

Statistics Using Technology: Rguroo Edition

Kathryn Kozak

John Doe

Jane Smith

2025-12-31

Table of contents

1 Home

Preface

This is test of my macros: Inline examples (should be colored): - The variable is x . - The defined term is *standard deviation*. - The descriptor is Frequency. - The dataset is Cars. - The dropdown selection is Mean. - The function is *lm*. - The dialog name is Descriptive Statistics. - The repository is Rguroo Datasets. - The answer is 42.

I hope you find this book useful in teaching statistics. When writing this book, I tried to follow the GAISE Standards (GAISE recommendations).

- Teach statistical thinking.
- Focus on conceptual understanding.
- Integrate real data with a context and a purpose.
- Foster active learning.
- Use technology to explore concepts and analyze data.
- Use assessments to improve and evaluate student learning

To this end, I ask students to interpret the results of their calculations. I incorporated the use of technology (R Studio) for most calculations. Because of that you will not find me using any of the computational formulas for standard deviations or correlation and regression since I prefer students understand the concept of these quantities. Also, because I utilize technology you will not find the standard normal table, Student's t-table, binomial table, chi-square distribution table, and F-distribution table in the book. Another difference between this book and other statistics books is the order of hypothesis testing and confidence intervals. Most books present confidence intervals first and then hypothesis tests. I find that presenting hypothesis testing first and then confidence intervals is more understandable for students. Lastly, I have de-emphasized the use of the z-test. In fact, I only use it to introduce hypothesis testing, and never utilize it again. Two samples should be emphasized over one sample test. Lastly, to aid student understanding and interest, most of the homework and examples utilize real data with multiple variables. The beauty of multiple variables, is that you can ask the students to investigate different analysis with different variables. This way students can work with data and come up with connections of asking questions and using data to answer the questions. Again, I hope you find this book useful for your introductory statistics class.

Mathematical Knowledge Assumed

I want to make a comment about the mathematical knowledge that I assumed the students possess. The course for which I wrote this book has a higher prerequisite than most introductory statistics books. However, I do feel that students can read and understand this book as long as they can read critically. I do not show how to create most of the graphs, but all graphs are created with R Studio. So I hope the mathematical level is appropriate for your course.

Technology Used

The technology that I utilized for creating the graphs and statistical analysis is R Studio. This is a statistical software that are used by statisticians and so using it gives students skills they may need in the future. Please feel free to use any other technology that is more appropriate for your students. Do make sure that you use some technology. I worked on the [StatPREP project](#) and there are Little Apps that can be used to explore data. There are also activities that can be used in your classes that utilize the Little Apps on the website.

Acknowledgments

I would like to thank the following people for taking their valuable time to review the book. Their comments and insights improved this book immensely.

- Daniel Kaplan, Macalester College
- Jane Tanner, Onondaga Community College
- Rob Farinelli, College of Southern Maryland
- Carrie Kinnison, retired engineer
- Sean Simpson, Westchester Community College
- Kim Sonier, Coconino Community College
- Jim Ham, Delta College
- Brian Birgen, Wartburg College
- Christopher Cunningham, Elgin Community College
- Kendra Feinstein, Tacoma Community College
- David Straayer, Tacoma Community College
- Students of Coconino Community College

- Students of Elgin Community College
- Students of Tacoma Community College
- Students of Wartburg College

I also want to thank Coconino Community College for granting me a sabbatical so that I would have the time to write the book. On a personal note, I wanted to thank my brother, John Matic, his wife Jenelle, and their children Hannah and Eli for their hospitality when writing the first edition. In addition to allowing my family access to their home, John provided numerous examples and data sets for business applications in this book. I inadvertently left this thank you out of the first edition of the book, His help and his family's hospitality were invaluable to me. Lastly, I want to thank my husband Rich and my son Dylan for supporting me in this project. Without their love and support, I would not have been able to complete the book.

New to the Fourth Edition

The additions to this edition mostly involve format changes and other edits to make the textbook more accessible for students with visual disabilities. Have a textbook that is accessible to all is very important to me, so please let me know if more changes need to be made. Minor changes and corrections were also made. One change is that every hypothesis test and confidence interval has assumptions that must be true to make the inference valid. Instead of calling them assumptions though, I decided to call them conditions to remove confusion about other assumptions.

Packages Needed for r studio

You will need the following packages installed and loaded in r Studio: arm, HNANES, MASS, mosaic, Weighted.Desc.Stat.

License

Creative Commons Attribution Sharealike.

2025 Kathryn Kozak



ISBN:

2 Statistical Basics

You are exposed to statistics regularly. If you are a sports fan, then you have the statistics for your favorite player. If you are interested in politics, then you look at the polls to see how people feel about certain issues or candidates. If you are an environmentalist, then you research arsenic levels in the water of a town or analyze the global temperatures. If you are in the business profession, then you may track the monthly sales of a store or use quality control processes to monitor the number of defective parts manufactured. If you are in the health profession, then you may look at how successful a procedure is or the percentage of people infected with a disease. There are many other examples from other areas. To understand how to collect data and analyze it, you need to understand what the field of statistics is and the basic definitions.

2.1 What is Statistics?

Statistics is the study of how to collect, organize, analyze, and interpret data collected from a group.

There are two branches of statistics. One is called descriptive statistics, which is where you collect and organize data. The other is called inferential statistics, which is where you analyze and interpret data. First you need to look at descriptive statistics since you will use the descriptive statistics when making inferences.

To understand how to create descriptive statistics and then conduct inferences, there are a few definitions that you need to look at. Note, many of the words that are defined have common definitions that are used in non-statistical terminology. In statistics, some have slightly different definitions. It is important that you notice the difference and utilize the statistical definitions.

The first thing to decide in a statistical study is whom you want to measure and what you want to measure. You always want to make sure that you can answer the question of whom you measured and what you measured. The who is known as the observation and the what is the variable(s).

observation, or simply observations: a person or object that you are interested in finding out information about.

Variable: the measurement or observation of the observation

Having the observation and the variables is part of picture of a **data set** or **data frame**. To make a data set or data frame into what is called tidy data, it should be organized in a way that each row of the data frame is an observation, and the variables should be well defined and are easily identified. An example of a data frame that is tidy data is:

Table 2.1: Example of a Data frame

name	children	fat	type	calories	protein	fat	sodium	fiber	carbo	sugar	potass	vitamin	self	weight	cup	rating
100%_Bran	N	N	C	70	4	1	130	10.0	5.0	6	280	25	3	1	0.33	68.40297
100%_Natural_Bran	N	N	C	120	3	5	15	2.0	8.0	8	135	0	3	1	1.00	33.98368
All-Bran	N	K	C	70	4	1	260	9.0	7.0	5	320	25	3	1	0.33	59.42551
All-Bran_with_Extra_Fiber	N	K	C	50	4	0	140	14.0	8.0	0	330	25	3	1	0.50	93.70491
Almond_Delight	N	R	C	110	2	2	200	1.0	14.0	8	-1	25	3	1	0.75	34.38484
Apple_Cinnamon_Cereal	N	C	C	110	2	2	180	1.5	10.5	10	70	25	1	1	0.75	29.50954

Collecting multiple variables from one observation makes sense. If you wanted to figure out the diameter of breast height of Ponderosa Pine trees in the Coconino National Forest, you need to physically measure a bunch of trees. While you are measuring the diameter, you might also want to measure the height of the tree, if the tree has a bark beetle infestation, the estimated age of the tree, the color of the bark, and how many branches it has. You may only want to estimate the average diameter at breast height, but now you have the ability to estimate other quantities too. No sense walking all over the forest and only measure one thing.

A large data frame is one that has at least 5 variables and at least 1000 units of observations. If a data frame only has 3 variables and 500 rows, that doesn't make it not usable. The 1000 observations and 5 variables is just a guideline to work with.

If you put the observation and the variable into one statement, then you obtain a population.

Population: set of all values of the variable for the entire group of units of observations

Notice, the population answers who you want to measure and what you want to measure. Make sure that your population always answers both of these questions. If it doesn't, then you haven't given someone who is reading your study the entire picture. As an example, if you just say that you are going to collect data from the senators in the U.S. Congress, you haven't told your reader what you are going to collect. Do you want to know their income, their highest degree earned, their voting record, their age, their political party, their gender, their marital status, or how they feel about a particular issue? Without telling what you want to measure, your reader has no idea what your study is actually about.

Sometimes the population is very easy to collect. Such as if you are interested in finding the average age of all of the current senators in the U.S. Congress, there are only 100 senators. This wouldn't be hard to find. However, if instead you were interested in knowing the average

age that a senator in the U.S. Congress first took office for all senators that ever served in the U.S. Congress, then this would be a bit more work. It is still doable, but it would take a bit of time to collect. But what if you are interested in finding the average diameter of breast height of all of the Ponderosa Pine trees in the Coconino National Forest? This would be impossible to actually collect. What do you do in these cases? Instead of collecting the entire population, you take a smaller group of the population, kind of a snap shot of the population. This smaller group is called a sample.

Sample: a subset from the population. It looks just like the population, but contains less data.

In today of big data, there is some confusion between really large data frames and populations. The population is a theoretical concept and even if you have a very large data frame, that doesn't mean you have the population. Most populations are not actually able to be collected. They are considered an ideal that you are trying to make decisions about.

How you collect your sample can determine how accurate the results of your study are. There are many ways to collect samples. Some of them create better samples than others. No sampling method is perfect, but some are better than others. Sampling techniques will be discussed later. For now, realize that every time you take a sample you will find different data values. The sample is a snapshot of the population, and there is more information than is in the picture. The idea is to try to collect a sample that gives you an accurate picture, but you will never know for sure if your picture is the correct picture. Unlike previous mathematics classes where there was always one right answer, in statistics there can be many answers, and you don't know which are right.

Once you have your data frame, either from a population or a sample, you need to know how you want to summarize the data. As an example, suppose you are interested in finding the proportion of people who like a candidate, the average height a plant grows to using a new fertilizer, or the variability of the test scores. Understanding how you want to summarize the data helps to determine the type of data you want to collect. Since the population is what we are interested in, then you want to calculate a number from the population. This is known as a parameter. As mentioned already, you can't really collect the entire population. Even though this is the number you are interested in, you can't really calculate it. Instead you use a number calculated from the sample, called a statistic, to estimate the parameter. Since no sample is exactly the same, the statistic values are going to be different from sample to sample. They estimate the value of the parameter, but again, you do not know for sure if your answer is correct.

Parameter: a number calculated from the population. Usually denoted with a Greek letter. This number is a fixed, unknown number that you want to find.

Statistic: a number calculated from the sample. Usually denoted with letters from the Latin alphabet, though sometimes there is a Greek letter with a $\hat{\cdot}$ (called a hat) above it. Since you can find samples, it is readily known, though it changes depending on the sample taken. It is used to estimate the parameter value.

One last concept to mention is that there are two different types of variables – qualitative (categorical) and quantitative (numerical). Each type of variable has different parameters and statistics that you find. It is important to know the difference between them.

Qualitative or categorical variable: answer is a word or name that describes a quality of the observation

Quantitative or numerical variable: answer is a number, something that can be counted or measured from the observation

2.1.1 Example: Stating Definitions for Qualitative Variable

In 2010, the Pew Research Center questioned 1500 adults in the U.S. to estimate the proportion of the population favoring marijuana use for medical purposes. It was found that 73% are in favor of using marijuana for medical purposes. State the observation, variable, population, and sample.

2.1.1.1 Solution

Observation: a U.S. adult

Variable: the response to the question “should marijuana be used for medical purposes?” This is qualitative data since you are recording a person’s response — yes or no.

Population: set of responses of all adults in the U.S.

Sample: set of responses of 1500 adults in the U.S.

Parameter: proportion of all U.S. Adults who favor marijuana for medical purposes

Statistic — proportion of 1500 U.S. Adults who favor marijuana for medical purposes

2.1.2 Example: Stating Definitions for Qualitative Variable

A parking control officer records the manufacturer of every 5th car in the college parking lot in order to determine the most common manufacturer. State the observation, variable, population, and sample.

2.1.2.1 Solution

Observation: a car in the college parking lot

Variable: the name of the manufacturer. This is qualitative data since you are recording a car type.

Population: set of names of the manufacturer of all cars in the college parking lot.

Sample: set of names of the manufacturer of the a particular number of cars in college parking lot

Parameter: proportion of each car type of all cars in the college parking lot

Statistic: proportion of each car type a particular number of cars in the college parking lot

2.1.3 Example: Stating Definitions for Quantitative Variable

A biologist wants to estimate the average height of a plant that is given a new plant food. She gives 10 plants the new plant food and measures the plant height on day 50. State the observation, variable, population, and sample.

2.1.3.1 Solution

Observation: a plant given the new plant food

Variable: the height of the plant on day 50 (Note: it is not the average height since you cannot measure an average – it is calculated from data.) This is quantitative data since you will have a number.

Population: set of heights on day 50 of all plants when the new plant food is used

Sample: set of heights on day 50 of 10 plants when the new plant food is used

Parameter: average height on day 50 of all plants when the new plant food is used

Statistic: average height on day 50 of 10 plants when the new plant food is used

Note: in Example: Stating Definitions for Qualitative Variable, you most likely will be comparing the new plant food to an old plant food. So you would have more units of observations, but the plants given the new plant food are what you are interested in in this case. You may also want to have measurements on other days after you give the plant food. In your data frame you would need to have many variables besides just the height of the plant on day 50. Examples of variables would be `plant_number`, `fertilizer` (yes or no), `height` on day 20, `height` on day 30, `height` on day 50, and so forth. One other comment, you variable names should make sense to your reader, and be one word for ease in analyzing by a computer program.

2.1.4 Example: Stating Definitions for Quantitative Variable

A doctor wants to see if a new treatment for cancer extends the life expectancy of a patient versus the old treatment. She gives one group of 25 cancer patients the new treatment and another group of 25 the old treatment. She then measures the life expectancy of each of the patients. State the units of observations, variables, populations, and samples.

2.1.4.1 Solution

In this example there are two observations, two variables, two populations, and two samples.

Observation 1: cancer patient given new treatment

Observation 2: cancer patient given old treatment

Variable 1: life expectancy when given new treatment. This is quantitative data since you will have a number.

Variable 2: life expectancy when given old treatment. This is quantitative data since you will have a number.

Population 1: set of life expectancies of all cancer patients given new treatment

Population 2: set of life expectancies of all cancer patients given old treatment

Sample 1: set of life expectancies of 25 cancer patients given new treatment

Sample 2: set of life expectancies of 25 cancer patients given old treatment

Parameter 1: average life expectancy of all cancer patients given new treatment

Parameter 2: average life expectancy of all cancer patients given old treatment

Statistic 1: average life expectancy of 25 cancer patients given new treatment

Statistic 2: average life expectancy of 25 cancer patients given old treatment

There are different types of quantitative variables, called discrete or continuous. The difference is in how many values can the data have. If you can actually count the number of data values (even if you are counting to infinity), then the variable is called discrete. If it is not possible to count the number of data values, then the variable is called continuous.

Discrete data can only take on particular values like integers. Discrete data are usually things you count.

Continuous data can take on any value. Continuous data are usually things you measure.

2.1.5 Example: Discrete or Continuous

Classify the quantitative variable as discrete or continuous.

- a. The weight of a cat.
- b. The number of fleas on a cat.
- c. The size of a shoe.

2.1.5.1 Solution

- a. The weight of a cat.

This is continuous since it is something you measure.

- b. The number of fleas on a cat.

This is discrete since it is something you count.

- c. The size of a shoe.

This is discrete since you can only be certain values, such as 7, 7.5, 8, 8.5, 9. You can't buy a 9.73 shoe.

There are also are four measurement scales for different types of data with each building on the ones below it. They are:

2.1.6 Measurement Scales:

Nominal: data is just a name or category. There is no order to any data and since there are no numbers, you cannot do any arithmetic on this level of data. Examples of this are gender, car name, ethnicity, and race.

Ordinal: data that is nominal, but you can now put the data in order, since one value is more or less than another value. You cannot do arithmetic on this data, but you can now put data values in order. Examples of this are grades (A, B, C, D, F), place value in a race (1st, 2nd, 3rd), and size of a drink (small, medium, large).

Interval: data that is ordinal, but you can now subtract one value from another and that subtraction makes sense. You can do arithmetic on this data, but only addition and subtraction. Examples of this are temperature and time on a clock.

Ratio: data that is interval, but you can now divide one value by another and that ratio makes sense. You can now do all arithmetic on this data. Examples of this are height, weight, distance, and length of time.

Nominal and ordinal data come from qualitative variables. Interval and ratio data come from quantitative variables.

Most people have a hard time deciding if the data are nominal, ordinal, interval, or ratio. First, if the variable is qualitative (words instead of numbers) then it is either nominal or ordinal. Now ask yourself if you can put the data in a particular order. If you can it is ordinal. Otherwise, it is nominal. If the variable is quantitative (numbers), then it is either interval or ratio. For ratio data, a value of 0 means there is no measurement. This is known as the absolute zero. If there is an absolute zero in the data, then it means it is ratio. If there is no absolute zero, then the data are interval. An example of an absolute zero is if you have \ \$0 in your bank account, then you are without money. The amount of money in your bank account is ratio data. **Word of caution:** sometimes ordinal data is displayed using numbers, such as 5 being strongly agree, and 1 being strongly disagree. These numbers are not really numbers. Instead they are used to assign numerical values to ordinal data. In reality you should not perform any computations on this data, though many people do. If there are numbers, make sure the numbers are inherent numbers, and not numbers that were assigned.

2.1.7 Example: Measurement Scale

State which measurement scale each is.

- a. Time of first class
- b. Hair color
- c. Length of time to take a test
- d. Age groupings (baby, toddler, adolescent, teenager, adult, elderly)

2.1.7.1 Solution

- a. Time of first class

This is interval since it is a number, but 0 o'clock means midnight and not the absence of time.

- b. Hair color

This is nominal since it is not a number, and there is no specific order for hair color.

- c. Length of time to take a test.

This is ratio since it is a number, and if you take 0 minutes to take a test, it means you didn't take any time to complete it.

- d. Age groupings (baby, toddler, adolescent, teenager, adult, elderly)

This is ordinal since it is not a number, but you could put the data in order from youngest to oldest or the other way around.

2.1.8 Homework for What is Statistics Section

1. Suppose you want to know how Arizona workers age 16 or older travel to work. To estimate the percentage of people who use the different modes of travel, you take a sample containing 500 Arizona workers age 16 or older. State the observation, variable, population, sample, parameter, and statistic.
2. You wish to estimate the mean cholesterol levels of patients two days after they had a heart attack. To estimate the mean you collect data from 28 heart patients. State the observation, variable, population, sample, parameter, and statistic.
3. Print-O-Matic would like to estimate their mean salary of all employees. To accomplish this they collect the salary of 19 employees. State the observation, variable, population, sample, parameter, and statistic.
4. To estimate the percentage of households in Connecticut which use fuel oil as a heating source, a researcher collects information from 1000 Connecticut households about what fuel is their heating source. State the observation, variable, population, sample, parameter, and statistic.
5. The U.S. Census Bureau needs to estimate the median income of males in the U.S., they collect incomes from 2500 males. State the observation, variable, population, sample, parameter, and statistic.
6. The U.S. Census Bureau needs to estimate the median income of females in the U.S., they collect incomes from 3500 females. State the observation, variable, population, sample, parameter, and statistic.
7. Eyeglassmatic manufactures eyeglasses and they would like to know the percentage of each defect type made. They review 25,891 defects and classify each defect that is made. State the observation, variable, population, sample, parameter, and statistic.
8. The World Health Organization wishes to estimate the mean density of people per square kilometer, they collect data on 56 countries. State the observation, variable, population, sample, parameter, and statistic
9. State the measurement scale for each.
 - a. Cholesterol level
 - b. Defect type
 - c. Time of first class

- d. Opinion on a 5 point scale, with 5 being strongly agree and 1 being strongly disagree
10. State the measurement scale for each.
- a. Temperature in degrees Celsius
 - b. Ice cream flavors available
 - c. Pain levels on a scale from 1 to 10, 10 being the worst pain ever
 - d. Salary of employees

2.2 Sampling Methods

As stated before, if you want to know something about a population, it is often impossible or impractical to examine the whole population. It might be too expensive in terms of time or money. It might be impractical — you can't test all batteries for their length of lifetime because there wouldn't be any batteries left to sell. You need to look at a sample. Hopefully the sample behaves the same as the population.

When you choose a sample you want it to be as similar to the population as possible. If you want to test a new painkiller for adults you would want the sample to include people who are fat, skinny, old, young, healthy, not healthy, male, female, etc.

There are many ways to collect a sample. None are perfect, and you are not guaranteed to collect a representative sample. That is unfortunately the limitations of sampling. However, there are several techniques that can result in samples that give you a semi-accurate picture of the population. Just remember to be aware that the sample may not be representative. As an example, you can take a random sample of a group of people that are equally males and females, yet by chance everyone you choose is female. If this happens, it may be a good idea to collect a new sample if you have the time and money. There are many sampling techniques, though only four will be presented here.

The simplest, and the type that is desired for is a **simple random sample**. This is where you pick the sample such that every sample has the same chance of being chosen. This type of sample is actually hard to collect, since it is sometimes difficult to obtain a complete list of all observations. There are many cases where you cannot conduct a truly random sample. However, you can get as close as you can.

Now suppose you are interested in what type of music people like. It might not make sense to try to find the most popular type of music preferred by everyone in the U.S. You probably don't like the same music as your parents. The answers vary so much you probably couldn't find an answer for everyone all at once. It might make sense to look at people in different age groups, or people of different ethnicities. This is called a **stratified sample**. The issue with this sample type is that sometimes people subdivide the population too much. It is best

to just have one stratification. Also, a stratified sample has similar problems that a simple random sample has.

If your population has some order in it, then you could do a **systematic sample**. This is popular in manufacturing. The problem is that it is possible to miss a manufacturing mistake because of how this sample is taken.

If you are collecting polling data based on location, then a **cluster sample** that divides the population based on geographical means would be the easiest sample to conduct. The problem is that if you are looking for opinions of people, and people who live in the same region may have similar opinions. As you can see each of the sampling techniques have pluses and minuses.

One last type of sample that is sometimes conducted is called a **convenience sample**. This sample is not one that should be conducted since the idea of a convenience sample is that the sample is collected using the most convenient process for the researcher. The researcher may ask people who they know or who are easy to get a hold of, and it is in no way representative of the population.

A **simple random sample (SRS)** of size **n** is a sample that is selected from a population in a way that ensures that every different possible sample of size **n** has the same chance of being selected. Also, every observation associated with the population has the same chance of being selected.

Ways to select a simple random sample:

- Put all names in a hat and draw a certain number of names out.
- Assign each observation a number and use a random number table or a calculator or computer to randomly select the observations that will be measured.

2.2.1 Example: Choosing a Simple Random Sample

Describe how to take a simple random sample from a classroom.

2.2.1.1 Solution

Give each student in the class a number. Using a random number generator you could then pick the number of students you want to pick.

2.2.2 Example: How Not to Choose a Simple Random Sample

You want to choose 5 students out of a class of 20. Give some examples of samples that are **not** simple random samples.

2.2.2.1 Solution

Choose 5 students from the front row. The people in the last row have no chance of being selected. Choose the 5 shortest students. The tallest students have no chance of being selected. Ask your friend to pick numbers that have been assigned to each student. Your friend may prefer certain numbers and picks those. This is not known by your friend, but this happens.

2.2.3 Example: How to Choose a Simple Random Sample using R

You want to take a simple random sample of size 10 from a data frame known as NHANES Table ??, use these steps:

```
library("NHANES") # turns on the package NHANES in R
sample_NHANES<- # gives the new sample a name
  NHANES |> # states the dataframe to collect from
  slice_sample(n=10) # creates a random sample and saves it as Sample_NHANES
options(width = 60)
knitr::kable(sample_NHANES) #displays the sample just created
```

Table 2.2: Random Sample of size 10 from I

[illegible]

Table 2.2: Random Sample of size 10 from I

[illegible]

Stratified sampling is where you break the population into groups called strata, then take a simple random sample from each strata.

For example:

- If you want to look at musical preference, you could divide the observations into age groups and then conduct simple random samples inside each group.
- If you want to calculate the average price of textbooks, you could divide the observations into groups by major and then conduct simple random samples inside each group.

2.2.4 Example: How to Choose a Stratified Sample using R

To take a stratified sample using rStudio of size 20 from NHANES Table ?? using race as the strata, use these steps:

```
library("NHANES") # turns on the package NHANES in R
sample_NHANES<- # gives the new sample a name
  NHANES |> # states the dataframe to collect from
  group_by(Race1) |> # tells what variable is the strata
  slice_sample(n=20) # takes the random sample within each strata
options(width = 60)
knitr::kable(sample_NHANES) #displays the sample just created
```

Table 2.3: Stratified Sample of size 100 from NHANES w

[illegible]

Table 2.3: Stratified Sample of size 100 from NHANES w

[illegible]

Table 2.3: Stratified Sample of size 100 from NHANES w

[illegible]

Table 2.3: Stratified Sample of size 100 from NHANES w

[illegible]

Table 2.3: Stratified Sample of size 100 from NHANES w

[illegible]

Convenience sample is one where the researcher picks observations to be included that are easy for the researcher to collect.

- An example of a convenience sample is if you want to know the opinion of people about the criminal justice system, and you stand on a street corner near the county court house, and questioning the first 10 people who walk by. The people who walk by the county court house are most likely involved in some fashion with the criminal justice system, and their opinion would not represent the opinions of all observations.

On a rare occasion, you do want to collect the entire population. In which case you conduct a census.

A **census** is when every observation is measured.

2.2.5 Example: Sampling type

1. Banner Health is a several state nonprofit chain of hospitals. Management wants to assess the incident of complications after surgery. They wish to use a sample of surgery patients. Several sampling techniques are described below. Categorize each technique as simple random sample, stratified sample, systematic sample, cluster sample, or convenience sampling.
 - a. Obtain a list of patients who had surgery at all Banner Health facilities. Divide the patients according to type of surgery. Draw simple random samples from each group.
 - b. Obtain a list of patients who had surgery at all Banner Health facilities. Number these patients, and then use a random number table to obtain the sample.
 - c. Randomly select some Banner Health facilities from each of the seven states, and then include all the patients on the surgery lists of the states.
 - d. At the beginning of the year, instruct each Banner Health facility to record any complications from every 100th surgery.
 - e. Instruct each Banner Health facilities to record any complications from 20 surgeries this week and send in the results.

2.2.5.1 Solution

- a. Obtain a list of patients who had surgery at all Banner Health facilities. Divide the patients according to type of surgery. Draw simple random samples from each group.

This is a stratified sample since the patients were separated into different stratum and then random samples were taken from each strata. The problem with this is that some

types of surgeries may have more chances for complications than others. Of course, the stratified sample would show you this.

- b. Obtain a list of patients who had surgery at all Banner Health facilities. Number these patients, and then use a random number table to obtain the sample.

This is a random sample since each patient has the same chance of being chosen. The problem with this one is that it will take a while to collect the data.

- c. Randomly select some Banner Health facilities from each of the seven states, and then include all the patients on the surgery lists of the states.

This is a cluster sample since all patients are questioned in each of the selected hospitals. The problem with this is that you could have by chance selected hospitals that have no complications.

- d. At the beginning of the year, instruct each Banner Health facility to record any complications from every 100th surgery.

This is a systematic sample since they selected every 100th surgery. The problem with this is that if every 90th surgery has complications, you wouldn't see this come up in the data.

- e. Instruct each Banner Health facilities to record any complications from 20 surgeries this week and send in the results.

This is a convenience sample since they left it up to the facility how to do it. The problem with convenience samples is that the person collecting the data will probably collect data from surgeries that had no complications.

2.2.6 Homework for Sampling Methods Section

1. Researchers want to collect cholesterol levels of U.S. patients who had a heart attack two days prior. The following are different sampling techniques that the researcher could use. Classify each as simple random sample, stratified sample, systematic sample, cluster sample, or convenience sample.
 - a. The researchers randomly select 5 hospitals in the U.S. then measure the cholesterol levels of all the heart attack patients in each of those hospitals.
 - b. The researchers list all of the heart attack patients and measure the cholesterol level of every 25th person on the list.
 - c. The researchers go to one hospital on a given day and measure the cholesterol level of the heart attack patients at that time.
 - d. The researchers list all of the heart attack patients. They then measure the cholesterol levels of randomly selected patients.

- e. The researchers divide the heart attack patients based on race, and then measure the cholesterol levels of randomly selected patients in each race grouping.
2. The quality control officer at a manufacturing plant needs to determine what percentage of items in a batch are defective. The following are different sampling techniques that could be used by the officer. Classify each as simple random sample, stratified sample, systematic sample, cluster sample, or convenience sample.
 - a. The officer lists all of the batches in a given month. The number of defective items is counted in randomly selected batches.
 - b. The officer takes the first 10 batches and counts the number of defective items.
 - c. The officer groups the batches made in a month into which shift they are made. The number of defective items is counted in randomly selected batches in each shift.
 - d. The officer chooses every 15th batch off the line and counts the number of defective items in each chosen batch.

The officer divides the batches made in a month into which day they were made. Then certain days are picked and every batch made that day is counted to determine the number of defective items.

3. You wish to determine the GPA of students at your school. Describe what process you would go through to collect a sample if you use a simple random sample.
4. You wish to determine the GPA of students at your school. Describe what process you would go through to collect a sample if you use a stratified sample.
5. You wish to determine the GPA of students at your school. Describe what process you would go through to collect a sample if you use a systematic sample.
6. You wish to determine the GPA of students at your school. Describe what process you would go through to collect a sample if you use a cluster sample.
7. You wish to determine the GPA of students at your school. Describe what process you would go through to collect a sample if you use a convenience sample.

2.3 Experimental Design

The section is an introduction to experimental design. This is how to actually design an experiment or a survey so that they are statistical sound. Experimental design is a very involved process, so this is just a small introduction.

2.3.1 Guidelines for planning a statistical study

1. Identify the observations that you are interested in. Realize that you can only make conclusions for these observations. As an example, if you use a fertilizer on a certain genus of plant, you can't say how the fertilizer will work on any other types of plants. However, if you diversify too much, then you may not be able to tell if there really is an improvement since you have too many factors to consider.
2. Specify the variable. You want to make sure this is something that you can measure, and make sure that you control for all other factors too. As an example, if you are trying to determine if a fertilizer works by measuring the height of the plants on a particular day, you need to make sure you can control how much fertilizer you put on the plants (which would be your treatment), and make sure that all the plants receive the same amount of sunlight, water, and temperature.
3. Specify the population. This is important in order for you know what conclusions you can make and what observations you are making the conclusions about.
4. Specify the method for taking measurements or making observations.
5. Determine if you are taking a census or sample. If taking a sample, decide on the sampling method.
6. Collect the data.
7. Use appropriate descriptive statistics methods and make decisions using appropriate inferential statistics methods.
8. Note any concerns you might have about your data collection methods and list any recommendations for future.

There are two types of studies:

An **observational study** is when the investigator collects data merely by watching or asking questions. Nothing is change or controlled

An **experiment** is when the investigator changes a variable or imposes a treatment to determine its effect.

2.3.2 Example: Observational Study or Experiment

State if the following is an observational study or an experiment.

- a. Poll students to see if they favor increasing tuition.
- b. Give some students a tutor to see if grades improve.

2.3.2.1 Solution

- a. Poll students to see if they favor increasing tuition.

This is an observational study. You are only asking a question.

- b. Give some students a tutor to see if grades improve.

This is an experiment. The tutor is the treatment.

2.3.3 Survey

Many observational studies involve surveys. A **survey** uses questions to collect the data and needs to be written so that there is no bias.

2.3.4 Experiment Options

In an experiment, there are different options.

Randomized two-treatment experiment: in this experiment, there are two treatments, and observations are randomly placed into the two groups. Either both groups get a treatment, or one group gets a treatment and the other gets either nothing or a placebo. The group getting either an old treatment, no treatment or a placebo is called the control group. The group getting the treatment is called the treatment group. The idea of the placebo is that a person thinks they are receiving a treatment, but in reality they are receiving a sugar pill or fake treatment. Doing this helps to account for the placebo effect, which is where a person's mind makes their body respond to a treatment because they think they are taking the treatment when they are not really taking the treatment. Note, not every experiment needs a placebo, such when using animals or plants. Also, you can't always use a placebo or no treatment. As an example, if you are testing a new blood pressure medication you can't give a person with high blood pressure a placebo or no treatment because of moral reasons.

Randomized Block Design: a block is a group of subjects that are similar, but the blocks differ from each other. Then randomly assign treatments to subjects inside each block. An example would be separating students into full-time versus part-time, and then randomly picking a certain number full-time students to get the treatment and a certain number part-time students to get the treatment. This way some of each type of student gets the treatment and some do not.

Rigorously Controlled Design: carefully assign subjects to different treatment groups, so that those given each treatment are similar in ways that are important to the experiment. An example would be if you want to have a full-time student who is male, takes only night classes, has a full-time job, and has children in one treatment group, then you need to have the same

type of student getting the other treatment. This type of design is hard to implement since you don't know how many differentiation you would use, and should be avoided.

Matched Pairs Design: the treatments are given to two groups that can be matched up with each other in some ways. One example would be to measure the effectiveness of a muscle relaxer cream on the right arm and the left arm of observations, and then for each observation you can match up their right arm measurement with their left arm. Another example of this would be before and after experiments, such as weight before and weight after a diet.

No matter which experiment type you conduct, you should also consider the following:

Replication: repetition of an experiment on more than one observation so you can make sure that the sample is large enough to distinguish true effects from random effects. It is also the ability for someone else to duplicate the results of the experiment.

Blind study is where the subject used in the study does not know which treatment they are getting or if they are getting the treatment or a placebo.

Double-blind study is where neither the subject used in the study nor the researcher knows who is getting which treatment or who is getting the treatment and who is getting the placebo. This is important so that there can be no bias created by either the subject or the researcher.

One last consideration is the time period that you are collecting the data over. There are three types of time periods that you can consider.

Cross-sectional study: data observed, measured, or collected at one point in time.

Retrospective (or case-control) study: data collected from the past using records, interviews, and other similar artifacts.

Prospective (or longitudinal or cohort) study: data collected in the future from groups sharing common factors.

2.3.5 Homework for Experimental Design Section

1. You want to determine if cinnamon reduces a person's insulin sensitivity. You give patients who are insulin sensitive a certain amount of cinnamon and then measure their glucose levels. Is this an observation or an experiment? Why?
2. You want to determine if eating more fruits reduces a person's chance of developing cancer. You watch people over the years and ask them to tell you how many servings of fruit they eat each day. You then record who develops cancer. Is this an observation or an experiment? Why?
3. A researcher wants to evaluate whether countries with lower fertility rates have a higher life expectancy. They collect the fertility rates and the life expectancies of countries around the world. Is this an observation or an experiment? Why?

4. To evaluate whether a new fertilizer improves plant growth more than the old fertilizer, the fertilizer developer gives some plants the new fertilizer and others the old fertilizer. Is this an observation or an experiment? Why?
5. A researcher designs an experiment to determine if a new drug lowers the blood pressure of patients with high blood pressure. The patients are randomly selected to be in the study and they randomly pick which group to be in. Is this a randomized experiment? Why or why not?
6. Doctors trying to see if a new stent works longer for kidney patients, asks patients if they are willing to have one of two different stents put in. During the procedure the doctor decides which stent to put in based on which one is on hand at the time. Is this a randomized experiment? Why or why not?
7. A researcher wants to determine if diet and exercise together helps people lose weight over just exercising. The researcher solicits volunteers to be part of the study, randomly picks which volunteers are in the study, and then lets each volunteer decide if they want to be in the diet and exercise group or the exercise only group. Is this a randomized experiment? Why or why not?
8. To determine if lack of exercise reduces flexibility in the knee joint, physical therapists ask for volunteers to join their trials. They then randomly select the volunteers to be in the group that exercises and to be in the group that doesn't exercise. Is this a randomized experiment? Why or why not?
9. You collect the weights of tagged fish in a tank. You then put an extra protein fish food in water for the fish and then measure their weight a month later. Are the two samples matched pairs or not? Why or why not?
10. A mathematics instructor wants to see if a computer homework system improves the scores of the students in the class. The instructor teaches two different sections of the same course. One section utilizes the computer homework system and the other section completes homework with paper and pencil. Are the two samples matched pairs or not? Why or why not?
11. A business manager wants to see if a new procedure improves the processing time for a task. The manager measures the processing time of the employees then trains the employees using the new procedure. Then each employee performs the task again and the processing time is measured again. Are the two samples matched pairs or not? Why or why not?
12. The prices of generic items are compared to the prices of the equivalent named brand items. Are the two samples matched pairs or not? Why or why not?
13. A doctor gives some of the patients a new drug for treating acne and the rest of the patients receive the old drug. Neither the patient nor the doctor knows who is getting which drug. Is this a blind experiment, double blind experiment, or neither? Why?

14. One group is told to exercise and one group is told to not exercise. Is this a blind experiment, double blind experiment, or neither? Why?
15. The researchers at a hospital want to see if a new surgery procedure has a better recovery time than the old procedure. The patients are not told which procedure that was used on them, but the surgeons obviously did know. Is this a blind experiment, double blind experiment, or neither? Why?
16. To determine if a new medication reduces headache pain, some patients are given the new medication and others are given a placebo. Neither the researchers nor the patients know who is taking the real medication and who is taking the placebo. Is this a blind experiment, double blind experiment, or neither? Why?
17. A new study is underway to track the eating and exercise patterns of people at different time periods in the future, and see who is afflicted with cancer later in life. Is this a cross-sectional study, a retrospective study, or a prospective study? Why?
18. To determine if a new medication reduces headache pain, some patients are given the new medication and others are given a placebo. The pain levels of a patient are then recorded. Is this a cross-sectional study, a retrospective study, or a prospective study? Why?
19. To see if there is a link between smoking and bladder cancer, patients with bladder cancer are asked if they currently smoke or if they smoked in the past. Is this a cross-sectional study, a retrospective study, or a prospective study? Why?
20. The Nurses Health Survey was a survey where nurses were asked to record their eating habits over a period of time, and their general health was recorded. Is this a cross-sectional study, a retrospective study, or a prospective study? Why?
21. Consider a question that you would like to answer. Describe how you would design your own experiment. Make sure you state the question you would like to answer, then determine if an experiment or an observation is to be done, decide if the question needs one or two samples, if two samples are the samples matched, if this is a randomized experiment, if there is any blinding, and if this is a cross-sectional, retrospective, or prospective study.

2.4 How Not to Do Statistics

Many studies are conducted and conclusions are made. However, there are occasions where the study is not conducted in the correct manner or the conclusion is not correctly made based on the data. There are many things that you should question when you read a study. There are many reasons for the study to have bias in it. Bias is where a study may have a certain slant or preference for a certain result. The following are a list of some of the questions or issues you should consider to help decide if there is bias in a study.

One of the first issues you should ask is who funded the study. If the entity that sponsored the study stands to gain either profits or notoriety from the results, then you should question the results. It doesn't mean that the results are wrong, but you should scrutinize them on your own to make sure they are sound. As an example if a study says that genetically modified foods are safe, and the study was funded by a company that sells genetically modified food, then one may question the validity of the study. Since the company funds the study and their profits rely on people buying their food, there may be bias.

An experiment could have **lurking or confounding variables** when you cannot rule out the possibility that the observed effect is due to some other variable rather than the factor being studied. An example of this is when you give fertilizer to some plants and no fertilizer to others, but the no fertilizer plants also are placed in a location that doesn't receive direct sunlight. You won't know if the plants that received the fertilizer grew taller because of the fertilizer or the sunlight. Make sure you design experiments to eliminate the effects of confounding variables by controlling all the factors that you can.

Over generalization is where you do a study on one group and then try to say that it will happen on all groups. An example is doing cancer treatments on rats. Just because the treatment works on rats does not mean it will work on humans. Another example is that until recently most FDA medication testing had been done on white males of a particular age. There is no way to know how the medication affects other genders, ethnic groups, age groups, and races. The new FDA guidelines stresses using subjects from different groups.

Cause and effect is where people decide that one variable causes the other just because the variables are related. Unless the study was done as an experiment where a variable was controlled, you cannot say that one variable caused the other. There is the possibility that another variable caused both to change. As an example, there is a relationship between number of drownings at the beach and ice cream sales. This does not mean that ice cream sales increasing causes people to drown. Most likely the cause for both increasing is the heat.

Sampling error: This is the difference between the sample results and the true population results. This is unavoidable, and results in the fact that samples are different from each other. As an example, if you take a sample of 5 people's height in your class, you will get 5 numbers. If you take another sample of 5 people's heights in your class, you will likely get 5 different numbers.

Non-sampling error: This is where the sample is collected poorly either through a biased sample or through error in measurements. Care should be taken to avoid this error.

Lastly, there should be care taken in considering the difference between **statistical significance versus practical significance**. This is a major issue in statistics. Something could be statistically significance, which means that a statistical test shows there is evidence to show what you are trying to prove. However, in practice it doesn't mean much or there are other issues to consider. As an example, suppose you find that a new drug for high blood pressure does reduce the blood pressure of patients. When you look at the improvement it actually doesn't amount to a large difference. Even though statistically there is a change, it may not

be worth marketing the product because it really isn't that big of a change. Another consideration is that you find the blood pressure medication does improve a person's blood pressure, but it has serious side effects or it costs a great deal for a prescription. In this case, it wouldn't be practical to use it. In both cases, the study is shown to be statistically significant, but practically you don't want to use the medication. The main thing to remember in a statistical study is that the statistics is only part of the process. You also want to make sure that there is practical significance. One more comment on statistical significance, the American Statistical Association (ASA) recently came out with a statement, "Based on our review of the articles in this special issue and the broader literature, we conclude that it is time to stop using the term 'statistically significant' entirely." (Advanced Solutions International, Inc, 2019) Though the ASA suggests not using this term anymore, there are many studies that have been done in the past that uses this term, so it is presented here. However, it is not a term that should be use and will be down played in the rest of this book.

Surveys have their own areas of bias that can occur. A few of the issues with surveys are in the wording of the questions, the ordering of the questions, the manner the survey is conducted, and the response rate of the survey.

The wording of the questions can cause **hidden bias**, which is where the questions are asked in a way that makes a person respond a certain way. An example is that a poll was done where people were asked if they believe that there should be an amendment to the constitution protecting a woman's right to choose. About 60% of all people questioned said yes. Another poll was done where people were asked if they believe that there should be an amendment to the constitution protecting the life of an unborn child. About 60% of all people questioned said yes. These two questions deal with the same issue, though giving different results, but how the question was asked affected the outcome.

The ordering of the question can also cause **hidden bias**. An example of this is if you were asked if there should be a fine for texting while driving, but proceeding that question is the question asking if you text while drive. By asking a person if they actually partake in the activity, that person now personalizes the question and that might affect how they answer the next question of creating the fine.

Non-response is where you send out a survey but not everyone returns the survey. You can calculate the response rate by dividing the number of returns by the number of surveys sent. Most response rates are around 30-50%. A response rate less than 30% is very poor and the results of the survey are not valid. To reduce non-response, it is better to conduct the surveys in person, though these are very expensive. Phones are the next best way to conduct surveys, emails can be effective, and physical mailings are the least desirable way to conduct surveys.

Voluntary response is where people are asked to respond via phone, email or online. The problem with these is that only people who really care about the topic are likely to call or email. These surveys are not scientific and the results from these surveys are not valid. Note: all studies involve volunteers. The difference between a voluntary response survey and a scientific

study is that in a scientific study the researchers ask the subjects to be involved, while in a voluntary response survey the subjects become involved on their own choosing.

2.4.1 Example: Bias in a Study

Suppose a mathematics department at a community college would like to assess whether computer-based homework improves students' test scores. They use computer-based homework in one classroom with one teacher and use traditional paper and pencil homework in a different classroom with a different teacher. The students using the computer-based homework had higher test scores. What is wrong with this experiment?

2.4.1.1 Solution

Since there were different teachers, you do not know if the better test scores are because of the teacher or the computer-based homework. A better design would be have the same teacher teach both classes. The control group would utilize traditional paper and pencil homework and the treatment group would utilize the computer-based homework. Both classes would have the same teacher, and the students would be split between the two classes randomly. The only difference between the two groups should be the homework method. Of course, there is still variability between the students, but utilizing the same teacher will reduce any other confounding variables.

2.4.2 Example: Cause and Effect

Determine if the one variable did cause the change in the other variable.

- a. Cinnamon was giving to a group of people who have diabetes, and then their blood glucose levels were measured a time period later. All other factors for each person were kept the same. Their glucose levels went down. Did the cinnamon cause the reduction?
- b. There is a link between spray on tanning products and lung cancer. Does that mean that spray on tanning products cause lung cancer?

2.4.2.1 Solution

- a. Cinnamon was giving to a group of people who have diabetes, and then their blood glucose levels were measured a time period later. All other factors for each person were kept the same. Their glucose levels went down. Did the cinnamon cause the reduction?

Since this was a study where the use of cinnamon was controlled, and all other factors were kept constant from person to person, then any changes in glucose levels can be attributed to the use of cinnamon.

- b. There is a link between spray on tanning products and lung cancer. Does that mean that spray on tanning products cause lung cancer?

