



STAT 216 Activity Coursepack

Fall 2020

Contents

Preface	3
Fall 2020 Calendar of In-Class Activities	4
1 Martian Alphabet	6
2 Study Design	12
3 Current Population Survey	17
4 IMDb Movie Reviews	23
5 Movie Profits	28
6 Handedness of Male Boxers	34
7 Winter Sports Helmet Use and Head Injuries	43
8 COVID-19 and Air Pollution	50
9 Weather Patterns and Record Snowfall	57
10 Hand Dexterity	63

Preface

This coursepack accompanies the textbook for STAT 216: Introduction to Statistics at Montana State University. Each of the activities in this workbook is designed to target specific learning outcomes of the course, giving you practice with important statistical concepts in a group setting with instructor guidance. Bring this workbook with you to class each week, and take notes in the workbook as you would your own notes. A well-written complete workbook will provide an optimal study guide for exams!

Fall 2020 Calendar of In-Class Activities

Week	Activity No.	Day	Date	Activity
1	1	M	8/17	Martian Alphabet
1	1	T	8/18	Martian Alphabet
1	1	W	8/19	Martian Alphabet
1	1	H	8/20	Martian Alphabet
1	1	F	8/21	Martian Alphabet
2	2	M	8/24	Study Design
2	2	T	8/25	Study Design
2	2	W	8/26	Study Design
2	2	H	8/27	Study Design
2	2	F	8/28	Study Design
3	3	M	8/31	Current Population Survey
3	3	T	9/1	Current Population Survey
3	3	W	9/2	Current Population Survey
3	3	H	9/3	Current Population Survey
3	3	F	9/4	Current Population Survey
4	-	M	9/7	No class – Labor Day
4	4	T	9/8	IMDb Movie Reviews
4	4	W	9/9	IMDb Movie Reviews
4	4	H	9/10	IMDb Movie Reviews
4	4	F	9/11	IMDb Movie Reviews
5	4	M	9/14	IMDb Movie Reviews
5	5	T	9/15	Movie Profits
5	5	W	9/16	Movie Profits
5	5	H	9/17	Movie Profits
5	5	F	9/18	Movie Profits
6	5	M	9/21	Movie Profits
6	—	T–F	9/22–9/25	Exam 1
7	—	M	9/28	Exam 1
7	6	T	9/29	Handedness of Male Boxers
7	6	W	9/30	Handedness of Male Boxers
7	6	H	10/1	Handedness of Male Boxers
7	6	F	10/2	Handedness of Male Boxers
8	6	M	10/5	Handedness of Male Boxers
8	7	T	10/6	Winter Sports Helmet Use and Head Injuries
8	7	W	10/7	Winter Sports Helmet Use and Head Injuries
8	7	H	10/8	Winter Sports Helmet Use and Head Injuries
8	7	F	10/9	Winter Sports Helmet Use and Head Injuries
9	7	M	10/12	Winter Sports Helmet Use and Head Injuries
9	8	T	10/13	COVID-19 and Air Pollution
9	8	W	10/14	COVID-19 and Air Pollution
9	8	H	10/15	COVID-19 and Air Pollution
9	8	F	10/16	COVID-19 and Air Pollution
10	8	M	10/19	COVID-19 and Air Pollution
10	9	T	10/20	Weather Patterns and Record Snowfall
10	9	W	10/21	Weather Patterns and Record Snowfall
10	9	H	10/22	Weather Patterns and Record Snowfall
10	9	F	10/23	Weather Patterns and Record Snowfall

Week	Activity No.	Day	Date	Activity
11	9	M	10/26	Weather Patterns and Record Snowfall
11	10	T	10/27	Hand Dexterity
11	10	W	10/28	Hand Dexterity
11	10	H	10/29	Hand Dexterity
11	10	F	10/30	Hand Dexterity
12	10	M	11/2	Hand Dexterity
12	—	T	11/3	No class — Election Day
12	—	W–F	11/4–11/6	Exam 2
13	—	M–T	11/9–11/10	Exam 2
13	—	W	11/11	No class — Veterans Day
13	—	H–W	11/12–11/18	Review

Martian Alphabet

1.1 Learning outcomes

- Describe the statistical investigation process
- Identify observational units, variables, and variable types in a statistical study

1.2 Terminology review

Statistics is the study of how best to collect, analyze, and draw conclusions from data. Today in class you will be introduced to the following terms:

- Observational units or cases
- Variables: categorical or quantitative
- Proportions
- Graphs: frequency bar plot and relative frequency bar plot
- Distribution

For more on these concepts, read Sections 1.2 and 2.1 in the textbook.

1.3 Can you read “Martian”?

How well can humans distinguish one “Martian” letter from another? In today’s activity, we’ll find out. When shown the two Martian letters, Kiki and Bumba, write down whether you think Bumba is on the left or on the right.

Steps of the statistical investigation process

Step 1: The first step of any statistical investigation is to *ask a research question*. In this study the research question is: can we as a class read Martian? (We will refine this later on!).

Step 2: To answer any research question, we must *design a study and collect data*. For our question, the study consists of each student being presented with two Martian letters and asking which was Bumba. Your responses will become our observed data that we will explore.

Observational units or **cases** are the subjects data are collected on. In a spreadsheet of the data set, each row will represent a single observational unit.

1. What are the observational units in this study?
2. How many students are in class today? This is the *sample size*.

A **variable** is information collected or measured on each observational unit or case. Each column in a data set will represent a different variable. Today we are only measuring one variable on each observational unit.

3. Identify the variable we are collecting on each observational unit in this study, i.e., what are we measuring on each student?

We will look at two types of variables: **quantitative** and **categorical** (see Figure 1.1).

Quantitative variables are numerical measurements that can be discrete (whole, non-negative numbers) or continuous (any value within an interval). The number of students in a class would be a discrete variable as you can not have a partial student. GPA would be a continuous variable ranging from 0 to 4.0.

Categorical variables are data that are in groups or categories such as eye color, state of residency, or whether or not a student lives on campus. Categorical variables with a natural ordering are considered ordinal variables while those without a natural ordering are considered a nominal variable. All categorical variables will be treated as nominal for analysis in this course.

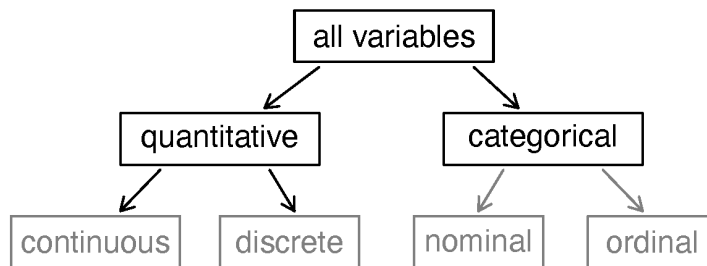


Figure 1.1: Types of variables.

4. Is the variable identified in question 3 categorical or quantitative?
5. Were you correct or incorrect in identifying Bumba?

Step 3: Once we have collected data, the next step is to *summarize and visualize the data*.

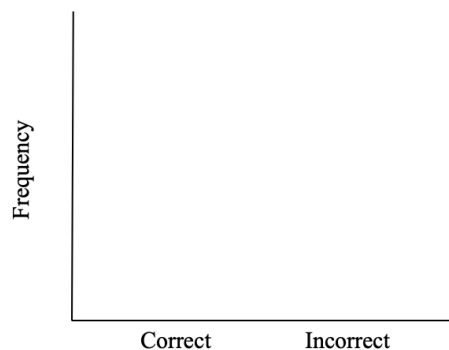
6. How many people in your class were correct in identifying Bumba? Using the class size from question 2, calculate the proportion of students who correctly identified Bumba.

$$\text{proportion} = \frac{\text{number of students who correctly identified Bumba}}{\text{total number of students}}$$

The proportion in question 6 is called a **summary statistic**—a single value that summarizes the data set. It is important to note that a variable is different than a summary statistic. A *variable* is measured on a *single observational unit* while a summary statistic is calculated from a group of observational units. For example, the variable “whether or not a student lives on campus” can be measured on each individual student. In a class of 50 students we can calculate the proportion of students who live on campus, the summary statistic. Look back and make sure you wrote the variable in question 3 as a variable, NOT a summary statistic.

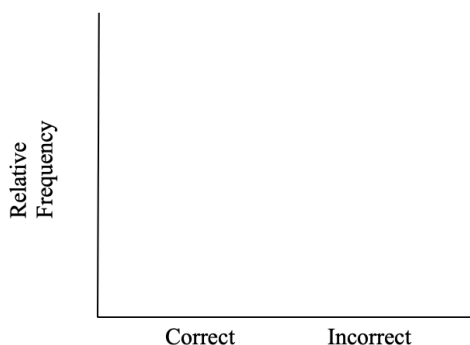
Looking at the data set and the summary statistics is only one way to display the data. We will also want to create a visualization or picture of the data. A **frequency bar plot** is used to display categorical data as a count or frequency. Since our variable has two levels, correct or incorrect, we will create two bars—one for each level.

7. Plot the observed class data using a frequency bar plot. Be sure to add a scale to the y-axis.



We can also visualize the data as a proportion in a **relative frequency bar plot**. Relative frequency is the proportion calculated for each level of the categorical variable.

8. Plot the observed class data using a relative frequency bar plot. Be sure to add a scale to the y-axis.



Step 4: The next step is to *use statistical analysis methods to draw inferences from the data*. To answer the research question, we will simulate what *could* have happened in our class given random chance, repeat that many times to understand the expected *variability* between different “randomly guessing” classes, then compare our class’s observed data to the simulation. This gives us an estimate of how often (or the probability of) our class’s result would occur if we were all merely guessing, allowing us to determine if the data provides evidence that we as a class can in fact read Martian.

9. If humans really don’t know Martian and are just guessing which is Bumba, what are the chances of getting it right?

How could we use a coin to simulate each student “just guessing” which Martian letter is Bumba?

How could we use coins to simulate the entire class “just guessing” which Martian letter is Bumba?

How many people in your class would you expect to choose Bumba correctly just by chance? Explain your reasoning.

10. Each of you will flip a coin one time to simulate your “guess”. Let Heads = correct, Tails = incorrect. What was the result of your simulation?

What was the result from your class’s simulation? What proportion of students “guessed” correctly in the simulation?

11. If students really don’t know Martian and are just guessing which is Bumba, which seems more unusual: the result from your class’s **simulation** or the observed proportion of students in your class that were correct (this is your summary statistic from question 6)? Explain your reasoning.

