Foundations for statistical inference - Confidence intervals

Title: CUNY SPS MDS DATA606_LAB5B"

Author: Charles Ugiagbe

Date: "10/9/2021"

Getting Started

In this lab, we will explore and visualize the data using the **tidyverse** suite of packages, and perform statistical inference using **infer**.

Let's load the packages.

library(tidyverse)

```
## Warning: package 'tidyverse' was built under R version 4.1.1
```

Warning: package 'tibble' was built under R version 4.1.1

Warning: package 'readr' was built under R version 4.1.1

library(openintro)

Warning: package 'openintro' was built under R version 4.1.1

library(infer)

Warning: package 'infer' was built under R version 4.1.1

The data

A 2019 Pew Research report states the following:

To keep our computation simple, we will assume a total population size of 100,000 (even though that's smaller than the population size of all US adults).

Roughly six-in-ten U.S. adults (62%) say climate change is currently affecting their local community either a great deal or some, according to a new Pew Research Center survey.

Source: Most Americans say climate change impacts their community, but effects vary by region

In this lab, you will assume this 62% is a true population proportion and learn about how sample proportions can vary from sample to sample by taking smaller samples from the population. We will first create our population assuming a population size of 100,000. This means 62,000 (62%) of the adult population think climate change impacts their community, and the remaining 38,000 does not think so.

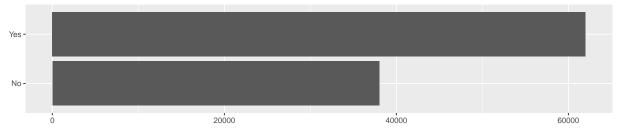
```
us_adults <- tibble(
  climate_change_affects = c(rep("Yes", 62000), rep("No", 38000))
)</pre>
```

The name of the data frame is us_adults and the name of the variable that contains responses to the question "Do you think climate change is affecting your local community?" is climate_change_affects.

We can quickly visualize the distribution of these responses using a bar plot.

```
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?



We can also obtain summary statistics to confirm we constructed the data frame correctly.

In this lab, you'll start with a simple random sample of size 60 from the population.

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

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.

Solution 1:

```
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
```

66.7% of US adults in the sample think that climate change affects their local community.

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?

Solution 2:

I wouldn't expect another student's sample proportion to be identical as mine because the each sample is unique and they are ramdomly chosen.

Confidence intervals

Return for a moment to the question that first motivated this lab: based on this sample, what can you infer about the population? With just one sample, the best estimate of the proportion of US adults who think climate change affects their local community would be the sample proportion, usually denoted as \hat{p} (here we are calling it p_hat). That serves as a good **point estimate**, but it would be useful to also communicate how uncertain you are of that estimate. This uncertainty can be quantified using a **confidence interval**.

One way of calculating a confidence interval for a population proportion is based on the Central Limit Theorem, as $\hat{p} \pm z^* S E_{\hat{p}}$ is, or more precisely, as

$$\hat{p} \pm z^* \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Another way is using simulation, or to be more specific, using **bootstrapping**. The term **bootstrapping** comes from the phrase "pulling oneself up by one's bootstraps", which is a metaphor for accomplishing an impossible task without any outside help. In this case the impossible task is estimating a population parameter (the unknown population proportion), and we'll accomplish it using data from only the given sample. Note that this notion of saying something about a population parameter using only information from an observed sample is the crux of statistical inference, it is not limited to bootstrapping.

In essence, bootstrapping assumes that there are more of observations in the populations like the ones in the observed sample. So we "reconstruct" the population by resampling from our sample, with replacement. The bootstrapping scheme is as follows:

This code will find the 95 percent confidence interval for proportion of US adults who think climate change affects their local community.

