

Lab 1: Hypothesis Testing

w203 Teaching Team

This is a team-based lab. Your instructor will divide you into teams of three, or possibly four students. To maximize learning, we would like all students to engage with every lab component, discussing strategy with teammates, reviewing solutions, and iterating on text and code.

The lab consists of two parts. Part one consists of foundation exercises similar to a homework. Part two is a written statistical analysis.

This lab is due June 28 at 2pm Berkeley time. You will find a separate place on Gradescope to submit each part. Please submit your `.Rmd` file, as well as your `.pdf` file. This is a team submission, so one person from your team should submit the teams' work, and associate the other members of the team with the submission (this is intuitive on Gradescope.)

Contents

| | | |
|----------|--|----------|
| 1 | Part 1: Foundational Exercises | 2 |
| 1.1 | Professional Magic | 2 |
| 1.2 | Wrong Test, Right Data | 3 |
| 1.3 | Test Assumptions | 4 |
| 1.3.1 | World Happiness | 4 |
| 1.3.2 | Legislators | 4 |
| 1.3.3 | Wine and health | 4 |
| 1.3.4 | Attitudes toward the religious | 4 |
| 2 | Part 2: Statistical Analysis | 6 |
| 2.1 | Data | 6 |
| 2.2 | The research question | 6 |
| 2.3 | Guidance from political scientists | 6 |
| 2.4 | Report guidelines | 6 |
| | General guidance | 7 |
| | Introduction | 7 |
| | Conceptualization and Operationalization | 7 |
| | Visual Design | 7 |
| | Data wrangling | 7 |
| | Hypothesis testing | 8 |
| | Test, results and interpretation | 8 |

1 Part 1: Foundational Exercises

1.1 Professional Magic

Your aunt (who is a professional magician), claims to have created a pair of magical coins that share a connection to each other that makes them land in the same way. The coins are always flipped at the same time. For a given flip $i \in \{1, 2, 3, \dots\}$, let X_i be a Bernoulli random variable representing the outcome of the first coin, and let Y_i be a Bernoulli random variable representing the outcome of the second coin. You assume that each flip of the pair is independent of all other flips of the pair. For all i , you also assume that X_i and Y_i have the joint distribution given in the following table.

| | $X_i = 0$ | $X_i = 1$ |
|-----------|-----------|-----------|
| $Y_i = 0$ | $p/2$ | $(1-p)/2$ |
| $Y_i = 1$ | $(1-p)/2$ | $p/2$ |

$p \in [0, 1]$ is a parameter.

Each flip of the pair is independent of all other flips of the pair. This means that whatever happens the first time that you flip both of the coins tells you nothing about the second time that you flip both of the coins, and so on. Essentially, this is a statement about the limits of your aunt's magic.

You design the following test to evaluate your aunt's claim: You flip the coins three times, and write down that your test statistic is the sum $X_1 + Y_1 + X_2 + Y_2 + X_3 + Y_3$. That is, your test statistic is essentially the number of heads that are shown.

Your null hypothesis is that $p = 1/2$, and you plan to reject the null if your test statistic is 0 or 6.

1. What is the type 1 error rate of your test?
2. What is the power of your test for the alternate hypothesis that $p = 3/4$?

1.2 Wrong Test, Right Data

Imagine that your organization surveys a set of customers to see how much they like your regular website, and how much they like your mobile website. Suppose that both of these preference statements are measured on 5-point Likert scales.

A Likert scale is one where a person is provided ordered categories that range from lowest to highest. You can read more about them in this seminal research design text by [Fowler](#), or this brief [overview](#). If you were to run a paired t-test using this data, what consequences would the violation of the metric scale assumption have for your interpretation of the test results? What would you propose to do to remedy this problem?

1.3 Test Assumptions

For the four following questions, your task is to evaluate the assumptions for the given test. It is not enough to say that an assumption is met or not met; instead, present your evidence in the form of background knowledge, visualizations, and numerical summaries. If you produce a histogram as part of your evaluation, be sure to consider what the most appropriate bin width is. The test that we ask you to evaluate may or may not be the most appropriate test for the scenario. Because the goal of this task is to evaluate whether the data satisfies the assumptions necessary for the test to provide meaningful results, you do not need perform the test (you may perform the test, but we will not be marking for the test results).