Since there is only a link, and not a study controlling the use of the tanning spray, then you cannot say that increased use causes lung cancer. You can say that there is a link, and that there could be a cause, but you cannot say for sure that the spray causes the cancer.

2.4.3 Example: Generalizations

- a. A researcher conducts a study on the use of ibuprofen on humans and finds that it is safe. Does that mean that all species can use ibuprofen?
- b. Aspirin has been used for years to bring down fevers in humans. Originally it was tested on white males between the ages of 25 and 40 and found to be safe. Is it safe to give to everyone?

2.4.3.1 Solution

- a. A researcher conducts a study on the use of ibuprofen on humans and finds that it is safe. Does that mean that all species can use ibuprofen?

No. Just because a drug is safe to use on one species doesn't mean it is safe to use for all species. In fact, ibuprofen is toxic to cats.

- b. Aspirin has been used for years to bring down fevers in humans. Originally it was tested on white males between the ages of 25 and 40 and found to be safe. Is it safe to give to everyone?

No. Just because one age group can use it doesn't mean it is safe to use for all age groups. In fact, there has been a link between giving a child under the age of 19 aspirin when they have a fever and Reye's syndrome.

2.4.4 Homework for How Not to Do Statistics Section

1. Suppose there is a study where a researcher conducts an experiment to show that deep breathing exercises helps to lower blood pressure. The researcher takes two groups of people and has one group to perform deep breathing exercises and a series of aerobic exercises every day and the other group was asked to refrain from any exercises. The researcher found that the group performing the deep breathing exercises and the aerobic exercises had lower blood pressure. Discuss any issue with this study.

2. Suppose a car dealership offers a low interest rate and a longer payoff period to customers or a high interest rate and a shorter payoff period to customers, and most customers choose the low interest rate and longer payoff period, does that mean that most customers want a lower interest rate? Explain.
3. Over the years it has been said that coffee is bad for you. When looking at the studies that have shown that coffee is linked to poor health, you will see that people who tend to drink coffee don't sleep much, tend to smoke, don't eat healthy, and tend to not exercise. Can you say that the coffee is the reason for the poor health or is there a lurking variable that is the actual cause? Explain.
4. When researchers were trying to figure out what caused polio, they saw a connection between ice cream sales and polio. As ice cream sales increased so did the incident of polio. Does that mean that eating ice cream causes polio? Explain your answer.
5. There is a positive correlation between having a discussion of gun control, which usually occur after a mass shooting, and the sale of guns. Does that mean that the discussion of gun control increases the likelihood that people will buy more guns? Explain.
6. There is a study that shows that people who are obese have a vitamin D deficiency. Does that mean that obesity causes a deficiency in vitamin D? Explain.
7. A study was conducted that shows that polytetrafluoroethylene (PFOA) (Teflon is made from this chemical) has an increase risk of tumors in lab mice. Does that mean that PFOA's have an increased risk of tumors in humans? Explain.
8. Suppose a telephone poll is conducted by contacting U.S. citizens via landlines about their view of gay marriage. Suppose over 50% of those called do not support gay marriage. Does that mean that you can say over 50% of all people in the U.S. do not support gay marriage? Explain.
9. Suppose that it can be shown to be statistically significant that a smaller percentage of the people are satisfied with your business. The percentage before was 87% and is now 85%. Do you change how you conduct business? Explain?
10. You are testing a new drug for weight loss. You find that the drug does in fact statistically show a weight loss. Do you market the new drug? Why or why not?
11. There was an online poll conducted about whether the mayor of Auckland, New Zealand, should resign due to an affair. The majority of people participating said he should. Should the mayor resign due to the results of this poll? Explain.
12. An online poll showed that the majority of Americans believe that the government covered up events of 9/11. Does that really mean that most Americans believe this? Explain.
13. A survey was conducted at a college asking all employees if they were satisfied with the level of security provided by the security department. Discuss how the results of this question could be biased.

14. An employee survey says, “Employees at this institution are very satisfied with working here. Please rate your satisfaction with the institution.” Discuss how this question could create bias.
15. A survey has a question that says, “Most people are afraid that they will lose their house due to economic collapse. Choose what you think is the biggest issue facing the nation today. a) Economic collapse, b) Foreign policy issues, c) Environmental concerns.” Discuss how this question could create bias.
16. A survey says, “Please rate the career of Roberto Clemente, one of the best right field baseball players in the world.” Discuss how this question could create bias.

3 Graphical Description of Data

In chapter 1, you were introduced to the concepts of population, which again is a collection of all the measurements from the individuals of interest. Remember, in most cases you can't collect the entire population, so you have to take a sample. Thus, you collect data either through a sample or a census. Now you have a large number of data values. What can you do with them? No one likes to look at just a set of numbers. One thing is to organize the data into a table or graph. Ultimately though, you want to be able to use that graph to interpret the data, to describe the distribution of the data set, and to explore different characteristics of the data. The characteristics that will be discussed in this chapter and the next chapter are:

1. Center: middle of the data set, also known as the average.
2. Variation: how much the data varies.
3. Distribution: shape of the data (symmetric, uniform, or skewed).
4. Qualitative data: analysis of the data
5. Outliers: data values that are far from the majority of the data.
6. Time: changing characteristics of the data over time.

This chapter will focus mostly on using the graphs to understand aspects of the data, and not as much on how to create the graphs. There is technology that will create most of the graphs, though it is important for you to understand the basics of how to create them.

This textbook uses RStudio to perform all graphical and descriptive statistics, and all statistical inference. When using RStudio, every command is performed the same way. You start off with a goal(explanatory variable ~ response variable, data=data frame_name,...)

RStudio uses packages to make calculations easier. For this textbook, you will mostly need the package mosaic. There will be others that you will need on occasion, but you will be told that at the time. Most likely, mosaic is already installed in your RStudio. If you wish to install other packages you use the command

```
install.packages("name of package")
```

where you replace the name of package with the package you wish to install.

Once the package is installed, then you will need to tell RStudio you want to use it every time you start RStudio. The command to tell RStudio you want to use a package is

```
library("name of package")
```


You will need to turn on the package mosaic. The NHANES package contains a data frame that is useful. Both are accessed by running the command `library("name of package")`.

Back to the basic command

```
goal(explanatory variable ~ response variable, data=data_frame_name,...)
```

The goal depends on what you want to do. If you want to create a graph then you would need

```
gf_graph_type(explanatory_variable ~ response_variable, data=data_frame_name, ...)
```

As an example if you want to create a density plot of cholesterol levels on day 2 from a data frame called Cholesterol, then your command would be

```
gf_density(~day2, data=Cholesterol)
```

You will see more on what the different commands are that you would use. A word about the ... at the end of the command. That means there are other things you can do, but that is up to you if you want to actually do them. They do not need to be used if you don't want to. The following sections will show you how to create the different graphs that are usually completed in an introductory statistics course.

3.1 Qualitative Data

Remember, qualitative data are words describing a characteristic of the individual. There are several different graphs that are used for qualitative data. These graphs include bar graphs, Pareto charts, and pie charts. Bar graphs can be created using a statistical program like RStudio.

Bar graphs or charts consist of the frequencies on one axis and the categories on the other axis. Drawing the bar graph using r is performed using the following command.

```
gf_bar(~explanatory variable, data=Dataframe)
```

3.1.1 Example: Drawing a Bar Chart

Data was collected for two semesters in a statistics class. The data frame is in Table ???. The command

```
head(data frame)
```

shows the variables and the first few lines of the data set. The data sets are usually larger than what is shown. The head command allows one to see the structure of the data frame.

```
Class<-read.csv( "https://krkozak.github.io/MAT160/class_survey.csv")
knitr::kable(head(Class))
```

Table 3.1: Head of Statistics Class Survey

vehicle	gender	distance_campus	ice_cream	rent	major	height	winter
None	Female	1.5	Cookie Dough	724	Environmental and Sustainability Studies	61	Liked it
Mercury	Female	14.7	Sherbet	200	Administrative Justice	60	Don't like it
Ford	Female	2.4	Chocolate Brownie.	600	Bio Chem	68	Liked it
Toyota	Female	5.2	coffee	0		66	Loved it
Jeep	Male	2.0	Cookie Dough	600	Pre-health Careers	71	Loved it
Subaru	Male	5.0	none	500	Finance	72	No opinion

Every data frame has a code book that describes the data set, the source of the data set, and a listing and description of the variables in the data frame.

Code book for data frame class

Description Survey results from two semesters of statistics classes at Coconino Community College in the years 2018-2019.

Format

This data frame contains the following columns:

vehicle: Type of car a student drives

gender: Self declared gender of a student

distance_campus: how far a student lives from the Lone Tree Campus of Coconino Community College (miles)

ice_cream: favorite ice cream flavor

rent: How much a student pays in rent

major: Students declared major

height: height of the student (inches)

winter: Student's opinion of winter (Love it, Like it, Don't like, No opinion)

Source

Kozak K (2019). Survey results form surveys collected in statistics class at Coconino Community College.

References

Kozak, 2019

Create a bar graph of vehicle type. To do this in RStudio, use the command

```
gf_bar(~variable, data=Data_Frame, ...)
```

where `gf_bar` is the goal, `vehicle` is the name of the response variable (there is no explanatory variable), the data frame is `Class`, and a title was added to the graph.

3.1.1.1 Solution

```
gf_bar(~vehicle, data=Class, title="Bar Chart of Cars driven by students in statistics class")
```

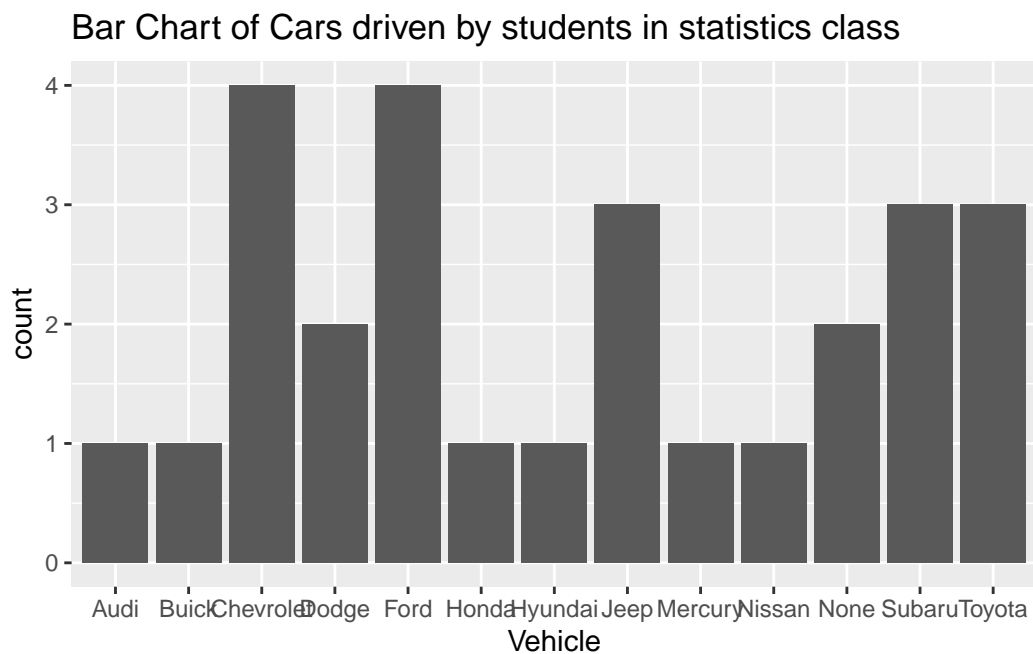


Figure 3.1: Cars driven by students in statistics class

Description of Figure ?? is a Bar graph with bars for Audi, Buick, Honda, Hyundai, Mercury, Nissan with height of 1, Dodge and None with height of 2, Jeep, Subaru, Toyota with heights of 3, and Chevrolet and Ford at height of 4.

Notice from Figure ??, you can see that Chevrolet and Ford are the more popular car, with Jeep, Subaru, and Toyota not far behind. Many types seems to be the lesser used, and tied for last place. However, more data would help to figure this out.

All graphs should have labels on each axis and a title for the graph.

The beauty of data frames with multiple variables is that you can answer many questions from the data. Suppose you want to see if gender makes a difference for the type of car a person drives. If you are a car manufacturer, if you knew that certain genders like certain cars, then you would advertise to the different genders. To create a bar graph that separates based on gender, perform the following command in RStudio.

```
gf_bar(~vehicle, fill=~gender, data=Class, title="Cars driving by students in statistics class")
```

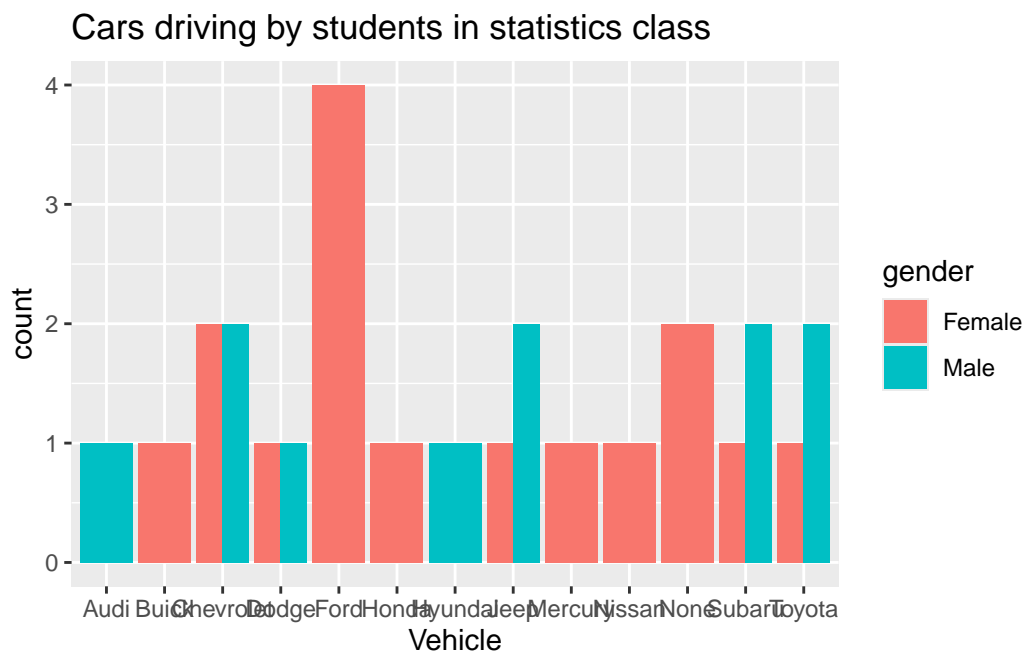


Figure 3.2: Bar graph of Cars driven by students in statistics class

Description of Figure ?? is a bar graph of number of vehicles separated by female and male. Audi and male has height of 1, Buick and female has a height of 1, Chevrolet and male and Chevrolet and female have heights of 2, Dodge and male and Dodge and female has heights of 1, Ford and female has a height of 4, Honda and female has a height of 1, Hyundai and male has a height of 1, Jeep and male has a height of 2 while Jeep and female has a height of 1, Mercury and female has a height of 1, Nissan and female has a height of 1, no car and female has a height of 2, Subaru and female has a height of 1, Subaru and male has a height of 2, Toyota and female has a height of 1, and Toyota and male has a height of 2.

Notice a Ford is driven by females more than any other car, while Chevrolet, Mercury, and Subaru cars are equally driven by males. Obviously a larger sample would be needed to make any conclusions from this data.

There are other types of graphs that can be created for quantitative variables. Another type is known as a dot plot. The command for this graph is as follows.

```
gf_dotplot(~vehicle, data=Class, title="Cars driven by students in statistics class", xlab="Vehicle")
```

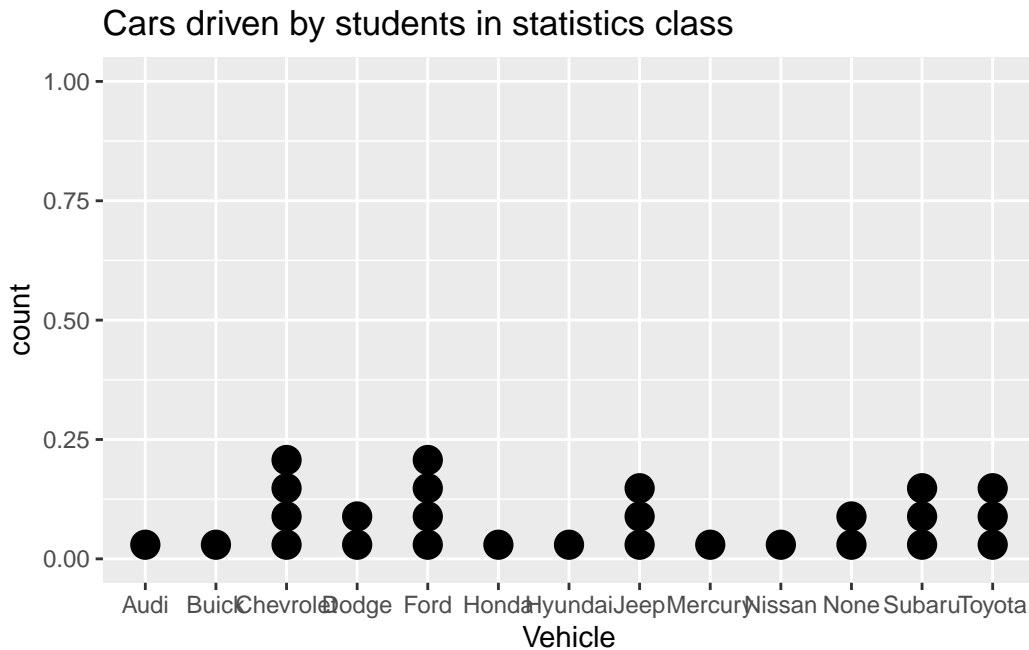


Figure 3.3: Cars driven by students in statistics class

Description of Figure ?? is a dot plot of number of vehicles with Audi, Buick, Honda, Hyundai, Mercury, Nissan with height of 1, Dodge and None with height of 2, Jeep, Subaru, Toyota with heights of 3, and Chevrolet and Ford at height of 4. Very similar to bar graph.

Notice a dot plot is like a bar chart. Both give you the same information. You can also divide a dot plot by gender.

Another type of graph that is also useful and similar to the dot plot is a point plot (scatter plot). In this plot you can graph the explanatory variable versus the response variable. The command for this in rStudio is as follows.

```
gf_point(vehicle~gender, data=Class, title="Cars driving by students in statistics class", xlab="Vehicle")
```

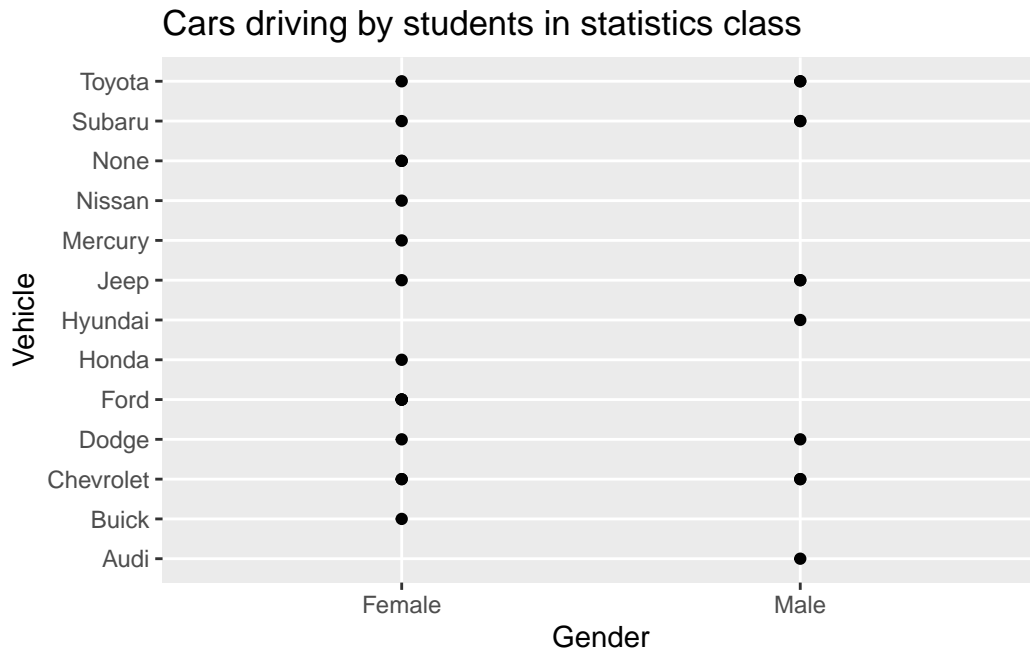


Figure 3.4: Cars driven by students in statistics class

Description of Figure ?? is a scatter plot of type of vehicles separated by female and male with females owning Toyota, Subaru, none, Nissan, Mercury, Jeep, Honda, Ford, Dodge, Chevrolet, and Buick, while males own Toyota, Subaru, Jeep, Hyundai, Dodge, Chevrolet, and Audi.

The problem with Figure ?? is that if there are multiple females who drive a Ford, only one dot is shown. So it is best to spread the dots out using a plot known as a jitter plot. In a jitter plot the dots are randomly moved off the center line. The command for a jitter plot is as follows:

```
gf_jitter(vehicle~gender, data=Class, title="Cars driving by students in statistics class",
```

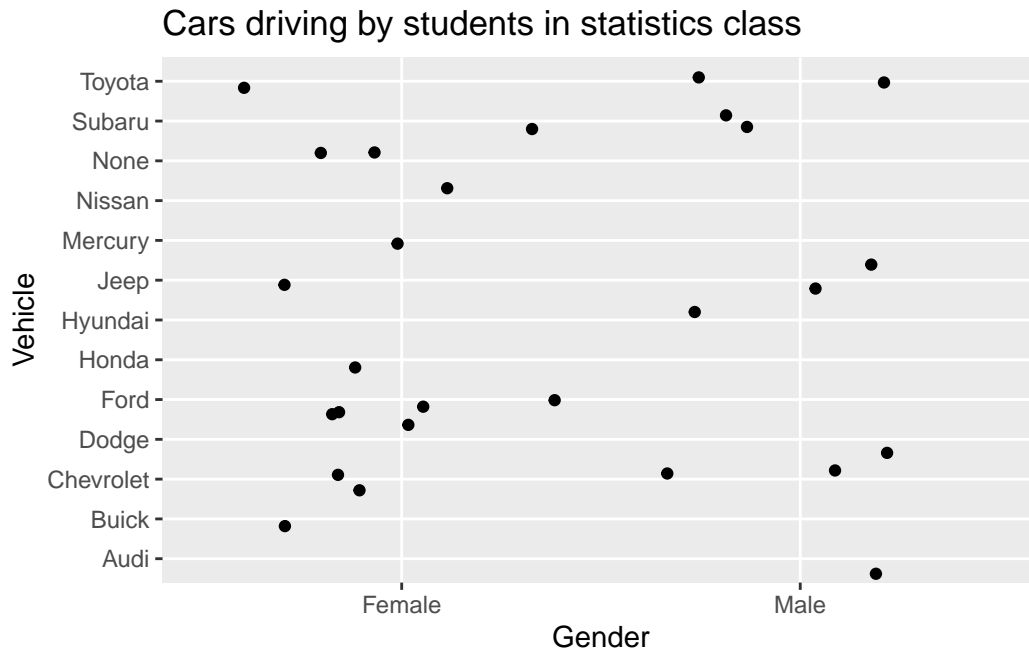


Figure 3.5: Cars driven by students in statistics class

Description of Figure ?? is a jitter plot of number of vehicles separated by female and male with females owning 1 Toyota, 1 Subaru, 2 with none, 1 Nissan, 1 Mercury, 1 Jeep, 1 Honda, 4 Fords, 1 Dodge, 2 Chevrolets, and 1 Buick, while males own 2 Toyotas, 2 Subarus, 2 Jeeps, 1 Hyundai, 1 Dodge, 1 Chevrolets, and 1 Audi.

Now you can observe that there are 4 females who drive a Ford. There is one female who drives a Honda. Other information about other cars and genders can be seen better than in the point plot and the bar graph. Jitter plots are useful to see how many data values are for each qualitative data values.

There are many other types of graphs that can be used on qualitative data. There are spreadsheet software packages that will create most of them, and it is better to look at them to see how to create them. It depends on your data as to which may be useful, but the bar, dot, and jitter plots are really the most useful.

3.1.2 Homework for Qualitative Data Section

1. Eyeglassomatic manufactures eyeglasses for different retailers. The number of lenses for different activities is in Table ??.

```
Eyeglasses<-read.csv( "https://krkozak.github.io/MAT160/eyglasses.csv")
knitr::kable(head(Eyeglasses))
```

Table 3.2: Head of Eyeglasses Data frame

activity
Grind
Grind
Grind
Grind
Grind
Grind

Code book for Data Frame Eyeglasses

Description Activities that an Eyeglass company performs when making eyeglasses, Grind means ground the lenses and put them in frames, multicoat means put tinting or coatings on lenses and then put them in frames, assemble means received frames and lenses from other sources and put them together, make frames means made the frames and put lenses in from other sources, receive finished means received glasses from other source unknown means do not know where the lenses came from.

Format

This data frame contains the following columns:

activity: The activity that is completed to make the eyeglasses by Eyeglassomatic

Source John Matic provided the data from a company he worked with. The company's name is fictitious, but the data is from an actual company.

References John Matic (2013)

Make a bar chart of this data. State any findings you can see from the graph.

2. Data was collected for two semesters in a statistics class drive. The data frame is in Table ??.

Code book for the Data Frame Class is found below Table ??.

Create a bar graph of the variable ice cream. State any findings you can see from the graphs.

3. The number of deaths in the US due to carbon monoxide (CO) poisoning from generators from the years 1999 to 2011 are in Table ?? (Hinaton, 2012). Create a bar chart of this data. State any findings you see from the graph.

```
Area<-read.csv( "https://krkozak.github.io/MAT160/area.csv")
knitr::kable(head(Area))
```


Table 3.3: Head of Area Data frame

deaths
Urban
Urban
Urban
Urban
Urban
Urban

4. Data was collected for two semesters in a statistics class drive. The data frame is in Table ???. Create a bar graph and dot plot of the variable major. Create a jitter plot of major and gender. State any findings you can see from the graphs.

Code book for the Data Frame Class is found below Table ??.

5. Eyeglassomatic manufactures eyeglasses for different retailers. They test to see how many defective lenses they made during the time period of January 1 to March 31. The table Table ?? gives the defect and the number of defects. Create a bar chart of the data and then describe what this tells you about what causes the most defects.

```
Defects<- read.csv( "https://krkozak.github.io/MAT160/defects.csv")
knitr::kable(head(Defects))
```

Table 3.4: Head of Defects Data frame

type
small
small
pd
flaked
scratch
spot

Code book for Data Frame Defects

Description Types of defects that an Eyeglass company sees in the lenses they make into eyeglasses.

Format

This data frame contains the following columns:

Source John Matic provided the data from a company he worked with. The company's name is fictitious, but the data is from an actual company.

6. American National Health and Nutrition Examination (NHANES) surveys is collected every year by the US National Center for Health Statistics (NCHS). The data frame is in Table ?? . Create a bar chart of MartialStatus. Create a jitter plot of MaritalStatus versus Education. Describe any findings from the graphs.

Table 3.5: NHANES Data frame

[illegible]

3.2 Quantitative Data

Histogram: a graph of frequencies (counts) on the vertical axis and classes on the horizontal axis. The height of the rectangles is the frequency and the width is the class width. The

width depends on how many classes (bins) are in the histogram. The shape of a histogram is dependent on the number of bins. In RStudio the command to create a histogram is

```
gf_histogram(~response variable, data=Data_Frame, title="title of the graph")
```

The last part of the command puts a title on the graph. You type in what ever you want for the title in the quotes.

Density Plot: Similar to a histogram, except smoothing is created to smooth out the graph. The shape is not dependent on the number of bins so the distribution is easier to determine from the density plot. In RStudio the command to create a density plot is

```
gf_density(~response variable, data=Data_Frame, title="title of the graph", xlab="Label", ylab="Label")
```

The last part of the command puts a title on the graph and labels on the axes. You type in what every you want for the title and labels in the quotes.

The last part of the command puts a title on the graph and labels on the axes. You type in what every you want for the title and labels in the quotes.

3.2.1 Example: Drawing a Histogram and Density plot

Data was collected for two semesters in a statistics class drive. The data frame is in Table ?? and the code book is below the data frame

Draw a histogram, density plot, and a dot plot for the variable the distance a student lives from the Lone Tree Campus of Coconino Community College. Describe the story the graphs tell.

3.2.1.1 Solution

```
gf_histogram(~distance_campus, data=Class, title="Distance in miles from the Lone Tree Campus")
```

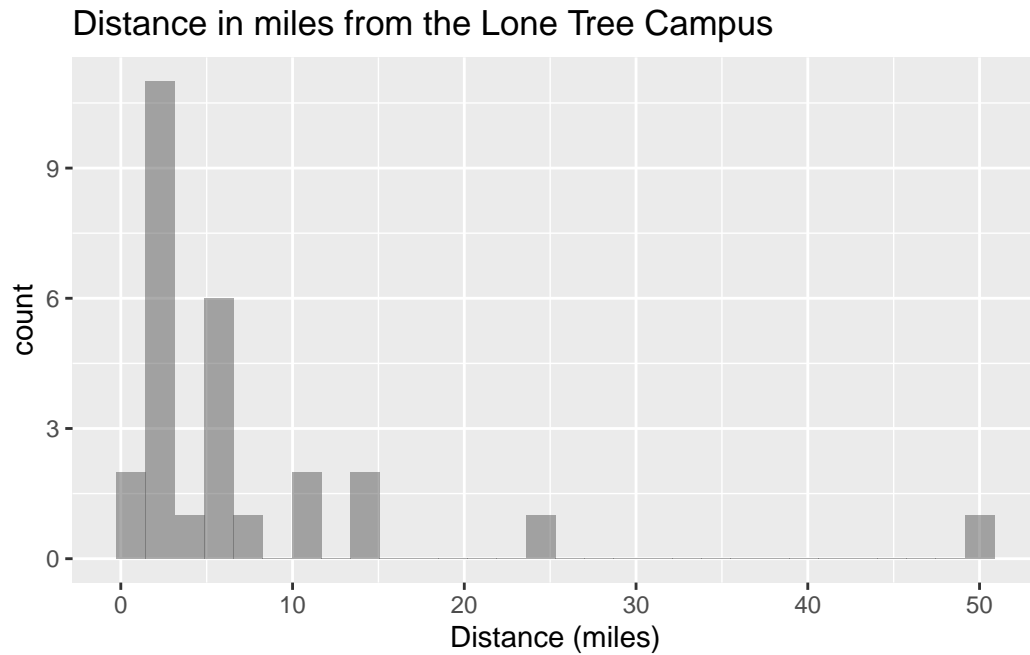


Figure 3.6: Distance in miles from the Lone Tree Campus

Description of the graph is histogram with high part on left and low part on right with several gaps. The graph contains bars.

```
gf_density(~distance_campus, data=Class, title="Distance in miles from the Lone Tree Campus")
```

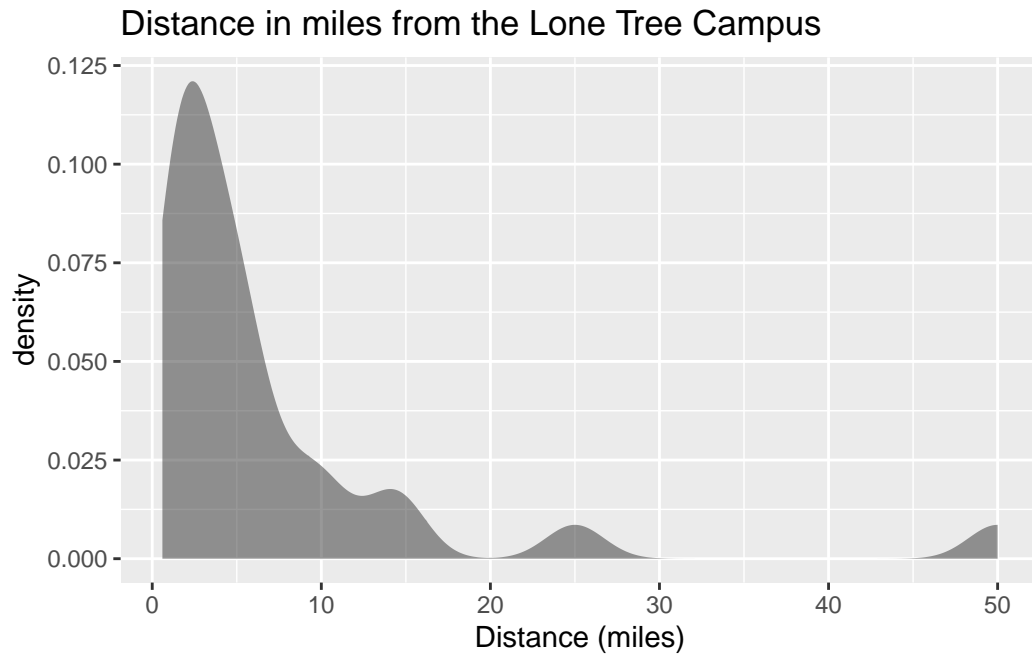


Figure 3.7: Distance in miles from the Lone Tree Campus

Description of the graph is density graph with high part on left and low part on right with several gaps. The graph is smooth.

```
gf_dotplot(~distance_campus, data=Class, title="Distance in miles from the Lone Tree Campus"
```

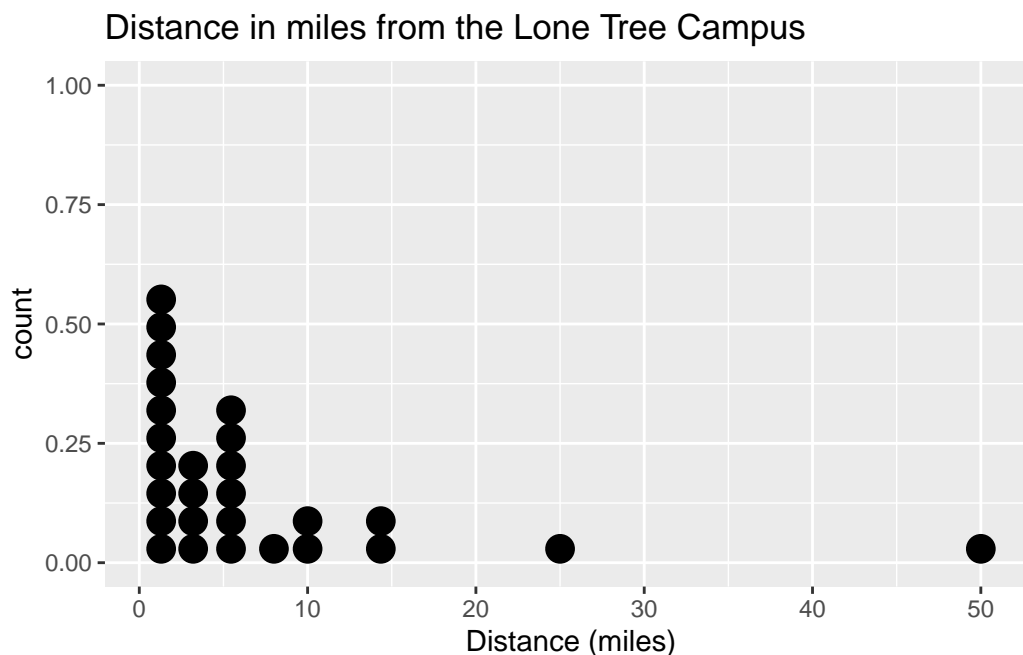


Figure 3.8: Distance in miles from the Lone Tree Campus

Description of the graph of dot plot with high part on left and low part on right with several gaps. The graph is with dots that represent each data value.

Notice the histogram, density plot, and dot plot are all very similar, but the density plot is smoother. They all tell you similar ideas of the shape of the distribution. Reviewing the graphs you can see that most of the students live within 10 miles of the Lone Tree Campus, in fact most live within 5 miles from the campus. However, there is a student who lives around 50 miles from the Lone Tree Campus. This is a great deal farther from the rest of the data. This value could be considered an outlier. An outlier is a data value that is far from the rest of the values. It may be an unusual value or a mistake. It is a data value that should be investigated. In this case, the student lived really far from campus, thus the value is not a mistake, and is just very unusual. The density plot is probably the best plot for most data frames.

There are other aspects that can be discussed, but first some other concepts need to be introduced.

3.2.2 Shapes of the distribution:

When you look at a distribution, look at the basic shape. There are some basic shapes that are seen in histograms. Realize though that some distributions have no shape. The common

shapes are symmetric, skewed, and uniform. Another interest is how many peaks a graph may have. This is known as modal.

Symmetric means that you can fold the graph in half down the middle and the two sides will line up. You can think of the two sides as being mirror images of each other. Skewed means one “tail” of the graph is longer than the other. The graph is skewed in the direction of the longer tail (backwards from what you would expect). A uniform graph has all the bars the same height.

Modal refers to the number of peaks. Unimodal has one peak and bimodal has two peaks. Usually if a graph has more than two peaks, the modal information is not longer of interest.

Other important features to consider are gaps between bars, a repetitive pattern, how spread out is the data, and where the center of the graph is.

3.2.3 Examples of graphs:

This graph is roughly symmetric and unimodal:

Graph: Symmetric Distribution

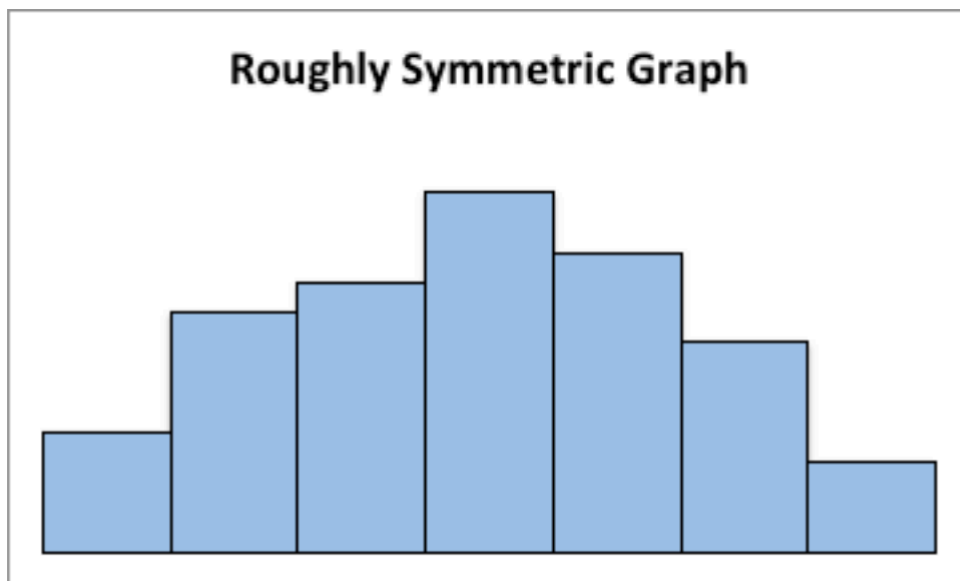


Figure 3.9: symmetric Graph

This graph is symmetric and bimodal:

Graph: Symmetric and Bimodal Distribution

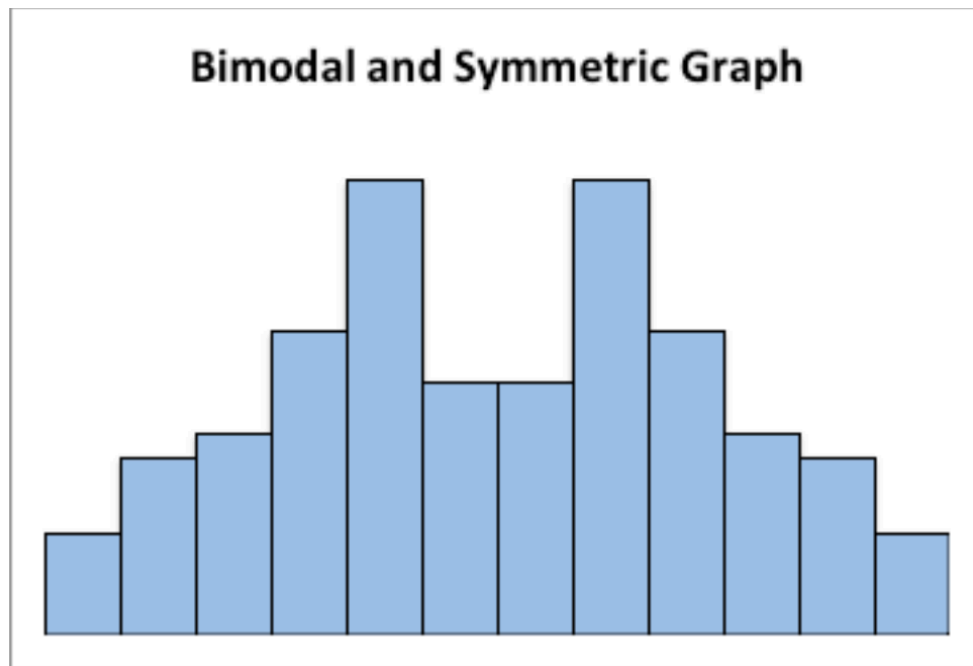


Figure 3.10: Bimodal and symmetric graph

This graph is skewed to the right:

Graph: Skewed Right Distribution

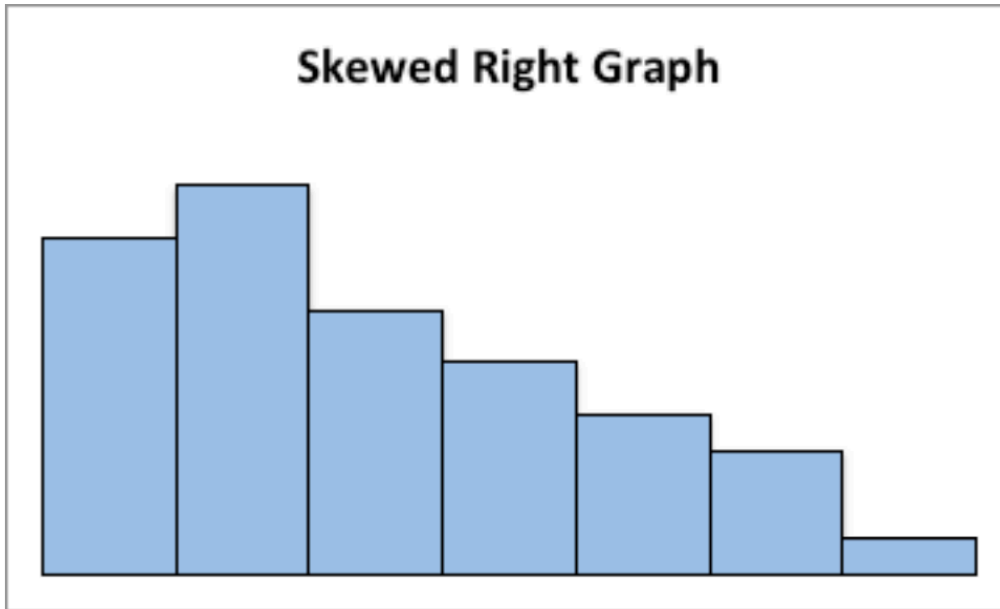


Figure 3.11: Skewed right graph

This graph is skewed to the left and has a gap:

Graph: Skewed Left Distribution

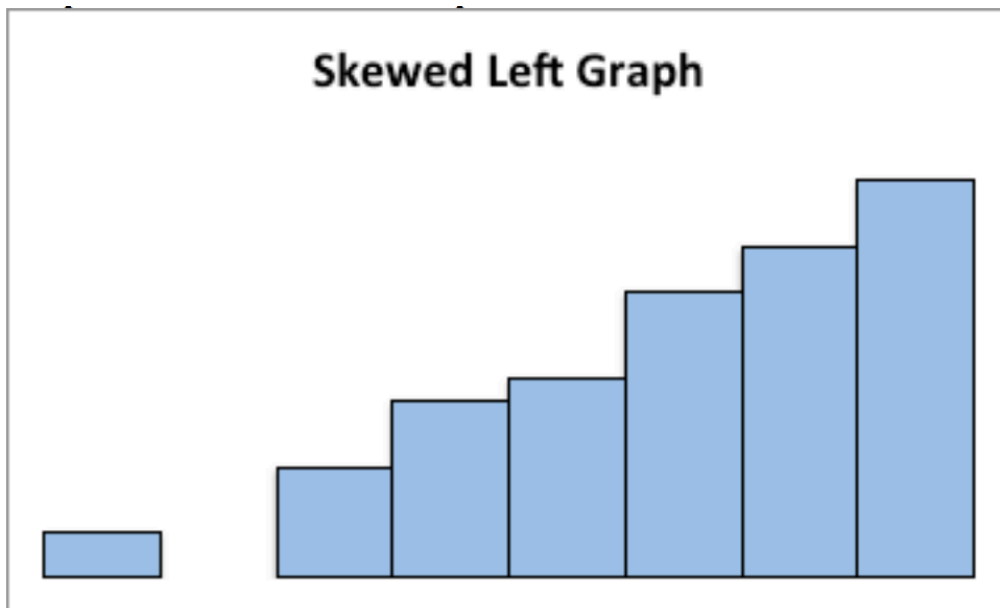


Figure 3.12: Skewed Left graph

This graph is uniform since all the bars are the same height:

Graph: Uniform Distribution

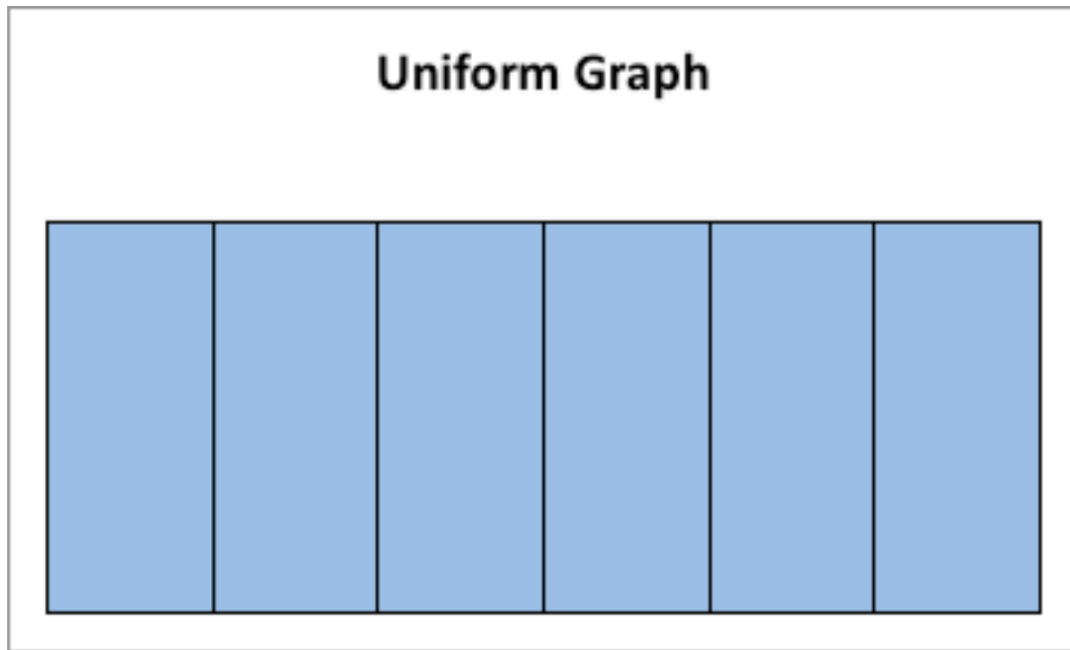


Figure 3.13: Uniform graph

3.2.4 Example: Drawing a Histogram and Density plot

Data was collected from the Chronicle of Higher Education for tuition from public four year colleges, private four year colleges, and for profit four year colleges. The data frame is in Table ???. Draw a density plot of instate tuition levels for all four year institutions, and then separate the density plot for instate tuition based on type of institution. Describe any findings from the graph.

```
Tuition<-read.csv( "https://krkozak.github.io/MAT160/Tuition_4_year.csv")
knitr::kable(head(Tuition))
```

Table 3.6: Head of Tuition Data Frame

INSTITUTION	TYPE	STATE	ROOM_INSTATE	ROOM_OUTSTATE	ROOM_TOT	INSTATE	OUTSTATE	TOT
University of Alaska Anchorage	Public_4 year	AK	12200	7688	19888	23858	36058	

Table 3.6: Head of Tuition Data Frame

INSTITUTION	TYPE	STATE	ROOM_BOARD	INSTATE_TUTION	OUTOFSTATE_TUTION	INSTATE_TOTAL	OUTOFSTATE_TOTAL
University of Alaska Fairbanks	Public_4_year	AK	8930	8087	17017	24257	33187
University of Alaska Southeast	Public_4_year	AK	9200	7092	16292	19404	28604
Alaska Bible College	Private_4_year	AK	5700	9300	15000	9300	15000
Alaska Pacific University	Private_4_year	AK	7300	20830	28130	20830	28130
Alabama Agricultural and Mechanical University	Public_4_year	AL	8379	9698	18077	17918	26297

Code book for Data Frame Tuition

Description Cost of four year institutions.

