

Introduction to Inference

Single Categorical Variable

Learning Objectives

- Identify the two possible explanations (one assuming the null hypothesis, and one assuming the alternative hypothesis) for a relationship seen in sample data
- Given a research question, construct the null and alternative hypotheses in words and using appropriate statistical symbols
- Describe and perform simulation-based hypothesis tests for a single proportion
- Interpret and evaluate a p-value
- Use bootstrapping to find a confidence interval for a single proportion

Terminology

These are a few of the terms we will cover in today's activity.

- Parameter of Interest
- Null Hypothesis
- Alternative Hypothesis
- Null distribution
- p-value
- Bootstrapping
- Confidence Interval

To review these concepts see Chapter 5 in your textbook

Steps of Statistical Investigation

We will work through a six step process to complete a hypothesis test for a single proportion.

- **Ask a research question** that can be addressed by collecting data. What are the researchers trying to show.
- **Design a study and collect data.** This step involves selecting the people or objects to be studied and how to gather relevant data on them.
- **Summarize and Visualize the data.** Calculate summary statistics and create graphical plots that best represent the research question.
- **Use statistical analysis methods to draw inferences from the data.** Choose an analysis technique appropriate for the data and identify the p-value. In this study, we will focus on using randomization.
- **Communicate the results and answer the research question.** Using the p-value and confidence interval from the analysis, determine whether the data provide statistical evidence against the null hypothesis.

- **Revisit and look forward** to point out limitations of the study and suggest new studies that could be performed to build on the findings of the study

Handedness of Male Boxers

Left-handedness is a trait that is found in about 10% of the population. The fighting claim states that left-handed men have an advantage in competition. Past studies have shown that left-handed men are over-represented among professional fighters. In this random sample of 500 male boxers we will see if there is an over-prevalence of left-handed fighters.

Summary Statistics Review

1. What are the observational units?
2. What variable are we testing? Is it categorical or quantitative?
3. What type of plot would be appropriate to visually display the data?
4. Write out in context the statistic will we calculate to summarize the data.

Ask a Research Question.

5. Identify the research question for this study.

Design a Study and Collect Data

6. What is the target population for this study?
7. Did the researchers use a biased or an unbiased method of selection? Explain your answer.

Summarize and Visualize the Data

```
# Counts for Handedness  
tally(~Stance, data=handedness_sub, margins=T)
```

```
## Stance  
## left-handed right-handed      Total  
##           81           419       500
```

```
#Tally creates a table with a count for each level of the variable
```

8. Calculate the appropriate summary statistic that represents the research question. Use appropriate notation.

Use statistical analysis methods to draw inferences from the data

When testing data we must first identify the null hypothesis. The null hypothesis is written about the parameter of interest, the true value of interest.

9. Write out the parameter of interest. (Hint: the true proportion of...)

10. We will assume that the true proportion of male boxers who are left handed is the same as the general population, 0.1. Using the parameter of interest in question 9, write out the null hypothesis in words.

The notation used for a true proportion is π . Since this summarizes a population, it is a parameter. When writing the null hypothesis in notation we set the parameter equal to the null value, $H_0 : \pi = \pi_0$

11. Write the null hypothesis in notation using the null value of 0.1.

The alternative hypothesis is the claim to be tested and the direction is based on the research question.

12. Based on the research question from question 5, are we testing that the parameter is greater than 0.1, less than 0.1 or different than 0.1?

13. Write out the alternative hypothesis in words.

14. Write out the alternative hypothesis in notation.

Remember that when utilizing a hypothesis test, we are evaluating two competing possibilities. For this study the **two possibilities** are either...

- The true proportion of male boxers who are left handed is 0.1 and our results just occurred by random chance or
- The true proportion of male boxers who are left handed is greater than 0.1 and our results reflect this

Notice that these two competing possibilities represent the null and alternative hypotheses.

The null distribution is created under the assumption the null hypothesis is true. In this case, we assume the true proportion of male boxers who are left handed is 0.1 so we will create 1000 different simulations of 500 boxers under this assumption.

Let's think about how to use cards to create one simulation of 500 boxers under the assumption the null hypothesis is true. Suppose blue cards represents left-handed and red cards represents right-handed.

15. How many cards total do we need? How many blue ones? How many red ones?

16. Next, we would mix the cards together and draw 1 card, write down if it's red or blue, and replace the card. How many times would we need to repeat this process to simulate our sample?

