Lab 3: Foundations for statistical inference - Confidence intervals

If you have access to data on an entire population, say the size of every house in Ames, Iowa, it's straight forward to answer questions like, "How big is the typical house in Ames?" and "How much variation is there in sizes of houses?". If you have access to only a sample of the population, as is often the case, the task becomes more complicated. What is your best guess for the typical size if you only know the sizes of several dozen houses? This sort of situation requires that you use your sample to make inference on what your population looks like.

Template for lab report

The data

We consider real estate data from the city of Ames, Iowa. The details of every real estate transaction in Ames is recorded by the City Assessor's office. Our particular focus for this lab will be all residential home sales in Ames between 2006 and 2010. This collection represents our population of interest. In this lab we would like to learn about these home sales by taking smaller samples from the full population. Let's load the data.

download.file("http://www.openintro.org/stat/data/ames.RData", destfile = "ames.RData")

load("ames.RData")

We see that there are quite a few variables in the data set, enough to do a very in-depth analysis. For this lab, we'll start with a simple random sample of size 60 from the population. Specifically, this is a simple random sample of size 60. Note that the data set has information on many housing variables, but for the first portion of the lab we'll focus on the size of the house, represented by the variable Gr.Liv.Area.

population <- ames\$Gr.Liv.Area

Before we do the sampling we'll once again set a seed. Choose a family member's birthday and enter that in the set.seed function below. For example, if the birthday is October 4, you would enter 104.

```
set.seed(104)  # make sure to change the seed
samp <- sample(population, 60)</pre>
```

Exercise 1 Describe the distribution of your sample. What would you say is the "typical" size within your sample? Also state precisely what you interpreted "typical" to mean. Make sure to include a visualization in your response.

Exercise 2 Now compare your distribution to another class mate's or the professor's. Do they look similar? Are they identical? Why, or why not?

This is a product of OpenIntro that is released under a Creative Commons Attribution-NonCommercial-NoDerivs 3.0 Unported (*http://creativecommons.org/licenses/by-nc-nd/3.0*/). This lab was written for OpenIntro by Andrew Bray and Mine Çetinkaya-Rundel.

Confidence intervals

One of the most common ways to describe the typical or central value of a distribution is to use the mean. In this case we can calculate the mean of the:

sample_mean <- mean(samp)</pre>

Return for a moment to the question that first motivated this lab: based on this sample, what can we infer about the population? Based only on this single sample, the best estimate of the average living area of houses sold in Ames would be the sample mean, usually denoted as \bar{x} (here we're calling it sample_mean). That serves as a good *point estimate* but it would be useful to also communicate how uncertain we are of that estimate. This can be captured by using a *confidence interval*.

We can calculate a 95% confidence interval for a sample mean by adding and subtracting 1.96 standard errors to the point estimate.^{\dagger}

```
se <- sd(samp)/sqrt(60)
lower <- sample_mean - 1.96 * se
upper <- sample_mean + 1.96 * se
c(lower, upper)</pre>
```

Exercise 3 Interpret this interval in context of the data. Make sure to include all relevant code for calculating the interval in your report.

This is an important inference that we've just made: even though we don't know what the full population looks like, we're 95% confident that the true average size of houses in Ames lies between the values lower and upper. There are a few conditions that must be met for this interval to be valid.

Exercise 4 For the confidence interval to be valid, the sample mean must be normally distributed and have standard error s/\sqrt{n} . What conditions must be met for this to be true? Explain why you believe your sample does/does not meet each of these conditions.

Confidence levels

Exercise 5 What does "95% confidence" mean? Be specific. If you're not sure, see Section 4.2.2.

In this case we have the luxury of knowing the true population mean since we have data on the entire population. This value can be calculated using the following command:

mean(population)

Exercise 6 Does your confidence interval capture the true average size of houses in Ames?

Exercise 7 Each student in your class should have gotten a slightly different confidence interval.

⁺See Section 4.2.3 if you are unfamiliar with this formula.

What proportion of those intervals would you expect to capture the true population mean? Why?

Using R, we're going to recreate many samples to learn more about how sample means and confidence intervals vary from one sample to another. *Loops* come in handy here.

Here is the rough outline:

- (1) Obtain a random sample.
- (2) Calculate the sample's mean and standard deviation.
- (3) Use these statistics to calculate a confidence interval.
- (4) Repeat steps (1)-(3) 50 times.