Format

This data frame contains the following columns:

INSTITUTION: Name of four year institution

TYPE: Type of four year institution, Public_4_year, Private_4_year, For_profit_4_year.

STATE: What state the institution resides

ROOM_BOARD: The cost of room and board at the institution (\\$)

INSTATE_TUTION: The cost of instate tuition (\\$)

INSTATE_TOTAL: The cost of room and board and instate tuition (\\$ per year)

OUTOFSTATE_TUTION: The cost of out of state tuition (\\$ per year)

OUTOFSTATE_TOTAL: The cost of room and board and out of state tuition (\\$ per year)

Source Tuition and Fees, 1998-99 Through 2018-19. (2018, December 31). Retrieved from <https://www.chronicle.com/interactives/tuition-and-fees>

References Chronicle of Higher Education *, December 31, 2018.

3.2.4.1 Solution

```
gf_density(~INSTATE_TUITION, data=Tuition, title="Instate Tuition at all Four Year institutions")
```

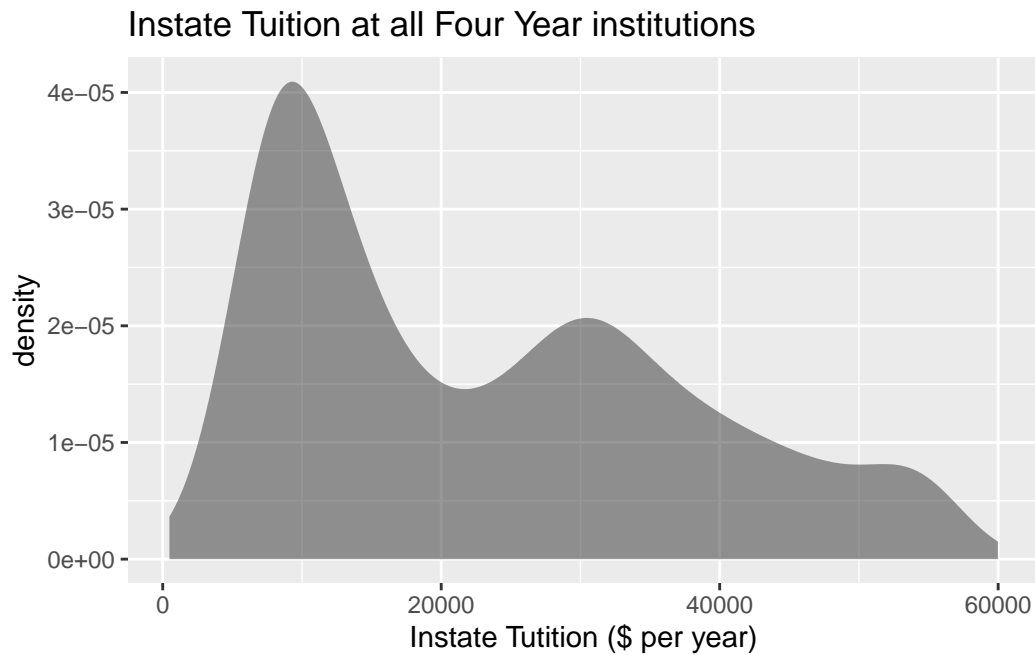


Figure 3.14: Density Plot for Instate Tuition Levels at all Four-Year Colleges

Description of the graph is a density with high part on left, then a dip and up to peak in the middle that is lower than the left peak and then the lowest peak on the right .

(ref:tuition-instate-type-cap) Density Plot for Instate Tuition Levels at all Four-Year Colleges

```
gf_density(~INSTATE_TUITION|TYPE, data=Tuition, title="Instate Tuition at all Four Year inst.")
```

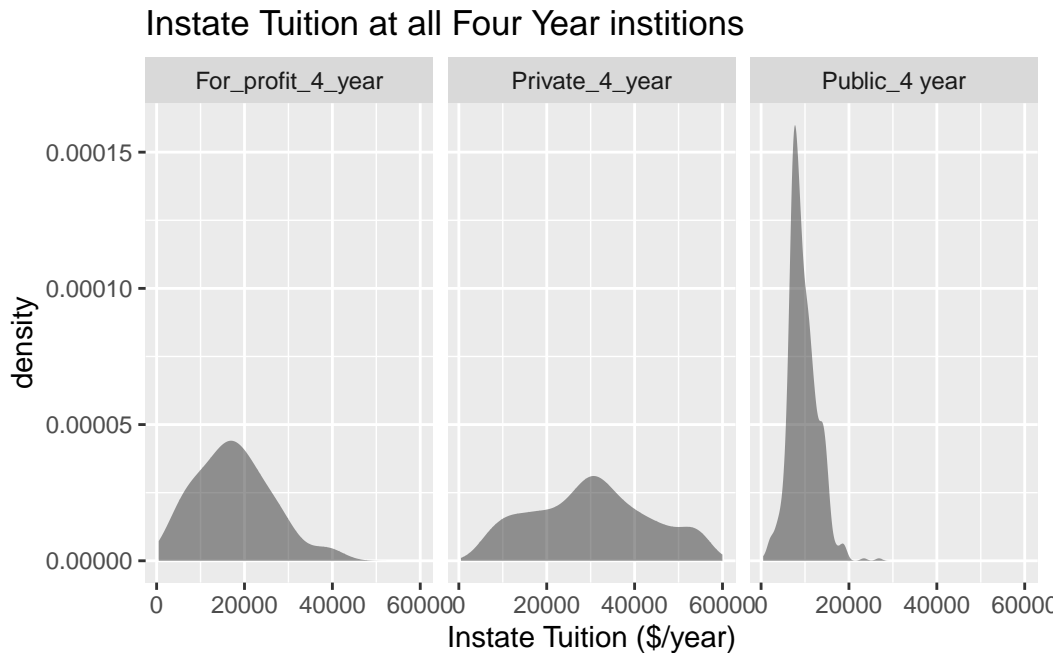


Figure 3.15: Instate Tuition at all Four Year institutions

Description of Figure ?? is a density plots separated by for profit 4 year with peak on left, private 4 year with peak in the middle, and public 4 year colleges with peak on the left. Public 4 year has the highest peak, with for profit 4 year is lower, and then private 4 year with the lowest peak.

The distribution is skewed right, with no gaps. Most institutions in state is less than \\$ 20,000 per year though some go as high as \\$ 60,000 per year. When separated by public versus private and for profit, most public are much less than \\$ 20,000 per year while private four year cost around \\$ 30,000 per year, and for profit are around \\$ 20,000 per year.

There are other types of graphs for quantitative data. They will be explored in the next section.

3.2.5 Homework for Quantitative Data Section

1. The weekly median incomes of males and females for specific occupations, are given in Table ?? (CPS News Releases. (n.d.). Retrieved July 8, 2019, from <https://www.bls.gov/cps/>). Create a density plot for males and females. Discuss any findings from the graph. Note: to put two graphs on the same axis, type the piping symbol `|>` (base r) or `%>%` (magrittr package) (Note: `|>` and `%>%` are piping symbols that can be thought of as “and then”) at the end of the first command and then type the command for the second graph on the next line. Also, use `fill=` “pick a color” in the

command to plot the graphs with different colors so the two graphs can be easier to distinguish.

```
Wages<- read.csv( "https://krkozak.github.io/MAT160/wages.csv")
knitr::kable(head(Wages))
```

Table 3.7: Head of Wages Data frame

Occupation	Numworkers	median_wage	male_wage	female_wage	male_wage	female_wage
Management, professional, and related occupations	48808	1246	23685	1468	25123	1078
Management, business, and financial operations occupations	19863	1355	10668	1537	9195	1168
Management occupations	13477	1429	7754	1585	5724	1236
Chief executives	1098	2291	790	2488	307	1736
General and operations managers	939	1338	656	1427	283	1139
Legislators	14	NA	10	NA	4	NA

Code book for Data Frame Wages

Description Median weekly earnings of full-time wage and salary workers by detailed occupation and sex. The Current Population Survey (CPS) is a monthly survey of households conducted by the Bureau of Census for the Bureau of Labor Statistics. It provides a comprehensive body of data on the labor force, employment, unemployment, persons not in the labor force, hours of work, earnings, and other demographic and labor force characteristics.

Format

This data frame contains the following columns:

Occupation: Occupations of workers.

Numworkers: The number of workers in each occupation (in thousands of workers)

median_wage: Median weekly wage (\\$)

male_worker: number of male workers (in thousands of workers)

male_wage: Median weekly wage of male workers (\\$)

female_worker: number of female workers (in thousands of workers)

female_wage: Median weekly wage of female workers (\\$)

Source CPS News Releases. (n.d.). Retrieved July 8, 2019, from <https://www.bls.gov/cps/>

References Current Population Survey (CPS) retrieved July 8, 2019.

2. The density of people per square kilometer for certain countries is in Table ?? (World Bank, 2019). Create density plot of density in 2018 for just Sub-Saharan Africa. Describe what story the graph tells.

```
Density<- read.csv( "https://krkozak.github.io/MAT160/density.csv")
knitr::kable(head(Density))
```

Table 3.8: Head of Density Data frame

[illegible]

Code book for Data Frame Density

Description Population density of all countries in the world

Format

Table 3.9: Head of Africa Data frame

Country	Region	Population	GDP	GDP_per_capita	Life expectancy	Infant mortality	Adult literacy	Urban population	Renewable energy	Public expenditure
Algeria	North Africa	34,129,817	215,456,539,299	6,315	75.2	23.6	67.1	71.3	0.8	1.3
Angola	Central Africa	24,663,832	253,245,705,713	10,267	53.7	120.4	33.4	36.9	0.3	0.3
Botswana	South Africa	2,354,916	21,508,780,023	9,135	54.7	43.8	64.4	59.3	0.2	0.2
Burkina Faso	West Africa	19,745,617	15,838,337,295	799	52.1	107.2	37.4	33.4	0.3	0.3
Burundi	East Africa	10,735,302	10,935,337,295	1,018	51.3	154.1	33.4	33.4	0.3	0.3
Cameroon	Central Africa	23,507,129	21,508,780,023	9,135	53.7	120.4	33.4	36.9	0.3	0.3
Cape Verde	West Africa	547,871	1,508,780,023	2,735	74.4	18.3	64.4	59.3	0.2	0.2
Cote d'Ivoire	West Africa	22,830,521	21,508,780,023	9,135	53.7	120.4	33.4	36.9	0.3	0.3
Egypt	North Africa	84,249,063	215,456,539,299	2,557	72.5	28.6	67.1	71.3	0.8	1.3
Ethiopia	East Africa	101,360,000	21,508,780,023	2,125	52.1	107.2	37.4	33.4	0.3	0.3
Ghana	West Africa	24,487,389	21,508,780,023	8,785	57.8	68.3	64.4	59.3	0.2	0.2
Guinea	West Africa	12,947,023	10,935,337,295	844	52.1	107.2	37.4	33.4	0.3	0.3
Guinea-Bissau	West Africa	1,912,847	10,935,337,295	572	52.1	107.2	37.4	33.4	0.3	0.3
Kenya	East Africa	44,299,834	21,508,780,023	4,855	53.7	120.4	33.4	36.9	0.3	0.3
Lesotho	South Africa	2,354,916	21,508,780,023	9,135	53.7	120.4	33.4	36.9	0.3	0.3
Liberia	West Africa	4,613,266	10,935,337,295	2,370	52.1	107.2	37.4	33.4	0.3	0.3
Madagascar	East Africa	22,830,521	10,935,337,295	478	52.1	107.2	37.4	33.4	0.3	0.3
Mali	West Africa	18,111,111	10,935,337,295	603	52.1	107.2	37.4	33.4	0.3	0.3
Mauritania	West Africa	3,427,746	10,935,337,295	319	52.1	107.2	37.4	33.4	0.3	0.3
Morocco	North Africa	33,240,000	215,456,539,299	6,481	74.4	28.6	67.1	71.3	0.8	1.3
Mozambique	East Africa	24,487,389	10,935,337,295	446	52.1	107.2	37.4	33.4	0.3	0.3
Niger	West Africa	19,745,617	10,935,337,295	554	52.1	107.2	37.4	33.4	0.3	0.3
Nigeria	West Africa	187,753,942	21,508,780,023	1,147	52.1	107.2	37.4	33.4	0.3	0.3
Rwanda	East Africa	11,533,857	10,935,337,295	948	52.1	107.2	37.4	33.4	0.3	0.3
Senegal	West Africa	15,427,023	21,508,780,023	1,394	52.1	107.2	37.4	33.4	0.3	0.3
Sierra Leone	West Africa	6,442,923	10,935,337,295	1,698	52.1	107.2	37.4	33.4	0.3	0.3
South Africa	South Africa	54,293,123	21,508,780,023	3,958	53.7	120.4	33.4	36.9	0.3	0.3
South Sudan	East Africa	11,533,857	10,935,337,295	948	52.1	107.2	37.4	33.4	0.3	0.3
Sudan	East Africa	43,849,063	10,935,337,295	249	52.1	107.2	37.4	33.4	0.3	0.3
Tanzania	East Africa	54,293,123	10,935,337,295	200	52.1	107.2	37.4	33.4	0.3	0.3
Togo	West Africa	7,854,746	10,935,337,295	1,389	52.1	107.2	37.4	33.4	0.3	0.3
Tunisia	North Africa	11,533,857	215,456,539,299	1,867	74.4	28.6	67.1	71.3	0.8	1.3
Zambia	East Africa	11,533,857	10,935,337,295	948	52.1	107.2	37.4	33.4	0.3	0.3
Zimbabwe	East Africa	11,533,857	10,935,337,295	948	52.1	107.2	37.4	33.4	0.3	0.3

- The Affordable Care Act created a market place for individuals to purchase health care plans. In 2014, the premiums for a 27 year old for the different levels health insurance are given in Table ?? (\“Health insurance marketplace,\” 2013). Create a density plot of bronze_lowest, then silver_lowest, and gold_lowest all on the same axes. Use `|> or %>%` at the end of each command. Describe the story the graphs tells.

```
Insurance<- read.csv( "https://krkozak.github.io/MAT160/insurance.csv")
knitr::kable(head(Insurance))
```

Table 3.10: Head of Insurance Data frame

state	average	bronze_lowest	silver_lowest	gold_lowest	platinum_lowest	bronze_posttax	silver_posttax	gold_posttax	platinum_posttax
AK	34	254	312	401	236	312	107	48	1131
AL	7	162	200	248	138	209	145	98	757
AR	28	181	231	263	135	241	145	85	873
AZ	106	141	164	187	107	166	145	120	600
DE	19	203	234	282	137	237	145	111	859
FL	102	169	200	229	132	218	145	96	789

Code book for Data Frame Insurance

Description The Affordable Care Act created a market place for individuals to purchase health care plans. The data is from 2014.

Format

This data frame contains the following columns:

state: state of insured.

average_QHP: The number of qualified health plans

bronze_lowest: premium for the lowest bronze level of insurance for a single person (\\$)

silver_lowest: premium for the lowest silver level of insurance for a single person (\\$)

gold_lowest: premium for the lowest gold level of insurance for a single person (\\$)

catastrophic: premium for the catastrophic level of insurance for a single person (\\$)

second_silver_pretax: premium for the second silver level of insurance for a single person pretax (\\$)

second_silver_posttax: premium for the second silver level of insurance for a single person posttax (\\$)

second_bronze_posttax: premium for the lowest bronze level of insurance for a single person posttax (\\$)

silver_family_pretax: premium for the silver level of insurance for a family pretax (\\$)

silver_family_posttax: premium for the silver level of insurance for a family posttax (\\$)

bronze_family_posttax: premium for the bronze level of insurance for a family posttax (\\$)

Source Health Insurance Market Place Retrieved from website: <http://aspe.hhs.gov/health/reports/2013/marketplace-premiums-for-2014>.

References Department of Health and Human Services, ASPE. (2013). Health insurance marketplace

4. Students in a statistics class took their first test. In Table ?? are the scores they earned. Create a density plot for grades. Describe the shape of the distribution.

```
Firsttest_1<- read.csv( "https://krkozak.github.io/MAT160/firsttest_1.csv")
knitr::kable(head(Firsttest_1))
```

Table 3.11: Head of First Test Data frame

grades
80
79
89

Table 3.11: Head of First Test Data frame

grades
74
73
67

5. Students in a statistics class took their first test. The scores they earned are in Table ??.
- Create a density plot for grades. Describe the shape of the distribution. Compare to the graph in question 4.

```
Firsttest_2<- read.csv( "https://krkozak.github.io/MAT160/firsttest_2.csv")
knitr::kable(head(Firsttest_2))
```

Table 3.12: Head of First Test Data frame

grades
67
67
76
47
85
70

3.3 Other Graphical Representations of Data

There are many other types of graphs. Some of the more common ones are the point plot (scatter plot), and a time-series plot. There are also many different graphs that have emerged lately for qualitative data. Many are found in publications and websites. The following is a description of the point plot (scatter plot), and the time-series plot.

3.3.1 Point Plots or Scatter Plot

Sometimes you have two different variables and you want to see if they are related in any way. A scatter plot helps you to see what the relationship would look like. A scatter plot is just a plotting of the ordered pairs.

3.3.2 Example: Scatter Plot

Is there a relationship between systolic blood pressure and weight? To answer this question some data is needed. The data frame NHANES contains this data, but given the size of the data frame, it may be not be very useful to look at the graph of all the data. It makes sense to take a sample from the data frame. A random sample is the better type of sample to take. Once the sample is taken, then a scatter plot can be created. The rStudio command for a scatter plot is

```
gf_point(response_variable ~ explanatory_variable, data= Data_Frame)
```

The sample is Table ??.

3.3.2.1 Solution

```
sample_NHANES <- NHANES |>
  sample_n(size = 100)
knitr::kable(head(sample_NHANES))
```

Table 3.13: Head of NHANES Sample Data

[illegible]

Preliminary: State the explanatory variable and the response variable

Let x =explanatory variable = Weight of a person (Weight)

y=response variable = Systolic blood pressure (BPSys1)

```
gf_point(BPSys1~Weight, data=sample_NHANES, xlab="Weight (kg)", ylab="Systolic Blood Pressure")
```

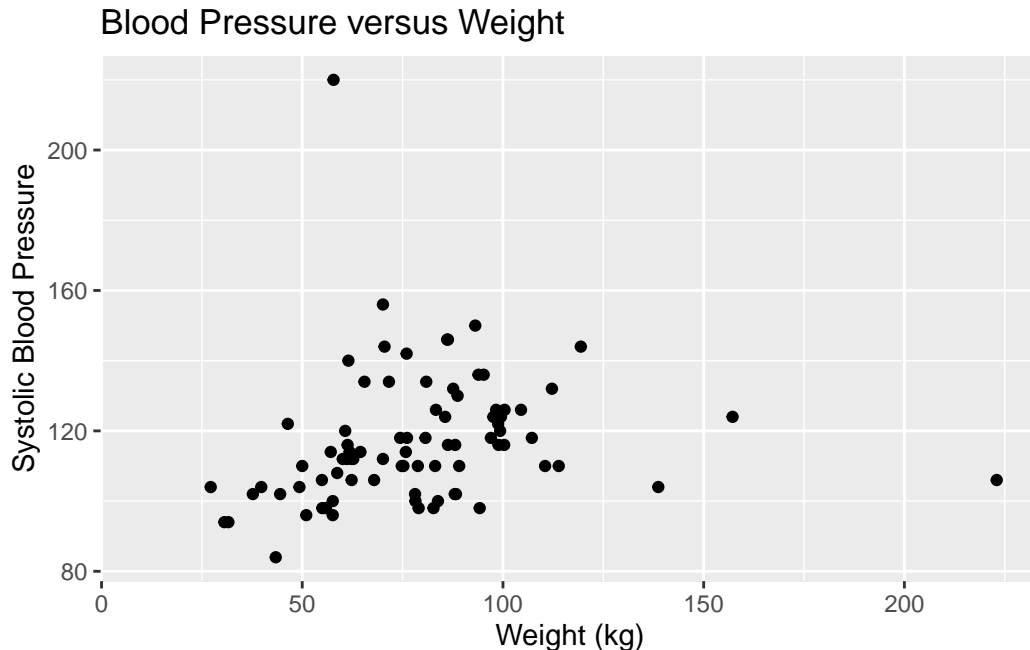


Figure 3.16: Blood Pressure versus Weight

Description of Figure ?? is a scatter plot with dots all over the plot though a line could be thought of fitting the dots with lower on the left and higher on the right.

Looking at the graph Figure ??, it appears that there is a linear relationship between weight and systolic blood pressure though it looks somewhat weak. It also appears to be a positive relationship, thus as weight increases, the systolic blood pressure increases.

3.3.3 Time-Series

A time-series plot is a graph showing the data measurements in chronological order, the data being quantitative data. For example, a time-series plot is used to show profits over the last 5 years. To create a time-series plot on RStudio, use the command

```
gf_line(response_variable ~ explanatory_variable, data=Data_Frame)
```

The purpose of a time-series graph is to look for trends over time. Caution, you must realize that the trend may not continue. Just because you see an increase, doesn't mean the

increase will continue forever. As an example, prior to 2007, many people noticed that housing prices were increasing. The belief at the time was that housing prices would continue to increase. However, the housing bubble burst in 2007, and many houses lost value, and haven't recovered.

3.3.4 Example: Time-Series Plot

The bank assets (in billions of Australia dollars (AUD)) of the Reserve Bank of Australia (RBA) and other financial organizations for the time period of September 1 1969, through March 1 2019, are contained in table Table ?? (Reserve Bank of Australia, 2019). Create a time-series plot of the total assets of Authorized Deposit-taking Institutions (ADIs) and interpret any findings.

```
Australian<- read.csv( "https://krkozak.github.io/MAT160/Australian_financial.csv")
knitr::kable(head(Australian))
```

Table 3.14: Head of Australian Data frame

Date	Day	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets	Assets
Sep69	9	2.7	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
Dec69	9	2.9	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
Mar70	1	80.0	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
Jun70	27	0.0	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
Sep70	3	0.0	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
Dec70	4	0.0	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA

Code book for Data frame Australian

Description The data is a range of economic and financial data produced by the Reserve Bank of Australia and other organizations.

Format

This data frame contains the following columns:

Date: quarters from September 1, 1969, to March 1, 2019

Day: The number of days since September 1, 1969, using 90 days between starts of a quarter. This column is to make it easier to graph in rStudio, and has no other purpose.

Assets_RBA: The assets for the Royal Bank of Australia

Assets_ADIs_Banks: The assets for Authorized Deposit-taking Institutions (ADIs), Banks

Assets_ADIs_Building: The assets for Authorized Deposit-taking Institutions (ADIs), Building societies

Assets_ADIs_CU: The assets for Authorized Deposit-taking Institutions (ADIs), Credit Unions

Assets_ADIs_Total: The assets for Authorized Deposit-taking Institutions (ADIs), total

Assets_RFCs_MM: The assets for Registered Financial Corporations (RFCs), Money Market Corporations

Assets_RFCs_Finance: The assets for Registered Financial Corporations (RFCs), Finance companies and general financiers

Assets_RFCs_Total: The assets for Registered Financial Corporations (RFCs) total

Assets_Life_offices: The Assets of Life offices and superannuation funds; Life insurance offices

Assets_Life_funds: The Assets of Life offices and superannuation funds; Superannuation funds

Assets_Life_Total: The Assets of Life offices and superannuation; Total

Assets_Other_Public_trusts: The Assets of Other managed funds; Public unit trusts

Assets_Other_Cash_trusts: The Assets of Other managed funds; Cash management trusts

Assets_Other_Common_funds: The Assets of Other managed funds; Common funds

Assets_Others_Friendly: The Assets of Other managed funds; Friendly societies

Assets_Other_General_insurance: The Assets of Other financial institutions; General insurance offices

Assets_Other_vehicles: The Assets Other financial institutions; Securitisation vehicles

Assets_Unconsolidated: The Assets of Unconsolidated; Statutory funds of life insurance offices; Superannuation

Source Reserve Bank of Australia. (2019, May 13). Statistical Tables. Retrieved July 10, 2019, from <https://www.rba.gov.au/statistics/tables/>

References Reserve Bank of Australia and other organizations

3.3.4.1 Solution

variable, x=total assets of Authorized Deposit-taking Institutions (ADIs)

Looking at the code book, one can see that the variable `Assets_ADIs_Total` is the variable in the data frame that is of interest here. With a time series plot, the other variable is time. In this case the variable in the data frame that represents time is `Date`. The problem with `Date` is that the units are every quarter. This is not easily interpreted by rStudio, so a column was created called `Day`. From the code book, this is the number of days since September 1, 1969, using 90 days between starts of a quarter. Even though this isn't perfect, it will work for determining trends. So create a time series plot of `Assets_ADIs_Total` versus `Day`. The command is:

```
gf_line(Assets_ADIs_Total~Day, data=Australian, title="Total Assets of Authorized Deposit-taking Institutions (ADIs)")
```

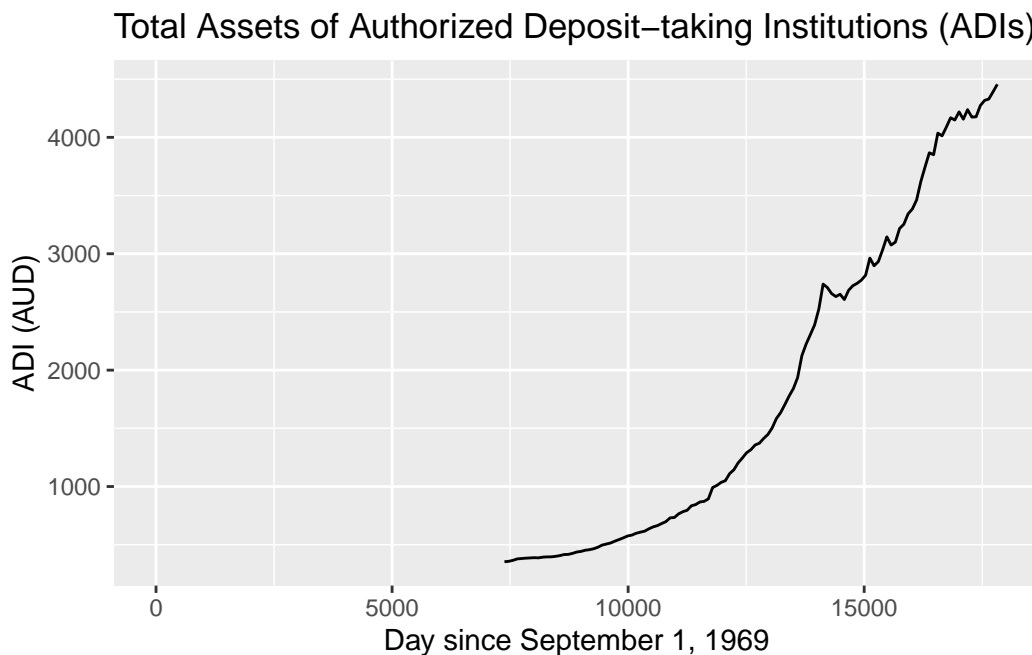


Figure 3.17: Total Assets of Authorized Deposit-taking Institutions

Description of Figure ?? is an increasing time series Graph of Total Assets of Authorized Deposit-taking Institutions from day 7500 to 17500. The first number starts at 0 and goes up to about 4500.

From the graph, total assets of Authorized Deposit-taking Institutions (ADIs) appear to be increasing with a slight dip around 14000 days since September 1, 1969. That would be around the year 2008 (14000 days /360 days per year + 1969).

Be careful when making a graph. If the vertical axis doesn't start at 0, then the change can look much more dramatic than it really is. For a graph to be useful to the reader, it needs to have a title that explains what the graph contains, the axes should be labeled so the reader knows what each axes represents, each axes should have a scale marked, and it is best if the vertical axis contains 0 to show the relationship.

3.3.5 Homework for Other Graphical Representations of Data Section

1. When an anthropologist finds skeletal remains, they need to figure out the height of the person. The height of a person (in cm) and the length of one of their metacarpal bone (in cm) were collected and are in Table ?? (Prediction of height, 2013). Create a scatter plot of length and height and state if there is a relationship between the height of a person and the length of their metacarpal.

```
Metacarpal<- read.csv( "https://krkozak.github.io/MAT160/metacarpal.csv")
knitr::kable(head(Metacarpal))
```

Table 3.15: Head of Metacarpal Data frame

length	height
45	171
51	178
39	157
41	163
48	172
49	183

Code book for Data frame Metacarpal

Description When anthropologists analyze human skeletal remains, an important piece of information is living stature. Since skeletons are commonly based on statistical methods that utilize measurements on small bones. The following data was presented in a paper in the American Journal of Physical Anthropology to validate one such method.

Format

This data frame contains the following columns:

length: length of Metacarpal I bone in mm

height: stature of skeleton in cm

Source Prediction of Height from Metacarpal Bone Length. (n.d.). Retrieved July 9, 2019, from <http://www.statsci.org/data/general/stature.html>

References Musgrave, J., and Harneja, N. (1978). The estimation of adult stature from metacarpal bone length. *Amer. J. Phys. Anthropology* 48, 113-120.

Devore, J., and Peck, R. (1986). *Statistics. The Exploration and Analysis of Data*. West Publishing, St Paul, Minnesota.

2. The value of the house and the amount of rental income in a year that the house brings in are in Table ?? (Capital and rental 2013). Create a scatter plot and state if there is a relationship between the value of the house and the annual rental income.

```
House<- read.csv( "https://krkozak.github.io/MAT160/house.csv")
knitr::kable(head(House))
```

Table 3.16: Head of House Data frame

capital	rental
61500	6656
67500	6864
75000	4992
75000	7280
76000	6656
77000	4576

Code book for Data frame House

Description The data show the capital value and annual rental value of domestic properties in Auckland in 1991.

Format

This data frame contains the following columns:

Capital: Selling price of house in Australian dollar (AUD)

rental: rental price of a house in Australian dollar (AUD)

Source Capital and rental values of Auckland properties. (2013, September 26). Retrieved from <http://www.statsci.org/data/oz/rentcap.html>

References Lee, A. (1994) *Data Analysis: An introduction based on R*. Auckland: Department of Statistics, University of Auckland. Data courtesy of Sage Consultants Ltd.

3. The World Bank collects information on the life expectancy of a person in each country (\“Life expectancy at,\” 2013) and the fertility rate per woman in the country (\“Fertility rate,\” 2013). The data for countries for the year 2011 are in Table ???. Create a scatter plot of the data and state if there appears to be a relationship between life expectancy and the number of births per woman in 2011.

```
Fertility<- read.csv( "https://krkozak.github.io/MAT160/fertility.csv")
knitr::kable(head(Fertility))
```

Table 3.17: Head of Fertility Data frame

country	lifexp_2011	fertilrate_2011	lifexp_2000	fertilrate_2000	lifexp_1990	fertilrate_1990
Macao SAR, China	79.91	1.03	77.62	0.94	75.28	1.69
Hong Kong SAR, China	83.42	1.20	80.88	1.04	77.38	1.27
Singapore	81.89	1.20	78.05	NA	76.03	1.87
Hungary	74.86	1.23	71.25	1.32	69.32	1.84
Korea, Rep.	80.87	1.24	75.86	1.47	71.29	1.59
Romania	74.51	1.25	71.16	1.31	69.74	1.84

Code book for Data frame Fertility

Description Data is from the World Bank on the life expectancy of countries and the fertility rates in those countries.

Format

This data frame contains the following columns:

Country: Countries in the World

lifexp_2011: Life expectancy of a person born in 2011

fertilrate_2011: Fertility rate in the country in 2011

lifexp_2000: Life expectancy of a person born in 2000

fertilrate_2000: Fertility rate in the country in 2000

lifexp_1990: Life expectancy of a person born in 1990

fertilrate_1990: Fertility rate in the country in 1990

Source Life expectancy at birth. (2013, October 14). Retrieved from <http://data.worldbank.org/indicator/SP.DY>

References Data from World Bank, Life expectancy at birth, total (years)

- The World Bank collected data on the percentage of gross domestic product (GDP) that a country spends on health expenditures (Current health expenditure (% of GDP), 2019), the fertility rate of the country (Fertility rate, total (births per woman), 2019), and the percentage of women receiving prenatal care (Pregnant women receiving prenatal care (%), 2019). The data for the countries where this information is available in Table ??.
- Create a scatter plot of the health expenditure and percentage of women receiving prenatal care in the year 2000, and state if there appears to be a relationship between percentage spent on health expenditure and the percentage of women receiving prenatal care.

```
Fert_prenatal<-read.csv( "https://krkozak.github.io/MAT160/fertility_prenatal.csv")
knitr::kable(head(Fert_prenatal))
```

Country	Health expenditure (% of GDP)	Fertility rate (births per woman)	Pregnant women receiving prenatal care (%)
Angola	7.47	52.45	39.27
Burkina Faso	6.16	46.36	46.56
Cameroon	6.06	76.86	97.07
Chad	9.77	17.27	37.47
Cote d'Ivoire	7.37	47.57	67.77
Egypt	7.77	87.87	98.98
Ghana	8.08	18.18	28.38
Guinea	4.48	58.48	58.68
Kenya	8.88	98.98	99.09
Madagascar	9.09	19.29	39.49
Mali	5.59	69.69	79.89
Morocco	3.39	29.29	25.49
Niger	6.76	76.76	86.86
Nigeria	7.77	87.87	97.97
Rwanda	6.66	66.66	66.66
Senegal	6.66	66.66	66.66
Sierra Leone	6.66	66.66	66.66
Tanzania	6.66	66.66	66.66
Togo	6.66	66.66	66.66
Tunisia	6.66	66.66	66.66
Zambia	6.66	66.66	66.66
Zimbabwe	6.66	66.66	66.66

Code book for Data frame Fert_prenatal

Description Data is from the World Bank on money spent on expenditure of countries and the percentage of women receiving prenatal care in those countries.

Format

This data frame contains the following columns:

Country.Name: Countries around the world

Country.Code: Three letter country code for countries around the world

Region: Location of a country around the world as classified by the World Bank

IncomeGroup: The income level of a country as classified by the World Bank

f1960-f2017: Fertility rate of a country from 1960-2017

p1986-p2018: Percentage of women receiving prenatal care in the country in 1986-2018

e200-2016: Expenditure amounts of the countries for medical care in 2000-2016 (% of GDP)

Source Fertility rate, total (births per woman). (n.d.). Retrieved July 8, 2019, from

<https://data.worldbank.org/indicator/SP.DYN.TFRT.IN> Pregnant women receiving prenatal care (%). (n.d.). Retrieved July 9, 2019, from <https://data.worldbank.org/indicator/SH.STA.ANVC.ZS>

Current health expenditure (% of GDP). (n.d.). Retrieved July 9, 2019, from <https://data.worldbank.org/indicator/SH.MVS.SRVS.CD>

References Data from World Bank, fertility rate, expenditure on health, and pregnant woman rate of prenatal care.

5. The Australian Institute of Criminology gathered data on the number of deaths (per 100,000 people) due to firearms during the period 1983 to 1997 (“Deaths from firearms,” 2013). The data is in Table ???. Create a time-series plot of the data and state any findings you can from the graph.

```
Firearm<- read.csv( "https://krkozak.github.io/MAT160/rate.csv")
knitr::kable(head(Firearm))
```

Table 3.19: Head of Firearm Data frame

year	rate
1983	4.31
1984	4.42
1985	4.52
1986	4.35
1987	4.39
1988	4.21

Code book for Data Frame Firearm

Description The data give the number of deaths caused by firearms in Australia from 1983 to 1997, expressed as a rate per 100,000 of population.

Format

This data frame contains the following columns:

Year: Years from 1983 to 1997

Rate: Rate of deaths caused by firearms in Australia per 100,000 population

Source Deaths from firearms. (2013, September 26). Retrieved from <http://www.statsci.org/data/oz/firearms.htm>

References Australian Institute of Criminology, 1999. The data was contributed by Rex Boggs, Glenmore State High School, Rockhampton, Queensland, Australia.

6. The economic crisis of 2008 affected many countries, though some more than others. Some people in Australia have claimed that Australia wasn't hurt that badly from the crisis. The bank assets (in billions of Australia dollars (AUD)) of the Reserve Bank of Australia (RBA) for the time period of September 1 1969, through March 1 2019, are contained in @bl-Australian (Reserve Bank of Australia, 2019). Create a time-series plot of the assets of the RBA and interpret any findings.

Code book for Data Frame Australian is below Table ??.

7. The consumer price index (CPI) is a measure used by the U.S. government to describe the cost of living. The cost of living for the U.S. from the years 1913 through 2019, with the year 1982 being used as the year that all others are compared (Consumer Price Index Data from 1913 to 2019, 2019) is given in Table ??. Create a time-series plot of the Average Annual CPI and interpret.

```
CPI<- read.csv( "https://krkozak.github.io/MAT160/CPI_US.csv")
knitr::kable(head(CPI))
```

Table 3.20: Head of CPI Data frame

Year	Jan	Feb	Mar	Apr	May	June	July	Aug	Sep	Oct	Nov	Dec	Annual_Avg	Dec_Delta	Dec_Avg
1913	9.8	9.8	9.8	9.8	9.7	9.8	9.9	9.9	10.0	10.0	10.1	10.0	9.9	—	—
1914	10.0	9.9	9.9	9.8	9.9	9.9	10.0	10.2	10.2	10.1	10.2	10.1	10.0	1	1
1915	10.1	10.0	9.9	10.0	10.1	10.1	10.1	10.1	10.1	10.2	10.3	10.3	10.1	2	1
1916	10.4	10.4	10.5	10.6	10.7	10.8	10.8	10.9	11.1	11.3	11.5	11.6	10.9	12.6	7.9
1917	11.7	12.0	12.0	12.6	12.8	13.0	12.8	13.0	13.3	13.5	13.5	13.7	12.8	18.1	17.4
1918	14.0	14.1	14.0	14.2	14.5	14.7	15.1	15.4	15.7	16.0	16.3	16.5	15.1	20.4	18

Code book for Data frame CPI

Description This table of Consumer Price Index (CPI) data is based upon a 1982 base of 100.

Format

This data frame contains the following columns:

Year: Year from 1913 to 2019

Jan, Feb, Mar, Apr, May, Jun, Jul, Aug, Sep, Oct, Nov, Dec: CPI for a particular month

Average_Avg: The average CPI for a particular year

PerDec_Dec: Percent change from December to December

Per_Avg_Avg: Percent change from Annual Average to Annual Average

Source Consumer Price Index Data from 1913 to 2019. (2019, June 12). Retrieved July 10, 2019, from <https://www.usinflationcalculator.com/inflation/consumer-price-index-and-annual-percent-changes-from-1913-to-2008/>

References US Inflation Calculator website, 2019.

8. The mean and median incomes income in current dollars is given in Table ???. Create a time-series plot and interpret.

```
US_income<- read.csv( "https://krkozak.github.io/MAT160/US_income.csv")
knitr::kable(head(US_income))
```

Table 3.21: Head of US_income Data frame

year	number	med_income_current	med_income_2017	mean_income_current	mean_income_2017
2017	127586	61372	61372	86220	86220
2016	126224	59039	60309	83143	84931
2015	125819	56516	58476	79263	82012
2014	124587	53657	55613	75738	78500
2013	122952	51939	54744	72641	76565
2012	122459	51017	54569	71274	76237

Code book for Data Frame US_income

Description This table is of US mean and median incomes in both current dollars and in 2017 dollars.

Format

This data frame contains the following columns:

Year: Year from 1975 to 2017

number: Households as of March of the following year. (in thousands)

med_income_current: median income of a US household in current dollars

med_income_2017: median income of a US household in 2017 CPI-U-RS adjusted dollars

mean_income_current: mean income of a US household in current dollars

mean_income_2017: mean income of a US household in 2017 CPI-U-RS adjusted dollars

Source US Census Bureau. (2018, March 06). Data. Retrieved July 21, 2019, from <https://www.census.gov/programs-surveys/cps/data-detail.html>

References U.S. Census Bureau, Current Population Survey, Annual Social and Economic Supplements.

4 Numerical Description of Data

Chapter 1 discussed what a population, sample, parameter, and statistic are, and how to take different types of samples. Chapter 2 discussed ways to graphically display data. There was also a discussion of important characteristics: center, variations, distribution, outliers, and changing characteristics of the data over time. Distributions and outliers can be answered using graphical means. Finding the center and variation can be done using numerical methods that will be discussed in this chapter. Both graphical and numerical methods are part of a branch of statistics known as **descriptive statistics**. Later descriptive statistics will be used to make decisions and/or estimate population parameters using methods that are part of the branch called **inferential statistics**.

4.1 Measures of Center

This section focuses on measures of central tendency. Many times you are asking what to expect on average. Such as when you pick a major, you would probably ask how much you expect to earn in that field. If you are thinking of relocating to a new town, you might ask how much you can expect to pay for housing. If you are planting vegetables in the spring, you might want to know how long it will be until you can harvest. These questions, and many more, can be answered by knowing the center of the data set. There are three measures of the “center” of the data. They are the mode, median, and mean. Any of the values can be referred to as the “average.”

The **mode** is the data value that occurs the most frequently in the data. To find it, you count how often each data value occurs, and then determine which data value occurs most often. The mode is not the most useful measure of center. This is because, a data set can have more than one mode. If there is a tie between two values for the most number of times then both values are the mode and the data is called bimodal (two modes). If every data point occurs the same number of times, there is no mode. If there are more than two numbers that appear the most times, then usually there is no mode.