1.3.1 World Happiness

The file [datasets/Happiness_WHR.csv](#) is subsetting from the World Happiness Report, a yearly publication that uses data from the Gallup World Poll surveys. The variable life ladder is a measure of happiness, described in the FAQ as follows:

This is called the Cantril ladder: it asks respondents to think of a ladder, with the best possible life for them being a 10, and the worst possible life being a 0. They are then asked to rate their own current lives on that 0 to 10 scale. The rankings are from nationally representative samples, for the years 2018-2020.

You would like to know whether people in countries with high GDP per capita (higher than the mean) are more happy or less happy than people in countries with low GDP (lower than the mean).

List all assumptions for a two-sample t-test. Then evaluate each assumption, presenting evidence based on your background knowledge, visualizations, and numerical summaries.

1.3.2 Legislators

The file [datasets/legislators-current.csv](#) is taken from the [congress-legislators project](#) on Github. You would like to test whether Democratic or Republican senators are older.

List all assumptions for a Wilcoxon rank-sum test (using the Hypothesis of Comparisons). Then evaluate each assumption, presenting evidence based on your background knowledge, visualizations, and numerical summaries.

1.3.3 Wine and health

The dataset `wine` can be accessed by installing the `wooldridge` package.

```
install.packages("wooldridge")
library(wooldridge)
?wine
wine
```

It contains observations of variables related to wine consumption for 21 countries. You would like to use this data to test whether countries have more deaths from heart disease or from liver disease.

List all assumptions for a signed-rank test. Then evaluate each assumption, presenting evidence based on your background knowledge, visualizations, and numerical summaries.

1.3.4 Attitudes toward the religious

The file [datasets/GSS_religion](#) is a subset of data from the 2004 General Social Survey (GSS).

The variables `prottemp` and `cathtemp` are measurements of how a respondent feels towards protestants and towards Catholics, respectively. The GSS questions are phrased as follows:

I'd like to get your feelings toward groups that are in the news these days. I will use something we call the feeling thermometer, and here is how it works:

I'll read the names of a group and I'd like you to rate that group using the feeling thermometer. Ratings between 50 degrees and 100 degrees mean that you feel favorable and warm toward the group. Ratings between 0 degrees and 50 degrees mean that you don't feel favorable toward the group and that you don't care too much for that group. If we come to a group whose name you Don't recognize, you don't need to rate that group. Just tell me and we'll move on to the next one. If you do recognize the name, but you don't feel particularly warm or cold toward the group, you would rate the group at the 50 degree mark.

How would you rate this group using the thermometer?

You would like to test whether the US population feels more positive towards Protestants or towards Catholics.

List all assumptions for a paired t-test. Then evaluate each assumption, presenting evidence based on your background knowledge, visualizations, and numerical summaries.

2 Part 2: Statistical Analysis

The American National Election Studies (ANES) conducts surveys of voters in the United States, with a flagship survey occurring immediately before and after each presidential election. In this part, you will use the ANES data to address a question about voters in the US. Your team will conduct a statistical analysis and generate a written report in pdf format.

This is an exercise in both statistics and professional communication. It is important that your techniques are properly executed; equally important is that your writing is clear and organized, and your argument well justified.

2.1 Data

Data for the lab should be drawn from the 2020 American National Election Studies (ANES). You can access this data at <https://electionstudies.org>. This is the official site of the ANES, a project that has [been ongoing since](#) 1948, and federally funded by the National Science Foundation since 1977.

To access the data, you will need to register for an account, confirm this account, and then login. The data that you need should come from the **2020 Time Series Study**.

While you're at the ANES website, you will also want to download the codebook, because all of the variables are marked as something like, V200002 – which isn't very descriptive without the codebook.

For a glimpse into some of the intricacies that go into the design of this study, take a look at the introduction to the codebook.

Like many modern surveys, the ANES includes survey weights, which are used to correct for situations in which members of one demographic group are more likely to respond to the survey than members of another demographic group. (The target proportions are ultimately based on US census data). The survey weights make it possible to generalize from the a population that represents people who take the survey to a population that represents the United States as a whole. These weights are beyond the scope of our class and you are not expected to utilize them. You will still be able to learn about a population model, even if applicability to the US population is limited.

2.2 The research question

Use the ANES data to address the following question:

Did Democratic voters or Republican voters experience more difficulty voting in the 2020 election?

2.3 Guidance from political scientists

Political identification in the US is a complex phenomenon that is the topic of a large academic literature. See `./background_literature/petrocik_2009.pdf` for some guidance about how stated political identity might not match with revealed political identity at the ballot box.