But before we do all of this, we need to first create empty vectors where we can save the means and standard deviations that will be calculated from each sample. And while we're at it, let's also store the desired sample size as n.

```
samp_mean <- rep(NA, 50)
samp_sd <- rep(NA, 50)
n <- 60</pre>
```

Now we're ready for the loop where we calculate the means and standard deviations of 50 random samples.

```
for(i in 1:50){
  samp <- sample(population, n) # obtain a sample of size n = 60 from the population
  samp_mean[i] <- mean(samp) # save sample mean in ith element of samp_mean
  samp_sd[i] <- sd(samp) # save sample sd in ith element of samp_sd
}</pre>
```

Interlude: The for loop

Let's take a break from the statistics for a moment to let that last block of code sink in. You have just run your first for loop, a cornerstone of computer programming. The idea behind the for loop is *iteration*: it allows you to execute code as many times as you want without having to type out every iteration. In the case above, we wanted to iterate the three lines of code inside the curly braces that take a random sample of size 60 from population then save the mean of that sample into the sample_mean vector and the standard deviation of the sample to the sample_sd vector. Without the for loop, this would be painful:

```
sample_mean <- rep(0, 50)
samp <- sample(population, n)
sample_mean[1] <- mean(samp)
sample_sd[1] <- sd(samp)
samp <- sample(population, n)
sample_mean[2] <- mean(samp)
sample_sd[2] <- sd(samp)</pre>
```

```
samp <- sample(population, n)
sample_mean[3] <- mean(samp)
sample_sd[3] <- sd(samp)
samp <- sample(population, n)
sample_mean[4] <- mean(samp)
sample_sd[4] <- sd(samp)</pre>
```

and so on ...

With the for loop, these many lines of code are compressed into a handful of lines. We've added one extra line to the code below, which prints the variable i during each iteration of the for loop. Run this code.

```
samp_mean <- rep(NA, 50)
samp_sd <- rep(NA, 50)
for (i in 1:50) {
   samp <- sample(population, n)
   samp_mean[i] <- mean(samp)
   samp_sd[i] <- sd(samp)
   print(i)
}</pre>
```

Let's consider this code line by line to figure out what it does. In the first two lines we *initialized a vector*. In this case, we created a vector of 50 NAs called sample_mean and another vector of 50 NAs called sample_sd. These vectors store values generated within the for loop.

The third line calls the for loop itself. The syntax can be loosely read as, "for every element i from 1 to 50, run the following lines of code". You can think of i as the counter that keeps track of which loop you're on. Therefore, more precisely, the loop will run once when i=1, then once when i=2, and so on up to i=50.

The body of the for loop is the part inside the curly braces, and this set of code is run for each value of i. Here, on every loop, we take a random sample of size 60 from population, take its mean, and store it as the ith element of sample_mean. We also take its standard deviation, and store it as the ith element of sample_sd.

In order to display that this is really happening, we asked R to print i at each iteration. This line of code is optional and is only used for displaying what's going on while the for loop is running.

The for loop allows us to not just run the code 50 times, but to neatly package the results, element by element, into the empty vector that we initialized at the outset.

Exercise 8 To make sure you understand what you've done in this loop, try running a smaller version. Initialize a vector of 100 zeros called sample_mean_practice. Run a loop that takes a sample of size 50 from population and stores the sample mean in sample_mean_practice, but only iterate from 1 to 100. Print the output to your screen (type sample_mean_practice into the console and press enter). How many elements are there in this object called sample_mean_practice? What does each element represent?

Back to confidence intervals

Using the previously created vectors sample_mean and sample_sd which store sample means and standard deviations from 50 random samples of size n = 60, we construct the confidence intervals.

lower <- samp_mean - 1.96 * samp_sd/sqrt(n)
upper <- samp_mean + 1.96 * samp_sd/sqrt(n)</pre>

Lower bounds of these 50 confidence intervals are stored in lower, and the upper bounds are in upper.

Exercise 9 View the first interval using the code below. Does this interval capture the true population mean? Note that you will need to include earlier code calculating lower and upper in your report for the below command to work.

c(lower[1], upper[1])

Exercise 10 Using the following function (which was downloaded with the data set), plot all intervals. What proportion of your confidence intervals include the true population mean? Is this proportion exactly equal to the confidence level? If no, explain why.[†]

plot_ci(lower, upper, mean(population))

Exercise 11 Pick a confidence level of your choosing, provided it is not 95%. What is the appropriate critical value (z^*) ?

Exercise 12 Calculate 50 confidence intervals at the confidence level you chose in Exercise 10. You do not need to obtain new samples, simply calculate new intervals based on the sample means and standard deviations you have already collected. Using the plot_ci function, plot all intervals and calculate the proportion of intervals that include the true population mean. How does this percentage compare to the confidence level level selected for the intervals?

Exercise 13 What concepts from the textbook are covered in this lab? What concepts, if any, are not covered in the textbook? Have you seen these concepts elsewhere, e.g. lecture, discussion section, previous labs, or homework problems? Be specific in your answer.

⁺This figure should look familiar, if not, see Section 4.2.2.