# Lab 5B Confidence intervals

## 2025-03-06

```r
install.packages('tidyverse')
```

```
## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.4'
## (as 'lib' is unspecified)
```

```r
install.packages('openintro')
```

```
## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.4'
## (as 'lib' is unspecified)
```

```r
install.packages('infer')
```

```
## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.4'
## (as 'lib' is unspecified)
```

```
## also installing the dependency 'ggplot2'
```

```r
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----------------------- tidyverse 2.0.0 --
## v dplyr     1.1.4     v readr     2.1.5
## v forcats   1.0.0     v stringr   1.5.1
## v ggplot2   3.5.2     v tibble    3.2.1
## v lubridate 1.9.4     v tidyr     1.3.1
## v purrr     1.0.4
## -- Conflicts ------------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become error
```

```r
library(openintro)
```

```
## Loading required package: airports
## Loading required package: cherryblossom
## Loading required package: usdata
```
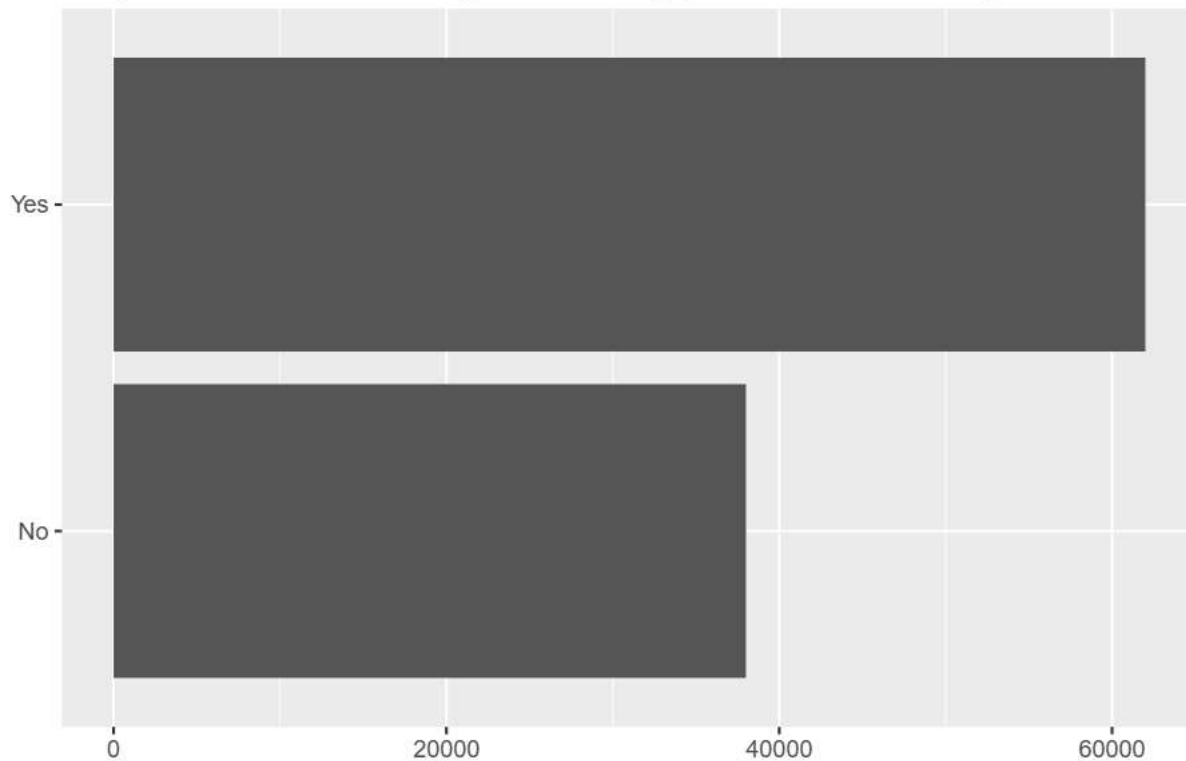
```r
library(infer)
```

```r
us_adults <- tibble(
  climate_change_affects = c(rep("Yes", 62000), rep("No", 38000))
)
```

```r
ggplot(us_adults, aes(x = climate_change_affects)) +
  geom_bar() +
  labs(
    x = "", y = "",
    title = "Do you think climate change is affecting your local community?"
  ) +
  coord_flip()
```

## Do you think climate change is affecting your local community?



```
us_adults %>%
  count(climate_change_affects) %>%
  mutate(p = n /sum(n))
```

```
## # A tibble: 2 x 3
##   climate_change_affects     n      p
##   <chr>                  <int>  <dbl>
## 1 No                     38000   0.38
## 2 Yes                    62000   0.62
```

```
n <- 60
samp <- us_adults %>%
  sample_n(size = n)
```