12. While your observed class data is likely far different from the simulated “just-guessing” class, comparing our class data to a single simulation does not seem to give enough information. The differences seen could just be due to that set of coin flips! Let’s simulate another class. Each student should flip your coin again. What was the result from your class’s second simulation? What proportion of students “guessed” correctly in the second simulation? Create a plot to compare the two simulated results with the observed class result.
13. We still unfortunately only have a couple of simulations to compare our class data to. It would be much better to be able to see how our class compared to hundreds or thousands of “just-guessing” classes. Since we don’t want to flip coins all class period, your instructor will use a computer simulation to get 1000 trials. Fill in the following blanks to describe how we would create a simulation of random guessing with 1000 trials.
- Probability of correct guesses: _____
- Sample size: _____
- Number of repetitions: _____
14. Sketch the distribution displayed by your instructor here, being sure to label each axis appropriately.
15. Is your class particularly good or bad at Martian? How can you use the plot in question 14 to tell?
16. Is it *possible* that we could see our class results just by chance if everyone was just guessing? Explain your reasoning.
17. Is it *likely* that we could see our class results just by chance if everyone was just guessing? Explain your reasoning.

Step 5: The next step in the statistical investigation process is to *communicate the results and answer the research question*.

18. Does this activity provide strong evidence that students were not just guessing at random? If so, what do you think is going on here? Can we as a class read Martian?¹

Step 6: The final step of any statistical investigation is to *revisit and look ahead*.

19. Can you think of any limitations of our study? Can you think of a new topic that might be of interest based on the results of our study?

1.4 Take home messages

1. In this course we will learn how to evaluate a claim by comparing observed results (classes' "guesses" when asked to identify Bumba) to a distribution of many simulated results under an assumption like "blind guessing."
2. Blind guessing between two outcomes will be correct only about half the time. We can create data (via computer simulation) to fit the assumption of blind guessing.
3. Unusual observed results will make us doubt the assumptions used to create the simulated distribution. A large number of correct "guesses" is evidence that a person was not just blindly guessing.

1.5 Additional notes

Use this space to summarize your thoughts and take additional notes on today's activity, and to write down the names and contact information of your team mates.

¹Reference for "Martian alphabet" is a TED talk given by Vilayanur Ramachandran in 2007. The synesthesia part begins at roughly 17:30 minutes: http://www.ted.com/talks/vilayanur_ramachandran_on_your_mind.

Study Design

2.1 Learning outcomes

- Explain the purpose of random sampling and its effect on scope of inference
- Explain the purpose of random assignment and its effect on scope of inference
- Identify whether a study is observational or an experiment
- Identify confounding variables in observational studies and explain why they are confounding
- Identify the types of bias present in a study

2.2 Terminology review

In today's activity, we will examine different types of sampling bias and study designs, confounding variables, and how to determine the scope of inference for a study. Some terms covered in this activity are:

- Population
- Sample
- Parameter
- Statistic
- Selection bias
- Response bias
- Non-response bias
- Scope of inference
- Explanatory variable
- Response variable
- Confounding variable
- Experiment
- Observational study

To review these concepts, see Sections 1.3 through 1.6 in the textbook.

2.3 Types of sampling bias

There are two parts to study design: how the sample was selected and how the study was conducted. First, we will look at sampling and types of bias (selection, non-response, or response).

In these next questions, identify the target population, the sample, the variable, and the type of bias present.

1. To determine if the proportion of out-of-state undergraduate students at Montana State University has increased in the last 10 years, a statistics instructor sent an email survey to 500 randomly selected current undergraduate students. One of the questions on the survey asked whether they had in-state or out-of-state residency. She only received 378 responses.

Target population:

Sample:

Variable:

Type(s) of bias:

2. Pew Research surveys US adults about many different topics. Recently, a survey was conducted to assess current presidential approval. A random sample of 6395 US adults was taken. Of those surveyed, 42% said they agree with President Trump on many or nearly all of the top issues facing the country today.

Target population:

Sample:

Variable:

Type(s) of bias:

3. A television station is interested in predicting whether or not a local referendum to legalize marijuana for adult use will pass. It asks its viewers to phone in and indicate whether they are in favor or opposed to the referendum. Of the 2241 viewers who phoned in, forty-five percent were opposed to legalizing marijuana.

Target population:

Sample:

Variable:

Type(s) of bias:

4. To gauge the interest in a new swimming pool, a local organization stood outside of the Bogart Pool during open hours. One of the questions they asked was, "Since the Bogart Pool is in such bad repair, don't you agree that the city should fund a new pool?"

Target population:

Sample:

Variable:

Type(s) of bias:

5. The Bozeman school district is interested in surveying parents of students about their opinions on returning to school this fall following the COVID-19 pandemic. They divided the school district into 10 divisions based on location and randomly surveyed 20 households within each division.

Target population:

Sample:

Variable:


Type(s) of bias:

2.4 Study design


The two main study designs we will cover are **observational studies** and **experiments**. Both the sampling method and the study design will help to determine the **scope of inference** for a study. Remember that only in a randomized experiment can we conclude a **causal** (cause and effect) relationship between the explanatory and response variable.

Scope of Inference: If evidence of an association is found in our sample, what can be concluded?

	Study Type		
Selection of cases	Randomized experiment	Observational study	
Random sample (and no other sampling bias)	Causal relationship, and can generalize results to population.	Cannot conclude causal relationship, but can generalize results to population.	→ Inferences to population can be made
No random sample (or other sampling bias)	Causal relationship, but cannot generalize results to a population.	Cannot conclude causal relationship, and cannot generalize results to a population.	→ Can only generalize to those similar to the sample due to potential sampling bias



Can draw cause-and-effect conclusions



Can only discuss association
due to potential confounding
variables

For the next exercises, identify the explanatory variable, the response variable, the study design (observational study or experiment), and the scope of inference.

- The pharmaceutical company Moderna Therapeutics is working in conjunction with the National Institutes of Health towards a vaccine for COVID-19 and has recently begun Phase 3 clinical trials. US clinical research sites will enroll 30,000 volunteers without COVID-19 to participate. Participants will be randomly assigned to receive either the candidate vaccine or a saline placebo. They will then be followed to assess vaccine-related symptoms and development of COVID-19. The trial is double-blind, so neither the investigators nor the participants will know who is assigned to which group.

Explanatory variable:

Response variable:

Study design:

What is the scope of inference for this study?

7. In another study, a local health department randomly selected 1000 US adults without COVID-19 to participate in a health survey. Each participant was assessed at the beginning of the study and then followed for one year. They were interested to see which participants elected to receive a vaccination for COVID-19 and whether any participants developed COVID-19.

Explanatory variable:

Response variable:

Study design:

What is the scope of inference for this study?

8. For each of the studies in questions 6 and 7, determine whether confounding variables could be an issue. If so, identify a potential confounding variable and explain how it meets the definition of a confounding variable.

2.5 Additional notes

Use this space to summarize your thoughts and take additional notes on today's activity.

Current Population Survey

3.1 Learning outcomes

- Identify and create appropriate summary statistics and plots given a data set or research question
- Plots for a single categorical variable: bar plot
- Plots for association between two categorical variables: segmented bar plot, mosaic plot
- Recognize and simulate probabilities as long-run frequencies
- Construct two-way tables to evaluate conditional probabilities

3.2 Terminology review

In today's activity, we will review summary measures and plots for categorical variables. Some terms covered in this activity are:

- Proportions
- Bar plots
- Segmented bar plots
- Probability
- Conditional probability
- Two-way tables

To review these concepts, see Sections 2.1 and 2.2 in the textbook.

3.3 “Current” Population Survey: 1985

The data set we will use for this activity is from the Current Population Survey (CPS) in 1985. The CPS is a survey sponsored by the Census Bureau and the Bureau of Labor Statistics to track labor force statistics for the United States population. The following table describes the variables in the data set:

Variable	Description
<code>educ</code>	Number of years of education
<code>south</code>	Whether lives in southern region of the US: S = lives in south, NS = does not live in south
<code>sex</code>	Sex: M = male, F = female
<code>exper</code>	Number of years of work experience (inferred from age and education)
<code>union</code>	Whether union member: Union or Not
<code>wage</code>	Wage (dollars per hour)
<code>age</code>	Age (years)
<code>race</code>	Race: W = white, NW = not white
<code>sector</code>	Sector of the economy: clerical , const (construction), management , manufacturing , professional , sales , service , other
<code>married</code>	Marital status: Married or Single

Vocabulary review

1. What are the observational units?
2. Which variables are categorical?
3. What types of plots can be used to display categorical data?

An important part of understanding data is to create visual pictures of what the data represent. In this activity, we will create graphical representations of categorical data.

R code

R is a free statistical analysis software program we will use in Stat 216. Please see D2L for instructions on how to download or access a version of R on your laptop, or plan to use the school computers for some parts of assigned out of class work for this course.

Throughout these activities, we will often include the R code you would use in order to produce output or plots. These “code chunks” appear in gray. In the code chunk below, we demonstrate how to read the data set into R using the `read.csv()` function, and tell R to treat the `sector` and `sex` variables as categorical variables (“factors”).

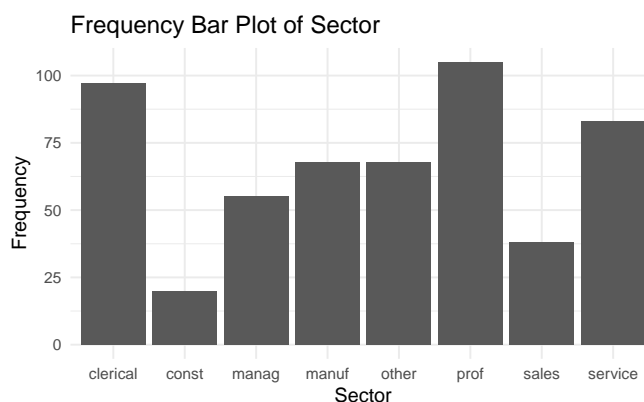
The `#` sign is not part of the R code. It is used by these authors to add comments to the R code and explain what each call is telling the program to do.

```
cps <- read.csv("data/cps.csv") #This will read in the data set
```

Displaying a single categorical variable

If we wanted to know how many people in our data set were in each sector, we would create a frequency bar plot of the variable `sector`.