As practical guidance:

1. Is it reasonable to use the vote that someone cast to identify their party preference in this case? What if someone had so difficult a time voting that they did not cast a ballot?
2. Please treat individuals who “lean” in one direction or another as members of that party. This means that someone who “Leans Democratic” should be classified as as Democrat; and someone who “Leans Republican” should be classified as a Republican.

2.4 Report guidelines

This section provides some guidance for you as you write your report. In [rubric.md](#) we provide you with specific statements of how we will evaluate your report.

General guidance

- You should knit an .Rmd file to create your pdf report.
- Your report should be no more than 3 pages in standard latex formatting (i.e. `output: pdf_document`)
- You should assume your reader is familiar with statistics, but has no special knowledge of the ANES survey.
- Your report should contain either a plot or a table that advances the argument.

Introduction

- Begin your report with an introduction to motivate the analysis.
- Introduce the topic area and explain why the research question is interesting.
- The introduction must “do work,” connecting the general topic to the specific techniques in the report.

Conceptualization and Operationalization

How do you get from a research question to data? First, ensure that the concepts in your question are clear.

- Who or what is a voter?
- Who is a “Republican” and who is a “Democrat”?
- What is difficulty voting?

Only after you have informed your reader of what these concepts are can you then describe how you are going to *measure* these concepts.

- What would be the best **possible** method of measuring this concept? Is this method possible? Why or why not?
- What is the best **available** method of measuring this concept? Why have you opted to use this measurement instead of other possibilities? Map the concept definitions that you have written down onto the variables that you are going to use. Describe, precisely, how the variables were generated, if they come from survey data, provide the text of the question that the respondent is reacting to, not the variable name.
- What, if any, changes have you made to the dataset from how it was provided? Why did you make those changes, how much data was affected, and what are the consequences for any estimates that you produce?

Visual Design

- Any plots or tables that you include must follow basic principles of visual design.
 - A plot/figure must have a title that is informative.
 - Variables must be labeled in plain language. As an example, `v20002` does not work for a label.
 - A plot should have a good ratio of information to ink / space on the page. Do not select a large or complicated plot when a simple table conveys the same information directly.
- Do not include any plot (or R output in general) that you do not discuss in your narrative.
- The code that makes your plot/figure should be included in your report .Rmd file, but should not be shown in your final report. To accomplish this, you can use an `echo=FALSE` argument in the code chunk that produces the plot/figure.

Data wrangling

To answer your research question, you will have to clean, tidy, and structure the data (A.K.A. wrangle).

- The code to wrangle data should be included with your deliverable somehow. If you choose to include it in your report .Rmd file, then it should not be shown in the PDF of your final report. To accomplish this, you can use an `echo=FALSE` argument for the code chunk that does the wrangling.

- A better practice – not strictly necessary for this lab – would be to write a function that loads and cleans *all* of the data that is being used by your team for its reports. This way, a single function can be run (and evaluated by your reader) for all the loading, cleaning, and manipulating.
- While we do not want to prohibit you from using additional tools for data manipulation, you should be able to complete this lab with no more than the base `stats` library, plus `dplyr` and `ggplot2` for data manipulation and plotting. Other tools within the tidyverse are available to use, but don't feel like you have to search them out.
- You will learn more by writing your own function than you would searching for a package that does one thing for your report.

Hypothesis testing

To answer your research question, you will have to execute one of the statistical tests from the course.

- The code that executes your test *should* be shown in your report, because it makes very clear the specific test that you're conducting.
- You need to argue, from the statistical principles of the course, why the test you are conducting is the *most appropriate* way to answer the research question.
- Although you might not do this for a report at your organization, for this class please list every assumption from your test, and evaluate it (assess whether the assumption is a reasonable reflection of the natural process that generated the data).
- If you identify problems with some assumptions for your test, that does not mean that you should abandon the analysis or hide the problem. If these “limitations” exist, please describe them honestly, and provide your interpretation of the consequences for your test.
- While you can choose to display the results of your test in the report, you also *certainly* need to write about these results. This should be accomplished using [inline code chunks](#), rather than by hard-coding / hard-writing output into your written report. An example of this is included in `lab_1_example_solution.Rmd`.

Test, results and interpretation

Please discuss whether any statistically significant results that you find are of *practical significance*. There are many ways to do this, but the best will provide your reader enough context to understand any measured differences in a scale appropriate to your variables. Explain the main takeaway of your analysis and how it relates to the broader context you identified in the introduction.