The **median** is the data value in the middle of a sorted list of data. To find it, you put the data in order, and then determine which data value is in the middle of the data set.

The **mean** is the arithmetic average of the numbers. This is the center that most people call the average, though all three -- mean, median, and mode -- really are averages.

There are no symbols for the mode and the median, but the mean is used a great deal, and statisticians gave it a symbol. There are actually two symbols, one for the population parameter and one for the sample statistic. In most cases you cannot find the population parameter, so you use the sample statistic to estimate the population parameter.

4.1.1 Population Mean

$$\mu = \frac{\sum x}{N}, \text{ pronounced mu}$$

N is the size of the population.

x represents a data value.

$\sum x$ means to add up all of the data values.

4.1.2 Sample Mean:

$$\bar{x} = \frac{\sum x}{n}, \text{ pronounced x bar.}$$

n is the size of the sample.

x represents a data value.

$\sum x$ means to add up all of the data values.

The value for \bar{x} is used to estimate μ since μ can't be calculated in most situations.

4.1.3 Example: Finding the Mean and Median using r

Suppose a vet wants to find the average weight of cats. The weights (in kg) of cats are in Table ??.

```
knitr::kable(head(cats))
```

Table 4.1: Head of Cats

Sex	Bwt	Hwt
F	2.0	7.0
F	2.0	7.4
F	2.0	9.5
F	2.1	7.2
F	2.1	7.3
F	2.1	7.6

The head command shows the variable names and the first few unit of observation rows
Find the mean and median of the weight of a cat.

4.1.3.1 Solution

Before starting any mathematics problem, it is always a good idea to define the unknown in the problem. In statistics, you want to define the variable. The symbol for the variable is x .

The variable is x = weight of a cat

Mean: To find with r Studio, perform the command:

```
df_stats(~Bwt, cats, mean)
```

```
      response      mean
1      Bwt 2.723611
```

The mean weight is 2.72 kg

Median: To find with r Studio, perform the command:

```
df_stats(~Bwt, cats, median)
```

```
      response median
1      Bwt      2.7
```

The median weight is 2.7 kg also. It appears the average weight is 2.7 kg of all cats.

4.1.4 Example: Finding Mean and Median with filtering

Looking at the data frame for cats weights Table ?? you see that there are several variables You may want to know what the other variables are. A Code Book describes the data set, explains what the variables are including the units, and the source of the data frame. To review the code book for a data frame, complete the following command.

```
??cats
```

Then click on MASS::cats.

The output looks like:

Image 3.1.1: Code book for cats data frame

cats {MASS}

R Documentation

Anatomical Data from Domestic Cats

Description

The heart and body weights of samples of male and female cats used for *digitalis* experiments. The cats were all adult, over 2 kg body weight.

Usage

```
cats
```

Format

This data frame contains the following columns:

Sex

sex: Factor with evels "F" and "M".

Bwt

body weight in kg.

Hwt

heart weight in g.

Source

R. A. Fisher (1947) The analysis of covariance method for the relation between a part and the whole, *Biometrics* **3**, 65–68.

References

Venables, W. N. and Ripley, B. D. (2002) *Modern Applied Statistics with S*. Fourth edition. Springer.

Figure 4.1: Code book for cats data frame

Suppose you want to know if male cats weigh more than female cats. Looking at the variables, you notice that there is a variable for the sex of the cat. You can look at the weights of males and females separately. This looks like:

4.1.4.1 Solution

To find the mean and median, separated by sex, use this command in R Studio:

```
df_stats(Bwt~Sex, data=cats, mean, median)
```

	response	Sex	mean	median
1	Bwt	F	2.359574	2.3
2	Bwt	M	2.900000	2.9

Notice that the female cats' mean weigh 2.4 kg and the male cats' mean weigh 2.9 kg. The median weight of female cats is 2.3 kg and for males is 2.9 kg. So it does appear that males cats weight a bit more than the female cats.

There are many different summary statistics that can be found. An example is the minimum and maximum value. In this example, you will see how to find the min and max values and then filter them out of a data set to see what effect they have on the mean and median.

4.1.5 Example: Affect of Extreme Values on Mean and Median

Find the minimum and maximum values of cats weights.

4.1.5.1 Solution

The command in RStudio for finding the minimum and maximum is very similar to how to find the mean and median. In fact all summary statistics start with

```
df_stats(~variable, data=Data Frame, desired statistics)
```

Here is the command in RStudio for the minimum and maximum of cat's body weight.

```
df_stats(~Bwt, data=cats, min, max)
```

	response	min	max
1	Bwt	2	3.9

The minimum weight of a cat in this data frame is 2 kg and the maximum weight of a cat is 3.9 kg.

Now create two new data sets. One data set will exclude the maximum value. You can call it anything you want, but it would make sense to call it something like `nomax`. The command to create the new data set is:

```
nomax <- filter(cats, Bwt<3.9)
```

Then create a data set that excludes the minimum value; call it `nomin`:

```
nomin<-filter(cats, Bwt>2)
```

The `<-` is the way to indicate to `r` what the data set `nomin` is equivalent to what follows the symbol. Notice that it doesn't look like anything happened, but new data sets were created in the background. Now you can find the mean and median of each new data set:

```
df_stats(~Bwt, data=nomax, mean, median)
```

	response	mean	median
1	Bwt	2.707042	2.7

The mean without the maximum value is 2.70 kg, and the median is 2.7 kg.

```
df_stats(~Bwt, data=nomin, mean, median)
```

	response	mean	median
1	Bwt	2.74964	2.7

The mean without the minimum value is 2.75 kg, and the median is 2.7 kg.

From Example: Affect of Extreme Values on Mean and Median, the mean of the data set with all the values is 2.72 kg where the median is 2.7 kg. Notice that when the maximum value was excluded from the data set, the mean decreased a little but the median didn't change, and when the minimum value was excluded from the data set, the mean increased a little but the median didn't change. The mean is much higher than the median. Why is this? This is because the mean is affected by extreme values, while the median is not. We say the median is a much more resistant measure of center because it isn't affected by extreme values as much.

An outlier is a data value that is very different from the rest of the data. It can be really high or really low. Extreme values may be an outlier if the extreme value is far enough from the

center. If there are extreme values in the data, the median is a better measure of the center than the mean. If there are no extreme values, the mean and the median will be similar so most people use the mean. The mean is not a resistant measure because it is affected by extreme values. The median is a resistant measure because it is not affected by extreme values.

As a consumer you need to be aware that people choose the measure of center that best supports their claim. When you read an article in the newspaper and it talks about the “average” it usually means the mean but sometimes it refers to the median. Some articles will use the word “median” instead of “average” to be more specific. If you need to make an important decision and the information says “average”, it would be wise to ask if the “average” is the mean or the median before you decide.

As an example, suppose that a company wants to use the mean salary as the average salary for the company. This is because the high salaries of the administrators will pull the mean higher. The company can say that the employees are paid well because the average is high. However, the employees want to use the median since it discounts the extreme values of the administration and will give a lower value of the average. This will make the salaries seem lower and that a raise is in order.

Why use the mean instead of the median? The reason is because when multiple samples are taken from the same population, the sample means tend to be more consistent than other measures of the center.

To understand how the different measures of center related to skewed or symmetric distributions, see Graph \#3.1.1. As you can see sometimes the mean is smaller than the median, sometimes the mean is larger than the median, and sometimes they are the same values.

Graph \#3.1.1: Mean, Median, Mode as Related to a Distribution

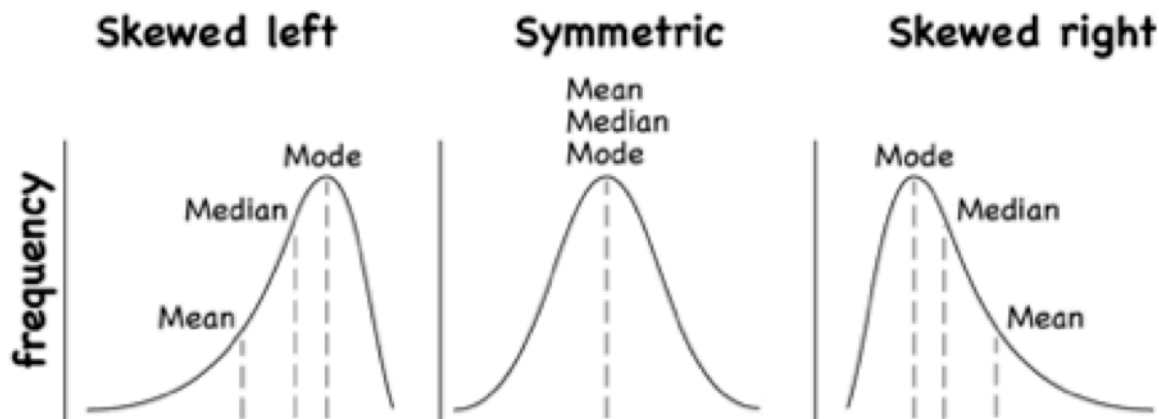


Figure 4.2: Mean, median, mode as related to distribution

One last type of average is a weighted average. **Weighted averages** are used quite often in different situations. Some teachers use them in calculating a student’s grade in the course,

or a grade on a project. Some employers use them in employee evaluations. The idea is that some activities are more important than others. As an example, a full time teacher at a community college may be evaluated on their service to the college, their service to the community, whether their paperwork is turned in on time, and their teaching. However, teaching is much more important than whether their paperwork is turned in on time. When the evaluation is completed, more weight needs to be given to the teaching and less to the paperwork. This is a weighted average.

4.1.6 Weighted Average

$$\text{weighted average} = \frac{\sum x*w}{\sum w}$$

where **w** is the weight of the data value, **x**.

4.1.7 Example: Weighted Average

In your biology class, your final grade is based on several things: a lab score, scores on two major tests, and your score on the final exam. There are 100 points available for each score. The lab score is worth 15% of the course, the two exams are worth 25% of the course each, and the final exam is worth 35% of the course. Suppose you earned scores of 95 on the labs, 83 and 76 on the two exams, and 84 on the final exam. Compute your weighted average for the course.

4.1.7.1 Solution

Variable: **x** = score

A weighted average can be found using technology. The commands for finding the weighted mean using RStudio is as follows:

```
x<-c(type in the scores with commas in between)
```

```
w<-c(type in the weights as decimals with commas in between)
```

```
weighted.mean(x,w)
```

The **x** and **w** represent the variables, <- means make the variables equivalent to what follows, the c(means combine all the values in the () as one combined variable.

For this example, the commands would be

```
x<-c(95, 83, 76, 84)
w<-c(.15, .25, .25, .35)
weighted.mean(x,w)
```


[1] 83.4

Your weighted mean in the biology class is 83.4%. Using the traditional grading scale, you have a B in the class.

4.1.8 Example: Weighted Average

The faculty evaluation process at John Jingle University rates a faculty member on the following activities: teaching, publishing, committee service, community service, and submitting paperwork in a timely manner. The process involves reviewing student evaluations, peer evaluations, and supervisor evaluation for each teacher and awarding him/her a score on a scale from 1 to 10 (with 10 being the best). The weights for each activity are 20 for teaching, 18 for publishing, 6 for committee service, 4 for community service, and 2 for paperwork.

- a) One faculty member had the following ratings: 8 for teaching, 9 for publishing, 2 for committee work, 1 for community service, and 8 for paperwork. Compute the weighted average of the evaluation.
- b) Another faculty member had ratings of 6 for teaching, 8 for publishing, 9 for committee work, 10 for community service, and 10 for paperwork. Compute the weighted average of the evaluation.
- c) Which faculty member had the higher average evaluation?

4.1.8.1 Solution

- a) One faculty member had the following ratings: 8 for teaching, 9 for publishing, 2 for committee work, 1 for community service, and 8 for paperwork. Compute the weighted average of the evaluation.

Variable: \mathbf{x} = rating, \mathbf{w} = weight

```
x<-c(8, 9, 2, 1, 8)
w<-c(20, 18, 6, 4, 2)
weighted.mean(x,w)
```

[1] 7.08

The weighted average is 7.08.

- b) Another faculty member had ratings of 6 for teaching, 8 for publishing, 9 for committee work, 10 for community service, and 10 for paperwork. Compute the weighted average of the evaluation.

```
x<-c(6, 8, 9, 10, 10)
w<-c(20, 18, 6, 4, 2)
weighted.mean(x,w)
```

```
[1] 7.56
```

The weighted average for this employee is 7.56.

c) Which faculty member had the higher average evaluation?

The second faculty member has a higher average evaluation.

The last thing to mention is which average is used on which type of data.

Mode can be found on nominal, ordinal, interval, and ratio data, since the mode is just the data value that occurs most often. You are just counting the data values.

Median can be found on ordinal, interval, and ratio data, since you need to put the data in order. As long as there is order to the data you can find the median.

Mean can be found on interval and ratio data, since you must have numbers to add together.

4.1.9 Homework for Measures of Center Section

Use Technology on all problems. State the variable on all problems.

1. Cholesterol levels were collected from patients certain days after they had a heart attack and are in Table ???. Find the mean and median for cholesterol levels 2 days after the heart attack.

```
Cholesterol<-read.csv( "https://krkozak.github.io/MAT160/cholesterol.csv")
knitr::kable(head(Cholesterol))
```

Table 4.2: Head of Cholesterol Levels of Patients After Heart Attack

patient	day2	day4	day14
1	270	218	156
2	236	234	NA
3	210	214	242
4	142	116	NA
5	280	200	NA
6	272	276	256

Code book for Data Frame Cholesterol

Description Cholesterol levels were collected from patients certain days after they had a heart attack

This data frame contains the following columns:

Patient: Patient number

day2: Cholesterol level of patient 2 days after heart attack. (mg/dL)

day4: Cholesterol level of patient 4 days after heart attack. (mg/dL)

day14: Cholesterol level of patient 14 days after heart attack. (mg/dL)

Source Ryan, B. F., Joiner, B. L., & Ryan, Jr, T. A. (1985). Cholesterol levels after heart attack. Retrieved from <http://www.statsci.org/data/general/cholest.html>

References Ryan, Joiner & Ryan, Jr, 1985

2. The lengths (in kilometers) of rivers on the South Island of New Zealand and what body of water they flow into are listed in Table ?? (Lee, 1994). Find the mean and median length of rivers that flow into the Pacific Ocean and the mean and median length of rivers that flow into the Tasman Sea.

```
Length<-read.csv( "https://krkozak.github.io/MAT160/length.csv")
knitr::kable(head(Length))
```

Table 4.3: Head of Length of New zealand rivers (km)

river	length	flowsto
Clarence	209	Pacific
Conway	48	Pacific
Waiau	169	Pacific
Hurunui	138	Pacific
Waipara	64	Pacific
Ashley	97	Pacific

Code book for data frame Length

Description Rivers in New Zealand, the lengths of river and what body of water the river flows into

This data frame contains the following columns:

River: Name of the river

length: how long the river is in kilometers

flowsto: what body of water the river flows into Pacific Ocean is Pacific and the Tasman Sea is Tasman

Source Lee, A. (1994). Data analysis: An introduction based on r. Auckland. Retrieved from <http://www.statsci.org/data/oz/nzrivers.html>

References Lee, A. (1994). Data analysis: An introduction based on r. Auckland.

3. Print-O-Matic printing company's employees have salaries that are contained in Table ??.

```
Pay<-read.csv( "https://krkozak.github.io/MAT160/pay.csv")
knitr::kable(head(Pay))
```

Table 4.4: Head of Salaries of Print-O-Matic Printing Company Employees

employee	salary
CEO	272500
Driver	58456
CD74	100702
CD65	57380
Embellisher	73877
Folder	65270

Code book for data frame Pay

Description Salaries of Print-O-Matic printing company's employees

This data frame contains the following columns:

employee: employees position in the company

salary: salary of that employee (Australian dollars (AUD))

Source John Matic provided the data from a company he worked with. The company's name is fictitious, but the data is from an actual company.

References John Matic (2013)

- a. Find the mean and median.
- b. Find the mean and median with the CEO's salary removed.
- c. What happened to the mean and median when the CEO's salary was removed? Why?

- d. If you were the CEO, who is answering concerns from the union that employees are underpaid, which average (mean or median) using the complete data set of the complete data set would you prefer? Why?
- e. If you were a platen worker, who believes that the employees need a raise, which average (mean or median) using the complete data set would you prefer? Why?
4. Print-O-Matic printing company spends specific amounts on fixed costs every month. The costs of those fixed costs are in a Table ??.

```
Cost<-read.csv( "https://krkozak.github.io/MAT160/cost.csv")
knitr::kable(head(Cost))
```

Table 4.5: Fixed Costs for Print-O-Matic Printing Company

	charges	cost
Bank charges		482
Cleaning		2208
Computer expensive		2471
Lease payments		2656
Postage		2117
Uniforms		2600

Code book for data frame Cost

Description fixed monthly charges for Print-0-Matic printing company

This data frame contains the following columns:

charges: Categories of monthly fixed charges

cost: fixed month costs (AUD)

Source John Matic provided the data from a company he worked with. The company's name is fictitious, but the data is from an actual company.

References John Matic (2013)

- a. Find the mean and median.
- b. Find the mean and median with the bank charges removed.
- c. What happened to the mean and median when the bank charges was removed? Why?
- d. If it is your job to oversee the fixed costs, which average (mean or median) using the complete data set would you prefer to use when submitting a report to administration to show that costs are low? Why?

- e. If it is your job to find places in the budget to reduce costs, which average (mean or median) using the complete data set would you prefer to use when submitting a report to administration to show that fixed costs need to be reduced? Why?
5. Looking at graph 3.1.2, state if the graph is skewed left, skewed right, or symmetric and then state which is larger, the mean or the median?

Graph 3.1.2: Skewed or Symmetric Graph

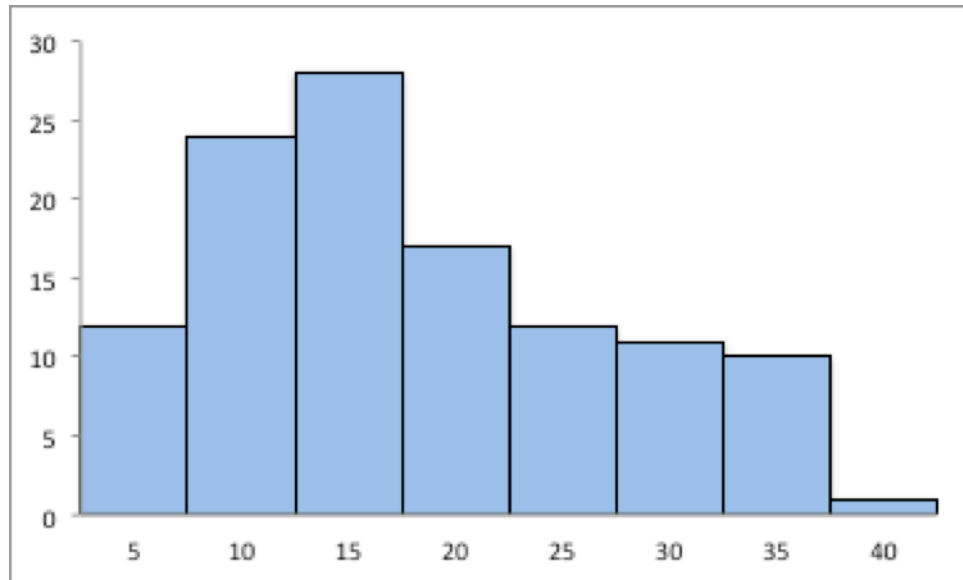


Figure 4.3: Graph #3.1.2

6. Looking at graph 3.1.3, state if the graph is skewed left, skewed right, or symmetric and then state which is larger, the mean or the median?

Graph 3.1.3: Skewed or Symmetric Graph

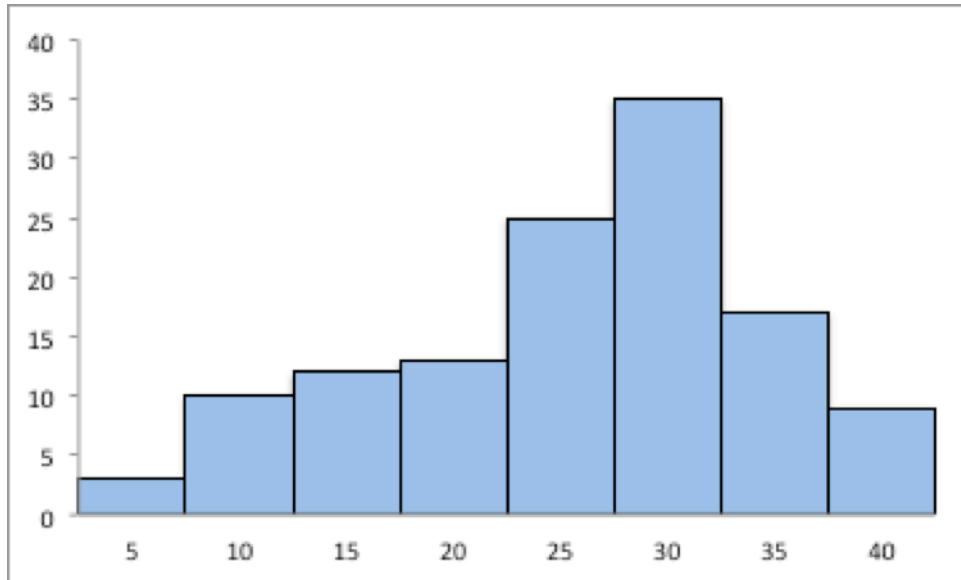


Figure 4.4: Graph 3.1.3

7. An employee at Coconino Community College (CCC) is evaluated based on goal setting and accomplishments toward the goals, job effectiveness, competencies, and CCC core values. Suppose for a specific employee, goal 1 has a weight of 30%, goal 2 has a weight of 20%, job effectiveness has a weight of 25%, competency 1 has a weight of 4%, competency 2 has a weight of 3%, competency 3 has a weight of 3%, competency 4 has a weight of 3%, competency 5 has a weight of 2%, and core values has a weight of 10%. Suppose the employee has scores of 3.0 for goal 1, 3.0 for goal 2, 2.0 for job effectiveness, 3.0 for competency 1, 2.0 for competency 2, 2.0 for competency 3, 3.0 for competency 4, 4.0 for competency 5, and 3.0 for core values. Find the weighted average score for this employee. If an employee has a score less than 2.5, they must have a Performance Enhancement Plan written. Does this employee need a plan?
8. An employee at Coconino Community College (CCC) is evaluated based on goal setting and accomplishments toward goals, job effectiveness, competencies, CCC core values. Suppose for a specific employee, goal 1 has a weight of 20%, goal 2 has a weight of 20%, goal 3 has a weight of 10%, job effectiveness has a weight of 25%, competency 1 has a weight of 4%, competency 2 has a weight of 3%, competency 3 has a weight of 3%, competency 4 has a weight of 5%, and core values has a weight of 10%. Suppose the employee has scores of 2.0 for goal 1, 2.0 for goal 2, 3.0 for goal 3, 2.0 for job effectiveness, 2.0 for competency 1, 3.0 for competency 2, 2.0 for competency 3, 3.0 for competency 4, and 4.0 for core values. Find the weighted average score for this employee. If an employee that has a score less than 2.5, they must have a Performance Enhancement Plan written. Does this employee need a plan?
9. A statistics class has the following activities and weights for determining a grade in the

course: test 1 worth 15% of the grade, test 2 worth 15% of the grade, test 3 worth 15% of the grade, homework worth 10% of the grade, semester project worth 20% of the grade, and the final exam worth 25% of the grade. If a student receives an 85 on test 1, a 76 on test 2, an 83 on test 3, a 74 on the homework, a 65 on the project, and a 79 on the final, what grade did the student earn in the course?

10. A statistics class has the following activities and weights for determining a grade in the course: test 1 worth 15% of the grade, test 2 worth 15% of the grade, test 3 worth 15% of the grade, homework worth 10% of the grade, semester project worth 20% of the grade, and the final exam worth 25% of the grade. If a student receives a 92 on test 1, an 85 on test 2, a 95 on test 3, a 92 on the homework, a 55 on the project, and an 83 on the final, what grade did the student earn in the course?

4.2 Measures of Spread

Variability is an important idea in statistics. If you were to measure the height of everyone in your classroom, every observation gives you a different value. That means not every student has the same height. Thus there is variability in people's heights. If you were to take a sample of the income level of people in a town, every sample gives you different information. There is variability between samples too. Variability describes how the data are spread out. If the data are very close to each other, then there is low variability. If the data are very spread out, then there is high variability. How do you measure variability? It would be good to have a number that measures it. This section will describe some of the different measures of variability, also known as variation.

In Example: Finding the Mean and Median using r , the average weight of a cat was calculated to be 2.72 kg. How much does this tell you about the weight of all cats? Can you tell if most of the weights were close to 2.72 kg or were the weights really spread out? The highest weight and the lowest weight are known, but is there more that you can tell? All you know is that the center of the weights is 2.72 kg.

You need more information.

The **range** of a set of data is the difference between the highest and the lowest data values (or maximum and minimum values). The **interval** is the lowest and highest values. The range is one value while the interval is two.

4.2.1 Example: Range

From Example: Affect of Extreme Values on Mean and Median, the maximum is 3.9 kg and the minimum is 2 kg. So the range is $3.9 - 2 = 1.9\text{kg}$. But what does that tell you? You don't know if the weights are really spread out, or if they are close together.

Unfortunately, range doesn't really provide a very accurate picture of the variability. A better way to describe how the data is spread out is needed. Instead of looking at the distance the highest value is from the lowest how about looking at the distance each value is from the mean. This distance is called the **deviation**. You might want to find the average of the deviation. Though the calculation for finding the average deviation is not very straight forward, you end up with a value called the **variance**. The symbol for the population variance is σ^2 , and it is the average squared distance from the mean. Statisticians like the variance, but many other people who work with statistics use a descriptive statistics which is the square root of the variance. This gives you the average distance from the mean. This is called the standard deviation, and is denoted with the letter σ .

The standard deviation is the average (mean) distance from a data point to the mean. It can be thought of as how much a typical data point differs from the mean.

The **sample variance** formula: $s^2 = \frac{\sum(x-\bar{x})^2}{n-1}$, where \bar{x} is the sample mean, n is the sample size, and \sum means to find the sum of the values. The $n-1$ on the bottom has to do with a concept called degrees of freedom. Basically, it makes the sample variance a better approximation of the population variance.

The **sample standard deviation** formula: $s = \sqrt{\frac{\sum(x-\bar{x})^2}{n-1}}$.

The **population variance** formula: $\sigma^2 = \frac{\sum(x-\mu)^2}{N}$, where σ is the Greek letter sigma and σ^2 represents the population variance, μ is the population mean, and N is the size of the population.

The **population standard deviation** formula: $\sigma = \sqrt{\frac{\sum(x-\mu)^2}{N}}$

Both the sample variance and sample standard deviation can be found using technology. If using rStudio, you would use

```
df_stats(~variable, data=data_frame, var, sd)
```

The next example will demonstrate this command.

4.2.2 Example: Finding the Standard Deviation

For the data frame Cats Table ?? find the variance and standard derivation for weight of cats. Then find the variance and standard deviation separated by sex of the cat.

4.2.2.1 Solution

The variance and standard deviation for all cats is found by performing the command:

```
df_stats(~Bwt, data=cats, var, sd)
```

	response		var	sd
1	Bwt		0.2355225	0.4853066

The variance for all cats is 0.24 kg^2 and the standard deviation is 0.49 kg.

To find out the mean, variance, and standard deviation for each sex of the cats, use the command:

```
df_stats(Bwt~Sex, data=cats, mean, var, sd)
```

	response	Sex	mean	var	sd
1	Bwt	F	2.359574	0.07506938	0.2739879
2	Bwt	M	2.900000	0.21854167	0.4674844

You can see that the mean weight of females cats is 2.36 kg, the variance is 0.075 kg^2 , and the standard deviation is 0.27 kg. For males cats, the mean is 2.9 kg, the variance is 0.22 kg^2 , and the standard deviation is 0.47 kg. This means that female cats weigh less than males and since the variance and standard deviations are much less for female cats than males cats, female cats' weights are more consistent than male cats.

In general a “small” variance and standard deviation means the data is close together (more consistent) and a “large” variance and standard deviation means the data is spread out (less consistent). Sometimes you want consistent data and sometimes you don't. As an example if you are making bolts, you want the lengths to be very consistent so you want a small standard deviation. If you are administering a test to see who can be a pilot, you want a large standard deviation so you can tell who are the good pilots and who are the not so good pilots.

What do “small” and “large” standard deviation mean? To a bicyclist whose average speed is 20 mph, $s = 20\text{mph}$ is huge. To an airplane whose average speed is 500 mph, $s = 20\text{mph}$ is nothing. The “size” of the variation depends on the size of the numbers in the problem and the mean. Another situation where you can determine whether a standard deviation is small or large is when you are comparing two different samples such as in Example: Finding the Standard Deviation. A sample with a smaller standard deviation is more consistent than a sample with a larger standard deviation.

Many other books and authors stress that there is a computational formula for calculating the standard deviation. However, this formula doesn't give you an idea of what standard deviation

is and what you are doing. It is only good for doing the calculations quickly. It goes back to the days when standard deviations were calculated by hand, and the person needed a quick way to calculate the standard deviation. It is an archaic formula that this author is trying to eradicate. It is not necessary anymore, computers will do the calculations for you with as much meaning as this formula gives. It is suggested that you never use it. If you want to understand what the standard deviation is doing, then you should use the definition formula. If you want an answer quickly, use a computer.

4.2.3 Use of Standard Deviation

One of the uses of the standard deviation is to describe how a population is distributed. This describes where much of the data is for most distributions. A general rule is that about 95% of the data is within 2 standard deviations of the mean. This is not perfect, but it works for many distributions. There are rules like the empirical rule and Chebyshev's theorem that give you more detailed percentages, but 95% in 2 standard deviations is a very good approximation.

4.2.4 Example: the general rule

The U.S. Weather Service has provided the information in Table ?? about the total monthly/annual number of reported tornadoes in Oklahoma for the years 1950 to 2018. (US Department of Commerce & NOAA, 2016)

```
Tornado<-read.csv("https://krkozak.github.io/MAT160/Tornado_OK.csv")
knitr::kable(head(Tornado))
```

Table 4.6: Monthly/Annual Number of tornadoes in Oklahoma

Year	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec	Annual
1950	0	1	1	5	12	1	0	0	2	1	0	0	23
1951	0	2	0	11	11	11	4	2	1	1	0	0	43
1952	0	0	0	7	5	5	4	1	0	0	0	0	22
1953	0	4	7	9	8	13	4	2	0	0	5	2	54
1954	0	0	7	13	19	4	4	2	3	1	0	0	53
1955	1	1	0	15	32	22	4	2	0	0	0	0	77

Code book for data frame Tornado

Description The U.S. Weather Service has collected data on the monthly and annual number of tornadoes in Oklahoma.

This data frame contains the following columns:

Year: Year from 1950-2018

Jan, Feb, Mar, Apr, May, Jun, Jul, Aug, Sep, Oct, Nov, Dec: Tornado numbers in each month of the year

Annual: Total number of tornadoes for each year

Source US Department of Commerce, & NOAA. (2016, November 15). 1950 Oklahoma Tornadoes. Retrieved from <https://www.weather.gov/oun/tornadodata-ok-1950>

References The data was supplied by The U.S. Weather Service

Find the general interval that contains about 95% of the data.

4.2.4.1 Solution

Variable: x = number of annual tornadoes in Oklahoma

Find the mean and standard deviation:

```
df_stats(~Annual, data=Tornado, mean, sd)
```

	response	mean	sd
1	Annual	56.02899	27.56061

The mean is $\mu = 56$ tornadoes and the standard deviation is $\sigma = 27.6$ tornadoes. The interval will be $\mu \pm 2 * \sigma = 56 \pm 2 * 27.6 = (0.8, 111.2)$

About 95% of the years have between 0.8 or 1 and 111 tornadoes in Oklahoma.

The general rule says that about 95% of the data is within two standard deviations of the mean. That percentage is fairly high. There isn't much data outside two standard deviations. A rule that can be followed is that if a data value is within two standard deviations, then that value is a common data value. If the data value is outside two standard deviations of the mean, either above or below, then the number is uncommon. It could even be called unusual. An easy calculation that you can do to figure it out is to find the difference between the data point and the mean, and then divide that answer by the standard deviation. As a formula this would be

$$z = \frac{x - \mu}{\sigma}$$

If you don't know the population mean, μ , and the population standard deviation, σ , then use the sample mean, \bar{x} , and the sample standard deviation, s , to estimate the population parameter values. Realize that using the sample standard deviation may not actually be very accurate.

4.2.5 Example: Determining If a Value Is Unusual

- a. In 1974, there were 45 tornadoes in Oklahoma. Is this value unusual? Why or why not?
- b. In 1999, there were 145 tornadoes in the Oklahoma. Is this value unusual? Why or why not?

4.2.5.1 Solution

- a. In 1974, there were 45 tornadoes in Oklahoma. Is this value unusual? Why or why not?

Variable: x = number of tornadoes in Oklahoma

To answer this question, first find how many standard deviations 45 is from the mean. From 3.2.4 example, we know $\mu = 56$ and $\sigma = 27.6$. For $x=45$, $z = \frac{45-56}{27.6} = -0.399$

Since this value is between -2 and 2, then it is not unusual to have 45 tornadoes in a year in Oklahoma. The z value is negative, so that means that 45 is less than the mean number of tornadoes.

- b. In 1999, there were 145 tornadoes in the Oklahoma. Is this value unusual? Why or why not?

Variable: x = number of tornadoes in Oklahoma

For this question the $x = 145$, $z = \frac{145-56}{27.6} = 3.22$

Since this value is more than 2, then it is unusual to have only 145 tornadoes in a year in Oklahoma.

4.2.6 Homework for Measures of Spread Section

Use Technology on all problems. State the variable on all problems.

1. Cholesterol levels were collected from patients certain days after they had a heart attack and are in Table ???. Find the mean, median, range, variance, and standard deviation for cholesterol levels 2 days after the heart attack.

Code book for Data Frame Cholesterol is below Table ??.

2. The lengths (in kilometers) of rivers on the South Island of New Zealand and what body of water they flow into are listed in Table ?? (Lee, 1994). Find the mean, median, range, variance, and standard deviation of the length of rivers that flow into the Pacific Ocean and the mean, median, range, variance, and standard deviation of the length of rivers that flow into the Tasman Sea. Compare and contrast the length of rivers that flow to

the Pacific Ocean versus the ones that flow into the Tasman Sea using both measures of spread and measures of variability.

Code book for data frame Length is below Table ??.

3. Print-O-Matic printing company's employees have salaries that are contained in Table ??. Find the mean, median, range, variance, and standard deviation for the salaries of all employees.

Code book for data frame Pay below Table ??.

4. Print-O-Matic printing company spends specific amounts on fixed costs every month. The costs of those fixed costs are in Table ??. Find the mean, median, range, variance, and standard deviation for the fixed costs.

Code book for Data frame Cost is below Table ??.

5. The data frame Pulse Table ?? contains various variables about a person including their pulse rates before the subject exercised and after the subject ran in place for one minute.

```
Pulse<-read.csv("https://krkozak.github.io/MAT160/pulse.csv")
knitr::kable(head(Pulse))
```

Table 4.7: Head of Pulse Rates of people Before and After Exercise

height	weight	age	gender	smokes	alcohol	exercise	ran	pulse_before	pulse_after	year
170	68	22	male	yes	yes	moderate	sat	70	71	93
182	75	26	male	yes	yes	moderate	sat	80	76	93
180	85	19	male	yes	yes	moderate	ran	68	125	95
182	85	20	male	yes	yes	low	sat	70	68	95
167	70	22	male	yes	yes	low	sat	92	84	96
178	86	21	male	yes	yes	low	sat	76	80	98

Code book for data frame Pulse

Description Students in an introductory statistics class (MS212 taught by Professor John Eccleston and Dr Richard Wilson at The University of Queensland) participated in a simple experiment. The students took their own pulse rate. They were then asked to flip a coin. If the coin came up heads, they were to run in place for one minute. Otherwise they sat for one minute. Then everyone took their pulse again. The pulse rates and other physiological and lifestyle data are given in the data.

Five class groups between 1993 and 1998 participated in the experiment. The lecturer, Richard Wilson, was concerned that some students would choose the less strenuous option of sitting

rather than running even if their coin came up heads, In the years 1995-1998 a different method of random assignment was used. In these years, data forms were handed out to the class before the experiment. The forms were pre-assigned to either running or non-running and there were an equal number of each. In 1995 and 1998 not all of the forms were returned so the numbers running and sitting was still not entirely controlled.

This data frame contains the following columns:

height: height of subject in cm

weight: weight of subject in kg

age: age of subject in years

gender: sex of subject, male, female

Smokes: whether a subject regularly smokes, yes means does smoke, no means does not smoke

alcohol: whether a subject regularly drinks alcohol, yes means the person does, no means the person does not

exercise: whether a subject exercises, low, moderate, high

ran: whether a subject ran one minute between pulse measurements (ran) or sat between pulse measurement (sat)

pulse_before: the pulse rate before a subject either ran or sat (bpm)

pulse_after: the pulse rate after a subject either ran or sat (bpm)

year: what year the data was collected (93-98)

Source Pulse rates before and after exercise. (2013, September 25). Retrieved from <http://www.statsci.org/data/oz/ms212.html>

References The data was supplied by Dr Richard J. Wilson, Department of Mathematics, University of Queensland.

Create a data frame that contains only males, who drink alcohol, but do not smoke. Then compare the pulse before and the pulse after using the mean and standard deviation. Discuss whether pulse before or pulse after has a higher mean and larger spread. The following command creates a new data frame with just males, who drink alcohol, but do not smoke, use the following command, where the new name is Males in Table ??.

```
Males<- Pulse |>
  filter(gender=="male", smokes == "no", alcohol == "yes")
knitr::kable(head(Males))
```

Table 4.8: Head of Pulse Rates of Nonsmoking Males Before and After Exercise

height	weight	age	gender	smokes	alcohol	exercise	ran	pulse_before	pulse_after	year
195	84	18	male	no	yes	high	sat	71	73	93
184	74	22	male	no	yes	low	ran	78	141	93
168	60	23	male	no	yes	moderate	ran	88	150	93
170	75	20	male	no	yes	high	ran	76	88	93
187	59	18	male	no	yes	high	sat	78	82	93
180	72	18	male	no	yes	moderate	sat	69	67	93

6. The data frame Pulse Table ?? contains various variables about a person including their pulse rates before the subject exercised and after the subject ran in place for one minute. Create a data frame that contains females, who do not smoke but do drink alcohol. Compare the pulse rate before and after exercise using the mean and standard deviation. Discuss whether pulse before or pulse after has a higher mean and larger spread.
7. To determine if Reiki is an effective method for treating pain, a pilot study was carried out where a certified second-degree Reiki therapist provided treatment on volunteers. Pain was measured using a visual analogue scale (VAS) and a likert scale immediately before and after the Reiki treatment (Olson & Hanson, 1997) and the data is in Table ??.

```
Reiki<- read.csv( "https://krkozak.github.io/MAT160/reki.csv")
knitr::kable(head(Reiki))
```

Table 4.9: Head of Pain Measurements Before and After Reiki Treatment

vas.before	vas.after	likert_before	likert_after
6	3	2	1
2	1	2	1
2	0	3	0
9	1	3	1
3	0	2	0
3	2	2	2

Code book for data frame Reiki

Description The purpose of this study was to explore the usefulness of Reiki as an adjuvant to opioid therapy in the management of pain. Since no studies in this area could be found, a pilot study was carried out involving 20 volunteers experiencing pain at 55 sites for a variety of reasons, including cancer. All Reiki treatments were provided by a certified second-degree Reiki therapist. Pain was measured using both a visual analogue scale (VAS) and a Likert

scale immediately before and after the Reiki treatment. Both instruments showed a highly significant ($p < 0.0001$) reduction in pain following the Reiki treatment.

This data frame contains the following columns:

vas.before: pain measured using a visual analogue scale (VAS) before Reiki treatment

vas.after: pain measured using a visual analogue scale (VAS) after Reiki treatment

likert_before: pain measured using a likert before Reiki treatment

likert_after: pain measured using a likert after Reiki treatment

Source Olson, K., & Hanson, J. (1997). Using reiki to manage pain: a preliminary report. Cancer Prev Control, 1(2), 108-13. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/9765732>

References** Using Reiki to manage pain: a preliminary report. Olson K1, Hanson J., Cancer Prev Control 1997, Jun; 1(2): 108-13.

Since the data was collected both before and after the treatment for all of the units of observations, you want to look at the effect size of the treatment. You want to find the difference between before and after for the pain scale. First you must create a new data frame that adds a column for the difference in before and after. This data is known as paired data. To create the new column in a new data frame called Newreiki use the following commands, Table ??.

```
Newreiki<-Reiki |>
  mutate(vas.diff=vas.before-vas.after)
knitr::kable(head(Newreiki))
```

Table 4.10: Head of Pain Measurements Before and After Reiki Treatment with Difference column

vas.before	vas.after	likert_before	likert_after	vas.diff
6	3	2	1	3
2	1	2	1	1
2	0	3	0	2
9	1	3	1	8
3	0	2	0	3
3	2	2	2	1

Now find the mean and standard deviation of the vas.diff variable in Newreiki. Perform similar commands to create the likert.diff variable. Then find the mean and standard deviation for likert.diff, and compare and contrast the vas and likert methods for describing pain.

8. Yearly rainfall amounts (in millimeters) in Sydney, Australia, are in Table ?? (Annual maximums of, 2013). a. Calculate the mean and standard deviation. b. Suppose Sydney, Australia received 300 mm of rainfall in a year. Would this be unusual?

```
Rainfall<-read.csv("https://krkozak.github.io/MAT160/rainfall.csv")
knitr::kable(head(Rainfall))
```

Table 4.11: Head of Yearly rainfall amounts in Sydney, Australia

amount
146.8
383.0
90.9
178.1
267.5
95.5

Code book for data frame Rainfall

Description Daily rainfall (in millimeters) was recorded over a 47-year period in Turramurra, Sydney, Australia. For each year, the wettest day was identified (that having the greatest rainfall). The data show the rainfall recorded for the 47 annual maxima.

This data frame contains the following columns:

amount: daily rainfall (mm)

Source Annual maximums of daily rainfall in Sydney. (2013, September 25). Retrieved from <http://www.statsci.org/data/oz/sydrain.html>

References Rayner J.C.W. and Best D.J. (1989) Smooth tests of goodness of fit. Oxford: Oxford University Press. Hand D.J., Daly F., Lunn A.D., McConway K.J., Ostrowski E. (1994). A Handbook of Small Data Sets. London: Chapman & Hall. Data set 157. Thanks to Jim Irish of the University of Technology, Sydney, for assistance in identifying the correct units for this data.

4.3 Ranking

Along with the center and the variability, another useful numerical measure is the ranking of a number. A **percentile** is a measure of ranking. It represents a location measurement of a data value to the rest of the values. Many standardized tests give the results as a percentile. Doctors also use percentiles to track a child's growth.

The k^{th} **percentile** is the data value that has $k\%$ of the data at or below that value.

4.3.1 Example: Interpreting Percentile

- a. What does a score of the 90th percentile mean?
- b. What does a score of the 70th percentile mean?

4.3.1.1 Solution

- a. What does a score of the 90th percentile mean?

This means that 90% of the scores were at or below this score. (A person did the same as or better than 90% of the test takers.)

- b. What does a score of the 70th percentile mean?

This means that 70% of the scores were at or below this score.

4.3.2 Example: Percentile Versus Score

If the test was out of 100 points and you scored at the 80th percentile, what was your score on the test?

4.3.2.1 Solution

You don't know! All you know is that you scored the same as or better than 80% of the people who took the test. If all the scores were really low, you could have still failed the test. On the other hand, if many of the scores were high you could have gotten a 95% or more.

There are special percentiles called **quartiles**. Quartiles are numbers that divide the data into fourths. One fourth (or a quarter) of the data falls between consecutive quartiles.

4.3.3 To find the quartiles:

The command in rStudio is

```
df_stats(~variable, data=data_frame, summary)
```

If you record the quartiles together with the maximum and minimum you have five numbers. This is known as the five-number summary. The five-number summary consists of the minimum, the first quartile ($Q1$), the median, the third quartile ($Q3$), and the maximum (in that order).

The interquartile range, IQR , is the difference between the first and third quartiles, $Q1$ and $Q3$. Half of the data (50%) falls in the interquartile range. If the IQR is “large” the data is spread out and if the IQR is “small” the data is closer together.

Interquartile Range (IQR)

Determining probable outliers from IQR : **fences**

A value that is less than $Q1 - 1.5 * IQR$ (this value is often referred to as a **low fence**) is considered an outlier.

Similarly, a value that is more than $Q3 + 1.5 * IQR$ (the **high fence**) is considered an outlier.

A boxplot (or box-and-whisker plot) is a graphical display of the five-number summary. It can be drawn vertically or horizontally. The basic format is a box from $Q1$ to $Q3$, a vertical line across the box for the median and horizontal lines as whiskers extending out each end to the minimum and maximum. The minimum and maximum can be represented with dots. Don’t forget to label the tick marks on the number line and give the graph a title.