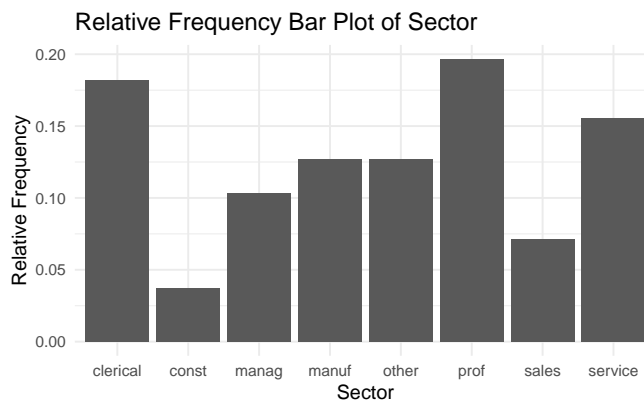
```
cps %>% #Data set piped into...
ggplot(aes(y = sector)) + #This specifies the variable
  geom_bar(stat = "count") + #Tell it to make a bar plot
  labs(title = "Frequency Bar Plot of Sector", #Give your plot a title
       x = "Frequency", #Label the x axis
       y = "Sector") + #Label the y axis
  coord_flip() #Turn the bars so they are vertical
```



4. Which sector of the economy has the largest number of people in it? Approximately how many people are in this sector?

We could also choose to display the data as a proportion in a relative frequency bar plot. To find the relative frequency, divide the count in each sector by the sample size. These are sample proportions.

```
cps %>% #Data set piped into...
ggplot(aes(x = sector)) + #This specifies the variable
  geom_bar(aes(y = ..prop.., group = 1)) + #Tell it to make a bar plot with proportions
  labs(title = "Relative Frequency Bar Plot of Sector", #Give your plot a title
       x = "Sector", #Label the x axis
       y = "Relative Frequency") #Label the y axis
```

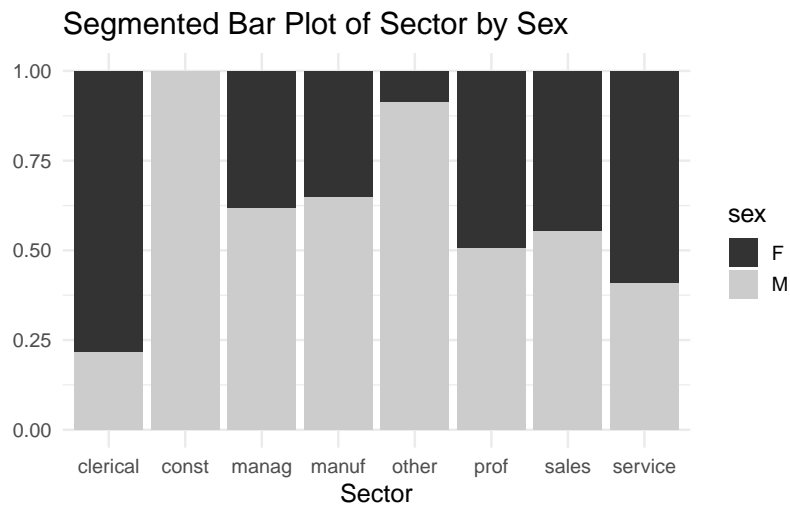


5. Which features in the relative frequency bar plot are the same as the frequency bar plot? Which are different?

Displaying two categorical variables

To examine the differences proportion of males and females across sectors, we would create a segmented bar plot of `sector` segmented by `sex`.

```
cps %>% #Data set piped into...
ggplot(aes(x = sector, fill = sex)) + #This specifies the variables
  geom_bar(stat = "count", position = "fill") + #Tell it to make a stacked bar plot
  labs(title = "Segmented Bar Plot of Sector by Sex", #Make sure to title your plot
       x = "Sector", #Label the x axis
       y = "") + #Remove y axis label
  scale_fill_grey() #Make figure black and white
```



6. Using the segmented bar plot, which sector has about the same proportion of males and females?
7. Which sector has the highest proportion of females?
8. Which variable is the bar plot treating as the explanatory variable? Which is the response variable?

3.4 Probability

9. A study was reported in which ninth grade Minnesota teens were asked whether they had gambled at least once a week in the past year. The sample consisted of 49.1% boys. The proportion of boys who had gambled at least once per week during the past year was 0.229, while among non-boys this proportion was only 0.045.

Let B = the event the person is a boy, and C = the event the person is a weekly gambler.

- a. Draw a segmented bar plot of sex segmented by gambling. Make sure to clearly label your axes and legend.

- b. Identify what each numerical value given in the problem represents in probability notation.

$$0.491 =$$

$$0.229 =$$

$$0.045 =$$

- c. Create a hypothetical two-way table to represent the situation. Recall that in a two-way table, the explanatory variable should be your column headers (similar to the x -axis in a segmented bar graph!) while the response variable becomes the row headers.

		Total
Total		100,000

- d. Find $P(B \text{ and } C)$. What does this probability represent in the context of the problem?

- e. Find the probability that a selected non-gambler is a non-boy. What is the notation used for this probability?

10. In a computer store, 30% of the computers in stock are laptops and 70% are desktops. Five percent of the laptops are on sale, while 10% of the desktops are on sale.

Let L = the event the computer is a laptop, and S = the event the computer is on sale.

- a. Identify what each numerical value given in the problem represents in probability notation.

$$0.30 =$$

$$0.70 =$$

$$0.05 =$$

$$0.10 =$$

- b. Create a hypothetical two-way table to represent the situation.

		Total
Total		100,000

- c. Calculate the probability that a randomly selected computer will be a desktop, given that the computer is on sale. What is the notation used for this probability?

- d. Find $P(S|L^C)$. What does this probability represent in context of the problem?

3.5 Additional notes

Use this space to summarize your thoughts and take additional notes on today's activity.

IMDb Movie Reviews

4.1 Learning objectives

- Identify and create appropriate summary statistics and plots given a data set or research question for quantitative data
- Interpret the following summary statistics in context: median, lower quartile, upper quartile, standard deviation, inter-quartile range
- Given a plot or set of plots, describe and compare the distribution(s) of a single quantitative variable (center, variability, shape, outliers)

4.2 Terminology review

In today's activity, we will review summary measures and plots for quantitative variables. Some terms covered in this activity are:

- Two measures of center: mean, median
- Two measures of spread (variability): standard deviation, inter-quartile range (IQR)
- Types of graphs: box plots, dot plots, histograms

To review these concepts, see Section 2.3 in the textbook.

4.3 Movies released in 2016

A data set was collected on movies released in 2016. Here is a list of some of the variables collected on these movies.

Variable	Description
budget_mil	Amount of money (in US \$ millions) budgeted for the production of the movie
revenue_mil	Amount of money (in US \$ millions) the movie made after release
duration	Length of the movie (in minutes)
content_rating	Rating of the movie (G, PG, PG-13, R, Not Rated)
imdb_score	IMDb user rating score from 1 to 10
genres	Categories the movie falls into (e.g., Action, Drama, etc.)
movie_facebook_likes	Number of likes a movie receives on Facebook

Vocabulary review

1. What are the observational units in this data set?
2. Which of the above listed variables are categorical?
3. Which of the above listed variables are quantitative?

Summarizing a single quantitative variable

The `favstats` function gives the summary statistics for a quantitative variable. Here we have the summary statistics for the variable `imdb_score`.

```
movies <- read.csv("data/Movies2016.csv") # Read in data set
movies %>% #Data set piped into...
  summarise(favstats(imdb_score)) #Apply favstats function to imdb_score
```

```
#>   min    Q1 median   Q3 max      mean      sd  n missing
#> 1 3.4 5.65    6.4 7.1 8.2 6.309783 1.086689 92      0
```

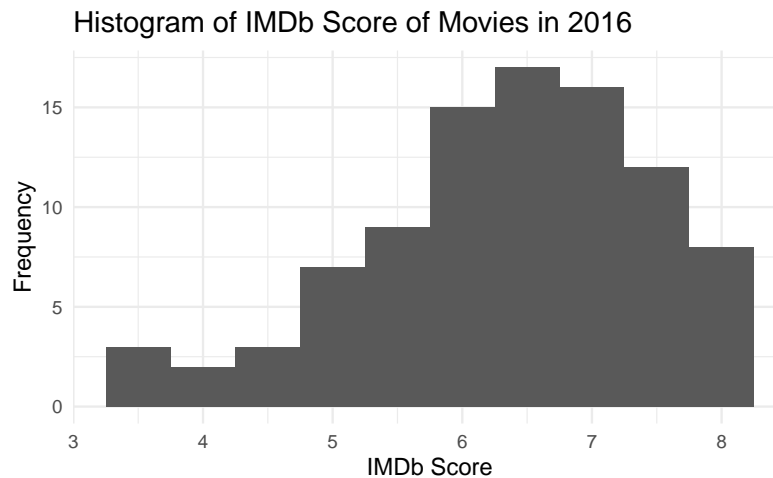
4. Give the values for the two measures of center.
5. Calculate the IQR.
6. Report the value of the standard deviation and interpret this value in context of the problem.

Displaying a single quantitative variable

7. What are the three types of plots used to plot a single quantitative variable?

A histogram of the variable 'IMDb Score' is shown below. Notice that the bin width is 0.5. For example the first bin consists of the number of movies in the data set with an IMDb score of 3.25 to 3.75. It is important to note that a movie with a IMDb score of 4.75 will fall into the bin for 4.75 - 5.25. Visually this shows us the range of IMDb scores for Movies released in 2016.

```
movies %>% #Data set piped into...
ggplot(aes(x = imdb_score)) + #Name variable to plot
  geom_histogram(binwidth = 0.5) + #Create histogram with specified binwidth
  labs(title = "Histogram of IMDb Score of Movies in 2016", #title for plot
        x = "IMDb Score", #Label for x axis
        y = "Frequency") #Label for y axis
```



8. Which range of IMDb scores have the highest frequency?
9. What is the shape of the distribution of IMDb scores?
10. Which five summary statistics are used in creating a box plot? *Hint:* Together they are called the **five-number summary** of the variable.
11. The three smallest IMDb scores in the data set are 3.4, 3.5, and 3.7 and the three largest IMDb scores are 8.5, 8.7, and 9.1:

```
movies %>% # Data set pipes into...
  select(imdb_score) %>% # Select imdb_score variable
  slice_min(imdb_score, n = 3) # Show 3 smallest values
```

```
#>   imdb_score
#> 1      3.4
#> 2      3.5
#> 3      3.7
```

```
movies %>% # Data set pipes into...
  select(imdb_score) %>% # Select imdb_score variable
  slice_max(imdb_score, n = 3) # Show 3 largest values
```

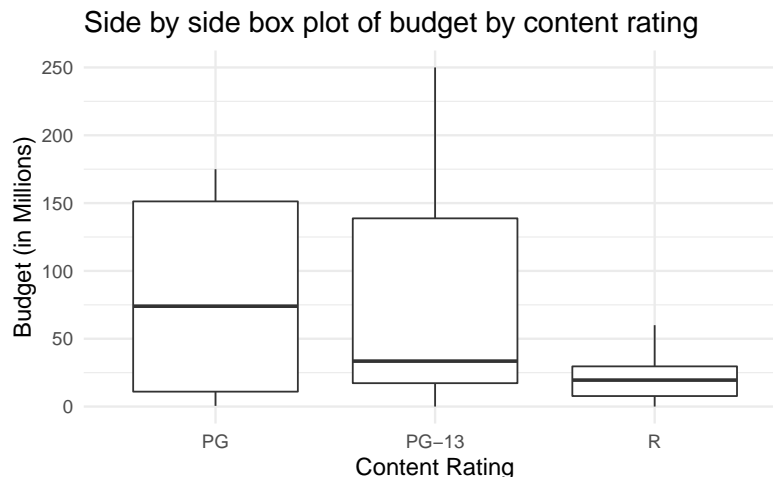
```
#>   imdb_score
#> 1      8.2
#> 2      8.1
#> 3      8.0
```

Using the summary statistics above, and the smallest and largest values of variable to check for outliers, sketch a box plot of IMDb Score. Be sure to label the axes.