```
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 0.667
```

- In specify we specify the response variable and the level of that variable we are calling a success.
- In generate we provide the number of resamples we want from the population in the reps argument (this should be a reasonably large number) as well as the type of resampling we want to do, which is "bootstrap" in the case of constructing a confidence interval.
- Then, we calculate the sample statistic of interest for each of these resamples, which is proportion.

Feel free to test out the rest of the arguments for these functions, since these commands will be used together to calculate confidence intervals and solve inference problems for the rest of the semester. But we will also walk you through more examples in future chapters.

To recap: even though we don't know what the full population looks like, we're 95% confident that the true proportion of US adults who think climate change affects their local community is between the two bounds reported as result of this pipeline.

Confidence levels

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

Solution 3:

A 95% confidence interval mean that we are 95% sure that a population parameter will fall between a set of values for a certain proportion of times. Hence, it means that we are 95% confident that the population parameter will fall within 2 standard deviations from the mean.

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?

Solution 4:

Yes. The confident interval captures the true population proportion of US adult. The interval is (0.55, 0.78); while the population that believe climate change affect is 0.62. The population proportion falls within the range.

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?

Solution 5:

Yes each confidence interval would have been slightly different. However, I would expect at least 95% of confidence intervals to capture the true population value because each interval was constructed with a 95% level of confidence to ensure that the interval captures the true population value.

In the next part of the lab, you will collect many samples to learn more about how sample proportions and confidence intervals constructed based on those samples vary from one sample to another.

- Obtain a random sample.
- Calculate the sample proportion, and use these to calculate and store the lower and upper bounds of the confidence intervals.
- Repeat these steps 50 times.

Doing this would require learning programming concepts like iteration so that you can automate repeating running the code you've developed so far many times to obtain many (50) confidence intervals. In order to keep the programming simpler, we are providing the interactive app below that basically does this for you and created a plot similar to Figure 5.6 on OpenIntro Statistics, 4th Edition (page 182).

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.

Solution 6:

98% of confidence intervals include the true population proportion. This is not exactly equal to the confidence interval, but it is within the interpretation of the confidence interval.

More Practice

7. Choose a different confidence level than 95%. Would you expect a confidence interval at this level to me wider or narrower than the confidence interval you calculated at the 95% confidence level? Explain your reasoning.

Solution 7:

When I choose 99% confidence level, the confidence interval became wider because it is higher than 95% level of confidence. Increasing the confidence level basically means widening the interval to be sure that it contains the population parameter.

8. Using code from the **infer** package and data fromt 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.

Solution 8:

```
samp %>%
  specify(response = climate_change_affects, success = "Yes") %>%
  generate(reps = 1000, type = "bootstrap") %>%
  calculate(stat = "prop") %>%
  get_ci(level = 0.90)

## # A tibble: 1 x 2
## lower_ci upper_ci
## <dbl> <dbl>
## 1 0.45 0.667
```

The confidence interval at 99% level of confidence for the samp data set is (0.566, 0.766). This means that we are 90% confident that the population proportion will fall within that interval.

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?

Solution 9:

When I use the app at 99% confidence level, 50 confidence intervals, and 100 bootstraps, the proportion of intervals that include the true population proportion is 98%. It is slightly lower than the confidence level.

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.

Solution 10:

1

0.417

0.7

```
samp %>%
  specify(response = climate_change_affects, success = "Yes") %>%
  generate(reps = 1000, type = "bootstrap") %>%
  calculate(stat = "prop") %>%
  get_ci(level = 0.98)

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

Using the app for 99% confidence level, the proportion of intervals that capture the true population proportion is 99%. Trying another confidence level of 90%, the confidence interval is (0.516, 0.8) which means that we are 90% confident that the population mean will lie within the interval.

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).

Solution 11:

Using the app, we discover that as sample sizes increase the width of the interval decreases i.e there is an inverse relationship between sample sizes and the width. Hence, the more we increase the sample size, the less the spread.

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?

Solution 12:

As i increase the number of bootstrap samples, the width of the interval reduces and the standard Error also reduces as the distribution tend to nearly normal.