An alternate form of a Boxplot, known as a modified box plot, only extends the left line to the smallest value greater than the **low fence**, and extends the right line to the largest value less than the **high fence**, and displays markers (dots, circles or asterisks) for each outlier.

If the data are **symmetrical**, then the box plot will be visibly symmetrical. If the data distribution has a left skew or a right skew, the line on that side of the box plot will be visibly long. If the plot is symmetrical, and the four quartiles are all about the same length, then the data are likely a near **uniform** distribution. If a box plot is symmetrical, and both outside lines are noticeably longer than the $Q1$ to median and median to $Q3$ distance, the distribution is then probably **bell-shaped**.

4.3.4 Example: Five-number Summary and Boxplot

Find the five-number summary, the interquartile range ($*IQR*$), and draw a box-and-whiskers plot for the weight of cats Table ??.

4.3.4.1 Solution

Variable: x = weight of cats To compute the five-number summary on RStudio, use the command:

```
df_stats(~Bwt, data=cats, summary)
```

	response	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1	Bwt	2	2.3	2.7	2.723611	3.025	3.9

Note rStudio also calculates the mean as part of the summary command, but the five-number summary is just the five numbers:

Minimum: 2 kg $Q1$: 2.3 kg Median: 2.7 kg $Q3$: 3.025 kg Maximum: 3.9 kg

To find the interquartile range, IQR find $Q3 - Q1$, so $IQR = 3.025 - 2.3 = 0.725kg$

To create a boxplot use the command

```
gf_boxplot(~variable, data=data_frame)
```

This is a modified boxplot which shows the outliers in the data.

```
gf_boxplot(~Bwt, data=cats, title="Weight of Cats", xlab="Body Weight (kg)")
```

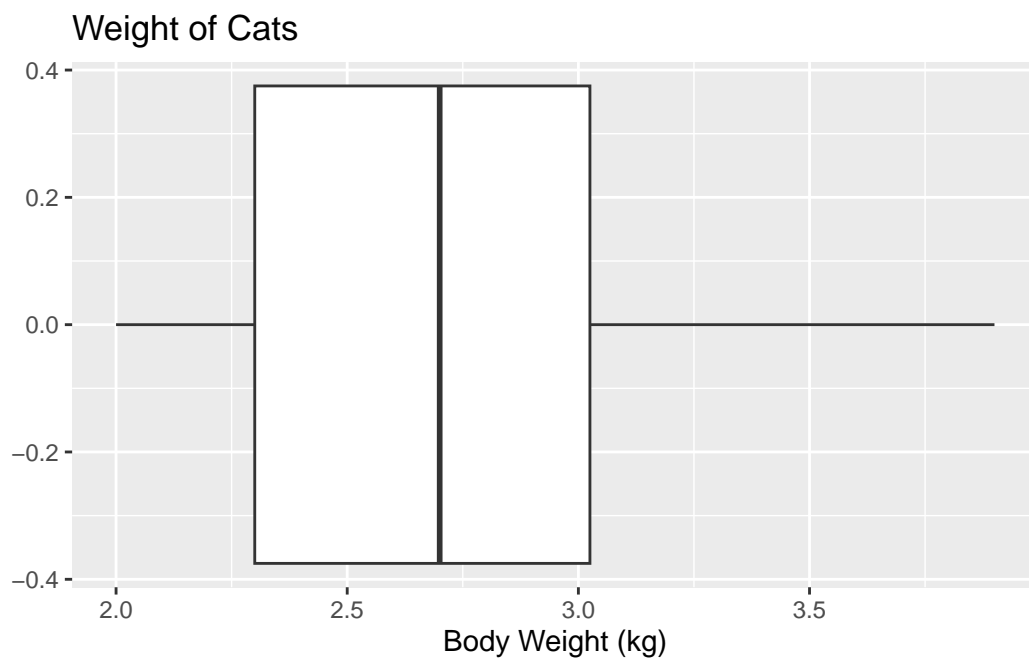


Figure 4.5: Weight of Cats

There are no outliers since there are no dots outside of the fences.

4.3.5 Example: Separating based on a factor

Find the five-number summary of the weights of cats separated by the sex of the cat. Then create a box plot of the weights of cats for each sex of the cat.

4.3.5.1 Solution

Variable: x_1 = weight of female cat

Variable: x_2 = weight of male cat

To find the five-number summary separated based on gender use the following command:

```
df_stats(~Bwt|Sex, data=cats, summary)
```

	response	Sex	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1	Bwt	F	2	2.15	2.3	2.359574	2.5	3.0
2	Bwt	M	2	2.50	2.9	2.900000	3.2	3.9

The five-number summary for female cats is (in kg)

Minimum: 2 Q1: 2.15 Median: 2.3 Q3: 2.5 Maximum: 3.0

The five-number summary for male cats is (in kg)

Minimum: 2 Q1: 2.50 Median: 2.9 Q3: 3.2 Maximum: 3.9

```
gf_boxplot(~Bwt|Sex, data=cats, title="Weights of Cats", xlab="Body Weight in (kg)")
```

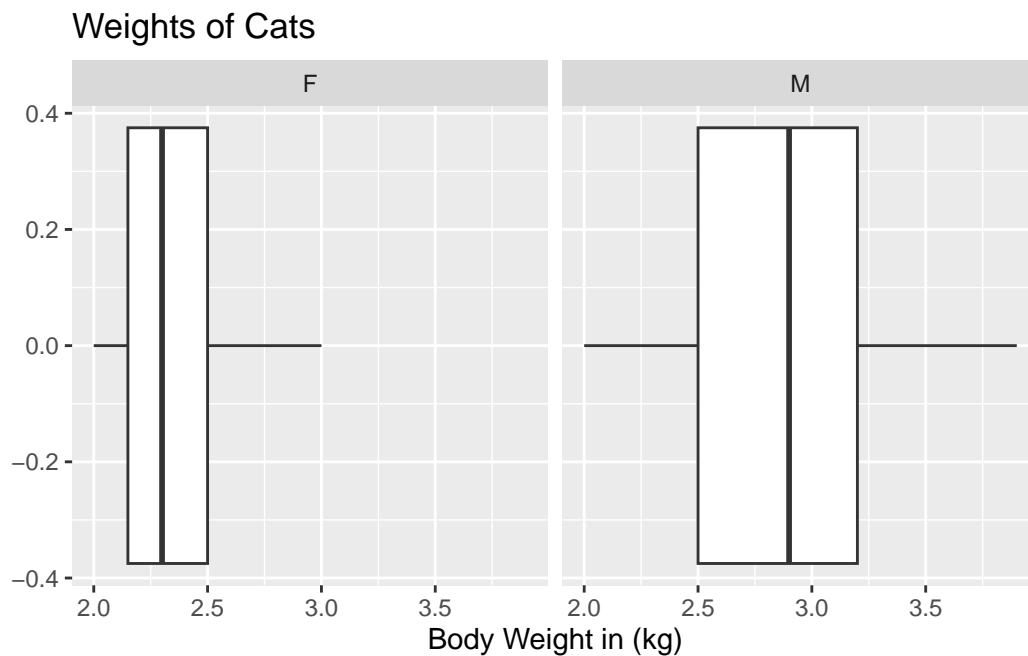


Figure 4.6: Weight of Cats Faceted by Sex

Notice that the weights of female cats has a median less than male cats, and in fact it can be seen that the $Q1$ to $Q3$ of the female cats is less than the $Q1$ to $Q3$ of the male cats.

4.3.6 Example: Putting it all together

The time (in 1/50 seconds) between successive pulses along a nerve fiber (“Time between nerve,” 2013) are given in Table ??.

```
Nerve<-read.csv( "https://krkozak.github.io/MAT160/Nerve_pulse.csv")
knitr::kable(head(Nerve))
```

Table 4.12: Head of Successive pulses along a nerve fiber

time
10.5
1.5
2.5
5.5
29.5
3.0

Code book for data frame Nerve

Description The data gives the time between 800 successive pulses along a nerve fiber. There are 799 observations rounded to the nearest half in units of 1/50 second.

This data frame contains the following columns:

time: time between successive Pulses along a nerve fiber, 1/50 second.

Source Time between nerve pulses. (2019, July 3). Retrieved from <<http://www.statsci.org/data/general/nerve>>.

References Fatt, P., and Katz, B. (1952). Spontaneous subthreshold activity at motor nerve endings. *Journal of Physiology* 117, 109-128.

Cox, D. R., and Lewis, P. A. W. (1966). *The Statistical Analysis of Series of Events*. Methuen, London.

Jorgensen, B. (1982). *The Generalized Inverse-Gaussian Distribution*. Springer-Verlag.

4.3.6.1 Solution

First, it might be useful to look at a visualization of the data, so create a density plot

```
gf_density(~time, data=Nerve, title="Time between Successive Nerve Pulses", xlab="Time (1/50
```

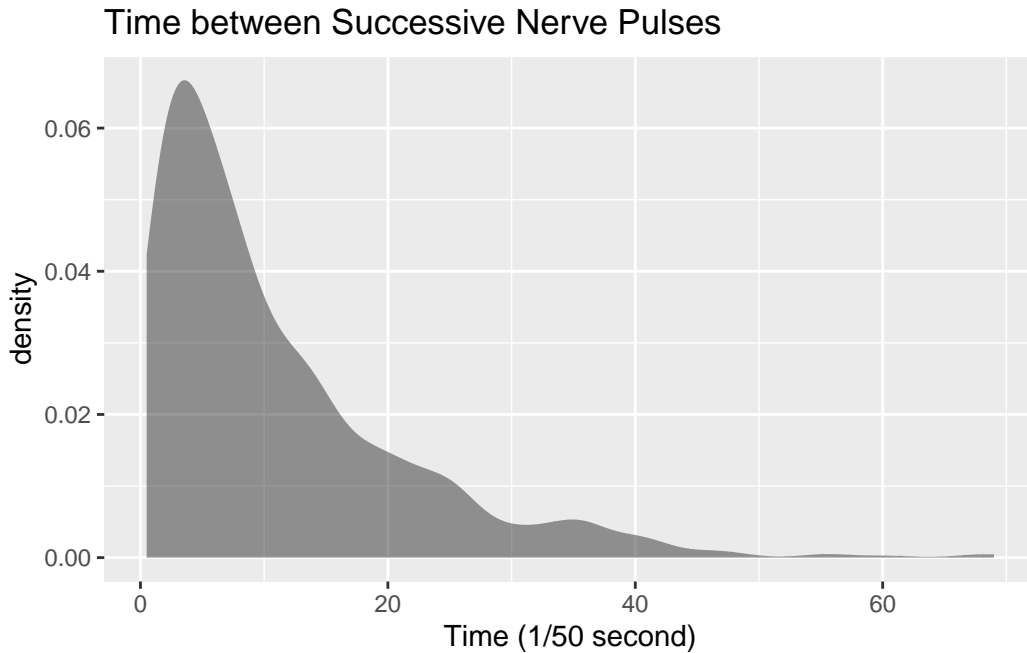


Figure 4.7: Weight of Cats Faceted by Sex

From the graph Figure ?? the data appears to be skewed right. Most of the time between successive nerve pulses appear to be around 5 or 10 1/50 second, but there are some times that are 60 1/50 second.

```
df_stats(~time, data=Nerve, mean, median, sd, summary)
```

response	mean	median	sd	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1 time	10.95119	7.5	10.45956	0.5	3.5	7.5	10.95119	15	69

Numerical descriptions might also be useful. Using technology, the mean is 11 1/50 second, the median is 7.5 1/50 second, the standard deviation is 10.5 1/50 second, and the five-number summary is minimum = 3.5, Q1 = 3.5, median = 7.5, Q3 = 15, and maximum = 69 1/50 second.

To visualize the five-number summary, create a box plot.


```
gf_boxplot(~time, data=Nerve, title="Nerve Pulses", xlab="Time (1/50 second)")
```

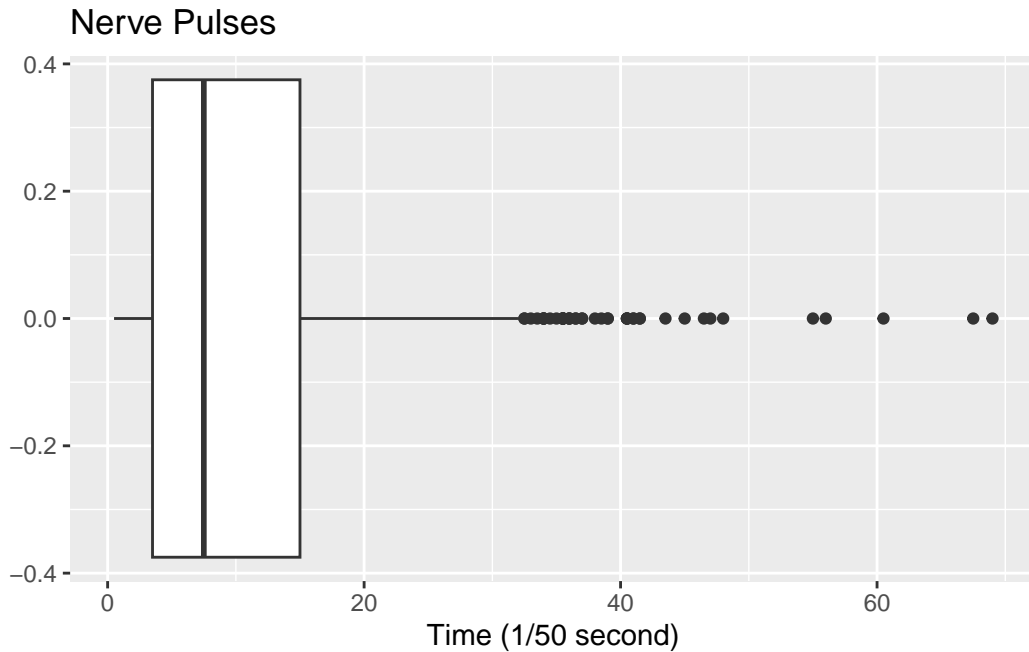


Figure 4.8: Boxplot of Nerve Pulses

Since there are many dots outside the upper fence the data has many outliers. From all of this information, one could say that nerve pulses between successive pulses is around 11 $\frac{1}{50}$ second, with a spread of 19.5 $\frac{1}{50}$ second. Most of the values are round 11 $\frac{1}{50}$ second, but they are not very consistent. The density plot and boxplot show that there is a great deal of spread of the data and it is skewed to the right. This means mostly the speed is around 11 $\frac{1}{50}$ second, but there is a great deal of variability in the values.

4.3.7 Homework for Ranking Section

Use Technology on all problems. State the variable on all problems.

1. Suppose you take a standardized test and you are in the 10th percentile. What does this percentile mean? Can you say that you failed the test? Explain.
2. Suppose your child takes a standardized test in mathematics and scores in the 96th percentile. What does this percentile mean? Can you say your child passed the test? Explain.

3. Suppose your child is in the 83rd percentile in height and 24th percentile in weight. Describe what this tells you about your child's stature.
4. Suppose your work evaluates the employees and places them on a percentile ranking. If your evaluation is in the 65th percentile, do you think you are working hard enough? Explain.
5. Cholesterol levels were collected from patients certain days after they had a heart attack and are in table Table ??.

Code book for Data Frame Cholesterol below Table ??.

Find the five-number summary and interquartile range (IQR) for the cholesterol level on day 2, and draw a boxplot

6. The lengths (in kilometers) of rivers on the South Island of New Zealand and what body of water they flow into are listed in table Table ?? (Lee, 1994).

Code book for data frame Length below Table ??.

Find the five-number summary and interquartile range (IQR) for the lengths of rivers that go to the Pacific Ocean and ones that go to the Tasman Sea, and draw a boxplot of both.

7. Print-O-Matic printing company's employees have salaries that are contained in Table ??
Find the five number summary and draw a boxplot for the salaries of all employees.

Code book for data frame Pay below Table ??.

8. The data frame Pulse Table ?? contains various variables about a person including their pulse rates before the subject exercised and after after the subject ran in place for one minute.

Code book for data frame Pulse below Table ??.

Create a data frame that contains only people who drink alcohol, but do not smoke. Then find the five number summary and draw a boxplot for both males and females separately.

9. To determine if Reiki is an effective method for treating pain, a pilot study was carried out where a certified second-degree Reiki therapist provided treatment on volunteers. Pain was measured using a visual analogue scale (VAS) and a likert scale immediately before and after the Reiki treatment (Olson & Hanson, 1997) and the data is in Table ??.

Code book for data frame Reiki below Table ??.

Find the five number summary for both the before and after VAS scores and draw boxplots of before and after VAS scores. To draw two boxplots at the same time, after the command to create the first box plot type the piping symbol `|>` (base r) or `%>%` (magrittr package) before pressing enter. (Note: `|>` and `%>%` are piping symbols that can be thought of as "and then.")

Then type the command for the second boxplot after the + symbol or on the next line in the r chunk if using an rmd or qmd file. Then press enter. You may want to graph each boxplot as a different color. To do this, the command would be

```
gf_boxplot(~variable, data=data_frame, color="red", xlab="type a label")
```

You can pick any color you want. Just replace the word red with the color you want to use. Now compare and contrast the before and after VAS scores.

5 Probability

Understanding probabilities are important in life. Examples of mundane questions that probability can answer for you are if you need to carry an umbrella or wear a heavy coat on a given day. More important questions that probability can help with are your chances that the car you are buying will need more maintenance, your chances of passing a class, your chances of winning the lottery, your chances of being in a car accident, and the chances that the U.S. will be attacked by terrorists. Most people do not have a very good understanding of probability, so they worry about being attacked by a terrorist but not about being in a car accident. The probability of being in a terrorist attack is much smaller than the probability of being in a car accident, thus it actually would make more sense to worry about driving. Also, the chance of you winning the lottery is very small, yet many people will spend their money on lottery tickets. Yet, if instead they saved the money that they spend on the lottery, they would have more money. In general, events that have a low probability (under 5%) are unlikely to occur. Whereas if an event has a high probability of happening (over 80%), then there is a good chance that the event will happen. This chapter will present some of the theory that you need to help make a determination of whether an event is likely to happen or not.

One story about how probability theory was developed is that a gambler wanted to know when to bet more and when to bet less. He talked to a couple of friends of his that happened to be mathematicians. Their names were Pierre de Fermat and Blaise Pascal. Since then many other mathematicians have worked to develop probability theory. There are actually two types of probability, **Empirical Probability** and **Theoretical Probability** that have been developed since the start of probability.

5.1 Empirical Probability

Empirical probabilities are found by actually conducting an experiment many times and counting the number of times the event happens. To understand how this is performed, first some definitions are needed.

Outcomes: the results of an experiment

Event: a set of certain outcomes of an experiment that you want to have happen

Sample Space: collection of all possible outcomes of the experiment. Usually denoted as SS .

Event space: the set of outcomes that make up an event. The symbol is usually a capital letter.

Frequency: how often an event happens

Relative Frequency: the frequency divided by the number of times the experiment is repeated

Start with an experiment. Suppose that the experiment is rolling a die. The sample space is $\{1, 2, 3, 4, 5, 6\}$. The event that you want is to get a 6, and the event space is $\{6\}$. To do this, roll a die 10 times. When you do that, you get a 6 two times. Based on this experiment, the probability of getting a 6 is 2 out of 10 or $1/5$. To get more accuracy, repeat the experiment more times. It is easiest to put this in a table, where n represents the number of times the experiment is repeated. When you put the number of 6s found over the number of times you repeat the experiment, this is the relative frequency.

Table 5.1: trials for Die Experiment

n	number_of_6s	relative_frequency
10	2	0.200
50	6	0.120
100	18	0.180
500	81	0.162
1000	163	0.163

Notice that as n increased, the relative frequency seems to approach a number. It looks like it is approaching 0.163. You can say that the probability of getting a 6 is approximately 0.163. If you want more accuracy, then increase n even more.

These probabilities are called **experimental probabilities** since they are found by actually doing the experiment. They come about from the relative frequencies and give an approximation of the true probability. The approximate probability of an event A , $P(A)$, is

5.1.1 Experimental Probabilities

$$P(A) = \frac{\text{number of times}}{\text{number of times experiment is conducted}}$$

For the event of getting a 6, the probability would be

$$P(6) = \frac{163}{1000} = 0.163.$$

You must do experimental probabilities whenever it is not possible to calculate probabilities using other means. An example is if you want to find the probability that a family has 5 children, you would have to actually look at many families, and count how many have 5

children. Then you could calculate the probability. Another example is if you want to figure out if a die is fair. You would have to roll the die many times and count how often each side comes up. Make sure you repeat an experiment many times, because otherwise you will not be able to estimate the true probability. This is due to the law of large numbers.

Law of large numbers: as n increases, the relative frequency tends towards the actual probability value.

Note: probability, relative frequency, percentage, and proportion are all different words for the same concept. Also, probabilities can be given as percentages, decimals, or fractions.

To find probabilities from data, you can take a data frame and count the number of values for each outcome.

5.1.2 Example: Statistics class survey

Data was collected for two semesters in a statistics class. The data frame in is in Table ??.

Find the probability (proportion) of people who like Cookie Dough ice cream.

5.1.2.1 Solution

To count the number of people who like cookie dough ice cream, use the following command in r Studio:

```
tally(~ice_cream, data=Class, margins=TRUE)
```

```
ice_cream
      Butter Pecan      Chocolate      Chocolate Brownie.
           2           2           1
      coffee      Cookie Dough      Cookies and Cream
           1           6           1
      Mint CC      Moose Tracks      none
           6           1           1
      Rocky Road      Sherbet Strawberry and banana
           2           2           1
      Vanilla      Total
           1          27
```

From this tally, it can be seen that 6 people like cookie dough ice cream. The probability that someone likes cookie dough is thus 6 divided by the number of people in the data frame, the Total. Instead of dividing, the following command will find the proportions for you. Proportions are just probabilities.

```
tally(~ice_cream, data=Class, format="proportion")
```

```
ice_cream
      Butter Pecan      Chocolate      Chocolate Brownie.
      0.07407407      0.07407407      0.03703704
      coffee      Cookie Dough      Cookies and Cream
      0.03703704      0.22222222      0.03703704
      Mint CC      Moose Tracks      none
      0.22222222      0.03703704      0.03703704
      Rocky Road      Sherbet Strawberry and banana
      0.07407407      0.07407407      0.03703704
      Vanilla
      0.03703704
```

So the probability that a person in the class likes cookie dough ice cream is 0.22.

5.1.3 Homework for Empirical Probability Section

1. The number of M&M's of each color that were found in a packet is in Table ?? (M&M's Color Distribution Analysis, 2019).

Table 5.2: M&M Distribution

color	type	pack
orange	plain	1
red	plain	1
green	plain	1
red	plain	1
yellow	plain	1
blue	plain	1

Code book for Data Frame MaM

Description An analysis of the colors in a case of M&M's to see if they match the published percentages

Usage MaM

Format

This data frame contains the following columns:

color: color of M&Ms

type: The type of M&M such as plain, peanut, peanut butter

pack: which pack the M&Ms came from.

Source M&M's Color Distribution Analysis. (n.d.). Retrieved July 11, 2019, from <https://joshmadison.com/2007/12/02/mms-color-distribution-analysis/>

References Josh Madison, 2019

Find the probability of choosing each color based on this data frame.

2. Eyeglassomatic manufactures eyeglasses for different retailers. They test to see how many defective lenses they made the time period of January 1 to March 31. The defect and the number of defects is in Table ??.

Code book for Data Frame Defects below Table ??.

Find the probability of each defect type based on this data.

3. In Australia in 1995, of the 2907 indigenous people in prison 17 of them died. In that same year, of the 14501 non-indigenous people in prison 42 of them died (\“Aboriginal deaths in,\” 2013). Find the probability that an indigenous person dies in prison and the probability that a non-indigenous person dies in prison. Compare these numbers and discuss what the numbers may mean.
4. A project conducted by the Australian Federal Office of Road Safety asked people many questions about their cars. One question was the reason that a person chooses a given car, and that data is in Table ?? (Car Preferences, 2019).

Table 5.3: Reason for Choosing a Car

ID	Age	Sex	Lic	Yic	Mch	Kids	Dr	Ref	Car	Real	Cont	Wl	Mi	Per	Fuel	Safety	AC	PS	Bark	Room	Doors	Prest	Colour
11018	mal	0	2	large	no	no	med	large	safety	sty	portant	very	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant
11125	fem	8	0	small	no	small	small	safety	very	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant
11263	mal	4	0	large	no	large	large	consid	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant
11351	fem	3	0	large	no	med	large	safety	little	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant
11419	fem	2	0	medium	no	med	small	looks	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant
11551	fem	3	0	medium	yes	med	large	consid	portant	very	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant	imptant

Code book for Data Frame Car_pref

Description These data were collected as part of a project for the Federal Office for Road Safety conducted by the Research Institute of Gender and Health at the University of Newcastle. There is evidence that women drivers who are involved in motor vehicle accidents are

more likely than men to be injured. A possible reason is that women often drive smaller cars that provide less protection in a collision. One of the aims of the project was to examine preferences for cars among men and women and investigate the extent to which safety was a factor in determining preferences. The survey was conducted by research assistants who asked people in car parks to participate and administered a structured questionnaire. They were instructed to obtain data from men and women with small, medium and large cars, with 50 people per group for a total of 300 respondents. (The sample size was based on power requirements for another part of the survey that involved anthropometric measurements.) The research assistants approached people in car parks of the University of Newcastle and nearby shopping centers during December 1997 and January 1998.

Usage Car_pref

Format

This data frame contains the following columns:

ID: Identification number of respondent

Age: Age of respondent (years)

Sex: female, male

LicYr: Time they have held a full driving licence, in years and months (years)

LicMth: Time they have held a full driving licence, in years and months (months)

ActCar: Make, model and year of car most often driven, coded to size of car small, medium, large

Kids5: Children under five, yes, no

Kids6: Children 6 to 16, yes, no

PrefCar: Preferred car, coded to size of car small, medium, large

Car15k: Preferred type of car if cost \ \$15000, small new car; large second-hand car

Reason: safety, reliability, cost, performance, comfort, looks

Cost: How important is cost when buying a car? not important, little importance, important, very important

Reliable: How important is reliability ...?

Perform: How important is performance ...?

Fuel: How important is fuel consumption ...?

Safety: How important is safety ...?

AC/PS: How important is air conditioning/power steering ...?

Park: How important is ease of parking ...?

Room: How important is space/roominess ...?

Doors: How important is the number of doors ...?

Prestige: How important is prestige/style ...?

Colour: How important is colour ...?

Source

Car Preferences. (n.d.). Retrieved July 11, 2019, from <http://www.statsci.org/data/oz/carprefs.html>

References

The data was contributed to OzDASL by Professor Annette Dobson, University of Queensland. Information on the data set was originally provided by Jenny Powers.

Find the probability a person chooses a car for each of the given reasons.

5.2 Theoretical Probability

It is not always feasible to conduct an experiment over and over again, so it would be better to be able to find the probabilities without conducting the experiment. These probabilities are called **Theoretical Probabilities**.

To be able to do theoretical probabilities, there is an assumption that you need to consider. It is that all of the outcomes in the sample space need to be ****equally likely outcomes****. This means that every outcome of the experiment needs to have the same chance of happening.

5.2.1 Example: Equally Likely Outcomes

Which of the following experiments have equally likely outcomes?

- Rolling a fair die.
- Flip a coin that is weighted so one side comes up more often than the other.
- Pull a ball out of a can containing 6 red balls and 8 green balls. All balls are the same size.
- Picking a card from a deck.
- Rolling a die to see if it is fair.

5.2.1.1 Solution

- a. Rolling a fair die.

Since the die is fair, every side of the die has the same chance of coming up. The outcomes are the different sides, so each outcome is equally likely

- b. Flip a coin that is weighted so one side comes up more often than the other.

Since the coin is weighted, one side is more likely to come up than the other side. The outcomes are the different sides, so each outcome is not equally likely

- c. Pull a ball out of a can containing 6 red balls and 8 green balls. All balls are the same size.

Since each ball is the same size, then each ball has the same chance of being chosen. The outcomes of this experiment are the individual balls, so each outcome is equally likely. Don't assume that because the chances of pulling a red ball are less than pulling a green ball that the outcomes are not equally likely. The outcomes are the individual balls and they are equally likely.

- d. Picking a card from a deck.

If you assume that the deck is fair, then each card has the same chance of being chosen. Thus the outcomes are equally likely outcomes. You do have to make this assumption. For many of the experiments you will do, you do have to make this kind of assumption.

- e. Rolling a die to see if it is fair.

In this case you are not sure the die is fair. The only way to determine if it is fair is to actually conduct the experiment, since you don't know if the outcomes are equally likely. If the experimental probabilities are fairly close to the theoretical probabilities, then the die is fair.

If the outcomes are not equally likely, then you must do experimental probabilities. If the outcomes are equally likely, then you can do theoretical probabilities.

Theoretical Probabilities: If the outcomes of an experiment are equally likely, then the probability of event A happening is

$$P(A) = \frac{\text{number of outcomes in event space}}{\text{number of outcomes in sample space}}$$

5.2.2 Example: Calculating Theoretical Probabilities

Suppose you conduct an experiment where you flip a fair coin twice

- What is the sample space?
- What is the probability of getting exactly one head?
- What is the probability of getting at least one head?
- What is the probability of getting a head and a tail?
- What is the probability of getting a head or a tail?
- What is the probability of getting a foot?
- What is the probability of each outcome? What is the sum of these probabilities?

5.2.2.1 Solution

- What is the sample space?

There are several different sample spaces you can do. One is $SS = \{0, 1, 2\}$ where you are counting the number of heads. However, the outcomes are not equally likely since you can get one head by getting a head on the first flip and a tail on the second or a tail on the first flip and a head on the second. There are 2 ways to get that outcome and only one way to get the other outcomes. Instead it might be better to give the sample space as listing what can happen on each flip. Let H = head and T = tail, and list which can happen on each flip.

$$SS = \{HH, HT, TH, TT\}$$

- What is the probability of getting exactly one head?

Let A = getting exactly one head. The event space is $A = \{HT, TH\}$. So $P(A) = \frac{2}{4}$

It may not be advantageous to reduce the fractions to lowest terms, since it is easier to compare fractions if they have the same denominator.

- What is the probability of getting at least one head?

Let B = getting at least one head. At least one head means get one or more. The event space is $B = \{HT, TH, HH\}$ and $P(B) = \frac{3}{4}$. Since $P(B)$ is greater than the $P(A)$, then event B is more likely to happen than event A .

- What is the probability of getting a head and a tail?

Let C = getting a head and a tail = $\{HT, TH\}$ and $P(C) = \frac{2}{4}$. This is the same event space as event A , but it is a different event. Sometimes two different events can give the same event space.

- e. What is the probability of getting a head or a tail?

Let D = getting a head or a tail. Since or means one or the other or both and it doesn't specify the number of heads or tails, then $D = \{HH, HT, TH, TT\}$ and $P(D) = \frac{3}{4}$

- f. What is the probability of getting a foot?

Let E = getting a foot. Since you can't get a foot, $E = \{\}$ or the empty set and $P(E) = \frac{0}{4} = 0$

- g. What is the probability of each outcome? What is the sum of these probabilities?

$P(HH) = P(HT) = P(TH) = P(TT) = \frac{1}{4}$. If you add all of these probabilities together you get 1.

This example had some results in it that are important concepts. They are summarized below:

5.2.3 Probability Properties

1. $0 \leq P(\text{event}) \leq 1$
2. If the $P(\text{event}) = 1$, then it will happen and is called the certain event
3. If the $P(\text{event}) = 0$, then it cannot happen and is called the impossible event
4. $\sum P(\text{all outcomes}) = 1$

5.2.4 Example: Calculating Theoretical Probabilities 2

Suppose you conduct an experiment where you pull a card from a standard deck.

- a. What is the sample space?
- b. What is the probability of getting a Spade?
- c. What is the probability of getting a Jack?
- d. What is the probability of getting an Ace?
- e. What is the probability of not getting an Ace?
- f. What is the probability of not getting an Ace?
- g. What is the probability of getting a Spade or an Ace?
- h. What is the probability of getting a Jack and an Ace?
- i. What is the probability of getting a Jack and an Ace?

5.2.4.1 Solution

- a. What is the sample space?

$SS = \{2S, 3S, 4S, 5S, 6S, 7S, 8S, 9S, 10S, JS, QS, KS, AS, 2C, 3C, 4C, 5C, 6C, 7C, 8C, 9C, 10C, JC, QC, KC, AC, 2D, 3D, 4D, 5D, 6D, 7D, 8D, 9D, 10D, JD, QD, KD, AD, 2H, 3H, 4H, 5H, 6H, 7H, 8H, 9H, 10H, JH, QH, KH, AH\}$

- b. What is the probability of getting a Spade?

Getting a spade = $\{2S, 3S, 4S, 5S, 6S, 7S, 8S, 9S, 10S, JS, QS, KS, AS\}$ so $P(spade) = \frac{13}{52}$

- c. What is the probability of getting a Jack?

Getting a Jack = $\{JS, JC, JH, JD\}$ so $P(jack) = \frac{4}{52}$

- d. What is the probability of getting an Ace?

Getting an Ace = $\{AS, AC, AH, AD\}$ so $P(ace) = \frac{4}{52}$

- e. What is the probability of not getting an Ace?

Not getting an Ace = $\{2S, 3S, 4S, 5S, 6S, 7S, 8S, 9S, 10S, JS, QS, KS, 2C, 3C, 4C, 5C, 6C, 7C, 8C, 9C, 10C, JC, QC, KC, 2D, 3D, 4D, 5D, 6D, 7D, 8D, 9D, 10D, JD, QD, KD, 2H, 3H, 4H, 5H, 6H, 7H, 8H, 9H, 10H, JH, QH, KH\}$ so $P(\text{not ace}) = \frac{48}{52}$

Notice, $P(ace) + P(\text{not ace}) = 1$, so you could have found the probability of not ace by doing 1 minus the probability of ace. $P(\text{not ace}) = 1 - P(ace) = 1 - \frac{4}{52} = \frac{48}{52}$

- f. What is the probability of getting a Spade and an Ace?

Getting a Spade and an Ace = $\{AS\}$ so $P(AS) = \frac{1}{52}$

- g. What is the probability of getting a Spade or an Ace?

Getting a Spade and an Ace = $\{2S, 3S, 4S, 5S, 6S, 7S, 8S, 9S, 10S, JS, QS, KS, AS, AC, AD, AH\}$ so $P(\text{spade and ace}) = \frac{16}{52}$

- h. What is the probability of getting a Jack and an Ace?

Getting a Jack and an Ace = $\{ \}$ since you can't do that when picking one card. So $P(\text{Jack and Ace}) = \frac{0}{52} = 0$

- i. What is the probability of getting a Jack or an Ace?

Getting a Jack or an Ace = $\{JS, JC, JD, JH, AS, AC, AD, AH\}$ so $P(\text{Jack or Ace}) = \frac{8}{52}$

5.2.5 Example: Calculating Theoretical Probabilities 3

Suppose you have an iPhone and playing iTunes with the following songs on it: 5 Rolling Stones songs, 7 Beatles songs, 9 Bob Dylan songs, 4 Faith Hill songs, 2 Taylor Swift songs, 7 U2 songs, 4 Mariah Carey songs, 7 Bob Marley songs, 6 Bunny Wailer songs, 7 Elton John songs, 5 Led Zeppelin songs, and 4 Dave Mathews Band songs. The different genre that you have are rock from the 60s which includes Rolling Stones, Beatles, and Bob Dylan; country includes Faith Hill and Taylor Swift; rock of the 90s includes U2 and Mariah Carey; Reggae includes Bob Marley and Bunny Wailer; rock of the 70s includes Elton John and Led Zeppelin; and bluegrass-rock includes Dave Mathews Band.

Suppose the iTunes is set to shuffle the songs, so it randomly picks the next song so you have no idea what the next song will be. Now you would like to calculate the probability that you will hear the type of music or the artist that you are interested in. The sample set is too difficult to write out, but you can figure it from looking at the number in each set and the total number. The total number of songs you have is 67.

- What is the probability that you will hear a Faith Hill song?
- What is the probability that you will hear a Bunny Wailer song?
- What is the probability that you will hear a song from the 60s?
- What is the probability that you will hear a Reggae song?
- What is the probability that you will hear a song from the 90s or a bluegrass-rock song?
- What is the probability that you will hear an Elton John or a Taylor Swift song?
- What is the probability that you will hear a country song or a U2 song?

5.2.5.1 Solution

- What is the probability that you will hear a Faith Hill song?

There are 4 Faith Hill songs out of the 67 songs, so $P(\text{Faith Hill}) = \frac{4}{67}$

- What is the probability that you will hear a Bunny Wailer song?

There are 6 Bunny Wailer songs, so $P(\text{Bunny Wailer}) = \frac{6}{67}$

- What is the probability that you will hear a song from the 60s?

There are 5, 7, and 9 songs that are classified as rock from the 60s, which is 21 total, so $P(\text{song from 60s}) = \frac{21}{67}$

- What is the probability that you will hear a Reggae song?

There are 6 and 7 songs that are classified as Reggae, which is 13 total, so $P(\text{Reggae}) = \frac{13}{67}$

- e. What is the probability that you will hear a song from the 90s or a bluegrass-rock song?

There are 7 and 4 songs that are songs from the 90s and 4 songs that are bluegrass-rock, for a total of 15, so $P(\text{song 90s or bluegrass-rock}) = \frac{15}{67}$

- f. What is the probability that you will hear an Elton John or a Taylor Swift song?

There are 7 Elton John songs and 2 Taylor Swift songs, for a total of 9, so $P(\text{Elton John or Taylor Swift}) = \frac{9}{67}$

- g. What is the probability that you will hear a country song or a U2 song?

There are 6 country songs and 7 U2 songs, for a total of 13, so $P(\text{country or U2}) = \frac{13}{67}$

Of course you can do any other combinations you would like.

Notice in Example: Calculating Theoretical Probabilities part e, it was mentioned that the probability of getting an ace plus the probability of not getting an ace was 1. This is because these two events have no outcomes in common, and together they make up the entire sample space. Events that have this property are called **complementary events**.

If two events are **complementary events** then to find the probability of one just subtract the probability of the other from one. Notation used for complement of A is not A or A^c .

$$P(A) + P(\text{not } A) = 1$$

5.2.6 Example: Complementary Events

- Suppose you know that the probability of it raining today is 0.45. What is the probability of it not raining?
- Suppose you know the probability of not getting the flu is 0.24. What is the probability of getting the flu?
- In an experiment of picking a card from a deck, what is the probability of not getting a card that is a Queen?

5.2.6.1 Solution

- Suppose you know that the probability of it raining today is 0.45. What is the probability of it not raining?

Since not raining is the complement of raining, then $P(\text{not raining}) = 1 - P(\text{raining}) = 1 - 0.45 = 0.55$

- b. Suppose you know the probability of not getting the flu is 0.24. What is the probability of getting the flu?

Since getting the flu is the complement of not getting the flu, then $P(\text{getting flu}) = 1 - P(\text{flu}) = 1 - 0.24 = 0.76$

- c. In an experiment of picking a card from a deck, what is the probability of not getting a card that is a Queen?

You could do this problem by listing all the ways to not get a queen, but that set is fairly large. One advantage of the complement is that it reduces the workload. You use the complement in many situations to make the work shorter and easier. In this case it is easier to list all the ways to get a Queen, find the probability of the Queen, and then subtract from one.

Queen = {QS, QC, QD, QH} so $P(\text{Queen}) = \frac{4}{52}$ and $P(\text{not Queen}) = 1 - P(\text{Queen}) = 1 - \frac{4}{52} = \frac{48}{52}$

The complement is useful when you are trying to find the probability of an event that involves the words at least or an event that involves the words at most. As an example of an at least event is suppose you want to find the probability of making at least \$50,000 when you graduate from college. That means you want the probability of your salary being greater than or equal to \$50,000. An example of an at most event is suppose you want to find the probability of rolling a die and getting at most a 4. That means that you want to get less than or equal to a 4 on the die. The reason to use the complement is that sometimes it is easier to find the probability of the complement and then subtract from 1.

5.2.7 Example: Using the Complement to Find Probabilities

- In an experiment of rolling a fair die one time, find the probability of rolling at most a 4 on the die.
- In an experiment of pulling a card from a fair deck, find the probability of pulling at least a 5 (ace is a high card in this example).

5.2.7.1 Solution

- In an experiment of rolling a fair die one time, find the probability of rolling at most a 4 on the die.

The sample space for this experiment is {1, 2, 3, 4, 5, 6}. You want the event of getting at most a 4, which is the same as thinking of getting 4 or less. The event space is {1, 2, 3, 4}. The probability is $P(\text{at most a 4}) = \frac{4}{6}$

Or you could have used the complement. The complement of rolling at most a 4 would be rolling number bigger than 4. The event space for the complement is {5, 6}. The

probability of the complement is $P(\text{more than } 4) = \frac{2}{6}$. The probability of at most 4 would be $P(\text{at most } 4) = 1 - P(\text{more than } 4) = 1 - \frac{2}{6} = \frac{4}{6}$

Notice you have the same answer, but the event space was easier to write out. For this example the complement probability wasn't that useful, but in the future there will be events where it is much easier to use the complement.

- b. In an experiment of pulling a card from a fair deck, find the probability of pulling at least a 5 (ace is a high card in this example).

The sample space for this experiment is $SS = \{2S, 3S, 4S, 5S, 6S, 7S, 8S, 9S, 10S, JS, QS, KS, AS, 2C, 3C, 4C, 5C, 6C, 7C, 8C, 9C, 10C, JC, QC, KC, AC, 2D, 3D, 4D, 5D, 6D, 7D, 8D, 9D, 10D, JD, QD, KD, AD, 2H, 3H, 4H, 5H, 6H, 7H, 8H, 9H, 10H, JH, QH, KH, AH\}$

Pulling a card that is at least a 5 would involve listing all of the cards that are a 5 or more. It would be much easier to list the outcomes that make up the complement. The complement of at least a 5 is less than a 5. That would be the event of 4 or less. The event space for the complement would be $\{2S, 3S, 4S, 2C, 3C, 4C, 2D, 3D, 4D, 2H, 3H, 4H\}$. The probability of the complement would be $\frac{12}{52}$. The probability of at least a 5 would be $P(\text{at least } 5) = 1 - P(\text{at most } 4) = 1 - \frac{12}{52} = \frac{40}{52}$

Another concept was shown in Example: Calculating Theoretical Probabilities 2 parts g and i. The problems were looking for the probability of one event or another. In part g, it was looking for the probability of getting a Spade or an Ace. That was equal to $\frac{16}{52}$. In part i, it was looking for the probability of getting a Jack or an Ace. That was equal to $\frac{8}{52}$. If you look back at the parts b, c, and d, you might notice the following result: $P(\text{Jack or Ace}) = P(\text{Jack}) + P(\text{Ace})$ but $P(\text{Spade or Ace}) \neq P(\text{Spade}) + P(\text{Ace})$.

Why does adding two individual probabilities together work in one situation to give the probability of one or another event and not give the correct probability in the other?

The reason this is true in the case of the Jack and the Ace is that these two events cannot happen together. There is no overlap between the two events, and in fact the $P(\text{Jack and Ace}) = 0$. However, in the case of the Spade and Ace, they can happen together. There is overlap, mainly the ace of spades. The $P(\text{Spade and Ace}) \neq 0$.

When two events cannot happen at the same time, they are called **mutually exclusive**. In the above situation, the events Jack and Ace are mutually exclusive, while the events Spade and Ace are not mutually exclusive.

5.2.8 Addition Rules:

If two events A and B are mutually exclusive, then $P(A \text{ and } B) = 0$ and $P(A \text{ or } B) = P(A) + P(B)$

If two events A and B are not mutually exclusive, then $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

5.2.9 Example: Using Addition Rules

Suppose your experiment is to roll two fair dice.

- What is the sample space?
- What is the probability of getting a sum of 5?
- What is the probability of getting the first die a 2?
- What is the probability of getting a sum of 7?
- What is the probability of getting a sum of 5 and the first die a 2?
- What is the probability of getting a sum of 5 or the first die a 2?
- What is the probability of getting a sum of 5 and sum of 7?
- What is the probability of getting a sum of 5 or sum of 7?

5.2.9.1 Solution

- What is the sample space?

As with the other examples you need to come up with a sample space that has equally likely outcomes. One sample space is to list the sums possible on each roll. That sample space would look like: $SS = \{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$. However, there are more ways to get a sum of 7 than there are to get a sum of 2, so these outcomes are not equally likely. Another thought is to list the possibilities on each roll. As an example you could roll the dice and on the first die you could get a 1. The other die could be any number between 1 and 6, but say it is a 1 also. Then this outcome would look like (1,1). Similarly, you could get (1, 2), (1, 3), (1,4), (1, 5), or (1, 6). Also, you could get a 2, 3, 4, 5, or 6 on the first die instead. Putting this all together, you get the sample space:

$SS = \{(1,1), (1,2), (1,3), (1,4), (1,5), (1,6),$
 $(2,1), (2,2), (2,3), (2,4), (2,5), (2,6),$
 $(3,1), (3,2), (3,3), (3,4), (3,5), (3,6),$
 $(4,1), (4,2), (4,3), (4,4), (4,5), (4,6),$
 $(5,1), (5,2), (5,3), (5,4), (5,5), (5,6),$
 $(6,1), (6,2), (6,3), (6,4), (6,5), (6,6)\}$

Notice that a (2,3) is different from a (3,2), since the order that you roll the die is important and you can tell the difference between these two outcomes. You don't need any of the doubles twice, since these are not distinguishable from each other in either order.