Displaying a single categorical and single quantitative variable

The box plot of 'Budget' in millions by 'Content rating' is plotted using the code below. This plot helps to compare the budget for different levels of content rating.

```
movies %>% #Data set piped into...
  filter(content_rating != "Not Rated") %>% # Remove Not Rated movies
  ggplot(aes(y = budget_mil, x = content_rating))+ #Identify variables
  geom_boxplot()+ #Tell it to make a box plot
  labs(title = "Side by side box plot of budget by content rating", #Title
       x = "Content Rating", #x-axis label
       y = "Budget (in Millions)") #y-axis label
```



12. Answer the following questions about the box plots above.

- Which content rating has the highest center?
- Which content rating has the largest spread?

- c. Which content rating is the most symmetric distribution?
- d. Fifty percent of movies in 2016 with a PG-13 content rating fall below what value?
- e. What is the value for the third quartile (Q3) for the PG-13 rating? Interpret this value in context.

4.4 Additional notes

Use this space to summarize your thoughts and take additional notes on today's activity.

Movie Profits

5.1 Learning objectives

- Identify and create appropriate summary statistics and plots given a data set with two quantitative variables
- Use scatterplots to assess the relationship between two quantitative variables
- Find the correlation coefficient
- Find the estimated line of regression using summary statistics and R linear model (`lm`) output
- Interpret the slope coefficient in context of the problem
- Interpret the coefficient of determination in context of the problem

5.2 Terminology review

In today's activity, we will review summary measures and plots for two quantitative variables. Some terms covered in this activity are:

- Scatterplot
- Correlation
- Slope
- Least-squares line of regression
- Coefficient of determination (r -squared)

To review these concepts, see Chapter 3 in the textbook.

5.3 Movies released in 2016

We will revisit the data set used last week collected on Movies released in 2016. Here is a reminder of the variables collected on these movies.

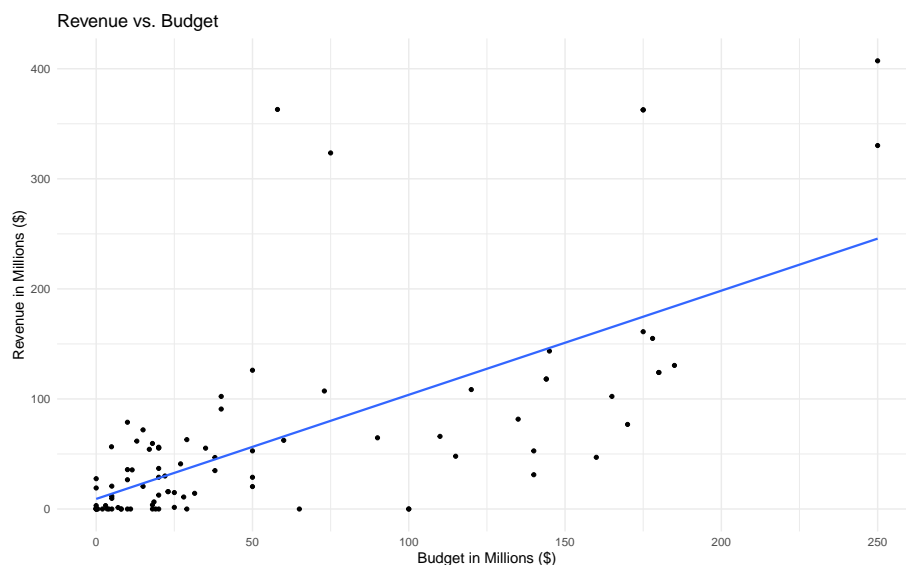
Variable	Description
<code>budget_mil</code>	Amount of money (in US \$ millions) budgeted for the production of the movie
<code>revenue_mil</code>	Amount of money (in US \$ millions) the movie made after release
<code>duration</code>	Length of the movie (in minutes)
<code>content_rating</code>	Rating of the movie (G, PG, PG-13, R, Not Rated)
<code>imdb_score</code>	IMDb user rating score from 1 to 10
<code>genres</code>	Categories the movie falls into (e.g., Action, Drama, etc.)
<code>movie_facebook_likes</code>	Number of likes a movie receives on Facebook

Vocabulary review

1. What type of plot is used to display two quantitative variables?
2. What summary statistics are used to describe the relationship between two quantitative variables?

We will look at the relationship between ‘Budget’ and ‘Revenue’ for movies released in 2016. This shows a scatterplot of ‘Budget’ as a predictor of ‘Revenue’ (note: both variables are measures in “millions of dollars”).

```
movies <- read.csv("data/Movies2016.csv") #Reads in data set
movies %>% #Data set pipes into...
ggplot(aes(x = budget_mil, y = revenue_mil))+ #Specify variables
  geom_point() + #Add scatterplot of points
  labs(x = "Budget in Millions ($)", #Label x-axis
       y = "Revenue in Millions ($)", #Label y-axis
       title = "Revenue vs. Budget") + #Be sure to title your plots
  geom_smooth(method = "lm", se = FALSE) #Add regression line
```



3. Assess the four features of the scatterplot that describe this relationship. Describe each feature using a complete sentence!
- Form (linear, non-linear)
 - Direction (positive, negative)
 - Strength
 - Unusual observations or outliers

4. Does there appear to be an association between 'Budget' and 'Revenue'? Explain.

Correlation

Correlation measures the strength and the direction between two quantitative variables. The closer the value of correlation to + or - 1 the stronger the linear relationship. Values close to zero indicate a very weak linear relationship between the two variables. The following output shows a correlation matrix between several pairs of quantitative variables.

```
# Take subset of variables
movies %>% #Data set pipes into...
  select(c("budget_mil", "revenue_mil", # Take subset of variables
           "duration", "imdb_score",
           "movie_facebook_likes")) %>%
  cor(use="pairwise.complete.obs") %>% # Calculate correlation matrix
  round(3) # Round to 3 decimals
```

```
#>
#>      budget_mil revenue_mil duration imdb_score
#> budget_mil      1.000      0.686    0.463     0.292
#> revenue_mil      0.686      1.000    0.227     0.398
#> duration         0.463      0.227    1.000     0.261
#> imdb_score       0.292      0.398    0.261     1.000
#> movie_facebook_likes 0.678      0.723    0.438     0.309
#>
#>      movie_facebook_likes
#> budget_mil              0.678
#> revenue_mil             0.723
#> duration                0.438
#> imdb_score              0.309
#> movie_facebook_likes    1.000
```

5. Using the output above, which two variables have the strongest correlation?
6. What is the value of correlation between 'Budget' and 'Revenue'?
7. Based on the value of correlation what would the sign of the slope be? Positive or negative? Explain.
8. Does your answer to question 7 match the direction you choose in question 3?
9. Explain why the correlation values on the diagonal are equal to 1.

Slope

The linear model function in R gives us the summary for the least squares regression line. The estimate for (Intercept) is the y-intercept for the line of least squares and the estimate for budget_mil is the value of b_1 , the slope.

```
# Fit linear model: y ~ x
revenueLM <- lm(revenue_mil ~ budget_mil, data=movies)
summary(revenueLM)$coefficients # Display coefficient summary
```

```
#>               Estimate Std. Error  t value    Pr(>|t|)
#> (Intercept)  9.1693054   9.0175499  1.016829 3.119606e-01
#> budget_mil   0.9460001   0.1056786  8.951670 4.339561e-14
```

You may remember from middle and high school that slope = $\frac{\text{rise}}{\text{run}}$.

Using b_1 to represent slope, we can write that as the fraction $\frac{b_1}{1}$.

Therefore, the slope predicts the how much the line will *rise* for each *run* of +1. In other words, as the x variable increases by 1 unit, the y variable is expected to change (increase/decrease) by the value of slope.

10. Write out the least squares line using the summary statistics provided in proper statistical notation.

11. Interpret the value of slope in context of the problem.

12. Using the least squares line from question 10, predict the revenue for a movie with a budget of 165 million.

Residuals

The model we are using assumes the relationship between the two variables follows a straight line. The residuals are the errors, or the part that hasn't been modeled by the line.

$$\text{Data} = \text{Model} + \text{Residual}$$

$$\text{Residual} = \text{Data} - \text{Model}$$

$$e_i = y_i - \hat{y}_i$$

13. The movie, *Independence Day: Resurgence*, had a budget of 165 million and revenue of 102.315 million. Find the residual for this movie.