17. What would we plot on the null distribution?

We will use the computer to simulate 1000 simulated proportions of male boxers who are left handed for a sample size of 500 based on the assumption that the true proportion of male boxers who are left handed is 0.1. This is called the null distribution because it is created based on the assumption that the null hypothesis is true.

To use the computer simulation, we will need to enter the "probability of success," "sample size," "number of repetitions," "as extreme as," and the "direction" (matches the direction of the alternative hypothesis).

18. What values should be entered into the simulation?

Probability of success:

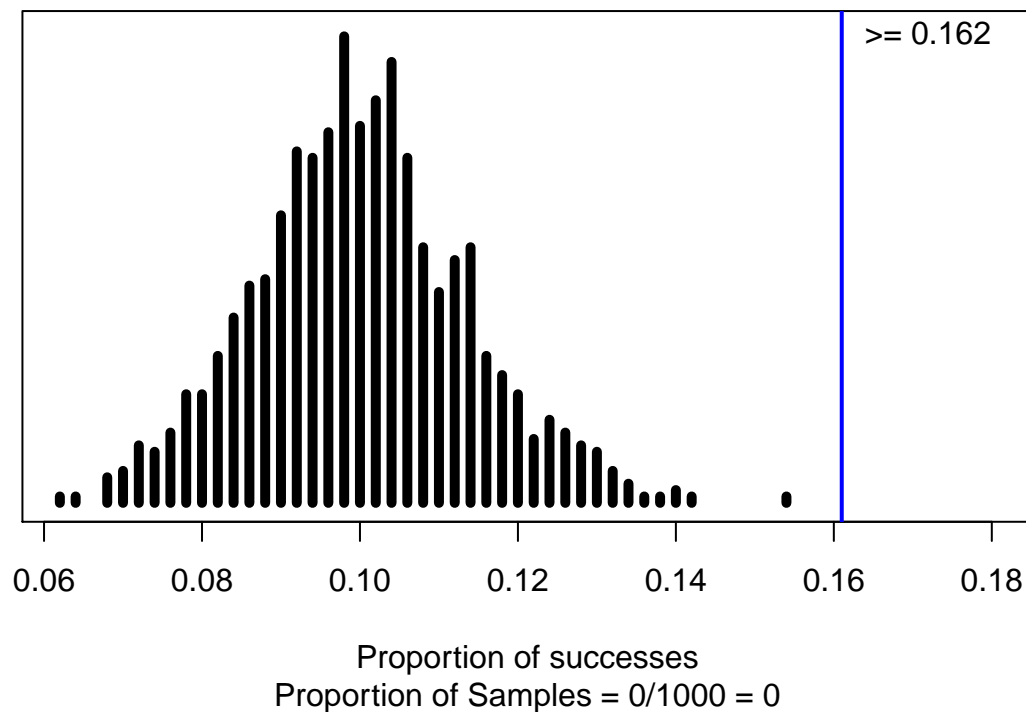
Sample size:

Number of repetitions:

As extreme as:

Direction:

```
one_proportion_test(probability_success = 0.1, #Null hypothesis value
                     sample_size = 500, #Enter sample size
                     number_repetitions = 1000, #Enter number of simulations
                     as_extreme_as = 81/500, #observed statistic
                     direction = "greater", #specify direction of alternative hypothesis
                     report_value = "proportion") #Reporting proportion or number of successes?
```



19. Around what value is the null distribution centered? Why does that make sense?

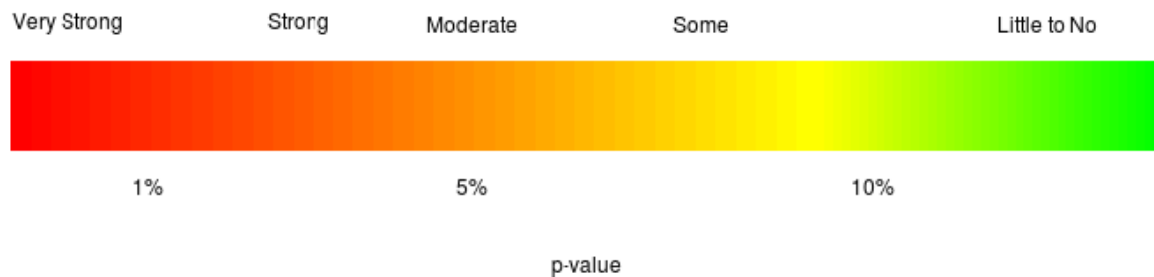
20. Where does the statistic (value from question 8) fall in the null distribution? Is it towards the center or in one of the tails?