This will always be the sample space for rolling two dice.

- b. What is the probability of getting a sum of 5?

Getting a sum of 5 = $\{(4,1), (3,2), (2,3), (1,4)\}$ so $P(\text{sum of 5}) = \frac{4}{36}$

- c. What is the probability of getting the first die a 2?

Getting first die a 2 = $\{(2,1), (2,2), (2,3), (2,4), (2,5), (2,6)\}$ so $P(\text{1st die 2}) = \frac{6}{36}$

- d. What is the probability of getting a sum of 7?

Getting a sum of 7 = $\{(6,1), (5,2), (4,3), (3,4), (2,5), (1,6)\}$ so $P(\text{sum of 7}) = \frac{6}{36}$

- e. What is the probability of getting a sum of 5 and the first die a 2?

This is events A and B which contains the outcome $\{(2,3)\}$ so $P(\text{sum of 5 and 1st die a 2}) = \frac{1}{36}$

- f. What is the probability of getting a sum of 5 or the first die a 2?

Notice from part e, that these two events are not mutually exclusive, so

$$P(\text{sum of 5 or 1st die a 2})$$

$$= P(\text{sum of 5}) + P(\text{1st die 2}) - P(\text{sum of 5 and 1st die a 2})$$

$$= \frac{4}{36} + \frac{6}{36} - \frac{1}{36} = \frac{9}{36}$$

- g. What is the probability of getting a sum of 5 and sum of 7?

These are the events parts a and c, which have no outcomes in common. Thus sum of 5 and sum of 7 = $\{ \}$ so $P(\text{sum of 5 and sum of 7}) = 0$

- h. What is the probability of getting a sum of 5 or sum of 7?

From part g, these two events are mutually exclusive, so $P(\text{sum of 5 or sum of 7}) = P(\text{sum of 5}) + P(\text{sum of 7})$

$$= \frac{4}{36} + \frac{6}{36} = \frac{10}{36}$$

5.2.10 Odds

Many people like to talk about the odds of something happening or not happening. Mathematicians, statisticians, and scientists prefer to deal with probabilities since odds are difficult to work with, but gamblers prefer to work in odds for figuring out how much they are paid if they win.

The **actual odds against** event A occurring are the ratio $\frac{P(\text{not } A)}{P(A)}$, usually expressed in the form $a : b$ or a to b , where a and b are integers with no common factors.

The **actual odds in favor** event A occurring are the ratio $\frac{P(A)}{P(\text{not } A)}$, which is the reciprocal of the odds against. If the odds against event A are $a : b$, then the odds in favor event A are $b : a$.

The **payoff odds** against event A occurring are the ratio of the net profit (if you win) to the amount bet.

payoff odds against event $A = (\text{net profit}) : (\text{amount bet})$

5.2.11 Example: Odds Against and Payoff Odds

In the game of Craps, if a shooter has a come-out roll of a 7 or an 11, it is called a natural and the pass line wins. The payoff odds are given by a casino as 1 : 1.

- Find the probability of a natural.
- Find the actual odds for a natural.
- Find the actual odds against a natural.
- If the casino pays 1:1, how much profit does the casino make on a \ \$10 bet?

5.2.11.1 Solution

- Find the probability of a natural.

A natural is a 7 or 11. The sample space is

$$SS = \{(1,1), (1,2), (1,3), (1,4), (1,5), (1,6), (2,1), (2,2), (2,3), (2,4), (2,5), (2,6), (3,1), (3,2), (3,3), (3,4), (3,5), (3,6), (4,1), (4,2), (4,3), (4,4), (4,5), (4,6), (5,1), (5,2), (5,3), (5,4), (5,5), (5,6), (6,1), (6,2), (6,3), (6,4), (6,5), (6,6)\}$$

The event space is $\{(1,6), (2,5), (3,4), (4,3), (5,2), (6,1), (5,6), (6,5)\}$

So $P(7 \text{ or } 11) = \frac{8}{36}$

- Find the actual odds for a natural.

odds of natural = $\frac{P(7 \text{ or } 11)}{P(\text{not } 7 \text{ or } 11)} = \frac{\frac{8}{36}}{1 - \frac{8}{36}} = \frac{\frac{8}{36}}{\frac{28}{36}} = \frac{8}{28} = \frac{2}{7}$

- c. Find the actual odds against a natural.

$$\text{odds of against a natural} = \frac{P(\text{not 7 or 11})}{P(7 \text{ or } 11)} = \frac{1 - \frac{8}{36}}{\frac{8}{36}} = \frac{\frac{28}{36}}{\frac{8}{36}} = \frac{28}{8} = \frac{3.5}{1}$$

- d. If the casino pays 1:1, how much profit does the casino make on a \\$10 bet?

The actual odds are 3.5 to 1 while the payoff odds are 1 to 1. The casino pays you \\$10 for your \\$10 bet. If the casino paid you the actual odds, they would pay \\$35.00 on every \\$1 bet, and on \\$10, they pay $3.5 \times 10 = \$35$. Their profit is $35 - 10 = \$25$.

5.2.12 Homework for Theoretical Probability Section

1. In Homework for Empirical Probability Section, the probabilities of each color of M&Ms in a packet were found. Use that information to answer the following questions.
 - a. Find the probability of choosing a green or red M&M.
 - b. Find the probability of choosing a blue, red, or yellow M&M.
 - c. Find the probability of not choosing a brown M&M.
 - d. Find the probability of not choosing a green M&M.
2. In Homework for Empirical Probability Section, the probabilities for defects in eyeglasses manufactured by Eyeglassomatic were calculated. Use that information to find the following probabilities.
 - a. Find the probability of picking a lens that is scratched or flaked.
 - b. Find the probability of picking a lens that is the wrong PD or was lost in lab.
 - c. Find the probability of picking a lens that is not scratched.
 - d. Find the probability of picking a lens that is not the wrong shape.
3. An experiment is to flip a fair coin three times.
 - a. State the sample space.
 - b. Find the probability of getting exactly two heads. Make sure you state the event space.
 - c. Find the probability of getting at least two heads. Make sure you state the event space.
 - d. Find the probability of getting an odd number of heads. Make sure you state the event space.
 - e. Find the probability of getting all heads or all tails. Make sure you state the event space.
 - f. Find the probability of getting exactly two heads or exactly two tails.
 - g. Find the probability of not getting an odd number of heads.
4. An experiment is rolling a fair die and then flipping a fair coin.
 - a. State the sample space.
 - b. Find the probability of getting a head. Make sure you state the event space.
 - c. Find the probability of getting a 6. Make sure you state the event space.

- d. Find the probability of getting a 6 or a head.
 - e. Find the probability of getting a 3 and a tail.
5. An experiment is rolling two fair dice.
 - a. State the sample space.
 - b. Find the probability of getting a sum of 3. Make sure you state the event space.
 - c. Find the probability of getting the first die is a 4. Make sure you state the event space.
 - d. Find the probability of getting a sum of 8. Make sure you state the event space.
 - e. Find the probability of getting a sum of 3 or sum of 8.
 - f. Find the probability of getting a sum of 3 or the first die is a 4.
 - g. Find the probability of getting a sum of 8 or the first die is a 4.
 - h. Find the probability of not getting a sum of 8.
 6. An experiment is pulling one card from a fair deck.
 - a. State the sample space.
 - b. Find the probability of getting a Ten. Make sure you state the event space.
 - c. Find the probability of getting a Diamond. Make sure you state the event space.
 - d. Find the probability of getting a Club. Make sure you state the event space.
 - e. Find the probability of getting a Diamond or a Club.
 - f. Find the probability of getting a Ten or a Diamond.
 7. An experiment is pulling a ball from an urn that contains 3 blue balls and 5 red balls.
 - a. Find the probability of getting a red ball.
 - b. Find the probability of getting a blue ball.
 - c. Find the odds for getting a red ball.
 - d. Find the odds for getting a blue ball.
 8. In the game of roulette, there is a wheel with spaces marked 0 through 36 and a space marked 00.
 - a. Find the probability of winning if you pick the number 7 and it comes up on the wheel.
 - b. Find the odds against winning if you pick the number 7.
 - c. The casino will pay you \$20 for every dollar you bet if your number comes up. How much profit is the casino making on the bet?

5.3 Conditional Probability

Suppose you want to figure out if you should buy a new car. When you first go and look, you find two cars that you like the most. In your mind they are equal, and so each has a 50% chance that you will pick it. Then you start to look at the reviews of the cars and realize that the first car has had 40% of them needing to be repaired in the first year, while the second

car only has 10% of the cars needing to be repaired in the first year. You could use this information to help you decide which car you want to actually purchase. Both cars no longer have a 50% chance of being the car you choose. You could actually calculate the probability you will buy each car, which is a conditional probability. You probably wouldn't do this, but it gives you an example of what a conditional probability is.

Conditional probabilities are probabilities calculated after information is given. This is where you want to find the probability of event A happening after you know that event B has happened. If you know that B has happened, then you don't need to consider the rest of the sample space. You only need the outcomes that make up event B . Event B becomes the new sample space, which is called the **restricted sample space, R** . If you always write a restricted sample space when doing conditional probabilities and use this as your sample space, you will have no trouble with conditional probabilities. The notation for conditional probabilities is $P(A, \text{ given } B) = P(A|B)$. The event following the vertical line is always the restricted sample space.

5.3.1 Example: Conditional Probabilities

- Suppose you roll two dice. What is the probability of getting a sum of 5, given that the first die is a 2?
- Suppose you roll two dice. What is the probability of getting a sum of 7, given the first die is a 4?
- Suppose you roll two dice. What is the probability of getting the second die a 2, given the sum is a 9?
- Suppose you pick a card from a deck. What is the probability of getting a Spade, given that the card is a Jack?
- Suppose you pick a card from a deck. What is the probability of getting an Ace, given the card is a Queen?

5.3.1.1 Solution

- Suppose you roll two dice. What is the probability of getting a sum of 5, given that the first die is a 2?

Since you know that the first die is a 2, then this is your restricted sample space, $R = \{(2,1), (2,2), (2,3), (2,4), (2,5), (2,6)\}$. Out of this restricted sample space, the way to get a sum of 5 is $\{(2,3)\}$. Thus

$$P(\text{sum of 5, given first die a 2}) = P(\text{sum of 5}|\text{first die 2}) = \frac{1}{6}$$

- Suppose you roll two dice. What is the probability of getting a sum of 7, given the first die is a 4?

Since you know that the first die is a 4, this is your restricted sample space, $R = \{(4,1), (4,2), (4,3), (4,4), (4,5), (4,6)\}$ Out of this restricted sample space, the way to get a sum of 7 is $\{(4,3)\}$. Thus

$$P(\text{sum of 7, given first die a 4}) = P(\text{sum of 7}|\text{first die 4}) = \frac{1}{6}$$

- c. Suppose you roll two dice. What is the probability of getting the second die a 2, given the sum is a 9?

Since you know the sum is a 9, this is your restricted sample space, $R = \{(3,6), (4,5), (5,4), (6,3)\}$. Out of this restricted sample space there is no way to get the second die a 2. Thus

$$P(\text{first die a 2, given sum is 9}) = P(\text{1st die a 2}|\text{sum of 9}) = \frac{0}{4}$$

- d. Suppose you pick a card from a deck. What is the probability of getting a Spade, given that the card is a Jack?

Since you know that the card is a Jack, this is your restricted sample space, $R = \{JS, JC, JD, JH\}$. Out of this restricted sample space, the way to get a Spade is $\{JS\}$. Thus

$$P(\text{spade, given card a Jack}) = P(\text{spade}|\text{Jack}) = \frac{1}{4}$$

- e. Suppose you pick a card from a deck. What is the probability of getting an Ace, given the card is a Queen?

Since you know that the card is a Queen, then this is your restricted sample space, $R = \{QS, QC, QD, QH\}$ Out of this restricted sample space, there is no way to get an Ace, thus

$$P(\text{Ace, given Queen}) = P(\text{Ace}|\text{Queen}) = \frac{0}{4}$$

If you look at the results of Example: Calculating Theoretical Probabilities 2 part d and Example: Calculating Theoretical Probabilities part b, you will notice that you get the same answer. This means that knowing that the first die is a 4 did not change the probability that the sum is a 7. This added knowledge did not help you in any way. It is as if that information was not given at all. However, if you compare example Example: Calculating Theoretical Probabilities 2 part b and Example: Calculating Theoretical Probabilities part a, you will notice that they are not the same answer. In this case, knowing that the first die is a 2 did change the probability of getting a sum of 5. In the first case, the events sum of 7 and first die is a 4 are called **independent events**. In the second case, the events sum of 5 and first die is a 2 are called **dependent events**.

Events A and B are considered **independent events** if the fact that one event happens does not change the probability of the other event happening. In other words, events A and B are independent if the fact that B has happened does not affect the probability of event A happening and the fact that A has happened does not affect the probability of event B happening. Otherwise, the two events are dependent.

In symbols, A and B are independent if $P(A|B) = P(A)$ or $P(B|A) = P(B)$

5.3.2 Example: Independent Events

- Suppose you roll two dice. Are the events “sum of 7” and “first die is a 3” independent?
- Suppose you roll two dice. Are the events “sum of 6” and “first die is a 4” independent?
- Suppose you pick a card from a deck. Are the events “Jack” and “Spade” independent?
- Suppose you pick a card from a deck. Are the events “Heart” and “Red” card independent?
- Suppose you have two children via separate births. Are the events “the first is a boy” and “the second is a girl” independent?
- Suppose you flip a coin 50 times and get a head every time, what is the probability of getting a head on the next flip?

5.3.2.1 Solution

- Suppose you roll two dice. Are the events “sum of 7” and “first die is a 3” independent?

To determine if they are independent, you need to see if $P(\text{sum of 7} | \text{first die a 3}) = P(\text{sum of 7})$ or the other way around. It doesn't matter which order these are calculated in, so pick whichever is easier. $\text{sum of 7} = \{(1,6), (2,5), (3,4), (4,3), (5,2), (6,1)\}$ $\text{first die is a 3} = \{(3,1), (3,2), (3,3), (3,4), (3,5), (3,6)\}$ $P(\text{sum of 7} | \text{first die a 3})$ means that you assume that first die is a 3 has happened. The restricted sample space is the first die is a 3, $R = \{(3,1), (3,2), (3,3), (3,4), (3,5), (3,6)\}$ In this restricted sample space, the way for a sum of 7 to happen is $\{(3,4)\}$, so $P(\text{sum of 7} | \text{first die a 3}) = \frac{1}{6}$ The $P(\text{sum of 7}) = \frac{6}{36} = \frac{1}{6}$. Since $P(\text{sum of 7} | \text{first die a 3}) = P(\text{sum of 7})$, the “sum of 7” and “first die is a 3” are independent events.

- Suppose you roll two dice. Are the events “sum of 6” and “first die is a 4” independent?

To determine if they are independent, you need to see if

$$P(\text{sum of 6} | \text{first die a 4}) = P(\text{sum of 6}).$$

Again it doesn't matter what order you do this in. Do which is easier. $\text{sum of 6} = \{(1,5), (2,4), (3,3), (4,2), (5,1)\}$ and $\text{first die is a 4} = \{(4,1), (4,2), (4,3), (4,4), (4,5), (4,6)\}$, if want $P(\text{sum of 6} | \text{first die a 4})$, the restricted sample space is 1st die is a 4, $R = \{(4,1), (4,2), (4,3), (4,4), (4,5), (4,6)\}$ In this restricted sample space, the way to get a sum of 6 is $\{(4,2)\}$, so $P(\text{sum of 6} | \text{first die a 4}) = \frac{1}{6}$. The $P(\text{sum of 6}) = \frac{5}{36}$ Notice $P(\text{sum of 6} | \text{first die a 4}) \neq P(\text{sum of 6})$, Thus “sum of 6” and “first die is a 4” are dependent.

- c. Suppose you pick a card from a deck. Are the events “Jack” and “Spade” independent?

To determine if they are independent, you need to see if

$$P(\text{Jack}|\text{Spade}) = P(\text{Jack}).$$

Remember, you can do this the other order if you wish. Jack = {JS, JC, JD, JH} and R = Spade {2S, 3S, 4S, 5S, 6S, 7S, 8S, 9S, 10S, JS, QS, KS, AS}. For $P(\text{Jack}|\text{Spade})$, the restricted sample space is Spade, $R = \{2S, 3S, 4S, 5S, 6S, 7S, 8S, 9S, 10S, JS, QS, KS, AS\}$. In this restricted sample space, the way to get a Jack is {JS}, so $P(\text{Jack}|\text{Spade}) = \frac{1}{13}$. The $P(\text{Jack}) = \frac{4}{52} = \frac{1}{13}$. Since $P(\text{Jack}|\text{Spade}) = P(\text{Jack})$, “Jack” and “Spade” are independent.

- d. Suppose you pick a card from a deck. Are the events “Heart” and “Red” card independent?

To determine if they are independent, you need to see if

$$P(\text{Heart}|\text{Red}) = P(\text{Heart}).$$

Heart = {2H, 3H, 4H, 5H, 6H, 7H, 8H, 9H, 10H, JH, QH, KH, AH} and Red card = {2D, 3D, 4D, 5D, 6D, 7D, 8D, 9D, 10D, JD, QD, KD, AD, 2H, 3H, 4H, 5H, 6H, 7H, 8H, 9H, 10H, JH, QH, KH, AH}. The restricted sample space is, red card, $R = \{2D, 3D, 4D, 5D, 6D, 7D, 8D, 9D, 10D, JD, QD, KD, AD, 2H, 3H, 4H, 5H, 6H, 7H, 8H, 9H, 10H, JH, QH, KH, AH\}$. In this restricted sample space, the way to get a heart is 13, and

$$P(\text{Heart}|\text{Red}) = \frac{13}{26}.$$

$$P(\text{Heart}) = \frac{13}{52}$$

Note $P(\text{Heart}|\text{Red}) \neq P(\text{Heart})$, so, “Heart” and “Red” card are dependent.

- e. Suppose you have two children via separate births. Are the events “the first is a boy” and “the second is a girl” independent?

In this case, you actually don’t need to do any calculations. The sex of one child does not affect the sex of the second child. The events are independent.

- f. Suppose you flip a coin 50 times and get a head every time, what is the probability of getting a head on the next flip?

Since one flip of the coin does not affect the next flip (the coin does not remember what it did the time before), the probability of getting a head on the next flip is still one-half.

5.3.3 Multiplication Rule:

Two more useful formulas:

If two events are dependent, then $P(A \text{ and } B) = P(A) * P(B|A)$

If two events are independent, then $P(A \text{ and } B) = P(A) * P(B)$

These two formulas are useful if the sample space is too large to write out, but if the sample space isn't too large, it is better to find probabilities of and statements using the sample space techniques.

If you solve the first equation for $P(B|A)$, you obtain $P(B|A) = \frac{P(A \text{ and } B)}{P(A)}$, which is a formula to calculate a conditional probability. However, it is easier to find a conditional probability by using the restricted sample space and counting unless the sample space is large.

5.3.4 Example: Multiplication Rule

- a. Suppose you pick three cards from a deck, what is the probability that they are all Queens if the cards are not replaced after they are picked?
- b. Suppose you pick three cards from a deck, what is the probability that they are all Queens if the cards are replaced after they are picked and before the next card is picked?

5.3.4.1 Solution

- a. Suppose you pick three cards from a deck, what is the probability that they are all Queens if the cards are not replaced after they are picked?

This sample space is too large to write out, so using the multiplication rule makes sense. Since the cards are not replaced, then the probability will change for the second and third cards. They are dependent events. This means that on the second draw there is one less Queen and one less card, and on the third draw there are two less Queens and 2 less cards.

$$\begin{aligned} P(3 \text{ Queens}) &= P(\text{Queen on 1st}) * P(\text{Queen on 2nd, given Queen on first}) * P(\text{Queen on third, given Queen on 1st and 2nd}) \\ &= \frac{4}{52} * \frac{3}{51} * \frac{2}{50} \end{aligned}$$

- b. Suppose you pick three cards from a deck, what is the probability that they are all Queens if the cards are replaced after they are picked and before the next card is picked?

Again, the sample space is too large to write out, so using the multiplication rule makes sense. Since the cards are put back, one draw has no affect on the next draw and they are all independent. $P(3 \text{ Queens}) = P(\text{Queen on 1st}) * P(\text{Queen on 2nd}) * P(\text{Queen on 3rd})$

$$= \frac{4}{52} * \frac{4}{52} * \frac{4}{52}$$

5.3.5 Example: Application Problem

A project conducted by the Australian Federal Office of Road Safety asked people many questions about their cars. One question was the reason that a person chooses a given car, and that data is in Table ?? (Car Preferences, 2019).

Code book for Data Frame `Car_pref` is below Table ??.

A contingency table, is a cross tabulation of the data into different categories. As an example, a contingency table of if the person has children under 5 and their current car is below, the following command in R Studio can be used to create this table.

```
tally(~Kids5+ActCar, data=Car_pref, margins=TRUE)
```

	ActCar			
Kids5	large	medium	small	Total
no	87	90	94	271
yes	13	10	6	29
Total	100	100	100	300

- Find the probability that a person questioned has kids under 5.
- Find the probability that a person questioned actually has a large car.
- Find the probability that a person questioned actually has a large car and has children under 5.
- Find the probability that a person questioned has a large car given that they have children under 5.
- Find the probability that a person has a large car or has children under 5.
- Find the probability that a person questioned has children under 5 given that they have a large car.
- Are the events that a person questioned has a “large car” and “kids under 5” independent events? Why or why not?

5.3.5.1 Solution

- Find the probability that a person questioned has kids under 5.

First, you need to find the number of people questioned. Add the first row, there are 150 people who did not have kids under 5. Adding the second row, there are 29 people who do have kids under 5. Also, you can add the columns. There are 100 people who have large cars, 100 people who have medium cars, and 100 who have small cars. Adding either the row totals or the columns totals, gives you 300 people in the study. Out of the 300 people, 29 people had kids under 5. So the

$$P(\text{kids under 5}) = \frac{29}{300}.$$

So 9.7% of the people questioned had children under 5.

- b. Find the probability that a person questioned actually has a large car.

There are 100 people with large cars out of 300 people. So,

$$P(\text{large car}) = \frac{100}{300}. \text{ There are 33\% of the people who have a large car.}$$

- c. Find the probability that a person questioned actually has a large car and has children under 5.

There are 13 people who have a large car and have children under 5, so the $P(\text{large car and children under 5}) = \frac{13}{300} = 0.043$. 4.3% of all people surveyed have a large car and children under 5.

- d. Find the probability that a person questioned has a large car given that they have children under 5.

In this case you know that the person had children under 5. You don't need to consider the people who don't. You only need to look at the row with people who have have children under 5. In that row, look to see how many people have a large car. There are 13 people with a large car out of the 29 people with kids under 5. So, $P(\text{large car}|\text{kids under 5}) = \frac{13}{29} = 0.45$

There is 45% chance that a person with a large car have children under 5.

- e. Find the probability that a person has a large car or has children under 5.

This problem can be done two ways. One is to use the addition formula, but a better way is to realize that there are 29 people who have kids under 5, and there are 100 people who have a large car. That is 129 people. But the 13 people who have large cars and kids under 5 were just counted twice. So subtract the 13 people from the 129. That give 116 people who have either kids under 5 or a large car. So

$$P(\text{large car or kids under 5}) = \frac{116}{300} = 0.39.$$

That means 39% of the people questioned has a large car or has children under 5.

- f. Find the probability that a person questioned has children under 5 given that they have a large car.

In this case you know that the person has a large car. You don't need to include the people who have medium or small cars. You only need to consider the column headed by large. In that column, there are 100 people who have large cars and out of those 100, 13 have children under 5. So, $P(\text{kids under 5}|\text{large}) = \frac{13}{100} = 0.13$. Thus 13% of people have children under 5 given that they have a large car.

- g. Are the events that a person questioned has a “large car” and “kids under 5” independent events? Why or why not?

In order for these events to be independent, either $P(\text{kids under 5}|\text{large car}) = P(\text{kids under 5})$ or $P(\text{large car}|\text{kids under 5}) = P(\text{large car})$ have to be true. Part (d) showed $P(\text{kids under 5}|\text{large car}) = 0.44$ and part (b) showed $P(\text{kids under 5}) = 0.33$. Since these are not equal, then these two events are dependent.

A big deal has been made about the difference between dependent and independent events while calculating the probability of “and” compound events. You must multiply the probability of the first event with the conditional probability of the second event.

Why do you care? Calculating probabilities when performing sampling is important, as this will be seen later. But here is a simplification that can make the calculations a lot easier: when the sample size is very small compared to the population size, you can assume that the conditional probabilities just don’t change very much over the sample.

For example, consider acceptance sampling. Suppose there is a big population of parts delivered to your factory, say 12,000 parts. Suppose there are 85 defective parts in the population. You decide to randomly select ten parts, and see if you should reject the shipment. What is the probability of rejecting the shipment?

There are many different ways you could reject the shipment. For example, maybe the first three parts are good, one is bad, and the rest are good. Or all ten parts could be bad, or maybe the first five. So many ways to reject! But there is only **one** way that you’d accept the shipment: if **all ten** parts are good. That would happen if the first part is good, **and** the second part is good, **and** the third part is good, and so on. Since the probability of the second part being good is (slightly) dependent on whether the first part was good, technically you should take this into consideration when you calculate the probability that all ten are good.

The probability of getting the first sampled part good is $\frac{1200-85}{1200} = \frac{1115}{1200}$. So the probability that all ten being good is

$$\frac{1115}{1200} * \frac{1114}{1200} * \frac{1113}{1200} * \dots * \frac{1106}{1200} = 0.931357.$$

If instead you assume that the probability doesn’t change much, you get $(\frac{1115}{1200})^{10} = 0.931382$. So as you can see, there is not much difference. So here is the rule: if the sample is very small compared to the size of the population, then you can assume that the probabilities are independent, even though they aren’t technically. By the way, the probability of rejecting the shipment is $1 - 0.9314 = 0.0686$.

5.3.6 Homework for Conditional Probability Section

1. Are owning a refrigerator and owning a car independent events? Why or why not?

2. Are owning a computer, tablet, or smart phone and paying for Internet service independent events? Why or why not?
3. Are passing your statistics class and passing your biology class independent events? Why or why not?
4. Are owning a bike and owning a car independent events? Why or why not?
5. An experiment is picking a card from a fair deck.
 - a. What is the probability of picking a Jack given that the card is a face card?
 - b. What is the probability of picking a heart given that the card is a three?
 - c. What is the probability of picking a red card given that the card is an ace?
 - d. Are the events Jack and face card independent events? Why or why not?
 - e. Are the events red card and ace independent events? Why or why not?
6. An experiment is rolling two dice.
 - a. What is the probability that the sum is 6 given that the first die is a 5?
 - b. What is the probability that the first die is a 3 given that the sum is 11?
 - c. What is the probability that the sum is 7 given that the first die is a 2?
 - d. Are the two events sum of 6 and first die is a 5 independent events? Why or why not?
 - e. Are the two events sum of 7 and first die is a 2 independent events? Why or why not?
7. You flip a coin four times. What is the probability that all four of them are heads?
8. You flip a coin six times. What is the probability that all six of them are heads?
9. You pick three cards from a deck with replacing the card each time before picking the next card. What is the probability that all three cards are kings?
10. You pick three cards from a deck without replacing a card before picking the next card. What is the probability that all three cards are kings?
11. A project conducted by the Australian Federal Office of Road Safety asked people many questions about their cars. One question was the reason that a person chooses a given car, and that data is in Table ?? (Car Preferences, 2019).

Code book for Data Frame Car__pref is below Table ??.

The contingency table for the sex of a person and the size car the person prefers is in table below

```
tally(~PreferCar+Sex, data=Car_pref, margins=TRUE)
```


	Sex		
PreferCar	female	male	Total
4	6	17	23
large	26	47	73
medium	75	61	136
small	43	25	68
Total	150	150	300

- What is the probability that a person questioned was female?
 - What is the probability that a person questioned prefers a medium car?
 - What is the probability that a person questioned prefers a medium car given that the person was female?
 - What is the probability that a person questioned was a female and prefers a medium car?
 - What is the probability that a person questioned was a female or prefers a medium car?
 - Are the events person questioned is a female and person questioned prefers a medium car mutually exclusive? Why or why not?
 - Are the events person questioned is a female and person questioned prefers a medium car independent? Why or why not?
12. Researchers watched groups of dolphins off the coast of Ireland in 1998 to determine what activities the dolphins partake in at certain times of the day (Activities of Dolphin Groups, 2019). The numbers in table \#4.3.5 represent the number of groups of dolphins that were partaking in an activity at certain times of days.

```
Dolphin<- read.csv( "https://krkozak.github.io/MAT160/dolphins.csv")
knitr::kable(head(Dolphin))
```

Table 5.4: Dolphin Activity

activity	period
Travel	Morning
Travel	Morning
Travel	Morning
Travel	Morning
Travel	Morning
Travel	Morning

Code book for Data Frame Dolphin

Description Groups of dolphins were observed off the coast of Iceland near Keflavik in 1998. The data here give the time of the day and the main activity of the group, whether travelling

quickly, feeding or socializing. The dolphin groups varied in size - usually feeding or socializing groups were larger than travelling groups.

Usage Dolphin

Format

This data frame contains the following columns:

Activity: Main activity of group: travelling (Travel), feeding (Feed) or socializing (Social)

Period: Time of the day: Morning, Noon, Afternoon or Evening

Source Activities of Dolphin Groups. (n.d.). Retrieved July 12, 2019, from <http://www.statsci.org/data/general/>

References Marianne Rasmussen, Department of Biology, University of Southern Denmark, Odense, Denmark.

- a. What is the probability that a dolphin group is partaking in travel?
- b. What is the probability that a dolphin group is around in the morning?
- c. What is the probability that a dolphin group is partaking in travel given that it is morning?
- d. What is the probability that a dolphin group is around in the morning given that it is partaking in socializing?
- e. What is the probability that a dolphin group is around in the afternoon given that it is partaking in feeding?
- f. What is the probability that a dolphin group is around in the afternoon and is partaking in feeding?
- g. What is the probability that a dolphin group is around in the afternoon or is partaking in feeding?
- h. Are the events dolphin group around in the afternoon and dolphin group feeding mutually exclusive events? Why or why not?
- i. Are the events dolphin group around in the morning and dolphin group partaking in travel independent events? Why or why not?

5.4 Counting Techniques

There are times when the sample space or event space are very large, that it isn't feasible to write it out. In that case, it helps to have mathematical tools for counting the size of the sample space and event space. These tools are known as counting techniques.

5.4.1 Multiplication Rule in Counting Techniques

If task 1 can be done ways, task 2 can be done ways, and so forth to task n being done ways. Then the number of ways to do task 1, 2,..., n together would be $m_1 * m_2 * \cdots * m_n$.

5.4.2 Example: Multiplication Rule in Counting

A menu offers a choice of 3 salads, 8 main dishes, and 5 desserts. How many different meals consisting of one salad, one main dish, and one dessert are possible?

5.4.2.1 Solution

There are three tasks, picking a salad, a main dish, and a dessert. The salad task can be done 3 ways, the main dish task can be done 8 ways, and the dessert task can be done 5 ways. The ways to pick a salad, main dish, and dessert are $3 * 8 * 5 = 120$.

5.4.3 Example: Multiplication Rule in Counting

How many three letter “words” can be made from the letters a, b, and c with no letters repeating? A “word” is just an ordered group of letters. It doesn’t have to be a real word in a dictionary.

5.4.3.1 Solution

There are three tasks that must be done in this case. The tasks are to pick the first letter, then the second letter, and then the third letter. The first task can be done 3 ways since there are 3 letters. The second task can be done 2 ways, since the first task took one of the letters. The third task can be done 1 way, since the first and second task took two of the letters. There are $3 * 2 * 1 = 6$

In Example: Multiplication Rule, the solution was found by find $3 * 2 * 1$. Many counting problems involve multiplying a list of decreasing numbers. This is called a **factorial**. There is a special symbol for this.

5.4.4 Factorial

$$n! = n(n-1)(n-2) * \dots * 2 * 1$$

As an example: $5! = 5 * 4 * 3 * 2 * 1 = 120$ $8! = 8 * 7 * 6 * 5 * 4 * 3 * 2 * 1 = 40320$

0 factorial is defined to be $0! = 1$ and **1 factorial** is defined to be $1! = 1$. In rStudio, the command for factorial is `factorial(number)`. As an example $7!$ using r Studio would be

```
factorial(7)
```

```
[1] 5040
```

Sometimes you are trying to select r objects from n total objects. The number of ways to do this depends on if the order you choose the r objects matters or if it doesn't. As an example if you are trying to call a person on the phone, you have to have their number in the right order. Otherwise, you call someone you didn't mean to. In this case, the order of the numbers matters. If however you were picking random numbers for the lottery, it doesn't matter which number you pick first. As long as you have the same numbers that the lottery people pick, you win. In this case the order doesn't matter. A **permutation** is an arrangement of items with a specific order. You use permutations to count items when the order matters. When the order doesn't matter you use combinations. A **combination** is an arrangement of items when order is not important. When you do a counting problem, the first thing you should ask yourself is "does order matter?"

5.4.5 Permutation Formula

Picking r objects from n total objects when order matters

$${}_nP_r = \frac{n!}{(n-r)!}$$

5.4.6 Combination Formula

Picking r objects from n total objects when order doesn't matter

$${}_nC_r = \frac{n!}{r!(n-r)!}$$

Most calculators have a factorial button on them, and many have the combination and permutation functions also.

5.4.7 Homework for Counting Techniques Sections

1. You are going to a benefit dinner, and need to decide before the dinner what you want for salad, main dish, and dessert. You have 2 different salads to choose from, 3 main dishes, and 5 desserts. How many different meals are available?
2. How many different phone numbers are possible in the area code 928?
3. You are opening a T-shirt store. You can have long sleeves or short sleeves, three different colors, five different designs, and four different sizes. How many different shirts can you make?
4. The California license plate has one number followed by three letters followed by three numbers. How many different license plates are there?
5. Find ${}_9P_4$
6. Find ${}_{10}P_6$
7. Find ${}_{10}C_5$
8. Find ${}_{20}P_4$
9. You have a group of twelve people. You need to pick a president, treasurer, and secretary from the twelve. How many different ways can you do this?
10. A baseball team has a 25-person roster. A batting order has nine people. How many different batting orders are there?
11. An urn contains five red balls, seven yellow balls, and eight white balls. How many different ways can you pick two red balls?
12. How many ways can you choose seven people from a group of twenty?

6 Discrete Probability Distribution

When computing probabilities, the sample space, which contains all the outcomes of the experiment, is listed. If the probabilities for all of the outcomes are also listed then these two together are called a probability distribution. With a probability distribution, the shape can be determined, the mean and standard deviation can be calculated, and the probability of events can be found. How to find all of these concepts depends on what type of quantitative variables are being considered. Remember there are different types of quantitative variables, called discrete or continuous. What is the difference between discrete and continuous data? **Discrete** data can only take on particular values in a range. **Continuous** data can take on any value in a range. Discrete data usually arises from counting while continuous data usually arises from measuring.

If you have a variable, and can find a probability associated with that variable, it is called a **random variable**. In many cases the random variable is what you are measuring, but when it comes to discrete random variables, it is usually what you are counting. So for the example of how tall is a plant given a new fertilizer, the random variable is the height of the plant given a new fertilizer. For the example of how many fleas are on prairie dogs in a colony, the random variable is the number of fleas on a prairie dog in a colony.

6.0.1 Examples of each:

How tall is a plant given a new fertilizer? Continuous. This is something you measure.

How many fleas are on prairie dogs in a colony? Discrete. This is something you count.

Now suppose you put all the values of the random variable together with the probability that the random variable would occur. You could then have a distribution like before, but now it is called a probability distribution since it involves probabilities. A **probability distribution** is an assignment of probabilities to the values of the random variable.

With the idea of a probability distribution, the next thing is to look at the basics of a probability distribution.

6.1 Basics of Probability Distributions

As a reminder, a variable or what will be called the random variable from now on, is represented by the letter x and it represents a quantitative (numerical) variable that is measured or observed in an experiment.

As with probabilities, probability distributions, have the properties, $0 \leq P(outcome) \leq 1$ and $\sum P(outcomes) = 1$

6.1.1 Example: Probability Distribution

The 2010 U.S. Census found the chance of a household being a certain size. The data is in Table ?? (“Households by age,” 2013). Note, the category 7 is really 7 or more people in the household. Draw the probability distribution and find the mean, variance, and standard deviation.

Table 6.1: Household Size from U.S. Census of 2010

size	prob
1	0.267
2	0.336
3	0.158
4	0.137
5	0.063
6	0.024
7	0.015

6.1.1.1 Solution

In this case, the random variable is $x = \text{size of household}$. This is a discrete random variable, since you are counting the number of people in a household.

It is a probability distribution since you have the x value and the probabilities that go with it, all of the probabilities are between zero and one, and the sum of all of the probabilities is one.

You can give a probability distribution in table form (as in Table ??) or as a graph. The graph looks like a histogram. To graph the histogram, use the following commands and process in rStudio.

First you need to load a few packages using the following commands. These packages are “arm” and “Weighted.Desc.Stat”. If these packages have not been installed, they need to be

installed before you can load them using library. Once you have installed them, they will always be available in /r Studio to be loaded. To load a package, use the command

```
library("name of package")
```

In this case the packages you need are arm and Weighted.Desc.Stat.

```
Loading required package: MASS
```

```
Loading required package: Matrix
```

```
Loading required package: lme4
```

```
arm (Version 1.14-4, built: 2024-4-1)
```

```
Working directory is /Users/mori/CSU Fullerton Dropbox/Mortaza Jamshidian/Statistics Text Us
```

To draw the probability distribution, use the following command. First you need to create variables for x , size, and the probability, $prob$, in r Studio. Then you can draw the distribution.

(ref:discrete-histogram-cap) Histogram of Size of Family

```
discrete.histogram(Household$size,Household$prob, bar.width = 1, main="Size of family", xlab=
```

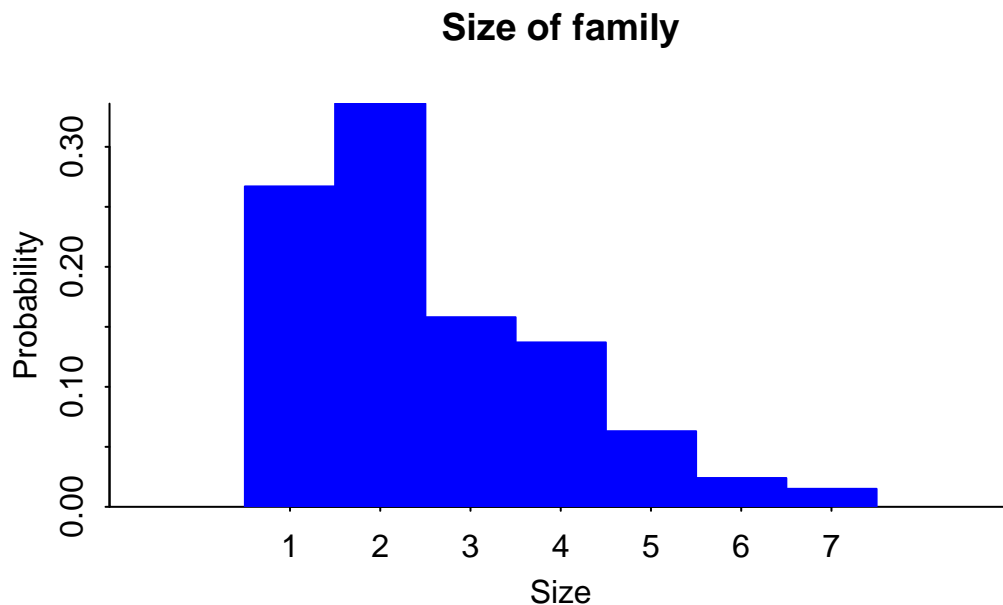


Figure 6.1: Histogram of Household Size from U.S. Census of 2010

This command is different than the commands used in the past, but is needed for discrete probability distributions. So putting a title on the graph uses the command `main="title you want"` instead of `title=` as before.

Notice this graph Figure ?? is skewed right, which means that most families have around 2 people in them and larger families become more and more rare.

To find the mean, variance, and standard deviation using r Studio, make sure that the package `Weighted.Desc.Stat` is loaded, then use the following commands.

```
w.mean(Household$size, Household$prob)
```

```
[1] 2.525
```

```
w.var(Household$size, Household$prob)
```

```
[1] 2.023375
```

```
w.sd(Household$size, Household$prob)
```

```
[1] 1.422454
```

The mean is 2.525 people, the variance is 2.02 *people*², and the standard deviation is 1.42 people.

When calculating the mean and standard deviation of a probability distribution, you can consider the population distribution the population even though it was most likely created from a large sample. Since a probability distribution is basically a population, the mean and standard deviation that are calculated are actually the population parameters and not the sample statistics. The notation used is the same as the notation for population mean, μ , and population standard deviation, σ , that was used in chapter 3. Note: the mean can also be thought of as the expected value. It is the value you expect to get if the trials were repeated infinite number of times. The mean or expected value does not need to be a whole number, even if the possible values of x are whole numbers. This means one can find what value they can expect to get in the long run for gambling or insurance including extended warranties using the mean of a probability distribution. First one needs to figure out the probability distribution, and then follow the process in example 5.1.1.

6.1.2 Example: Calculating the Expected Value

In the Arizona lottery game called Pick 3, a player pays \$1 and then picks a three-digit number. If those three numbers are picked in that specific order the person wins \$500. What is the expected value in this game?

6.1.2.1 Solution

To find the expected value, you need to first create the probability distribution. In this case, the random variable x = winnings. If you pick the right numbers in the right order, then you win \$500, but you paid \$1 to play, so you actually win \$499. If you didn't pick the right numbers, you lose the \$1, the x value is -\$1. You also need the probability of winning and losing. Since you are picking a three-digit number, and for each digit there are 10 numbers you can pick with each independent of the others, you can use the multiplication rule. To win, you have to pick the right numbers in the right order. The first digit, you pick 1 number out of 10, the second digit you pick 1 number out of 10, and the third digit you pick 1 number out of 10. The probability of picking the right number in the right order is $\frac{1}{1000}$. The probability of losing (not winning) would be

$$1 - \frac{1}{1000} = \frac{999}{1000}.$$

Putting this information into a table will help to organize the information and find the expected value.

Table 6.2: Probability Distribution of Lottery

outcome	amount	probability
win	499	0.001
lose	-1	0.999

Now type the values into R using the following command:

Now to find the expected value, it is the same as finding the mean, though the command is a little different since you don't have a data frame for this data.

```
weighted.mean(amount, probability)
```

```
[1] -0.5
```

The expected value (or mean) is -0.5. That is -\$0.50. Since it is negative, that means you lose \$0.50 every time you play the Pick 3. It seems you would be better off putting the \$1 every week into a savings account then playing the Pick 3 lottery.

The reason probability is studied in statistics is to help in making decisions in inferential statistics. To understand how that is done the concept of a rare event is needed.

6.1.3 Rare Event Rule for Inferential Statistics

If, under a given assumption, the probability of a particular observed event is extremely small, then you can conclude that the assumption is probably not correct.

An example of this is suppose you roll an assumed fair die 1000 times and get a six 600 times, when you should have only rolled a six around 160 times, then you should believe that your assumption about it being a fair die is untrue.

6.1.4 Determining if an event is unusual

If you are looking at a value of x for a discrete variable, and the P (the variable has a value of x or more) is less than 0.05, then you can consider the x an unusually high value. Another way to think of this is if the probability of getting such a high value is less than 0.05, then the event of getting the value x is unusual.

Similarly, if the P (the variable has a value of x or less) is less than 0.05, then you can consider this an unusually low value. Another way to think of this is if the probability of getting a value as small as x is less than 0.05, then the event x is considered unusual.

Why is it “ x or more” or “ x or less” instead of just “ x ” when you are determining if an event is unusual? Consider this example: you and your friend go out to lunch every day. Instead of Going Dutch (each paying for their own lunch), you decide to flip a coin, and the loser pays for both. Your friend seems to be winning more often than you’d expect, so you want to determine if this is unusual before you decide to change how you pay for lunch (or accuse your friend of cheating). The process for how to calculate these probabilities will be presented in the next section on the binomial distribution. If your friend won 6 out of 10 lunches, the probability of that happening turns out to be about 20.5%, not unusual. The probability of winning 6 or more is about 37.7%. But what happens if your friend won 501 out of 1,000 lunches? That doesn’t seem so unlikely! The probability of winning 501 or more lunches is about 47.8%, and that is consistent with your hunch that this isn’t so unusual. But the probability of winning exactly 501 lunches is much less, only about 2.5%. That is why the probability of getting exactly that value is not the right question to ask: you should ask the probability of getting that value or more (or that value or less on the other side).

The value 0.05 will be explained later, and it is not the only value you can use for unusual events.

6.1.5 Example: Is the Event Unusual

The 2010 U.S. Census found the chance of a household being a certain size. The data is in the table (“Households by age,” 2013).

The 2010 U.S. Census found the chance of a household being a certain size. The data is in Table ?? (“Households by age,” 2013). Note, the category 7 is really 7 or more people in the household.