14. Did the line of regression overestimate or underestimate the revenue for this movie?

Coefficient of determination (squared correlation)

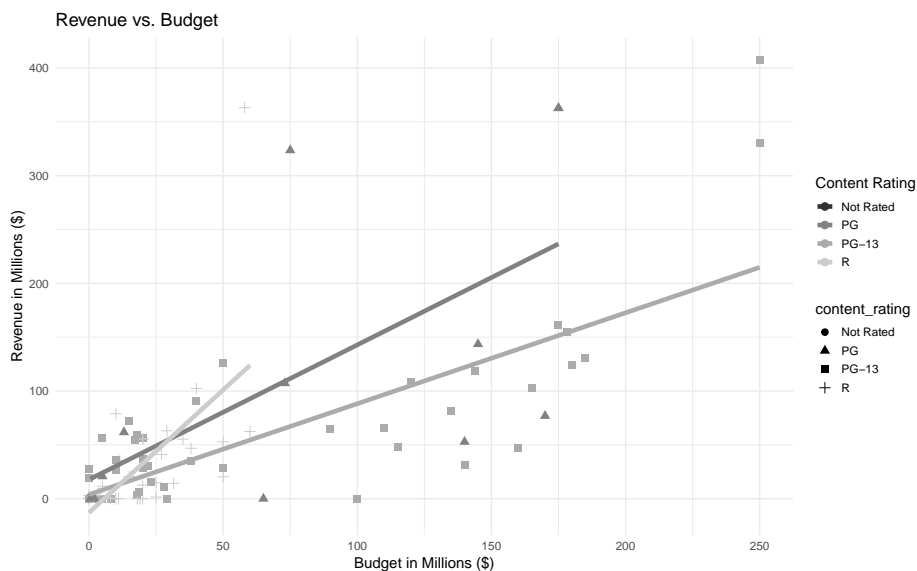
The coefficient of determination, r^2 , can also be used to describe the strength of the linear relationship between two quantitative variables. r^2 measures the proportion of variation in the response that is explained by the least squares line with the explanatory variable.

15. Use the correlation, r , to calculate the coefficient of determination between 'Budget' and 'Revenue', r^2 .
16. Interpret the coefficient of determination in context of the problem.

Multivariate plots

What if we wanted to see if the relationship between 'Budget' and 'Revenue' differs if we add another variable into the picture? The following plot visualized three variables, creating a **multivariate** plot.

```
movies %>% #Data set pipes into...
ggplot(aes(x = budget_mil, y = revenue_mil, color = content_rating)) + #Specify variables
  geom_point(aes(shape = content_rating), size = 3) + #Add scatterplot of points
  labs(x = "Budget in Millions ($)", #Label x-axis
       y = "Revenue in Millions ($)", #Label y-axis
       color = "Content Rating", #Label legend
       title = "Revenue vs. Budget") + #Be sure to tile your plots
  geom_smooth(method = "lm", se = FALSE, lwd = 2) + #Add regression lines
  scale_color_grey() #Make black and white
```



25. Identify the three variables plotted in this graph.
26. Does the relationship between 'Budget' and 'Revenue' differ among the different content ratings? Explain.

5.4 Additional notes

Use this space to summarize your thoughts and take additional notes on today's activity.

Handedness of Male Boxers

6.1 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
- Interpret a confidence interval

6.2 Terminology review

In today's activity, we will introduce simulation hypothesis testing and confidence intervals for a single categorical variable. Some terms covered in this activity are:

- Parameter of interest
- Null hypothesis
- Alternative hypothesis
- Simulation
- Null distribution
- p-value
- Bootstrapping
- Confidence interval

To review these concepts, see Chapter 5 in your textbook, focusing on Sections 5.1 through 5.3.

6.3 Steps of the statistical investigation process

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

6.4 Handedness of male boxers

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

```
handedness <- read.csv("data/Male_boxers_sample.csv") # Read in data set
handedness_sub <- handedness %>%
  select(Stance) # Select Stance variable
dim(handedness_sub) # Check dimensions of data set are 500 rows x 1 col
```

```
#> [1] 500 1
```

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

```
handedness_sub %>% count(Stance) # Count number in each Stance category
```

```
#>      Stance    n  
#> 1 left-handed  81  
#> 2 right-handed 419
```

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, or the true value of interest.

9. Write out the parameter of interest for this study. *Hint:* The parameter of interest is the true proportion of....
10. 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 in place of π_0 .

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. Once we have one simulated sample, what would we calculate and plot on the null distribution?

We will use the computer to simulate a null distribution of 1000 different samples of 500 male boxers, plotting the proportion who are left handed in each sample, based on the assumption that the true proportion of male boxers who are left handed is 0.1 (or that the null hypothesis is true).

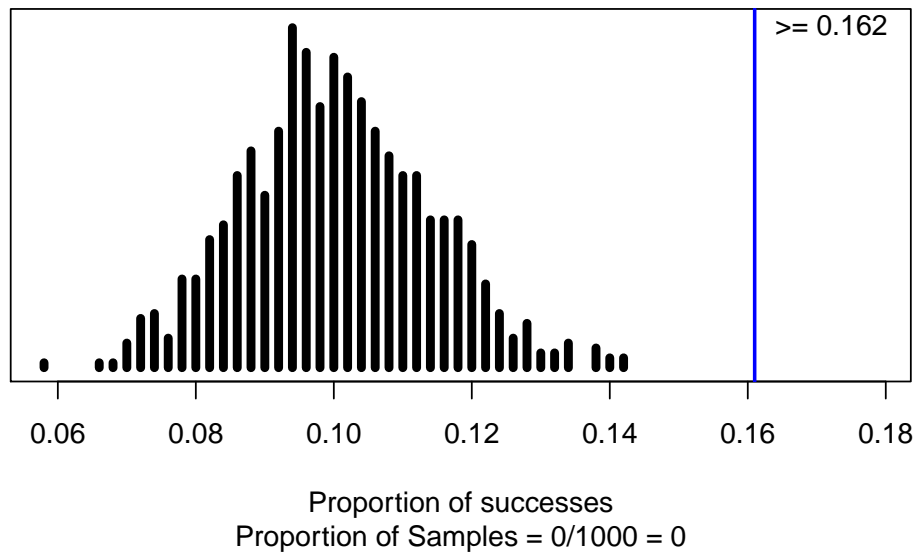
To use the computer simulation, we will need to enter the “probability of success” (π_0), “sample size” (the number of observational units in the sample), “number of repetitions” (the number of samples to be generated), “as extreme as” (the observed statistic), and the “direction” (matches the direction 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 ("greater", "less", or "two-sided"):

The following R code produced the null distribution with 1000 simulations that follows. Check that your answers to question 18 match the code inputs.

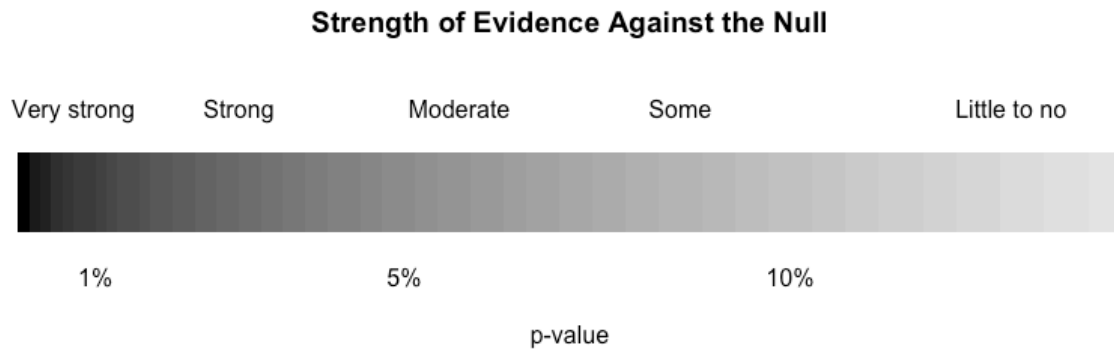
```
one_proportion_test(probability_success = 0.1, #Null hypothesis value
                     sample_size = 500, #Enter sample size
                     number_repetitions = 1000, #Enter number of simulations
                     as_extreme_as = 81/500, #observed statistic
                     direction = "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? (Circle one.)



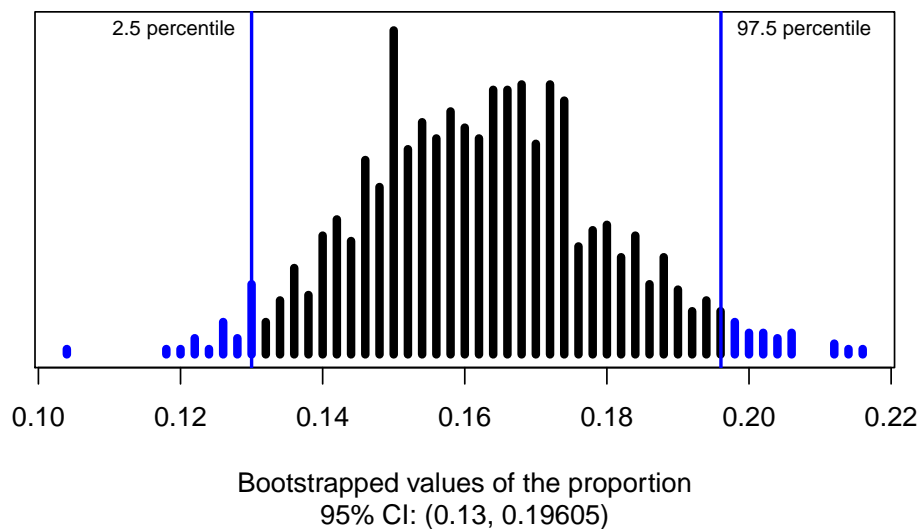
A **point estimate** provides a single plausible value for a parameter. However, a point estimate is rarely perfect; usually there is some error in the estimate. In addition to supplying a point estimate of a parameter, a next logical step would be to provide a plausible range of values for the parameter. 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. In your own words, explain the bootstrapping process.

The following R code produced the following bootstrap distribution with 1000 simulations.

```
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. What is the value at the center of this distribution? Why does this make sense?

26. Explain why the two vertical lines are at the 2.5th percentile and the 97.5th percentile.

27. Report the 95% bootstrapped confidence interval for π . Use interval notation: (lower value, upper value).

28. Interpret the 95% confidence interval in context.

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.

29. Write a paragraph summarizing the results. Be sure to describe:

- Summary statistic
- P-value and interpretation
- Conclusion (written to answer the research question)
- Confidence interval and interpretation
- Generalization - to what group do the results apply?

Revisit and look forward

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

6.5 Additional notes

Use this space to summarize your thoughts and take additional notes on today's activity.

Winter Sports Helmet Use and Head Injuries

7.1 Learning objectives

- Write out the null and alternative hypothesis for two categorical variables
- Assess the conditions to use the standard normal distribution for a difference in proportions
- Calculate the Z test statistic for the difference in proportions
- Find the p-value and assess the strength of evidence
- Create and interpret a confidence interval for the difference in proportions

7.2 Terminology review

In today's activity, we will use theory-based methods to analyze two categorical variables. Some terms covered in this activity are:

- Conditional proportion
- Z test
- z^* multiplier
- Null hypothesis
- Alternative hypothesis
- Test statistic
- Standard normal distribution
- Independence and success-failure conditions
- Type 1 and Type 2 errors
- Decisions

To review these concepts, see Chapter 5 in your textbook.

