.31 | 2.90 | 3.36 |
| 9 | 1.38 | 1.83 | 2.26 | 2.82 | 3.25 |
| 10 | 1.37 | 1.81 | 2.23 | 2.76 | 3.17 |

15

## Finding the critical $t$ ($t^\star$)



$df = 9$

95%

0    $t^\star = ?$

$n = 10$, $df = 10 - 1 = 9$, $t^\star$ is at the intersection of row $df = 9$ and two tail probability 0.05.

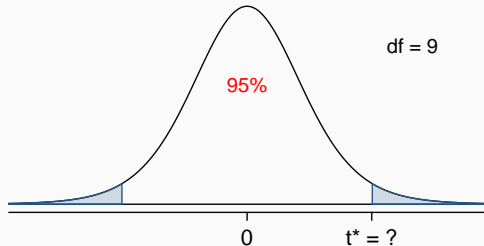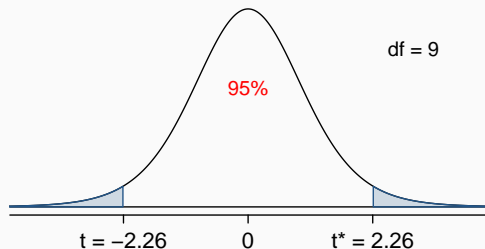| one tail | 0.100 | 0.050 | 0.025 | 0.010 | 0.005 |
|---|---|---|---|---|---|
| two tails | 0.200 | 0.100 | 0.050 | 0.020 | 0.010 |
| df    6 | 1.44 | 1.94 | 2.45 | 3.14 | 3.71 |
| 7 | 1.41 | 1.89 | 2.36 | 3.00 | 3.50 |
| 8 | 1.40 | 1.86 | 2.31 | 2.90 | 3.36 |
| 9 | 1.38 | 1.83 | 2.26 | 2.82 | 3.25 |
| 10 | 1.37 | 1.81 | 2.23 | 2.76 | 3.17 |

## Finding the critical $t$ ($t^\star$)



$n = 10$, $df = 10 - 1 = 9$, $t^\star$ is at the intersection of row $df = 9$ and two tail probability 0.05.

| one tail | 0.100 | 0.050 | 0.025 | 0.010 | 0.005 |
|---|---|---|---|---|---|
| two tails | 0.200 | 0.100 | 0.050 | 0.020 | 0.010 |
| df    6 | 1.44 | 1.94 | 2.45 | 3.14 | 3.71 |
| 7 | 1.41 | 1.89 | 2.36 | 3.00 | 3.50 |
| 8 | 1.40 | 1.86 | 2.31 | 2.90 | 3.36 |
| 9 | 1.38 | 1.83 | 2.26 | 2.82 | 3.25 |
| 10 | 1.37 | 1.81 | 2.23 | 2.76 | 3.17 |

15

# Finding the critical $t$ ($t^\star$)



$n = 10$, $df = 10 - 1 = 9$, $t^\star$ is at the intersection of row $df = 9$ and two tail probability 0.05.

| one tail | 0.100 | 0.050 | 0.025 | 0.010 | 0.005 |
|---|---|---|---|---|---|
| two tails | 0.200 | 0.100 | 0.050 | 0.020 | 0.010 |
| df     6 | 1.44 | 1.94 | 2.45 | 3.14 | 3.71 |
| 7 | 1.41 | 1.89 | 2.36 | 3.00 | 3.50 |
| 8 | 1.40 | 1.86 | 2.31 | 2.90 | 3.36 |
| 9 | 1.38 | 1.83 | 2.26 | 2.82 | 3.25 |
| 10 | 1.37 | 1.81 | 2.23 | 2.76 | 3.17 |

**Constructing a CI for a small sample mean**

Which of the following is the correct calculation of a 95% confidence interval for the difference between the traffic flow between Friday $6^{th}$ and $13^{th}$?

$$\bar{x}_{diff} = 1836 \qquad s_{diff} = 1176 \qquad n = 10 \qquad SE = 372$$

A) $1836 \pm 1.96 \times 372$

B) $1836 \pm 2.26 \times 372$

C) $1836 \pm -2.26 \times 372$

D) $1836 \pm 2.26 \times 1176$

## Constructing a CI for a small sample mean

Which of the following is the correct calculation of a 95% confidence interval for the difference between the traffic flow between Friday $6^{th}$ and $13^{th}$?

$$\bar{x}_{diff} = 1836 \qquad s_{diff} = 1176 \qquad n = 10 \qquad SE = 372$$

A) $1836 \pm 1.96 \times 372$

B) $1836 \pm 2.26 \times 372 \rightarrow (995, 2677)$

C) $1836 \pm -2.26 \times 372$

D) $1836 \pm 2.26 \times 1176$

## Interpreting the CI

Which of the following is the best interpretation for the confidence interval we just calculated?

$$\mu_{diff:6th-13th} = (995, 2677)$$

We are 95% confident that ...

A) the difference between the average number of cars on the road on Friday $6^{th}$ and $13^{th}$ is between 995 and 2,677.

B) on Friday $6^{th}$ there are 995 to 2,677 fewer cars on the road than on the Friday $13^{th}$, on average.

C) on Friday $6^{th}$ there are 995 fewer to 2,677 more cars on the road than on the Friday $13^{th}$, on average.

D) on Friday $13^{th}$ there are 995 to 2,677 fewer cars on the road than on the Friday $6^{th}$, on average.

Which of the following is the best interpretation for the confidence interval we just calculated?

$$\mu_{diff:6th-13th} = (995, 2677)$$

We are 95% confident that ...

A) the difference between the average number of cars on the road on Friday 6th and 13th is between 995 and 2,677.

B) on Friday 6th there are 995 to 2,677 fewer cars on the road than on the Friday 13th, on average.

C) on Friday 6th there are 995 fewer to 2,677 more cars on the road than on the Friday 13th, on average.

D) on Friday 13th there are 995 to 2,677 fewer cars on the road than on the Friday 6th, on average.

Does the conclusion from the hypothesis test agree with the findings of the confidence interval?

Do you think the findings of this study suggests that people believe Friday 13th is a day of bad luck?

**Synthesis**

Does the conclusion from the hypothesis test agree with the findings of the confidence interval?

- Yes, the hypothesis test found a significant difference, and the CI does not contain the null value of 0.

Do you think the findings of this study suggests that people believe Friday 13th is a day of bad luck?

## Synthesis

Does the conclusion from the hypothesis test agree with the findings of the confidence interval?

- Yes, the hypothesis test found a significant difference, and the CI does not contain the null value of 0.

Do you think the findings of this study suggests that people believe Friday 13th is a day of bad luck?

- No, this is an observational study. We have just observed a significant difference between the number of cars on the road on these two days. We have not tested for people's beliefs.

## Recap: Inference using the *t*-distribution

- If $\sigma$ is unknown, use the *t*-distribution with $SE = \frac{s}{\sqrt{n}}$.
- Conditions:
    - independence of observations (often verified by a random sample, and if sampling without replacement, $n < 10\%$ of population)
    - no extreme skew
- Hypothesis testing:

$$T_{df} = \frac{\text{point estimate} - \text{null value}}{SE}, \text{ where } df = n - 1$$

- Confidence interval:

$$\text{point estimate} \pm t_{df}^{\star} \times SE$$

———————

*Note: The example we used was for paired means (difference between dependent groups). We took the difference between the observations and used only these differences (one sample) in our analysis, therefore the mechanics are the same as when we are working with just one sample.*