State random variable:

Solution

State random variable

- a. Is it unusual for a household to have six people in the family?

size = number of people in a household

- a. Is it unusual for a household to have six people in the family?
- b. If you did come upon many families that had six people in the family, what would you think?
- c. Is it unusual for a household to have four people in the family?
- d. If you did come upon a family that has four people in it, what would you think?

6.1.5.1 Solution

To determine this, you need to look at probabilities. However, you cannot just look at the probability of six people. You need to look at the probability of x being six or less people or the probability of x being six or more people. The

$$P(x \leq 6) = P(1) + P(2) + P(3) + P(4) + P(5) + P(6) = 0.267 + 0.336 + 0.158 + 0.137 + 0.063 + 0.024 = 0.985$$

Since this probability is more than 5%, then six is not an unusually low value.

$$\text{The } P(x \geq 6) = P(6) + P(7) = 0.024 + 0.015 = 0.039$$

Since this probability is less than 5%, then six is an unusually high value. It is unusual for a household to have six people in the family.

- b. If you did come upon many families that had six people in the family, what would you think?

Since it is unusual for a family to have six people in it, then you may think that either the size of families is increasing from what it was or that you are in a location where families are larger than in other locations.

- c. Is it unusual for a household to have four people in the family?

To determine this, you need to look at probabilities. Again, look at the probability of x being four or less or the probability of x being four or more. The

$$P(x \leq 4) = P(0) + P(1) + P(2) + P(3) + P(4) = 0.267 + 0.336 + 0.158 + 0.137 = 0.898$$

Since this probability is more than 5%, four is not an unusually low value.

The

$$P(\geq 4) = P(4) + P(5) + P(6) + P(7) = 0.137 + 0.063 + 0.024 + 0.015 = 0.239$$

Since this probability is more than 5%, four is not an unusually low value. Thus, four is not an unusual size of a family.

- d. If you did come upon a family that has four people in it, what would you think?

Since it is not unusual for a family to have four members, then you would not think anything is amiss.

6.1.6 Homework for Basics of Probability Distributions Section

1. Eyeglassomatic manufactures eyeglasses for different retailers. The number of days it takes to fix defects in an eyeglass and the probability that it will take that number of days are in Table ??.

```
Days<- read.csv(
  "https://krkozak.github.io/MAT160/table_5_1_3.csv")
knitr::kable(Days)
```

Table 6.3: Nuumber of Days to fix Eyeglasses

days	prob
1	0.249
2	0.108
3	0.091
4	0.123
5	0.133
6	0.114
7	0.070
8	0.046
9	0.019
10	0.013
11	0.010
12	0.008

Table 6.3: Nuumber of Days to fix Eyeglasses

days	prob
13	0.006
14	0.004
15	0.002
16	0.002
17	0.001
18	0.001

- State the random variable.
 - Draw a histogram of the number of days to fix defects
 - Find the mean number of days to fix defects.
 - Find the variance for the number of days to fix defects.
 - Find the standard deviation for the number of days to fix defects.
 - Find probability that a lens will take at least 16 days to make a fix the defect.
 - Is it unusual for a lens to take 16 days to fix a defect?
 - If it does take 16 days for eyeglasses to be repaired, what would you think?
- Suppose you have an experiment where you flip a coin three times. You then count the number of heads.
 - State the random variable.
 - Write the probability distribution for the number of heads.
 - Draw a histogram for the number of heads.
 - Find the mean number of heads.
 - Find the variance for the number of heads.
 - Find the standard deviation for the number of heads.
 - Find the probability of having two or more number of heads.
 - Is it unusual for to flip two heads?
 - The Ohio lottery has a game called Pick 4 where a player pays \\$1 and picks a four-digit number. If the four numbers come up in the order you picked, then you win \\$2,500. What is your expected value?
 - An LG Dishwasher, which costs \\$800, has a 20% chance of needing to be replaced in the first 2 years of purchase. A two-year extended warranty costs \\$112.10 on a dishwasher. What is the expected value of the extended warranty assuming it is replaced in the first 2 years?

6.2 Binomial Probability Distribution

Section 5.1 introduced the concept of a probability distribution. The focus of the section was on discrete probability distributions. To find the probability distribution for a situation, you usually needed to actually conduct the experiment and collect data. Then you can calculate the experimental probabilities. Normally you cannot calculate the theoretical probabilities. However, there are certain types of experiment that allow you to calculate the theoretical probability. One of those types is called a **Binomial Experiment**.

Properties of a **binomial experiment** (or Bernoulli trial):

1. Fixed number of trials, n , which means that the experiment is repeated a specific number of times.
2. The n trials are independent, which means that what happens on one trial does not influence the outcomes of other trials.
3. There are only two outcomes, which are called a success and a failure.
4. The probability of a success doesn't change from trial to trial, where p = probability of success and $q = 1 - p$ = probability of failure.

If you know you have a binomial experiment, then you can calculate binomial probabilities. This is important because binomial probabilities come up often in real life. Examples of binomial experiments are:

Toss a fair coin ten times, and find the probability of getting two heads.

Question twenty people in class, and look for the probability of more than half being women?

Shoot five arrows at a target, and find the probability of hitting it five times?

6.2.1 Formula for the probabilities for a Binomial experiment

First, the random variable in a binomial experiment is x = number of successes.

Be careful, a success is not always a good thing. Sometimes a success is something that is bad, like finding a defect. A success just means you observed the outcome you wanted to see happen.

Binomial Formula for the probability of r successes in n trials is $P(X = r) = {}_nC_r * p^r * q^{n-r}$

where ${}_nC_r$ is the number of combinations of n things taking r at a time. It is read “ n choose r ”.

When solving problems, make sure you define your random variable and state what n, p , and r are. Without doing this, the problems are a great deal harder.

The command to find a binomial probability in r Studio is

$P(X = r) =$

`dbinom(r, n, p)`

$P(x \leq r) =$

`pbinom(r, n, p, lower.tail=TRUE)`

$P(x \geq r) =$

`pbinom(r-1, n, p, lower.tail = FALSE)`

6.2.2 Example: Calculating Binomial Probabilities

When looking at a person's eye color, it turns out that 1% of people in the world has green eyes ("What percentage of," 2013). Consider a group of 20 people.

- State the random variable.
- Argue that this is a binomial experiment
- Find the probability that none of the 20 people have green eyes.
- Find the probability that nine have green eyes.
- Find the probability that at most three have green eyes.
- Find the probability that at most two have green eyes.
- Find the probability that at least four have green eyes.
- In Europe, four people out of twenty have green eyes. Is this unusual? What does that tell you?

6.2.2.1 Solution

- State the random variable.

x = number of people with green eyes

- Argue that this is a binomial experiment.
 - There are 20 people, and each person is a trial, so there are a fixed number of trials. In this case, $n = 20$.
 - If you assume that each person in the group is chosen at random the eye color of one person doesn't affect the eye color of the next person, thus the trials are independent.
 - Either a person has green eyes or they do not have green eyes, so there are only two outcomes. In this case, the success is a person has green eyes.

4. The probability of a person having green eyes is 0.01. This is the same for every trial since each person has the same chance of having green eyes.
- c. Find the probability that none of the 20 people have green eyes.

If none have green eyes, then $r = 0$.

Probability that none have green eyes is $P(X = 0) = 0.818$, using the command:

```
dbinom(0,20,0.01)
```

```
[1] 0.8179069
```

- d. Find the probability that nine have green eyes.

If nine have green eyes, then $r = 9$.

Probability that 9 have green eyes is

$P(X = 9) = 1.50 \times 10^{-13}$. Notice that r gives the answer as 1.50391e-13. This is the way many computer programs write a number in scientific notation. It isn't possible for a computer to write it as 1.50381×10^{-13} , but it is possible for humans to write it correctly. So make sure the answer is written in the correct scientific notation.

```
dbinom(9,20,0.01)
```

```
[1] 1.50381e-13
```

- e. Find the probability that at most three have green eyes.

At most three means that three is the highest value you will have. Find the probability of x is less than or equal to three.

Since this is less than, then the lower tail of the probability distribution is being used, so $P(X \leq 3) = 0.99996$ using the command in r Studio of

```
pbinom(3,20,0.01, lower.tail=TRUE)
```

```
[1] 0.9999574
```

The reason the answer is written to more decimal places is because when it is rounded to three decimal places the rounding makes the answer 1. But 1 means that the event will happen, when in reality there is a slight chance that it won't happen. It is best to write the answer to more decimal places or it can be written as > 0.999 to represent that the number is very close to 1, but isn't 1.

- f. Find the probability that at most two have green eyes.

At most 2 means 2 or less. So find the probability that there are less than or equal to 2. $P(X \leq 2) = 0.999$, and again, this is the lower tail of the probability distribution, so use `lower.tail=TRUE` in the `r` command:

```
pbinom(2,20,0.01, lower.tail=TRUE)
```

```
[1] 0.9989964
```

- g. Find the probability that at least four have green eyes.

At least four means four or more. Find the probability of x being greater than or equal to four. Since it is greater than or equal to, this is the right tail of the probability distribution. However, if you just use `lower.tail=FALSE`, then the 4 is not included in `r` calculations. You want all numbers from 4 on up, so you need to use

$r = 4 - 1 = 3$ in the `r` command. This will include 4 in the calculation. $P(X \geq 4) = 4.26 \times 10^{-5}$

```
pbinom(4-1,20,0.01, lower.tail=FALSE)
```

```
[1] 4.262093e-05
```

- h. In Europe, four people out of twenty have green eyes. Is this unusual? What does that tell you?

Since the probability of finding four or more people with green eyes is much less than 0.05, it is unusual to find four people out of twenty with green eyes. That should make you wonder if the proportion of people in Europe with green eyes is more than the 1% for the general population. If this is true, then you may want to ask why Europeans have a higher proportion of green-eyed people. That of course could lead to more questions.

6.2.3 Example: Calculating Binomial Probabilities

According to the Center for Disease Control (CDC), about 1 in 88 children in the U.S. have been diagnosed with autism (“CDC-data and statistics,” 2013). Suppose you consider a group of 10 children.

- State the random variable.
- Argue that this is a binomial experiment
- Find the probability that none have autism.
- Find the probability that seven have autism.

- e. Find the probability that at least five have autism.
- f. Find the probability that at most two have autism.
- g. Suppose five children out of ten have autism. Is this unusual? What does that tell you?

6.2.3.1 Solution

- a. State the random variable.

x = number of children with autism.

- b. Argue that this is a binomial experiment

1. There are 10 children, and each child is a trial, so there are a fixed number of trials. In this case, $n = 10$.
2. If you assume that each child in the group is chosen at random, then whether a child has autism does not affect the chance that the next child has autism. Thus the trials are independent.
3. Either a child has autism or they do not have autism, so there are two outcomes. In this case, the success is a child has autism.
4. The probability of a child having autism is $\frac{1}{88}$. This is the same for every trial since each child has the same chance of having autism.

- c. Find the probability that none have autism.

$$P(X = 0) = 0.892$$

```
dbinom(0,10, 1/88)
```

```
[1] 0.892002
```

- d. Find the probability that seven have autism.

$$P(X = 7) = 2.84 \times 10^{-12}$$

```
dbinom(7,10, 1/88)
```

```
[1] 2.837346e-12
```

- e. Find the probability that at least five have autism.

$P(X \geq 5) = 4.553 \times 10^{-8}$. Again, this is the upper tail of the probability distribution, so use `lower.tail=FALSE` and

$r = 5 - 1 = 4$ to make sure that `r` calculates for 5 and on up.

```
pbinom(5-1, 10, 1/88, lower.tail=FALSE)
```

```
[1] 4.553416e-08
```

f. Find the probability that at most two have autism.

$P(X \leq 2) = 0.9998$. This is using the lower tail of the probability distribution.

```
pbinom(2, 10, 1/88, lower.tail=TRUE)
```

```
[1] 0.9998341
```

g. Suppose five children out of ten have autism. Is this unusual? What does that tell you?

Since the probability of five or more children in a group of ten having autism is much less than 5%, it is unusual to happen. If this does happen, then one may think that the proportion of children diagnosed with autism is actually more than $\frac{1}{88}$.

6.2.4 Homework for Binomial Probability Distribution Section

1. Approximately 10% of all people are left-handed (“11 little-known facts,” 2013). Consider a grouping of fifteen people.
 - a. State the random variable.
 - b. Argue that this is a binomial experiment
 - c. Find the probability that none are left-handed.
 - d. Find the probability that seven are left-handed.
 - e. Find the probability that at least two are left-handed.
 - f. Find the probability that at most three are left-handed.
 - g. Find the probability that at least seven are left-handed.
 - h. Seven of the last 15 U.S. Presidents were left-handed. Is this unusual? What does that tell you?
2. According to an article in the American Heart Association’s publication **Circulation**, 24% of patients who had been hospitalized for an acute myocardial infarction did not fill their cardiac medication by the seventh day of being discharged (Ho, Bryson & Rumsfeld, 2009). Suppose there are twelve people who have been hospitalized for an acute myocardial infarction.

- a. State the random variable.
 - b. Argue that this is a binomial experiment
 - c. Find the probability that all filled their cardiac medication.
 - d. Find the probability that seven did not fill their cardiac medication.
 - e. Find the probability that none filled their cardiac medication.
 - f. Find the probability that at most two did not fill their cardiac medication.
 - g. Find the probability that at least three did not fill their cardiac medication.
 - h. Find the probability that at least ten did not fill their cardiac medication.
 - i. Suppose of the next twelve patients discharged, ten did not fill their cardiac medication, would this be unusual? What does this tell you?
3. Eyeglassomatic manufactures eyeglasses for different retailers. In March 2010, they tested to see how many defective lenses they made, and there were 16.9% defective lenses due to scratches. Suppose Eyeglassomatic examined twenty eyeglasses.
 - a. State the random variable.
 - b. Argue that this is a binomial experiment
 - c. Find the probability that none are scratched.
 - d. Find the probability that all are scratched.
 - e. Find the probability that at least three are scratched.
 - f. Find the probability that at most five are scratched.
 - g. Find the probability that at least ten are scratched.
 - h. Is it unusual for ten lenses to be scratched? If it turns out that ten lenses out of twenty are scratched, what might that tell you about the manufacturing process?
 4. The proportion of brown M&M's in a milk chocolate packet is approximately 14% (Madison, 2013). Suppose a package of M&M's typically contains 52 M&M's.
 - a. State the random variable.
 - b. Argue that this is a binomial experiment
 - c. Find the probability that six M&M's are brown.
 - d. Find the probability that twenty-five M&M's are brown.
 - e. Find the probability that all of the M&M's are brown.
 - f. Would it be unusual for a package to have only brown M&M's? If this were to happen, what would you think is the reason?

6.3 Mean and Standard Deviation of Binomial Distribution

If you list all possible values of x in a Binomial distribution, you get the **Binomial Probability Distribution**. You can draw a histogram of the probability distribution and find the mean (expected value), variance, and standard deviation of it. To have r Studio calculate the binomial values and save them to a variable, use the command

```
x<-c(0:n) p<-dbinom(0:n, n, p)
```

6.3.1 Example: Finding the Probability Distribution, Mean, Variance and Standard Deviation of a Binomial Distribution

When looking at a person's eye color, it turns out that 1% of people in the world has green eyes ("What percentage of," 2013). Consider a group of 20 people.

- State the random variable.
- Write the probability distribution.
- Draw a histogram.
- Find the mean, variance, and standard deviation.

6.3.1.1 Solution

- State the random variable.

x = number of people who have green eyes

- Write the probability distribution.

In this case you need to write each value of x and its corresponding probability. It is easiest to do this by using the `r` Command:

```
green<-c(0:20)
probability_green<-dbinom(0:20,20, 0.01)
```

It looks like nothing happened, but `r` save the values as variables. To see what is in each of those values, type

```
green
```

```
[1] 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
```

```
probability_green
```

```
[1] 8.179069e-01 1.652337e-01 1.585576e-02 9.609552e-04 4.125313e-05
[6] 1.333434e-06 3.367259e-08 6.802543e-10 1.116579e-11 1.503810e-13
[11] 1.670900e-15 1.534344e-17 1.162381e-19 7.225371e-22 3.649177e-24
[16] 1.474415e-26 4.654088e-29 1.106141e-31 1.862190e-34 1.980000e-37
[21] 1.000000e-40
```

These can now be typed into a table if desired.

c. Draw a histogram.

On `r`, this is like what was done in Section 5.1. Makes sure that the packages “`arm`” and “`Weighted.Desc.Stat`” are loaded. Then perform the command to get:

```
discrete.histogram(green, probability_green, bar.width = 1, main="Number of People with Green Eyes")
```

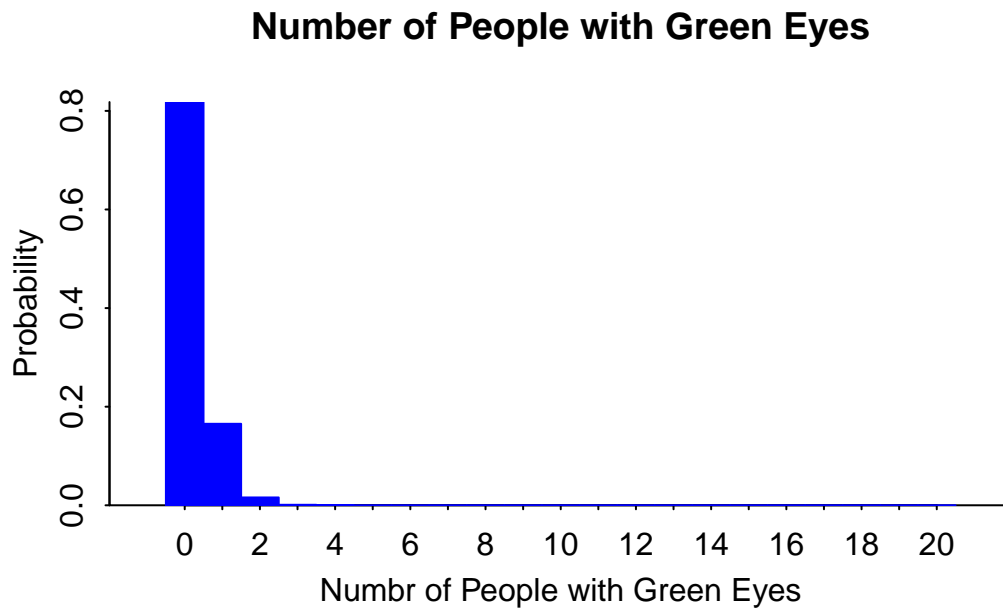


Figure 6.2: Histogram of Number of People with Green Eyes

Notice this graph Figure ?? is skewed right.

d. Find the mean, variance, and standard deviation

Using `r` Studio command such as those in Section 5.1:

```
w.mean(green, probability_green)
```

```
[1] 0.2
```

```
w.var(green, probability_green)
```

```
[1] 0.198
```

```
w.sd(green, probability_green)
```

```
[1] 0.4449719
```

You expect on average that out of 20 people, less than 1 person would have green eyes, with a variance of 0.198 people^2 and a standard deviation of 0.44 people.

6.3.2 Homework for Mean and Standard Deviation of Binomial Distribution Section

1. Suppose a random variable, x , arises from a binomial experiment. Suppose $n = 6$, and $p = 0.13$.
 - a. Write the probability distribution.
 - b. Draw a histogram.
 - c. Describe the shape of the histogram.
 - d. Find the mean.
 - e. Find the variance.
 - f. Find the standard deviation.
2. Suppose a random variable, x , arises from a binomial experiment. Suppose $n = 10$, and $p = 0.81$.
 - a. Write the probability distribution.
 - b. Draw a histogram.
 - c. Describe the shape of the histogram.
 - d. Find the mean.
 - e. Find the variance.
 - f. Find the standard deviation.
3. Suppose a random variable, x , arises from a binomial experiment. Suppose $n = 7$, and $p = 0.50$.
 - a. Write the probability distribution.
 - b. Draw a histogram.
 - c. Describe the shape of the histogram.
 - d. Find the mean.
 - e. Find the variance.
 - f. Find the standard deviation.
4. Approximately 10% of all people are left-handed. Consider a grouping of fifteen people.
 - a. State the random variable.

- b. Write the probability distribution.
 - c. Draw a histogram.
 - d. Describe the shape of the histogram.
 - e. Find the mean.
 - f. Find the variance.
 - g. Find the standard deviation.
5. According to an article in the American Heart Association's publication **Circulation**, 24% of patients who had been hospitalized for an acute myocardial infarction did not fill their cardiac medication by the seventh day of being discharged (Ho, Bryson & Rumsfeld, 2009). Suppose there are twelve people who have been hospitalized for an acute myocardial infarction.
- a. State the random variable.
 - b. Write the probability distribution.
 - c. Draw a histogram.
 - d. Describe the shape of the histogram.
 - e. Find the mean.
 - f. Find the variance.
 - g. Find the standard deviation.
6. Eyeglassomatic manufactures eyeglasses for different retailers. In March 2010, they tested to see how many defective lenses they made, and there were 16.9% defective lenses due to scratches. Suppose Eyeglassomatic examined twenty eyeglasses.
- a. State the random variable.
 - b. Write the probability distribution.
 - c. Draw a histogram.
 - d. Describe the shape of the histogram.
 - e. Find the mean.
 - f. Find the variance.
 - g. Find the standard deviation.
7. The proportion of brown M&M's in a milk chocolate packet is approximately 14% (Madison, 2013). Suppose a package of M&M's typically contains 52 M&M's.
- a. State the random variable.
 - b. Find the mean.
 - c. Find the variance.
 - d. Find the standard deviation.

7 Continuous Probability Distribution

Chapter 5 dealt with probability distributions arising from discrete random variables. Mostly that chapter focused on the binomial experiment. There are many other experiments from discrete random variables that exist but are not covered in this book.

Chapter 6 deals with probability distributions that arise from continuous random variables. The focus of this chapter is a distribution known as the normal distribution, though realize that there are many other distributions that exist. A few others are examined in future chapters.

Looking at the density plot of a quantitative variable, one can guess what the distribution of that variable is. As an example, consider the NHANES data frame. One variable to consider is Weight. The density plot of Weight is

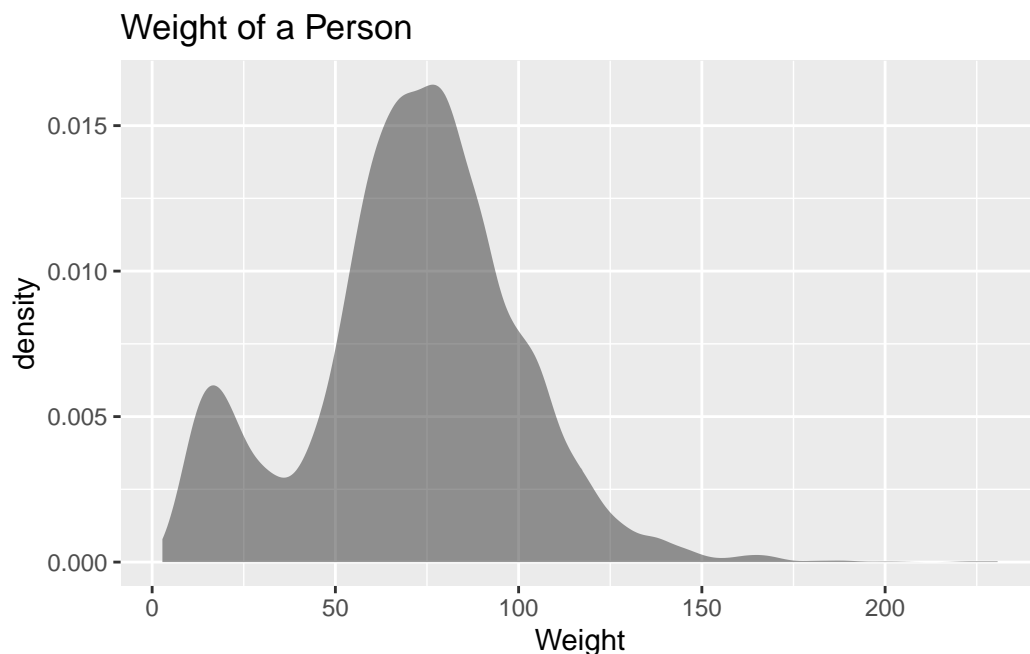


Figure 7.1: Density plot of the Weight of a Person

Figure ?? looks somewhat symmetric, and maybe bell shaped.

Consider, the variable head circumference (HeadCirc) in the NHANES data frame. The density plot for this variable is Figure ??

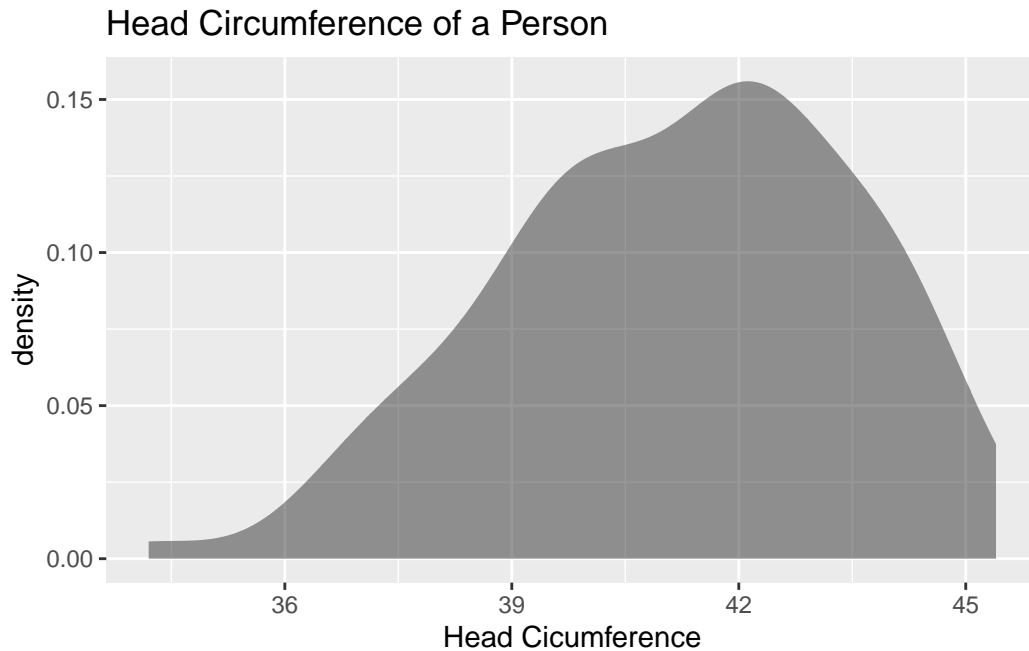


Figure 7.2: Density plot of head circumference of a person

Figure ?? looks somewhat skewed left.

Now consider the variable BMI from the NHANES data frame. The density plot is

This density plot Figure ?? appears to be skewed right

Now consider the variable SmokeAge. Its density plot is Figure ??

This distribution appears to be bimodal.

lastly, consider the variable Pulse. The density plot is Figure ??

This density plot appears to be symmetric and could almost be considered bell shaped.

The reason that one considers the density plots to understand the distribution of the population, is that in some cases the distribution can be approximated with a known distribution that has certain properties. There are many known distribution. Some examples are the Uniform distribution, the Chi-Squared distribution, the Student's T distribution, and the normal distribution. The normal distribution is one of the more common distributions to use as a model, and it will be explored in this chapter. But do realize that there are many other distributions that one can use.

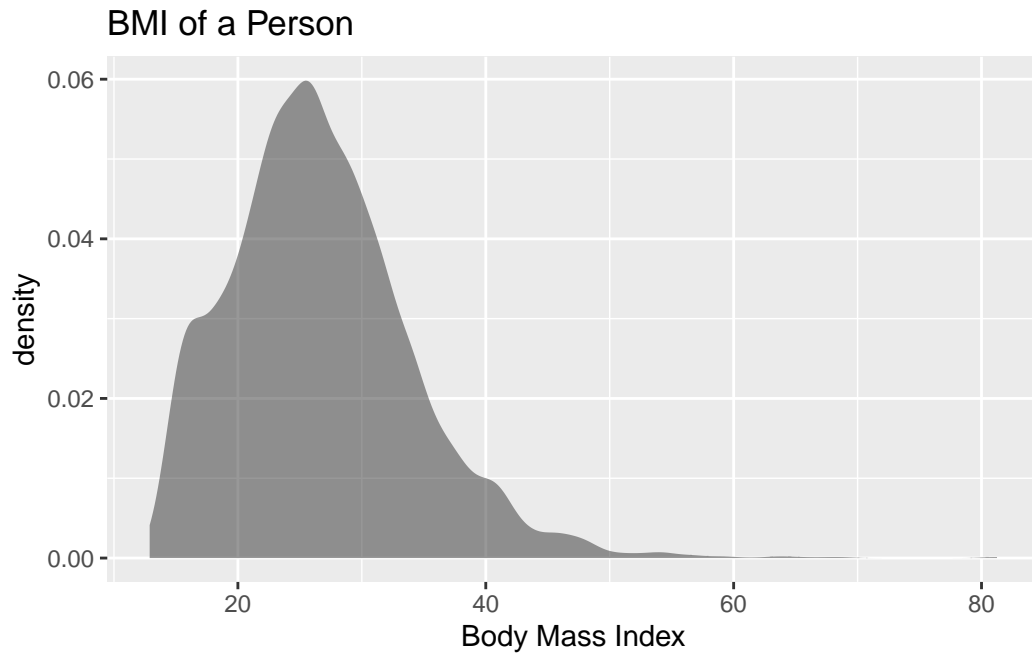


Figure 7.3: Density Plot of BMI of a Person

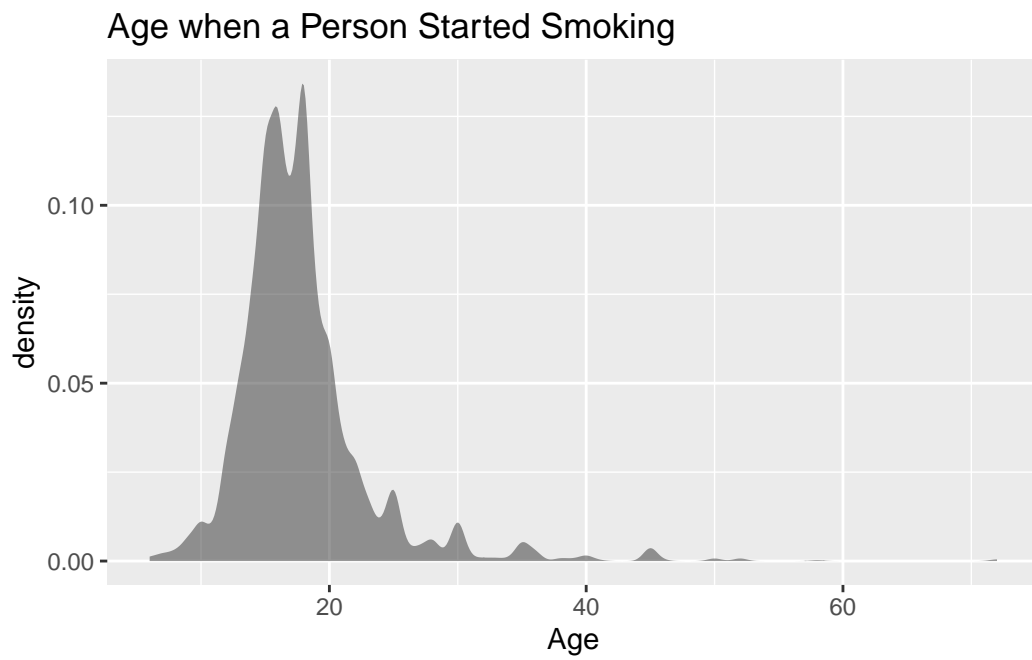


Figure 7.4: Density Plot of Age when Person Started Smoking

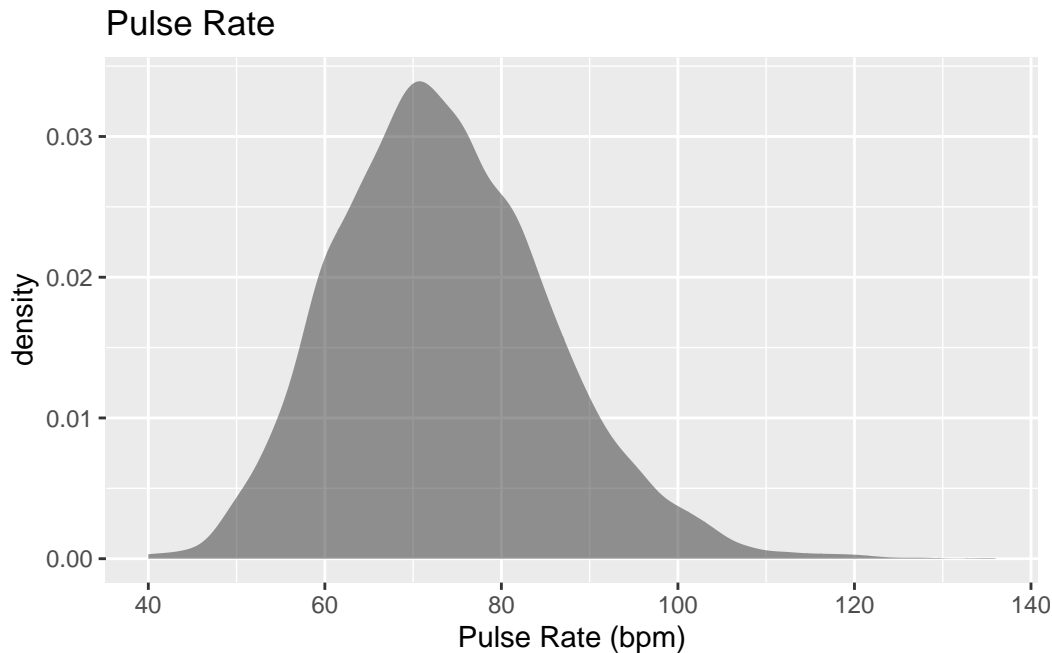


Figure 7.5: Density Plot of Pulse Rate of a person

7.1 Normal Distribution

Many populations have a distribution that is a symmetric, unimodal, and bell-shaped. For example: height, blood pressure, and cholesterol level. However, not every bell shaped curve is a normal curve. In a normal curve, there is a specific relationship between its “height” and its “width.” Normal curves can be tall and skinny or they can be short and fat. They are all symmetric, unimodal, and centered at μ , the population mean.

Figure ?? and Figure ?? show two different normal curves drawn on the same scale. Both have $\mu = 2$ but the one in Figure ?? has a standard deviation of 1 and the one in Figure ?? has a standard deviation of 4. Notice that the larger standard deviation makes the graph wider (more spread out) and shorter.

Every normal curve has common features.

- The center, or the highest point, is at the population mean, μ .
- The transition points are the places where the curve changes from a “hill” to a “valley”. The distance from the mean to the transition point is one standard deviation.
- The area under the whole curve is exactly 1. Therefore, the area under the half below or above the mean is 0.5.

Just as in a discrete probability distribution, the object is to find the probability of an event occurring. However, unlike in a discrete probability distribution where the event can be a

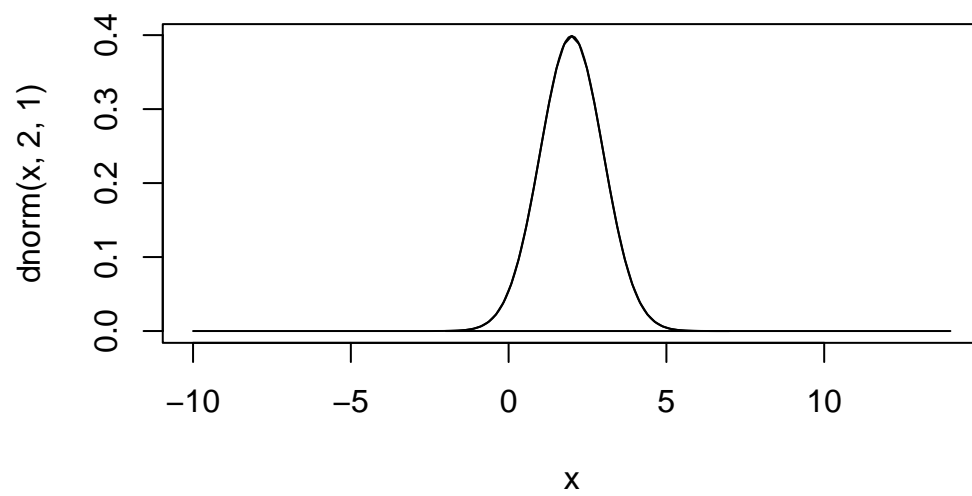


Figure 7.6: Normal curve with mean 2 and standard deviation 1

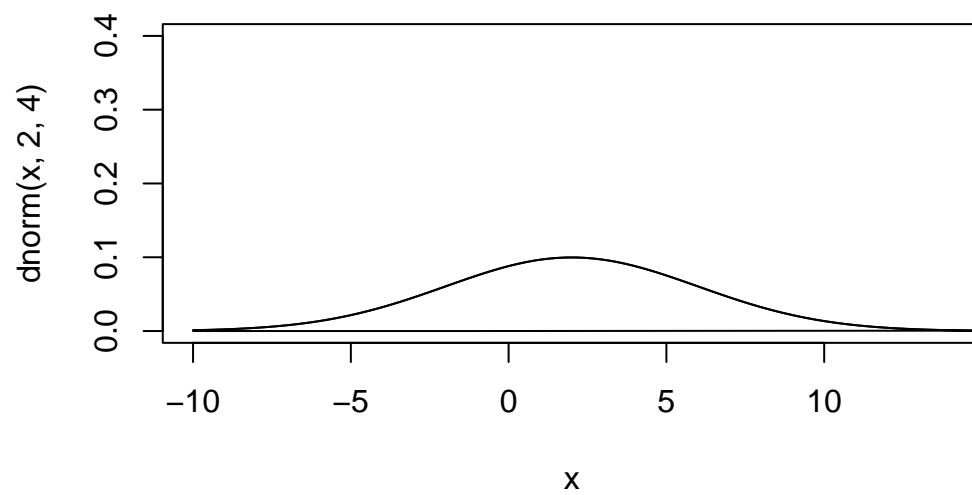


Figure 7.7: Normal curve with mean 2 and standard deviation 4

single value, in a continuous probability distribution the event must be a range. You are interested in finding the probability of x occurring in the range between a and b , or $P(a \leq x \leq b) = P(a < x < b)$. Calculus tells us this probability is the area under the curve in the interval from a to b .

Before looking at the process for finding the probabilities under the normal curve, it is somewhat useful to look at the **Empirical Rule** that gives approximate values for these areas. The Empirical Rule is just an approximation and it will only be used in this section to give you an idea of what the size of the probabilities is for different shadings. A more precise method for finding probabilities for the normal curve will be demonstrated in the next section. Please do not use the empirical rule except for real rough estimates.

The Empirical Rule for any normal distribution: Approximately 68% of the data is within one standard deviation of the mean. Approximately 95% of the data is within two standard deviations of the mean. Approximately 99.7% of the data is within three standard deviations of the mean.

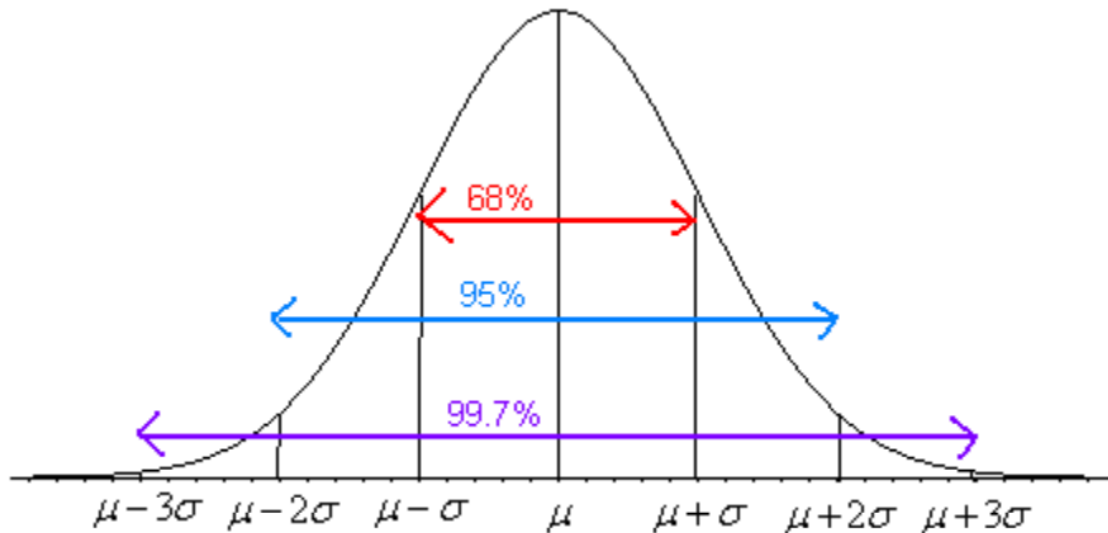


Figure 7.8: Empirical Rule Graph

Be careful, there is still some area left over in each end. Remember, the maximum a probability can be is 100%, so if you calculate you will see that for both ends together there is 0.3% of the curve. Because of symmetry, you can divide this equally between both ends and find that there is 0.15% in each tail beyond the 3rd standard deviations.

7.2 Finding Probabilities for the Normal Distribution

The Empirical Rule is just an approximation and only works for certain values. What if you want to find the probability for x values that are not integer multiples of the standard deviation? The probability is the area under the curve. To find areas under the curve, you need calculus. Before technology, you needed to convert every x value to a standardized number, called the z -score or z -value or simply just z . The z -score is a measure of how many standard deviations an x value is from the mean. To convert from a normally distributed x value to a z -score, you use the following formula.

$$z - score = \frac{x - \mu}{\sigma}$$

where μ = mean of the population of the x value and σ = standard deviation for the population of the x value

The z -score is normally distributed, with a mean of 0 and a standard deviation of 1. It is known as the standard normal curve. The z -score is a measure of how many standard deviations a data value is from its mean. If the z - score is positive, the data value is above the mean. If the z -score is negative, the data value is below the mean. The farther the z -value is from 0, the farther the data value is from the mean.

These days technology can find probabilities without converting to the z -score and looking the probabilities up in a table. There are many programs available that will calculate the probability for a normal curve. The command on r to find the area to the left $P(x < value)$ is

```
pnorm(value, mean, standard_deviation, lower.tail=TRUE)
```

The command on r to find the area to the right, $P(x > value)$ is

```
pnorm(value, mean, standard_deviation, lower.tail=FALSE)
```

7.2.1 Example: General Normal Distribution

The length of a human pregnancy is normally distributed with a mean of 272 days with a standard deviation of 9 days (Bhat & Kushtagi, 2006).

- State the random variable
- Find the probability of a pregnancy lasting more than 280 days.
- Find the probability of a pregnancy lasting less than 250 days.
- Find the probability that a pregnancy lasts between 265 and 280 days.
- Find the length of pregnancy that 10% of all pregnancies last less than.
- Suppose you meet a woman who says that she was pregnant for less than 250 days. Would this be unusual and what might you think?

7.2.1.1 Solution

a. State the random variable.

x = length of a human pregnancy

b. Find the probability of a pregnancy lasting more than 280 days.

First translate the statement into a mathematical statement. $P(x > 280)$

Now, draw a picture Figure ??.

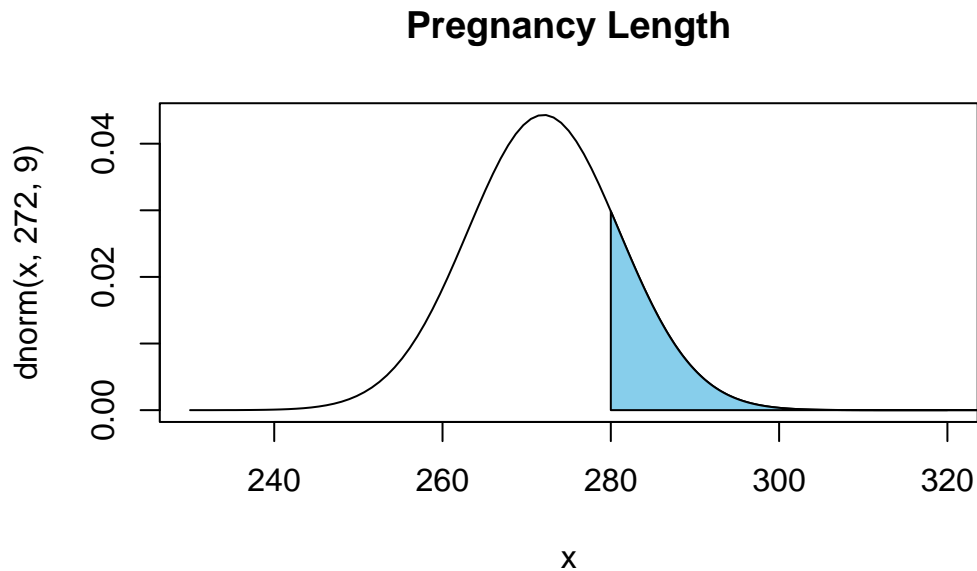


Figure 7.9: Normally distributed with mean 272 and standard deviation 9, and $P(x > 280)$

The probability of a pregnancy lasting longer than 280 days is $P(x > 280) = 0.187$. The command in rStudio is

```
pnorm(280, 272, 9, lower.tail=FALSE)
```

```
[1] 0.1870314
```

Thus 18.7% of all pregnancies last more than 280 days. This is not unusual since the probability is greater than 5%.

c. Find the probability of a pregnancy lasting less than 250 days.