7.3 Helmet use and head injuries

In “Helmet Use and Risk of Head Injuries in Alpine Skiers and Snowboarders” by Sullheim et. al., in the *Journal of the American Medical Association*, Vol. 295, No. 8 (2006), we can see the summary results from a random sample 3562 skiers and snowboarders involved in accidents in the two-way table below. Is there evidence that safety helmet use reduces the risk of head injury for skiers and snowboarders?

	Helmet Use	No Helmet Use	Total
Head Injury	96	480	576
No Head Injury	656	2330	2986
Total	752	2810	3562

These counts can be found in R by using the `count()` function:

```
injury <- read.csv("data/head_injury.csv") # Read data set in
injury <- # Write over original data with the following
  injury %>% # Pipe data set into
  mutate(Helmet <- factor(Helmet),
         Injury <- factor(Injury)) # Convert to factors

injury %>% group_by(Helmet) %>% count(Injury)
```

```
#> # A tibble: 4 x 3
#> # Groups:   Helmet [2]
#>   Helmet      Injury      n
#>   <chr>      <chr>    <int>
#> 1 No_Helmet Head_Injury    480
#> 2 No_Helmet No_Head_Injury 2330
#> 3 Wore_Helmet Head_Injury     96
#> 4 Wore_Helmet No_Head_Injury  656
```

Vocabulary review

1. What is the explanatory variable?
2. What is the response variable?
3. Is this an experiment or observational study? Justify your answer.
4. Put an X in the box that represents the appropriate scope of inference for this study.

		Study Type	
		Randomized Experiment	Observational Study
Selection of Cases	Random Sample		
	No Random Sample		

5. What is the conditional proportion of skiers/snowboarders with a head injury that wore a helmet?
6. What is the conditional proportion of skiers/snowboarders with a head injury that did not wear a helmet?

Ask a research question

In this study we are looking at the relationship between two groups or two parameters (π_1 and π_2). Remember we define the parameter for a categorical variable as the true proportion of observational units that are labeled as a “success” in the response variable.

7. What is the variable of interest in this study?
8. Write the two parameters of interest for this study. Let 1 = skier/snowboarder wore helmet, 2 = skier/snowboarder did not wear helmet.

π_1 -

π_2 -

When comparing two groups, we assume the two parameters are equal in the null hypothesis—there is no association between the variables.

9. Write the null hypothesis out in words using your answers to question 8.

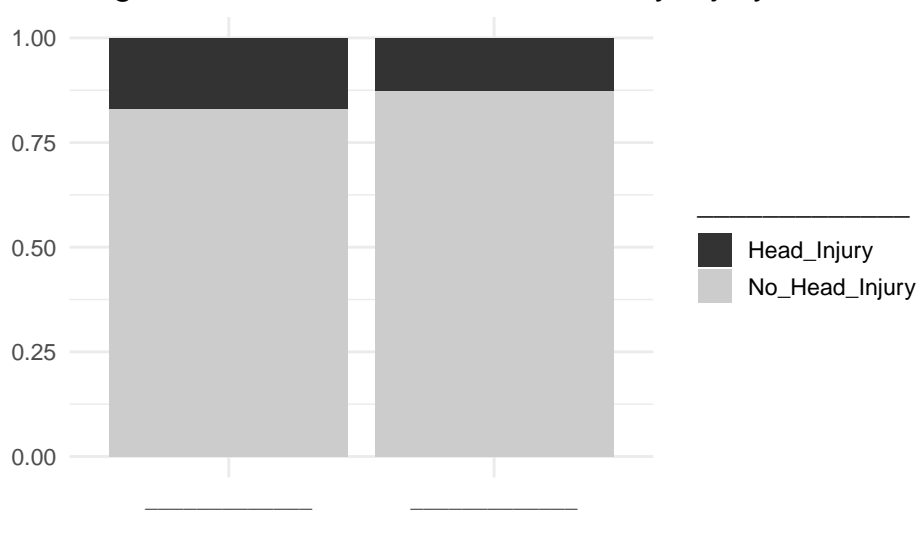
10. What is the research question?

11. Based on the research question fill in the appropriate sign for the alternative hypothesis:

$$H_A : \pi_1 - \pi_2 \text{ _____ } 0$$

Summarize and visualize the data

Segmented Bar Plot of Helmet Use by Injury



12. Fill in the blanks on the graph with the appropriate variables and values to complete the segmented bar plot showing the proportion of head injuries between those who use helmets and those who do not use helmets. *Hint:* Use the conditional proportions from questions 5 and 6.
13. Based on the segmented bar plot, Does there appear to be an association between helmet use and head injury? Explain.
14. Calculate the point estimate for this study. Use helmet use minus no helmet use as the order of subtraction.
15. What is the notation used for the value calculated in question 14?

Use statistical analysis methods to draw inferences from the data

To test the null hypothesis we could use simulation methods as we did with a single categorical variable. In this activity we will focus on theory-based methods. Like with a single proportion, the difference in proportions can be mathematically modeled using the normal distribution if certain conditions are met.

Conditions for the sample distribution of $\hat{p}_1 - \hat{p}_2$:

- Independence: The data are independent within and between the two groups.
- Success-Failure Condition: The success-failure condition holds for each group.

16. Is the independence condition met? Explain your answer.

17. Is the success-failure condition met for each group? Explain your answer.

To calculate the test statistic we use:

$$Z = \frac{\text{point estimate} - \text{null value}}{SE}$$

where the standard error is calculated using the pooled proportion of successes.

$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\hat{p}_{pool}(1 - \hat{p}_{pool})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}, \text{ where}$$
$$\hat{p}_{pool} = \frac{\text{number of "successes"}}{\text{number of cases}} = \frac{\hat{p}_1 n_1 + \hat{p}_2 n_2}{n_1 + n_2}$$

18. Calculate the $SE(\hat{p}_1 - \hat{p}_2)$.

19. Calculate the test statistic.

We will use the `pnorm` function in R to find the p-value. Use the provided R markdown file and enter the value of the test statistic at `xx`.

```
pnorm(xx, # Enter value of test statistic
      m=0, s=1 # Using the standard normal mean = 0, sd = 1
      lower.tail=TRUE) # Gives a p-value less than the test statistic
```

20. Report the p-value.

21. How much evidence does the p-value provide against the null hypothesis?

To find a confidence interval for the difference in proportions we will add and subtract the margin of error from the point estimate to find the two endpoints.

$$\hat{p}_1 - \hat{p}_2 \pm z^* SE(\hat{p}_1 - \hat{p}_2), \text{ where}$$

$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\left(\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2} \right)}$$

Note that the formula changes when calculating the variability around the statistic in order to calculate a confidence interval! Here, use the sample proportions for each group to calculate the standard error for the difference in proportions.

22. Calculate the standard error for a difference in proportions to create a 95% confidence interval.

The z^* multiplier is found under the standard normal distribution. We find the values that encompass the middle 95% of the distribution. If 95% of the standard normal distribution should be in the middle, that leaves 5% in the tails, or 2.5% in each tail. The `qnorm` function in R will tell us the z^* value for the desired percentile (in this case, 95% + 2.5% = 97.5% percentile).

```
qnorm(0.975) # Multiplier for 95% confidence interval
```

```
#> [1] 1.959964
```

23. Using the multiplier of $z^* = 1.96$, calculate the 95% confidence interval for the difference in true proportion of head injuries for those that used helmets minus those who did not.

24. Interpret the confidence interval found in question 23 in context of the problem.

		Test conclusion	
		Fail to reject H_0	Reject H_0
	H_0 true	good decision	Type 1 Error
Truth	H_A true	Type 2 Error	good decision

25. Write a paragraph summarizing the results of the study. Be sure to include:

- Summary statistic
- Test statistic and interpretation
- P-value and interpretation
- Conclusion (written to answer the research question)
- Confidence interval and interpretation
- Scope of inference

Types of errors

Hypothesis tests are not flawless. In a hypothesis test, there are two competing hypotheses: the null and alternative. We make a decision about which might be true, but we may choose incorrectly.

A **Type 1 Error** is rejecting the null hypothesis when H_0 is actually true. A **Type 2 Error** is failing to reject the null hypothesis when the alternative is actually true.

26. Using a significance level of 0.05, what decision do you make in regards to the null hypothesis?

27. What type of error could we have made?

28. Write this error in context of the problem.

7.4 Additional notes

Use this space to summarize your thoughts and take additional notes on today's activity.

COVID-19 and Air Pollution

8.1 Learning outcomes

- Given a research question, construct the null and alternative hypotheses in words and using appropriate statistical symbols
- Describe and perform a simulation-based hypothesis test for paired quantitative data
- Interpret and evaluate a p-value
- Find a confidence interval for the mean difference using bootstrapping
- Interpret a confidence interval
- Use a confidence interval to determine the conclusion of a hypothesis test

8.2 Terminology review

In today's activity, we will analyze paired quantitative data using simulation-based methods. Some terms covered in this activity are:

- Mean difference
- Paired data
- Independent groups
- Shifted null distribution

To review these concepts, see Section 6.2 in the textbook.

8.3 COVID-19 and air pollution

In June 2020, the social distancing efforts and stay-at-home directives to help combat the spread of COVID-19 appeared to help 'flatten the curve' across the United States, albeit at a high cost to many individuals and businesses. The impact of these measures, though, goes far beyond the infection and death rates from the disease. You may have seen images comparing air quality in large international cities like Rome, Milan, Wuhan, and New Delhi such as the one pictured below which seem to indicate, perhaps unsurprisingly, that fewer people driving and factories being shut down have reduced air pollutants.

Have high population-density U.S. cities seen the same improved air quality conditions? To study this question, data was gathered from the U.S. Environmental Protection Agency (EPA) AirData website which records the ozone (O₃) and fine particulate matter (PM_{2.5}) values for cities across the U.S. These measures are used to calculate an air quality index (AQI) score for each city each day of the year. Thirty-three of the most densely populated U.S. cities were selected and the AQI score recorded for April 20, 2020 as well as the five-year median AQI score for April 20th (2015 - 2019). Note that higher AQI scores indicate worse air quality.



Figure 8.1: The India Gate in New Delhi, India

	Mean	Standard deviation	Sample size
Current	$\bar{x}_1 = 47.394$	$s_1 = 14.107$	$n_1 = 33$
5 Year Median	$\bar{x}_2 = 51.545$	$s_2 = 17.447$	$n_2 = 33$
Differences	$\bar{x}_d = -4.152$	$s_d = 17.096$	$n_d = 33$

Vocabulary review

1. What is the sample size?
2. Identify the variables in this study. What role do each have?
3. Why is this treated as a paired study design and not two independent samples?
4. Is this an experiment or observational study? Justify your answer.