21. Is the statistic likely to happen or unlikely to happen if the true proportion of male boxers is 0.1? Explain your answer.

22. Using the simulation, what is the proportion of samples at this summary statistic or greater, if the true proportion of male boxers is 0.1? *Hint: Look under the simulation.*

This is the p-value. The smaller the p-value the more evidence we have against the null hypothesis.

23. Using the following guidelines for the strength of evidence, how much evidence do the data provide against the null hypothesis?

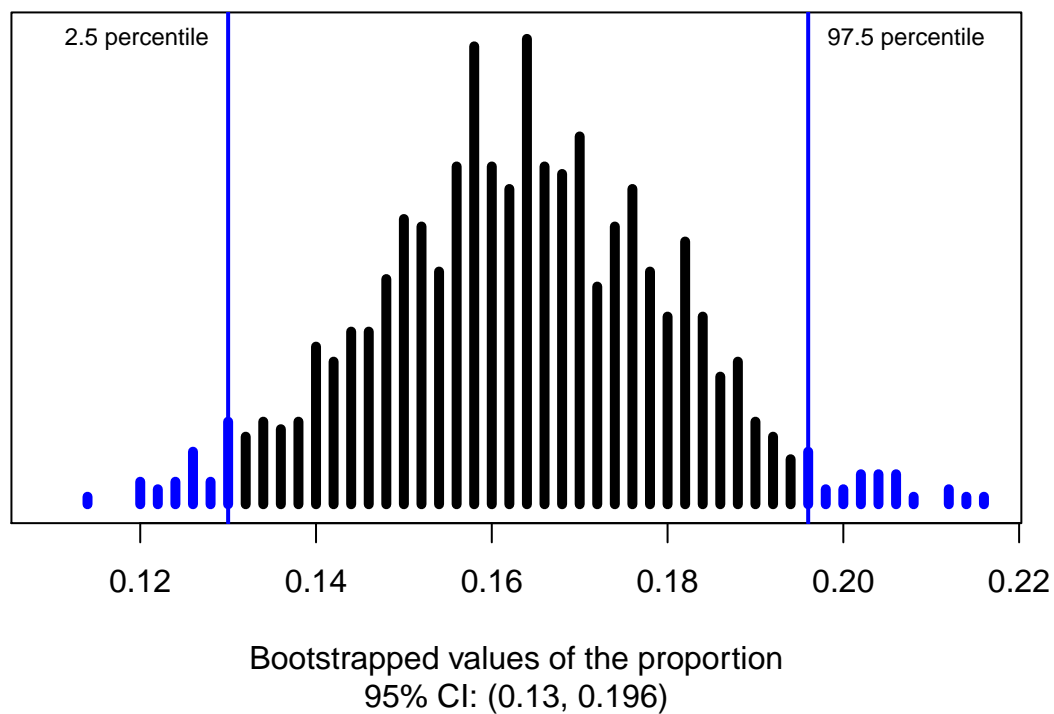


A point estimate provides a single plausible value for a parameter. However, a point estimate is rarely perfect; usually there is some error in the estimate. In addition to supplying a point estimate of a parameter, a next logical step would be to provide a plausible range of values for the parameter. This plausible range of values for the population parameter is called a confidence interval.

We will use bootstrapping to find the 95% confidence interval.

24. What is bootstrapping?

```
one_proportion_bootstrap_CI(sample_size = 500, #Sample size
                             number_successes = 81, #Observed number of successes
                             number_repetitions = 1000, #Number of bootstrap samples to use
                             confidence_level = 0.95) #Confidence level as a decimal
```



25. Explain why the blue lines are at the 2.5th percentile and the 97.5th percentile.
26. Report the 95% bootstrapped confidence interval for π . Use interval notation: (lower value, upper value).
27. What are we 95% confident is contained within this interval?

Communicate the results and answer the research question

When we write a conclusion we answer the research question by stating how much evidence there is for the alternative hypothesis.

28. Write a paragraph summarizing the results. Be sure to include:

- Summary statistic
- P-value
- Conclusion in context
- Confidence Interval
- Interpretation of the Confidence Interval

Revisit and look forward

29. Suggest a new research question that you might investigate, building on what you learned in this study.