Exercise 1

What percent of the adults in your sample think climate change affects their local community? Hint: Just like we did with the population, we can calculate the proportion of those in this sample who think climate change affects their local community.

```
n <- 60
samp <- us_adults %>%
  sample_n(size = n)
```

```
samp %>%
count(climate_change_affects) %>%
mutate(p = n /sum(n))
```

```
## # A tibble: 2 x 3
##   climate_change_affects     n      p
##   <chr>                  <int>  <dbl>
```

```
## 1 No                     27   0.45
## 2 Yes                    33   0.55
```

Exercise 2

Would you expect another student's sample proportion to be identical to yours? Would you expect it to be similar? Why or why not?

No i don't expect another student's proportion to be identical because the sample of 60 people might not be exactly the same proportion. It might be close.

```
samp %>%
  specify(response = climate_change_affects, success = "Yes") %>%
  generate(reps = 1000, type = "bootstrap") %>%
  calculate(stat = "prop") %>%
  get_ci(level = 0.95)
```

```
## # A tibble: 1 x 2
##   lower_ci upper_ci
##      <dbl>    <dbl>
## 1    0.433    0.683
```

Exercise 3

In the interpretation above, we used the phrase "95% confident". What does "95% confidence" mean?

It means that if we were to repeatedly take a sample from the same Us_adult population. We expect the true population parameter to fall within the 95% confidence interval as specified 'get_ci(level = 0.95)' for the proportion of US adults who think climate change affects their local community.

Exercise 4

Does your confidence interval capture the true population proportion of US adults who think climate change affects their local community? If you are working on this lab in a classroom, does your neighbor's interval capture this value?

Because the confidence level is 95% it captures a true population proportion. If i was in a classroom and other students used the same sample size our intervals would be close but not exactly the same.

Exercise 5

Each student should have gotten a slightly different confidence interval. What proportion of those intervals would you expect to capture the true population mean? Why?

The larger sample size provides a more precise estimate of the population mean because it narrows the confidence interval and margin of error.

Exercise 6

Given a sample size of 60, 1000 bootstrap samples for each interval, and 50 confidence intervals constructed (the default values for the above app), what proportion of your confidence intervals include the true population proportion? Is this proportion exactly equal to the confidence level? If not, explain why. Make sure to include your plot in your answer.

With a 95% confidence level if we repeat the sampling process 95% of those intervals contain the true population proportion.

Exercise 7

Using code from the infer package and data from the one sample you have (samp), find a confidence interval for the proportion of US Adults who think climate change is affecting their local community with a confidence level of your choosing (other than 95%) and interpret it.

```
samp %>%
  specify(response = climate_change_affects, success = "Yes") %>%
  generate(reps = 1000, type = "bootstrap") %>%
  calculate(stat = "prop") %>%
  get_ci(level = 0.99)
```

```
## # A tibble: 1 x 2
##   lower_ci upper_ci
##      <dbl>    <dbl>
## 1      0.4    0.717
```

I set the confidence level that change affects their local community to 99%.The lower confidence interval which is a narrower range that is less certain is 0.46 and the upper confidence interval that contains the true estimate of the population is 0.76.

Exercise 9

Using the app, calculate 50 confidence intervals at the confidence level you chose in the previous question, and plot all intervals on one plot, and calculate the proportion of intervals that include the true population proportion. How does this percentage compare to the confidence level selected for the intervals?

This percentage compared to the confidence level selected for the intervals is higher.

Exercise 10

Lastly, try one more (different) confidence level. First, state how you expect the width of this interval to compare to previous ones you calculated. Then, calculate the bounds of the interval using the infer package and data from samp and interpret it. Finally, use the app to generate many intervals and calculate the proportion of intervals that are capture the true population proportion.

The lower_Ci is slighly higher with 79% interval while the upper_ci is higher with a higher 99% confidence level.

```
samp %>%
  specify(response = climate_change_affects, success = "Yes") %>%
  generate(reps = 1000, type = "bootstrap") %>%
  calculate(stat = "prop") %>%
  get_ci(level = 0.79)
```

```
## # A tibble: 1 x 2
##   lower_ci upper_ci
##      <dbl>    <dbl>
## 1    0.467    0.633
```

Exercise 11

Using the app, experiment with different sample sizes and comment on how the widths of intervals change as sample size changes (increases and decreases).

As long as the confidence level is the same different sample sizes lead to the close and similar widths of intervals.

Exercise 12

Finally, given a sample size (say, 60), how does the width of the interval change as you increase the number of bootstrap samples. Hint: Does changing the number of bootstap samples affect the standard error?

Yes. As the bootstrap increases the standard error decreases.