Ask a research question

5. What are the two competing possibilities to run a hypothesis test for this study?
6. Write the null hypothesis in words.
7. What is the research question?
8. Write the alternative hypothesis in notation.

Summarize and visualize the data

9. Report the summary statistic for the data.
10. What notation is used for the value in question 9?

Use statistical inferential methods to draw inferences from the data

To simulate the null distribution we will use a bootstrapping method. Recall that the null distribution must be created under the assumption that the null hypothesis is true. Therefore, before bootstrapping we will need to shift each data point by the difference $\mu_0 - \bar{x}$. This will ensure that the mean of the shifted data is μ_0 and that the simulated null distribution will be centered at the null value.

11. Calculate the difference $\mu_0 - \bar{x}$. Will we need to shift the data up or down?
12. Use the provided R markdown file and enter the calculated value from question 11 for xx to simulate the null distribution and enter the summary statistic from question 9 for yy to find the p-value.

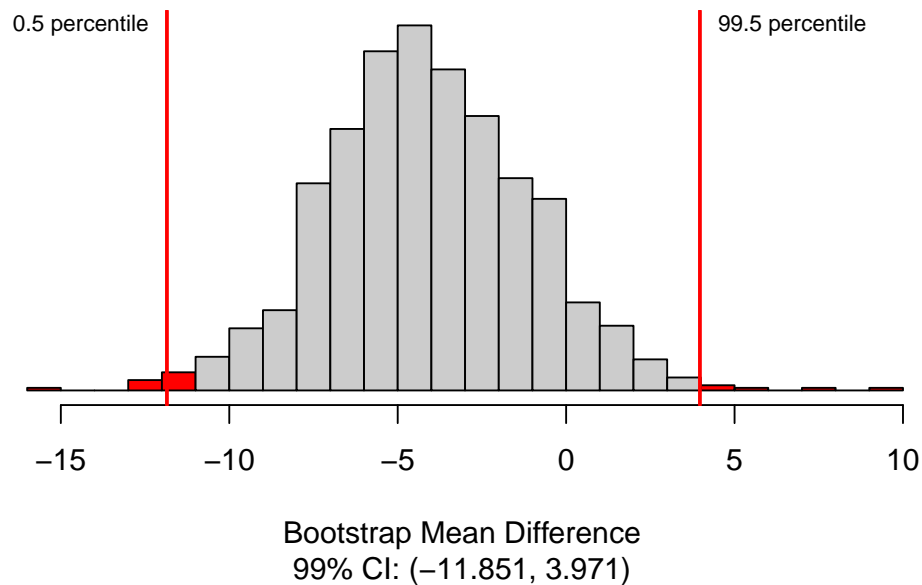
```
paired_test(data = Air$Difference,    #Vector of differences or data set with column for each group
             shift = xx,              #Shift needed for bootstrap hypothesis test
             as_extreme_as = yy,      #Observed statistic
             direction = "less",      #Direction of alternative
             number_repetitions = 1000, #Number of simulated samples for null distribution
             which_first = 1)         #Not needed when using calculated differences
```

13. Sketch the null distribution created in Question 12 here.
14. Explain why the null distribution is centered at zero.
15. What proportion of samples are at or less than the sample mean difference in AQI Scores for current scores minus 5 year median scores?

16. Interpret the p-value in the context of the problem.
17. How much evidence does this provide for improved air quality in US cities?
18. Write out the parameter of interest in context of the study.

The following R code creates a bootstrap distribution showing 1000 simulations of the mean difference.

```
paired_bootstrap_CI(data = Air$Difference, #Enter vector of differences
  number_repetitions = 1000, #Number of bootstrap samples for CI
  confidence_level = 0.99, #Confidence level in decimal form
  which_first = 1) #Not needed when entering vector of differences
```



19. Use the bootstrapped distribution above to find a 99% confidence interval for the parameter of interest. Report the confidence interval in interval notation.

Communicate the results and answer the research question.

20. Interpret the 99% confidence interval in the context of the problem.

21. Do the results of your confidence interval and hypothesis test agree? What does each tell you about the null hypothesis?

22. Write a paragraph summarizes the results of this study. Be sure to include:

- Summary statistic
- P-value and interpretation
- Conclusion (written to answer the research question)
- Confidence interval and interpretation
- Scope of inference

Revisit and look forward

23. Would it be possible to design an experiment to determine if the changed human behavior due to the COVID-19 pandemic causes a decrease in air pollution? Explain.

8.4 Additional notes

Use this space to summarize your thoughts and take additional notes on today's activity.

Weather Patterns and Record Snowfall

9.1 Learning objectives

- Write out the null and alternative hypothesis for one categorical and one quantitative variable
- Calculate and carry-out simulation based hypothesis test for a difference in means
- Interpret and evaluate a p-value
- Find a bootstrap confidence interval for a difference in means
- Interpret a confidence interval
- Use a confidence interval to determine the conclusion of a hypothesis test

9.2 Terminology review

In today's activity, we will use simulation-based methods to analyze one categorical and one quantitative variable, where the groups formed by the categorical variable are independent. Some terms covered in this activity are:

- Independent groups
- Difference in means

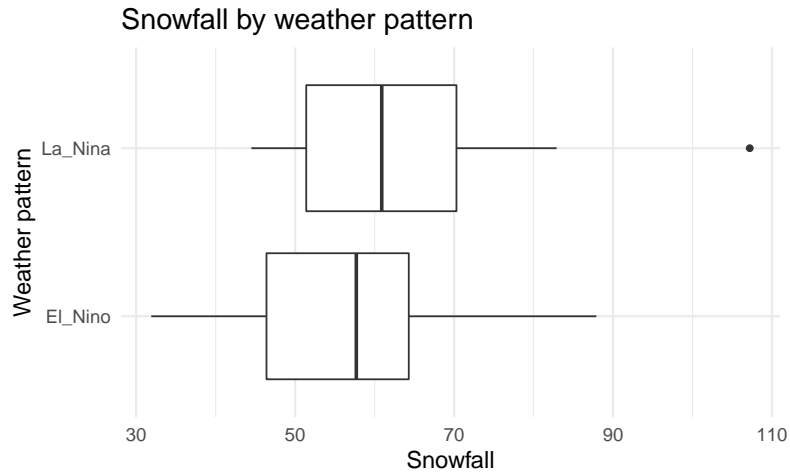
To review these concepts, see Section 6.3 in the textbook.

9.3 Weather patterns and record snowfall

In the winter of 2018-2019, Bozeman had a record snowfall which resulted in the collapse of two flat-roofed buildings on the MSU campus. A writer for the Washington Post predicted the heavy snowfall for 2018-2019 due to the El Niño weather pattern that occurred in that season. A meteorologist in Montana wanted to see if the weather pattern really was associated with total snowfall. She obtained historical data from 44 years on the weather pattern (El Niño or La Niña) and snowfall (in inches) at the Billings Weather Station.

```
Snow <- read.csv("data/SnowfallbyWeatherPattern.csv") # Read in data set
# Code categorical variables as factors
Snow <- # Write over original data with the following
  Snow %>% # Pipe data set into
  mutate(WeatherPattern = factor(WeatherPattern)) # Convert to factor
```

```
# Side-by-side box plots
Snow %>%
ggplot(aes(x = WeatherPattern, y = Snowfall)) +
  geom_boxplot() +
  labs(title = "Snowfall by weather pattern",
       x = "Weather pattern") +
  coord_flip()
```



```
# Summary statistics
Snow %>%
  group_by(WeatherPattern) %>%
  summarise(favstats(Snowfall))
```

```
#> # A tibble: 2 x 10
#>   WeatherPattern   min     Q1 median     Q3    max  mean   sd    n missing
#>   <fct>          <dbl> <dbl>  <dbl> <dbl> <dbl> <dbl> <dbl> <int>    <int>
#> 1 El_Nino       31.9  46.4  57.7  64.3  87.9  56.2  13.0   23      0
#> 2 La_Nina       44.5  51.4  60.9  70.3 107.   63.1  15.5   21      0
```

Quantitative variables review

1. The two variables assessed in this study are the type of weather pattern and snowfall. Identify the role for each variable (explanatory, response).
2. Which group (El Niño or La Niña) has the highest center? Explain which measure you are using.
3. Using the side-by-side boxplots, which group has the largest spread? How did you make that choice?

4. Is this an experiment or an observational study? Justify your answer.

5. Is this a paired data set or two independent groups? Explain your reasoning.

Ask a research question

6. Write out the parameter of interest in context of the study. Use proper notation and be sure to define your subscripts. Use El Niño minus La Niña as the order of subtraction.

7. What are the two competing possibilities we will evaluate in this study?

8. Identify which of your answers in question 7 is the null hypothesis and which is the alternative hypothesis.

Summarize and visualize the data

9. Calculate the summary statistic. Use El Niño minus La Niña as the order of subtraction. What is the appropriate notation for the statistic?

Use statistical inferential methods to draw inferences from the data

Remember that the null distribution is created based on the assumption the null hypothesis is true. In this study, we assume there is no association between variables. This means that a snowfall value could be in either an El Niño year or a La Niña year.

To demonstrate this your instructor will use cards to represent the sample.

10. How many cards will we start with?
11. What will we write on each card?
12. Next we will mix the cards together and shuffle into two piles. How many cards will go into each pile? What should we label the piles?
13. What value is calculated from the cards and plotted on the null distribution?
14. Once we create a null distribution of 1000 simulations, at what value do you expect the distribution to be centered? Explain your reasoning.

Simulation method

15. Using the provided R markdown file, enter the values for the variables, data set, first in subtraction, number of simulations, observed statistic, and direction of the alternative hypothesis.

```
two_mean_test(RESPONSE~PREDICTOR, data = DATASET, #Variables and data
               first_in_subtraction = "VALUE", #First value in order of subtraction
               number_repetitions = ###, #Number of simulations
               as_extreme_as = ###, #Observed statistic
               direction = "?") #Direction of alternative: "greater", "less", or "two-sided"
```

16. Report the p-value. How much evidence does the p-value provide against the null hypothesis?

17. Using bootstrapping find a 90% confidence interval. Use the provided R markdown file. Enter the variables, first in subtraction, number of repetitions, and the confidence level.

```
two_mean_bootstrap_CI(RESPONSE~EXPLANATORY, data = DATASET, #Variables and data
  first_in_subtraction = "VALUE", #First value in order of subtraction
  number_repetitions = ###, #Number of simulations
  confidence_level = ##)
```

18. Interpret the interval you calculated in Question 17.

Communicate the results and answer the research question

19. Write a paragraph summarizing the results of the study. Be sure to include:

- Summary statistic
- P-value and interpretation
- Conclusion (written to answer the research question)
- Confidence interval and interpretation
- Scope of inference

Revisit and look forward

20. Would the results from the theory-based test match the results we saw with the simulation? Explain why or why not.

21. If we had data on 45 La Niña years and 47 El Niño years and found a similar summary statistic, what would happen to the p-value? The width of the confidence interval? The power of the test?

9.4 Additional notes

Use this space to summarize your thoughts and take additional notes on today's activity.

Hand Dexterity

10.1 Learning outcomes

- Given a research question, construct the null and alternative hypotheses in words and using appropriate statistical symbols
- Describe and perform theory-based hypothesis tests for the slope
- Interpret and evaluate a p-value
- Construct and interpret a theory-based confidence interval for slope
- Use a confidence interval to determine the conclusion of a hypothesis test

10.2 Terminology review

In today's activity, we will use theory-based hypothesis tests and confidence intervals for a linear regression slope. Some terms covered in this activity are:

- Correlation
- Slope
- Regression line

To review these concepts, see Chapters 3 and 7 in the textbook.

10.3 Hand dexterity

Physical therapists often evaluate manual (hand) dexterity by having patients complete simple tasks, such as moving pegs on a board or threading objects through holes. Researchers want to examine the manual dexterity of children as part of a follow-up study of a test originally designed for adults to see how manual dexterity changes with age. In this test, 174 participants were given a board with 16 pegs, each in their own hole, arranged in a 4x4 grid. Participants were instructed to pick up the peg with one hand, flip it over by rotating their wrist, then reinsert it in the same hole. Using this test, researchers want to know if as people age the speed at which they can flip all 16 pegs increases.

The variables in this data set¹ consist of the following:

Variable	Description
time	Recorded time to flip all 16 pegs (seconds)
speed	Average speed to flip a peg for each participant (seconds per peg)
age	Age of the participants (years)
dominant	Whether the participant's dominant hand was used, coded as 0 for no, 1 for yes
gender	The participant's gender, recorded as a binary variable, 0 for male, 1 for female

¹Data source: Hand Dexterity in Children: Administration and Normative Values of the Functional Dexterity Test (FDT), Gogola, G., et al., 2013

Variable	Description
HD	The dominant hand of the participant, recorded as R for right hand, L for left hand
handUsed	Which hand the participant used to complete the test, recorded as R for right hand, L for left hand

```
# Read in data set
hands <- read.csv("data/hands.csv")

# Rename variables (odd original coding)
colnames(hands) <- c("time", "speed", "age", "dominant", "gender",
                    "HD", "handUsed")

# Code categorical variables as factors
hands <- # Write over original data with the following
hands %>% # Pipe data set into
mutate(dominant = factor(dominant), # Recode categorical variables as factors
       gender = factor(gender),
       HD = factor(HD),
       handUsed = factor(handUsed))
```

Vocabulary review

1. Explain why regression methods are appropriate to use to address the researchers' question. Make sure you clearly define the variables of interest in your explanation and their roles.
2. What is the scope of inference for this study? Explain your answer.
3. Use the provided R markdown file to create a scatterplot to examine the relationship between the speed at which a participant can flip a peg and the age of the participant by filling in the variable names ('speed' and 'age') for xx and yy. Provide this plot. Based on your plot, does it appear that there is a relationship between age and speed? Note: age should be on the x-axis.

```
hands %>% # Pipe data set into...
ggplot(aes(x = xx, y = yy)) + #Specify variables
  geom_point() + #Add scatterplot of points
  labs(x = "Age (yrs)", #Label x-axis
       y = "Speed (sec/peg)", #Label y-axis
       title = "Scatterplot of Age vs. Speed") + #Be sure to tile your plots
  geom_smooth(method = "lm", se = FALSE) #Add regression line
```

4. Describe the features of the plot you created in Question 3.

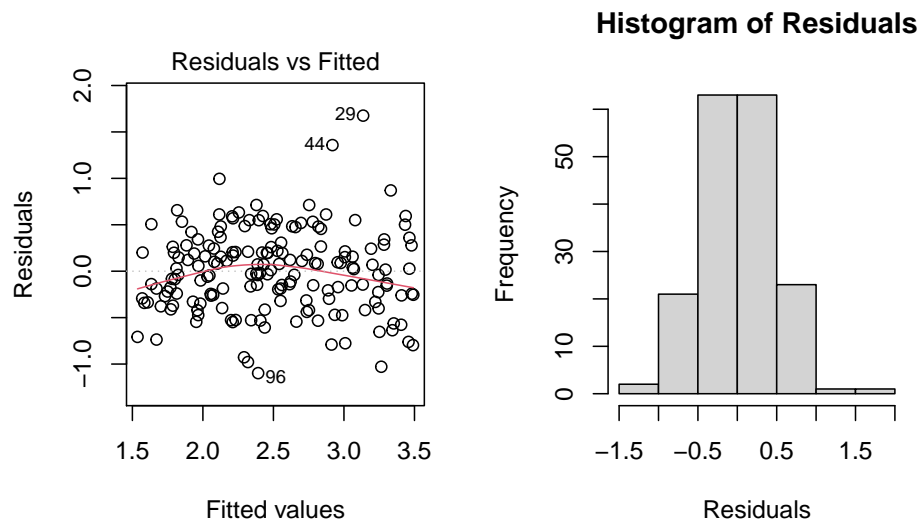
If you indicated there are potential outliers, which points are they?

Conditions for the least squares line

When performing inference on a least squares line, the follow conditions are generally required

- Linearity: the data should follow a linear trend
- Nearly normal residuals: residuals must be nearly normal
- Constant variability: the variability of points around the least squares line remains roughly constant
- Independent observations: individual data points must be independent

The scatterplot and the residual plots will be used to assess the conditions for approximating the data with the t -distribution.



5. Are the conditions met to approximate the t -distribution?

Ask a research question

6. Write out the null hypothesis in words.
7. Using the research question, write the alternative hypothesis in notation.

Summarize and visualize the data

Using the provided R markdown file, enter the response variable into the linear model function for xx and the explanatory variable for yy to get the linear model output.

```
lm.hand <- lm(xx~yy, data=hands) #lm(response~explanatory)
summary(lm.hand)$coefficients
```

8. Using the output from the evaluated R code above, write the equation of the regression line.
9. Interpret the slope in context of the problem.
10. Using your estimated line of best fit, predict the per peg speed for a participant who was 9.18 years old. Show all work.
11. Calculate the residual associated with the observation (9.18, 2.95), using your estimated regression line from question 8.

Use statistical inferential methods to draw inferences from the data

To find the value of the test statistic to test the slope we will use,

$$T = \frac{\text{slope estimate}}{SE} = \frac{b_1}{SE(b_1)}$$

We will use the linear model output above to get the estimate for slope and standard error.

12. Calculate the test statistic for slope. Identify where this calculated value is in the linear model output.

13. Interpret the test statistic in context of the problem.

14. Using the linear model output, report the p-value for the test of significance.

15. Based on the p-value, how much evidence is there against the null hypothesis?

Recall that a confidence interval is calculated by adding and subtracting the margin of error to the point estimate.

$$\text{point estimate} \pm t^* SE(\text{estimate})$$

$$b_1 \pm t^* SE(b_1)$$

The t^* multiplier comes from the t -distribution with $n - 2$ df. Recall for a 95% confidence interval, use the 97.5% percentile (95% of the distribution is in the middle, leaving 2.5% in each tail).

```
qt(0.95+0.025, 172) #95% t* multiplier
```

```
#> [1] 1.973852
```

16. Calculate the 95% confidence interval for the true slope.

Communicate the results and answer the research question

17. Based on the p-value, write a conclusion in context of the problem.
18. Does the p-value agree with the 95% confidence interval? What does each tell you about the null hypothesis?
19. Summarize the results of the study in a written paragraph. Be sure to include.
 - Summary statistic
 - Test statistic and interpretation
 - P-value and interpretation
 - Confidence interval and interpretation
 - Conclusion (written to answer the research question)
 - Scope of inference

Revisit and look forward

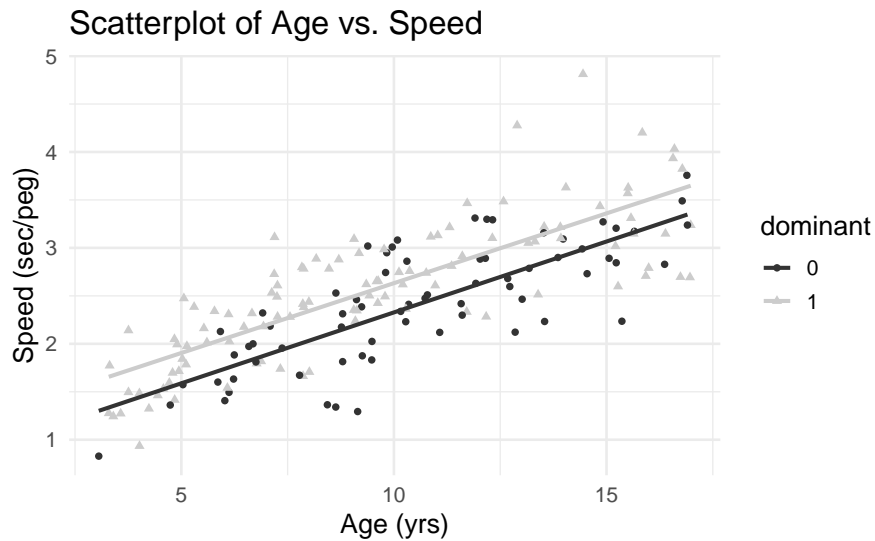
20. Is there an effect due to gender on the linear relationship between age and speed? Explain your answer using the scatterplot below.

```
hands %>% # Pipe data set into...
ggplot(aes(x = age, y = speed, color = dominant))+ #Specify variables
```

```

geom_point(aes(pch = dominant)) + #Add scatterplot of points
labs(x = "Age (yrs)", #Label x-axis
     y = "Speed (sec/peg)", #Label y-axis
     legend = "Dominant hand", #Label your legend
     title = "Scatterplot of Age vs. Speed") + #Be sure to tile your plots
geom_smooth(method = "lm", se = FALSE) + #Add regression line
scale_color_grey() #Make greyscale for printing

```



10.4 Additional notes

Use this space to summarize your thoughts and take additional notes on today's activity.