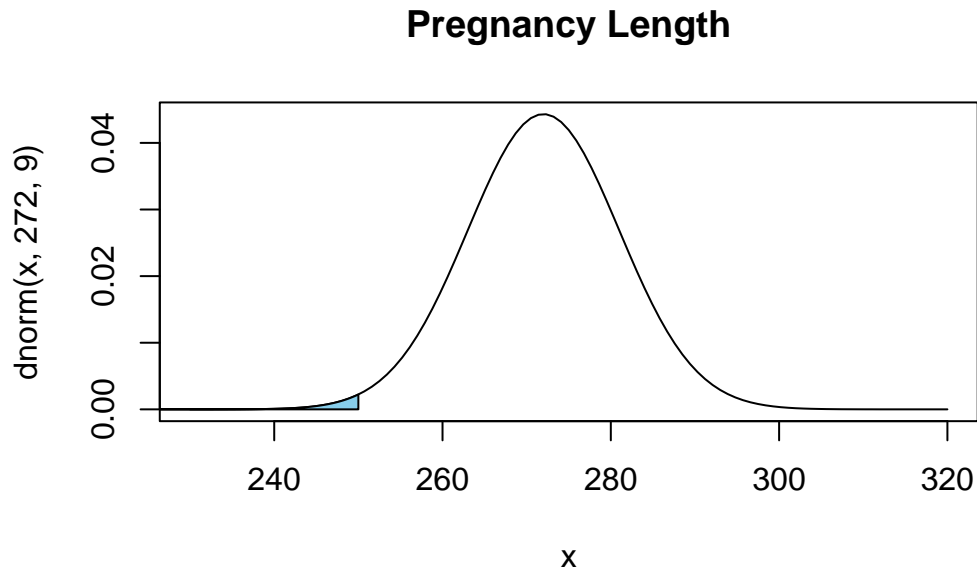


Figure 7.10: Density plot of pregnancy length. Normally distributed with mean 272 and standard deviation 9, and $P(x < 250)$

First translate the statement into a mathematical statement. $P(x < 250)$

Now, draw a picture Figure ??.

The probability of a pregnancy lasting longer than 250 days is $P(x < 250) = 0.0073$. The command in r Studio is

```
pnorm(250, 272, 9, lower.tail=TRUE)
```

```
[1] 0.007253771
```

Thus 0.73% of all pregnancies last less than 250 days. This is unusual since the probability is less than 5%.

d. Find the probability that a pregnancy lasts between 265 and 280 days.

First translate the statement into a mathematical statement. $P(265 < x < 280)$

Now draw a picture Figure ??.

The probability of a pregnancy lasting between 265 days and 280 days is $P(265 < x < 280) = 0.187$. To find the area between two values on the normal distribution, first, find the area to the left of the lower value, Graphically, this looks like Figure ??

Now find the area less than 280. Graphically this looks like Figure ??

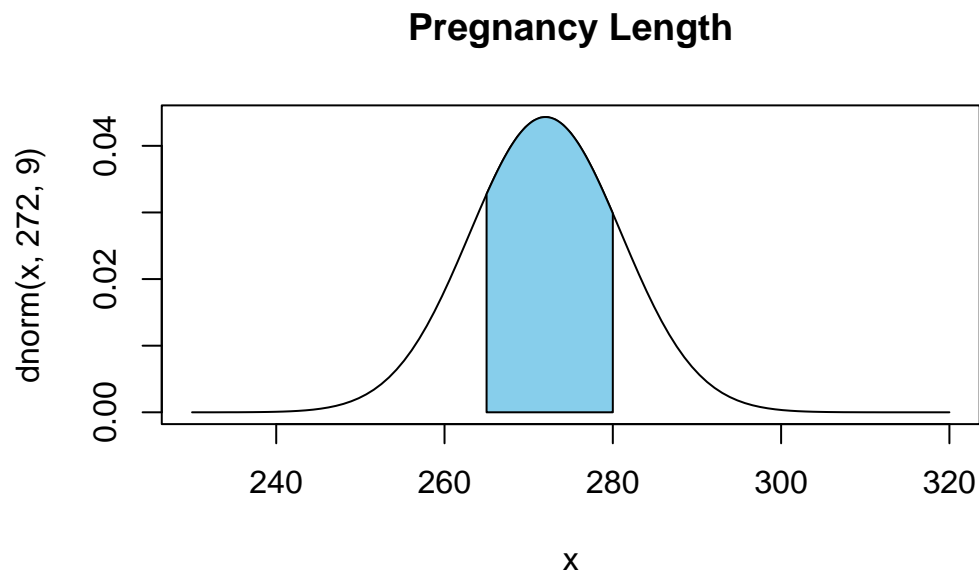


Figure 7.11: Density plot of pregnancy length. Normally distributed with mean 272 and standard deviation 9, and $P(265 < x < 280)$

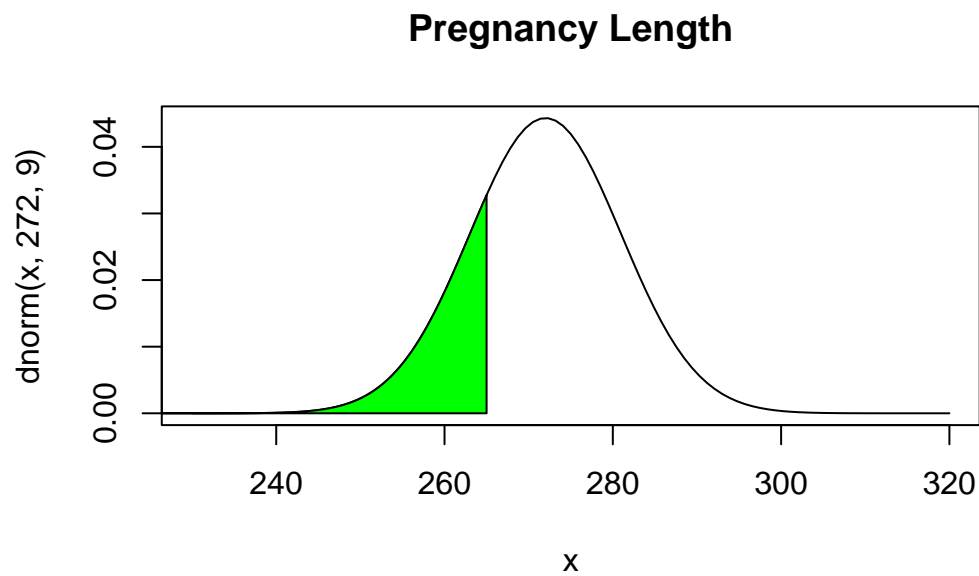


Figure 7.12: Density plot of pregnancy length. Normally distributed with mean 272 and standard deviation 9, and $P(x < 265)$

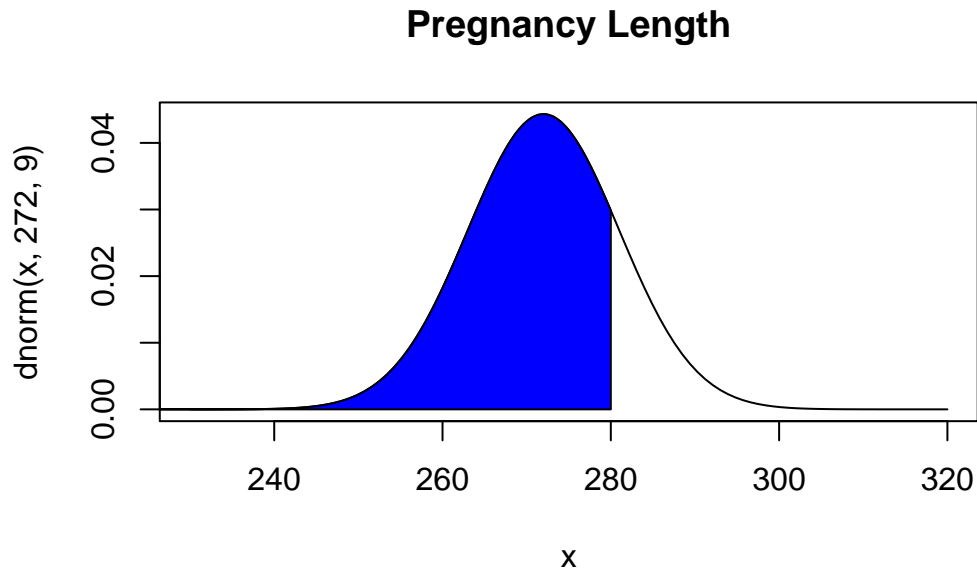


Figure 7.13: Density plot of pregnancy length. Normally distributed with mean 272 and standard deviation 9, and $P(x < 280)$

Looking at the three figures, if you take the area in Figure ?? and subtract the area in Figure ?? you get the area in Figure ??. In rStudio, the way to find the probability the probability of a pregnancy lasting between 265 days and 280 days, $P(265 < x < 280) = 0.595$ use the following command

```
pnorm(280, 272, 9, lower.tail=TRUE)-pnorm(265, 272, 9, lower.tail=TRUE)
```

```
[1] 0.5946186
```

Thus 59.5% of all pregnancies last between 265 and 280 days.

- e. Find the length of pregnancy that 10% of all pregnancies last less than.

This problem is asking you to find an x value from a probability. You want to find the x value that has 10% of the length of pregnancies to the left of it. In this case, you are given the probability. In r, the command is

```
qnorm(area, mean, standard_deviation, lower.tail=TRUE or FALSE)
```

For this example since you know the area in the lower tail, then use `lower.tail=TRUE`. So the command is

```
qnorm(0.1, 272, 9, lower.tail = TRUE)
```

```
[1] 260.466
```

Thus 10% of all pregnancies last less than approximately 260 days.

- f. Suppose you meet a woman who says that she was pregnant for less than 250 days. Would this be unusual and what might you think?

From part (c) you found the probability that a pregnancy lasts less than 250 days is 0.73%. Since this is less than 5%, it is very unusual. You would think that either the woman had a premature baby, or that she may be wrong about when she actually became pregnant.

7.2.2 Example: General Normal Distribution

The mean mathematics SAT score in 2012 was 514 with a standard deviation of 117 (“Total group profile,” 2012). Assume the mathematics SAT score is normally distributed.

- State the random variable.
- Find the probability that a person has a mathematics SAT score over 700.
- Find the probability that a person has a mathematics SAT score of less than 400.
- Find the probability that a person has a mathematics SAT score between a 500 and a 650.
- Find the mathematics SAT score that represents the top 1% of all scores.

7.2.2.1 Solution

- State the random variable.

x = mathematics SAT score

- Find the probability that a person has a mathematics SAT score over 700.

First translate the statement into a mathematical statement. $P(x > 700)$

Now, draw a picture Figure ??

To find $P(x > 700) = 0.0559$, the command in R would be

```
pnorm(700, 514, 117, lower.tail = FALSE)
```

```
[1] 0.05594631
```

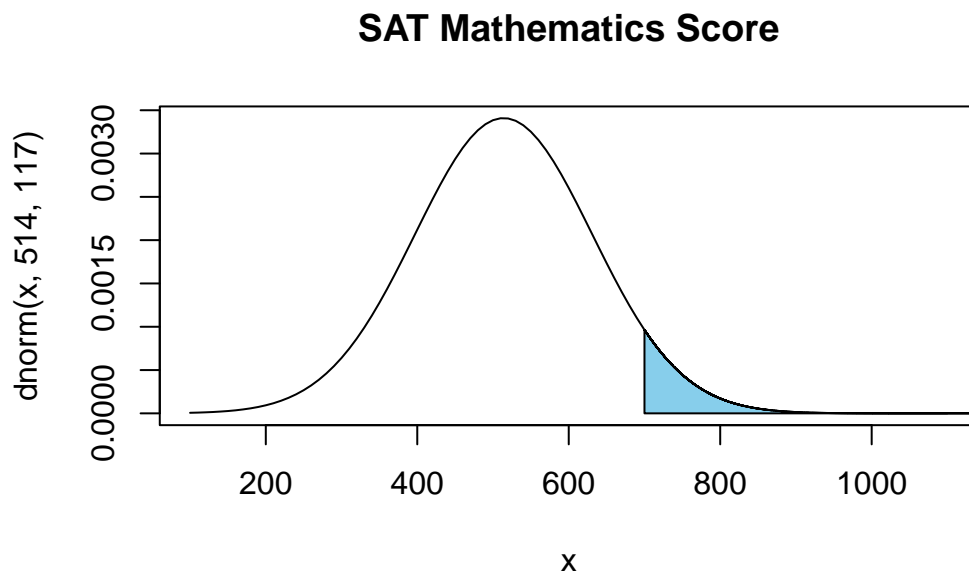


Figure 7.14: Density plot of SAT mathematics score. Normally distributed with mean 514 and standard deviation 117 and $P(x > 700)$

There is a 5.6% chance that a person scored above a 700 on the mathematics SAT test. This is not unusual.

- c. Find the probability that a person has a mathematics SAT score of less than 400.

First translate the statement into a mathematical statement. $P(x < 400)$

Now, draw a picture Figure ??

To find $P(x < 400) = 0.165$, the command in r would be

```
pnorm(400, 514, 117, lower.tail = TRUE)
```

```
[1] 0.1649392
```

So, there is a 16.5% chance that a person scores less than a 400 on the mathematics part of the SAT.

- d. Find the probability that a person has a mathematics SAT score between a 500 and a 650.

First translate the statement into a mathematical statement $P(500 < x < 650)$

Now, draw a picture Figure ??

To find $P(500 < x < 650) = 0.425$, the command in r would be

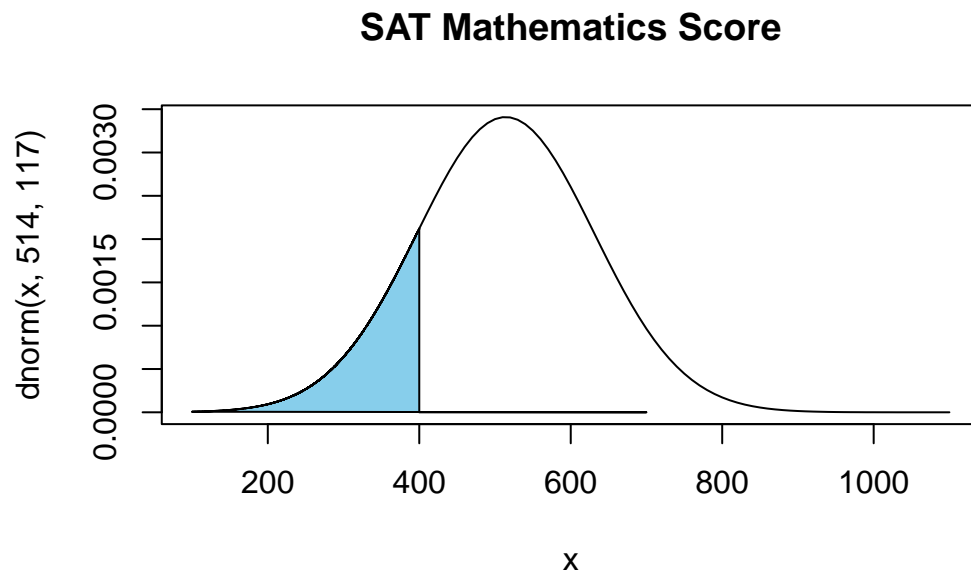


Figure 7.15: Density plot of SAT mathematics score. Normally distributed with mean 514 and standard deviation 117 and $P(x < 400)$

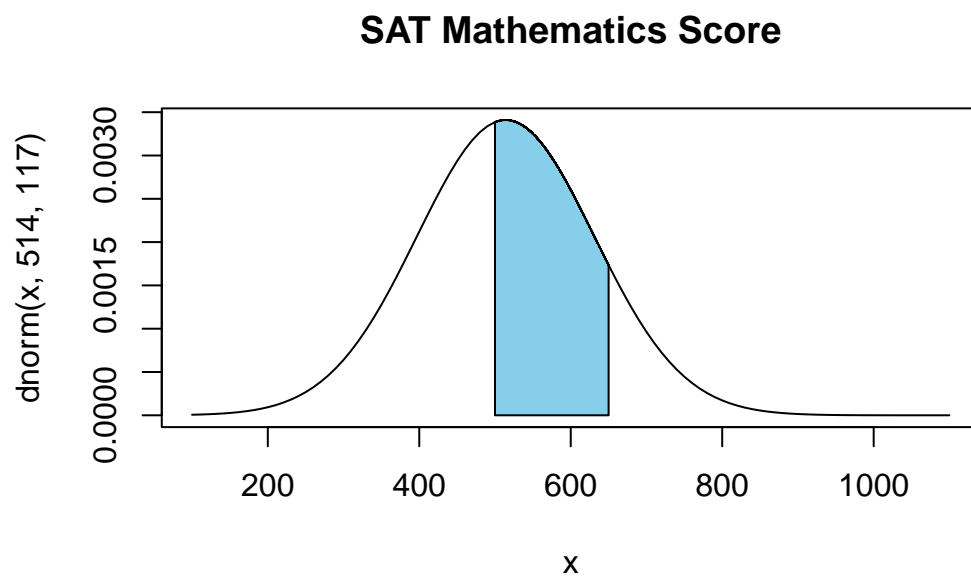


Figure 7.16: Density plot of SAT mathematics score. Normally distributed with mean 514 and standard deviation 117 and $P(514 < x < 650)$

```
pnorm(650, 514, 117, lower.tail = TRUE)-pnorm(500, 514, 117, lower.tail=TRUE)
```

```
[1] 0.4250851
```

So, there is a 42.5% chance that a person has a mathematical SAT score between 500 and 650.

- e. Find the mathematics SAT score that represents the top 1% of all scores.

This problem is asking you to find an x value from a probability. You want to find the x value that has 1% of the mathematics SAT scores to the right of it. In this case you are using the upper tail of the curve. To find this x value on rStudio, use the command

```
qnorm(0.01, 514, 117, lower.tail=FALSE)
```

```
[1] 786.1827
```

So, 1% of all people who took the SAT scored over about 786 points on the mathematics SAT.

7.2.3 Homework for Normal Distribution Section

- Find each of the probabilities, where z is a z -score from the standard normal distribution with mean of $\mu = 0$ and standard deviation $\sigma = 1$. It helps to draw a picture for each problem.
 - $P(z < 2.36)$
 - $P(z > 0.67)$
 - $P(0 < x < 2.11)$
 - $P(-2.78 < z < 1.97)$
- Find the z -score corresponding to the given area. Remember, z is distributed as the standard normal distribution with mean of $\mu = 0$ and standard deviation $\sigma = 1$.
 - The area to the left of z is 15%.
 - The area to the right of z is 65%.
 - The area to the left of z is 10%.
 - The area to the right of z is 5%.

- e. The area between $-z$ and z is 95%. (Hint draw a picture and figure out the area to the left of $-z$.)
 - f. The area between $-z$ and z is 99%.
3. If a random variable that is normally distributed has a mean of 25 and a standard deviation of 3, convert the given value to a z -score.
 - a. $x = 23$
 - b. $x = 33$
 - c. $x = 19$
 - d. $x = 45$
 4. According to the WHO MONICA Project the mean blood pressure for people in China is 128 mmHg with a standard deviation of 23 mmHg (Kuulasmaa, Hense & Tolonen, 1998). Assume that blood pressure is normally distributed.
 - a. State the random variable.
 - b. Find the probability that a person in China has blood pressure of 135 mmHg or more.
 - c. Find the probability that a person in China has blood pressure of 141 mmHg or less.
 - d. Find the probability that a person in China has blood pressure between 120 and 125 mmHg.
 - e. Is it unusual for a person in China to have a blood pressure of 135 mmHg? Why or why not?
 - f. What blood pressure do 90% of all people in China have less than?
 5. The size of fish is very important to commercial fishing. A study conducted in 2012 found the length of Atlantic cod caught in nets in Karlskrona to have a mean of 49.9 cm and a standard deviation of 3.74 cm (Ovegard, Berndt & Lunneryd, 2012). Assume the length of fish is normally distributed.
 - a. State the random variable.
 - b. Find the probability that an Atlantic cod has a length less than 52 cm.
 - c. Find the probability that an Atlantic cod has a length of more than 74 cm.
 - d. Find the probability that an Atlantic cod has a length between 40.5 and 57.5 cm.
 - e. If you found an Atlantic cod to have a length of more than 74 cm, what could you conclude?
 - f. What length are 15% of all Atlantic cod longer than?
 6. The mean cholesterol levels of women age 45-59 in Ghana, Nigeria, and Seychelles is 5.1 mmol/l and the standard deviation is 1.0 mmol/l (Lawes, Hoorn, Law & Rodgers, 2004). Assume that cholesterol levels are normally distributed.
 - a. State the random variable.

- b. Find the probability that a woman age 45-59 in Ghana, Nigeria, or Seychelles has a cholesterol level above 6.2 mmol/l (considered a high level).
 - c. Find the probability that a woman age 45-59 in Ghana, Nigeria, or Seychelles has a cholesterol level below 5.2 mmol/l (considered a normal level).
 - d. Find the probability that a woman age 45-59 in Ghana, Nigeria, or Seychelles has a cholesterol level between 5.2 and 6.2 mmol/l (considered borderline high).
 - e. If you found a woman age 45-59 in Ghana, Nigeria, or Seychelles having a cholesterol level above 6.2 mmol/l, what could you conclude?
 - f. What value do 5% of all woman ages 45-59 in Ghana, Nigeria, or Seychelles have a cholesterol level less than?
7. In the United States, males between the ages of 40 and 49 eat on average 103.1 g of fat every day with a standard deviation of 4.32 g (“What we eat,” 2012). Assume that the amount of fat a person eats is normally distributed.
- a. State the random variable.
 - b. Find the probability that a man age 40-49 in the U.S. eats more than 110 g of fat every day.
 - c. Find the probability that a man age 40-49 in the U.S. eats less than 93 g of fat every day.
 - d. Find the probability that a man age 40-49 in the U.S. eats less than 65 g of fat every day.
 - e. If you found a man age 40-49 in the U.S. who says he eats less than 65 g of fat every day, would you believe him? Why or why not?
 - f. What daily fat level do 5% of all men age 40-49 in the U.S. eat more than?
8. A dishwasher has a mean life of 12 years with an estimated standard deviation of 1.25 years (“Appliance life expectancy,” 2013). Assume the life of a dishwasher is normally distributed.
- a. State the random variable.
 - b. Find the probability that a dishwasher will last more than 15 years.
 - c. Find the probability that a dishwasher will last less than 6 years.
 - d. Find the probability that a dishwasher will last between 8 and 10 years.
 - e. If you found a dishwasher that lasted less than 6 years, would you think that you have a problem with the manufacturing process? Why or why not?
 - f. A manufacturer of dishwashers only wants to replace free of charge 5% of all dishwashers. How long should the manufacturer make the warranty period?
9. The mean starting salary for nurses is \$67,694 nationally (“Staff nurse -,” 2013). The standard deviation is approximately \$10,333. Assume that the starting salary is normally distributed.
- a. State the random variable.

- b. Find the probability that a starting nurse will make more than \\$80,000.
 - c. Find the probability that a starting nurse will make less than \\$60,000.
 - d. Find the probability that a starting nurse will make between \\$55,000 and \\$72,000.
 - e. If a nurse made less than \\$50,000, would you think the nurse was under paid? Why or why not?
 - f. What salary do 30% of all nurses make more than?
10. The mean yearly rainfall in Sydney, Australia, is about 137 mm and the standard deviation is about 69 mm (“Annual maximums of,”2013). Assume rainfall is normally distributed.
- a. State the random variable.
 - b. Find the probability that the yearly rainfall is less than 100 mm.
 - c. Find the probability that the yearly rainfall is more than 240 mm.
 - d. Find the probability that the yearly rainfall is between 140 and 250 mm.
 - e. If a year has a rainfall less than 100mm, does that mean it is an unusually dry year? Why or why not?
 - f. What rainfall amount are 90% of all yearly rainfalls more than?

7.3 Assessing Normality

The distributions you have seen up to this point have been assumed to be normally distributed, but how do you determine if it is normally distributed. One way is to take a sample and look at the sample to determine if it appears normal. If the sample looks normal, then most likely the population is also. Here are some guidelines that are use to help make that determination.

1. **Density Plot:** Make a density plot. For a normal distribution, the density plot should be roughly bell-shaped. For small samples, this is not very accurate, and another method is needed. A distribution may not look normally distributed from the density plot, but it still may be normally distributed.
2. **Normal quantile plot (or normal probability plot):** This plot is provided through statistical software on a computer. If the points lie close to a line, the data comes from a distribution that is approximately normally distributed. If the points do not lie close to a line or they show a pattern that is not a line, the data are likely to come from a distribution that is not normally distributed.

7.3.1 To create a density plot on rStudio:

Read the Data Frame into r Studio. The command for density is

```
gf_density(~variable, data=Data_Frame)
```

See chapter 2 for more examples of this.

7.3.2 To create a normal quantile plot on rStudio

Read the Data Frame into rStudio. The command for normal quantile plot is

```
gf_qqnorm(~variable, data=Data_Frame)
```

Realize that your random variable may be normally distributed, even if the sample fails the two tests. However, if the density plot definitely doesn't look symmetric and bell shaped, and the normal probability plot doesn't look linear, then you can be fairly confident that the data set does not come from a population that is normally distributed.

7.3.3 Example: Is It Normal?

In Kiama, NSW, Australia, there is a blowhole. The data in Table ?? are times in seconds between eruptions ("Kiama blowhole eruptions," 2013). Do the data come from a population that is normally distributed?

```
Eruption<-read.csv( "https://krkozak.github.io/MAT160/Blowhole_eruptions.csv")
knitr::kable(head(Eruption))
```

Table 7.1: Time (in Seconds) Between Kiama Blowhole Eruptions

Interval
83
51
87
60
28
95

Code book for Data Frame Eruption

Description The ocean swell produces spectacular eruptions of water through a hole in the cliff at Kiama, about 120km south of Sydney, known as the Blowhole. The times at which 65 successive eruptions occurred from 1340 hours on 12 July 1998 were observed using a digital watch.

Format This data frame contains the following columns:

Interval: Waiting time between eruptions (seconds)

Source Kiama Blowhole Eruptions. (n.d.). Retrieved from <http://www.statsci.org/data/oz/kiama.html>

References The data was collected and contributed by Jim Irish, Faculty of Engineering, University of Technology, Sydney.

- State the random variable
- Draw a Density plot
- Draw the normal quantile plot.
- Do the data come from a population that is normally distributed?

7.3.3.1 Solution

- State the random variable

x = time in seconds between eruptions of Kiama Blowhole

- Draw a Density plot

The density plot produced is in Figure ??

```
gf_density(~Interval, data=Eruption, title="Eruption times for Kiama Blowhole", xlab="Time (seconds)", ylab="density")
```

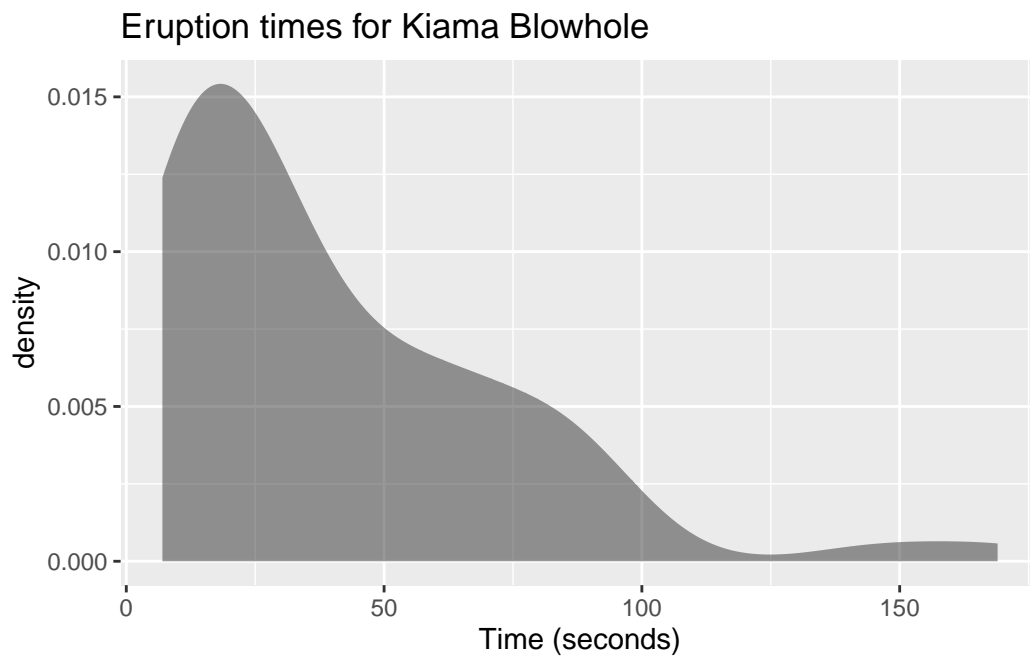


Figure 7.17: Density Plot of Eruption Times for Kiama Blowhole

This looks skewed right and not symmetric.

c. Draw the normal quantile plot.

The normal quantile plot is in Figure ??

```
gf_qq(~Interval, data=Eruption, title="Eruption times for Kiama Blowhole")
```

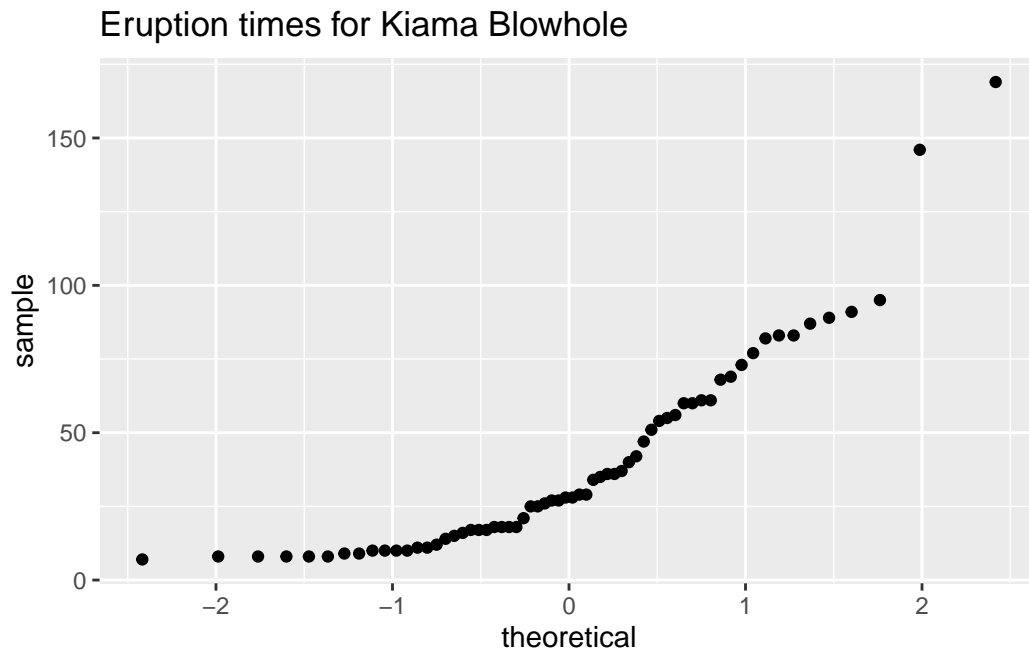


Figure 7.18: Normal Quantile Plot of Eruption Times for Kiama blowhole.

Figure ?? looks more like an exponential growth than linear.

d. Do the data come from a population that is normally distributed?

Considering the density plot is skewed right, and the normal probability plot does not look linear, then the conclusion is that this sample is not from a population that is normally distributed.

7.3.4 Example: Is It Normal?

The US National Center for Health Statistics (NCHS) conducted a series of health and nutrition surveys called NHANES. One of the many variables in NHANES is pulse. Determine if pulse is a normally distributed variable. The NHANES data frame is Table ??.

a. State the random variable

- b. Draw a density plot
- c. Draw the normal quantile plot.
- d. Do the data come from a population that is normally distributed?

7.3.4.1 Solution

- a. State the random variable

$x = \text{pulse}$

- b. Draw a density plot

The density plot is in Figure ??

```
gf_density(~Pulse, data=NHANES, title="Pulse Rate", xlab="Pulse Rate (bpm)")
```

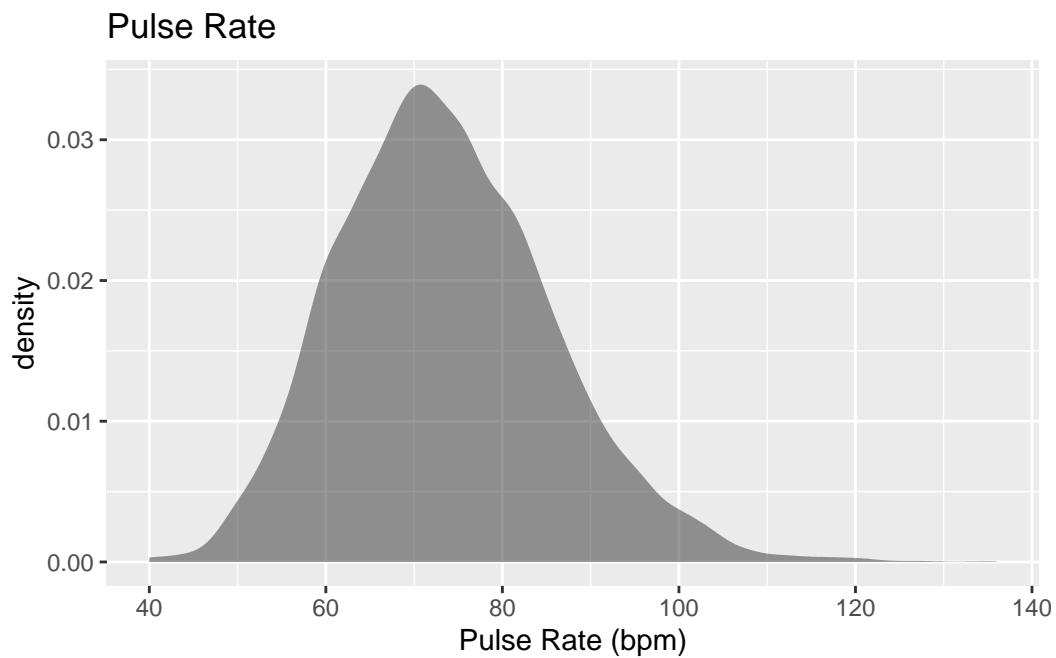


Figure 7.19: Density plot of Pulse Rate (bpm)

This looks somewhat symmetric and bell shaped.

- c. Draw the normal quantile plot.

The normal quantile plot is in Figure ??

```
gf_qq(~Pulse, data=NHANES, title="Pulse Rate")
```

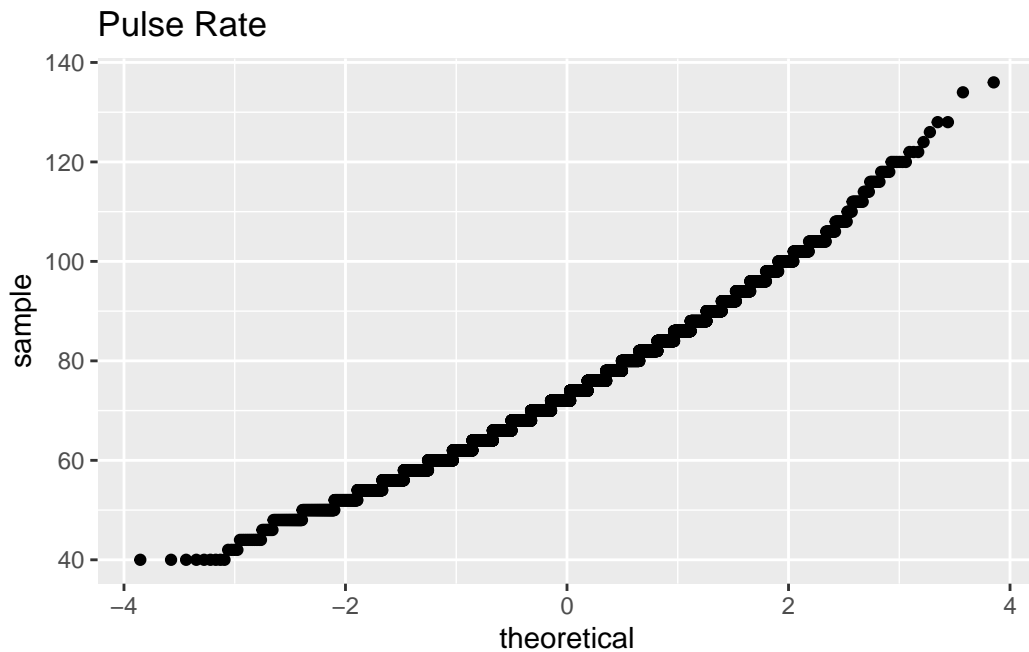


Figure 7.20: Normal Quantile Plot of Pulse Rate (bmp)

Figure ?? looks fairly linear.

- d. Do the data come from a population that is normally distributed?

Considering the density plot is bell shaped and the normal probability plot looks linear. The conclusion is that this sample is from a population that is normally distributed.

7.3.5 Homework for Assessing Normality Section

1. Cholesterol data was collected on patients four days after having a heart attack. The data is in Table ?. Assess if the data is from a population that is normally distributed.

Code Book for Cholesterol See is below Table ?.

2. The size of fish is very important to commercial fishing. A study conducted in 2012 collected the lengths of Atlantic cod caught in nets in Karlskrona (Ovegard, Berndt & Lunneryd, 2012). Data based on information from the study is in Table ?. Determine if the data is from a population that is normally distributed.


```
Cod<-read.csv( "https://krkozak.github.io/MAT160/cod.csv")
knitr::kable(head(Cod))
```

Table 7.2: Atlantic Cod Lengths

length
48
50
50
55
53
50

3. The WHO MONICA Project collected blood pressure data for people in China (Kuulasmaa, Hense & Tolonen, 1998). Data based on information from the study is in Table ?? . Determine if the data is from a population that is normally distributed.

```
BP<-read.csv( "https://krkozak.github.io/MAT160/bp.csv")
knitr::kable(head(BP))
```

Table 7.3: Blood Pressure Values for People in China

pressure
114
141
154
137
131
132

4. Annual rainfalls for Sydney, Australia are given in Table ?? (“Annual maximums of,” 2013). Can you assume rainfall is normally distributed?

```
Annual<-read.csv( "https://krkozak.github.io/MAT160/annual.csv")
knitr::kable(head(Annual))
```

Table 7.4: Annual Rainfall in Sydney, Australia

amount
146.8
383.0
90.9
178.1
267.5
95.5

7.4 Sampling Distribution and the Central Limit Theorem

You now have most of the skills to start statistical inference, but you need one more concept.

First, it would be helpful to state what statistical inference is in more accurate terms.

Statistical Inference: to make accurate decisions about parameters from statistics

When it says “accurate decision,” you want to be able to measure how accurate. You measure how accurate using probability. In both binomial and normal distributions, you needed to know that the random variable followed either distribution. You need to know how the statistic is distributed and then you can find probabilities. In other words, you need to know the shape of the sample mean or whatever statistic you want to make a decision about.

How is the statistic distributed? This is answered with a sampling distribution.

Sampling Distribution: how a sample statistic is distributed when repeated trials of size n are taken.

7.4.1 Example: Sampling Distribution

The NHANES data frame has the pulse rates for approximately 50,000 individuals. The random variable is $x = \text{pulse rate}$. The probability distribution of this random variable is presented in Figure ???. Although pulse rates from 50,000 individuals isn’t the entire population, the sample is most likely a good representation of the population. Thus, it is safe to assume the population is normally distributed. An estimate for the population mean is 73.6 bpm, and the population standard deviation estimate is 12.2 bpm.

```
gf_density(~Pulse, data=NHANES, title = "Pulse Rate", xlab="Pulse (bpm)")
df_stats(~Pulse, data=NHANES, mean, sd)
```

```

  response      mean      sd
1    Pulse 73.55973 12.15542

```

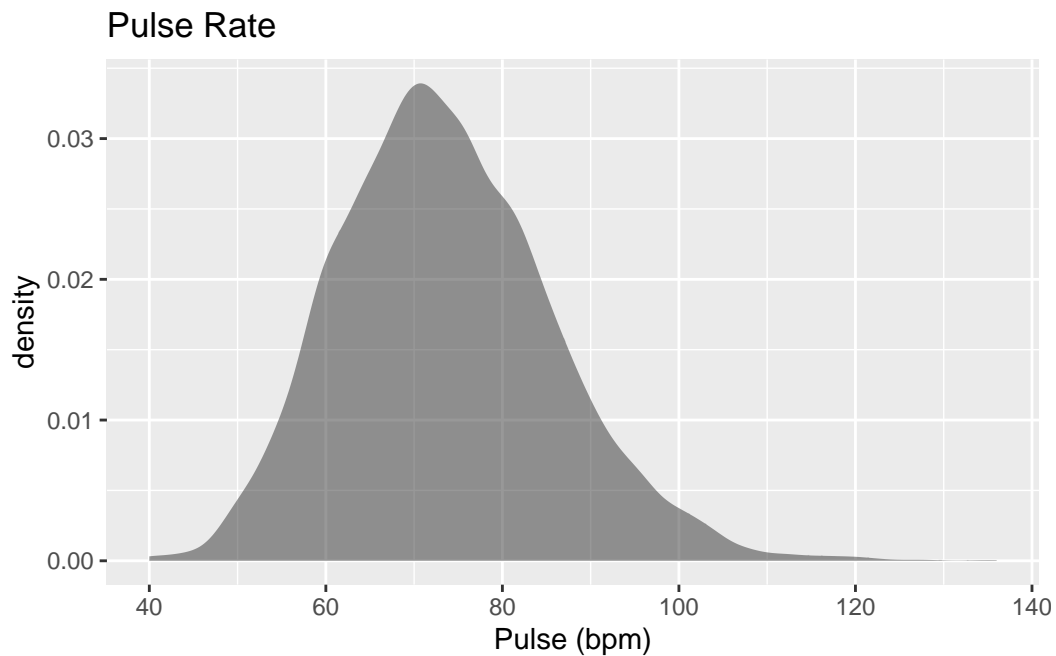


Figure 7.21: Distribution of Pulse Rate

Suppose you take a random sample of 10 pulse rates from those 50,000 individuals. A random sample of data from 10 individuals is:

```

NHANES|>
  sample_n(size=10)

```

A tibble: 10 x 76

	ID	SurveyYr	Gender	Age	AgeDecade	AgeMonths	Race1	Race3	Education
	<int>	<fct>	<fct>	<int>	<fct>	<int>	<fct>	<fct>	<fct>
1	62602	2011_12	female	80	<NA>	NA	White	White	Some Colle~
2	59756	2009_10	female	11	" 10-19"	142	White	<NA>	<NA>
3	54020	2009_10	female	20	" 20-29"	248	White	<NA>	9 - 11th G~
4	71879	2011_12	male	24	" 20-29"	NA	Black	Black	Some Colle~
5	69671	2011_12	female	58	" 50-59"	NA	Hispanic	Hispanic	8th Grade
6	52933	2009_10	female	12	" 10-19"	153	Mexican	<NA>	<NA>
7	57287	2009_10	female	33	" 30-39"	398	White	<NA>	Some Colle~
8	64629	2011_12	male	14	" 10-19"	NA	White	White	<NA>

```

 9 64870 2011_12 female    80 <NA>          NA Other    Asian    High School
10 64570 2011_12 female    65 " 60-69"      NA White    White    Some Colle~
# i 67 more variables: MaritalStatus <fct>, HHIncome <fct>, HHIncomeMid <int>,
#   Poverty <dbl>, HomeRooms <int>, HomeOwn <fct>, Work <fct>, Weight <dbl>,
#   Length <dbl>, HeadCirc <dbl>, Height <dbl>, BMI <dbl>,
#   BMICatUnder20yrs <fct>, BMI_WHO <fct>, Pulse <int>, BPSysAve <int>,
#   BPDiaAve <int>, BPSys1 <int>, BPDia1 <int>, BPSys2 <int>, BPDia2 <int>,
#   BPSys3 <int>, BPDia3 <int>, Testosterone <dbl>, DirectChol <dbl>,
#   TotChol <dbl>, UrineVol1 <int>, UrineFlow1 <dbl>, UrineVol2 <int>, ...

```

It might be useful to find the mean pulse rate from a random sample of size 10.

```

response mean
1 Pulse 68.75

```

Now suppose you took another random sample of size 10 and found the mean pulse rate for that sample. Repeat this process 100 times. At this point you would basically have a new sample of 100 mean pulse rates. You could assess how this sample is distributed by creating a density plot Figure ??

```

Trials <- do(100) * { NHANES |>
  sample_n(size = 10) |>
  df_stats( ~Pulse, means = mean)}
gf_density( ~means, data = Trials, title = "Density plot of sample mean when n=10")

df_stats(~means, data=Trials, mean, sd)

```

```

response mean sd
1 means 73.83725 3.82866

```

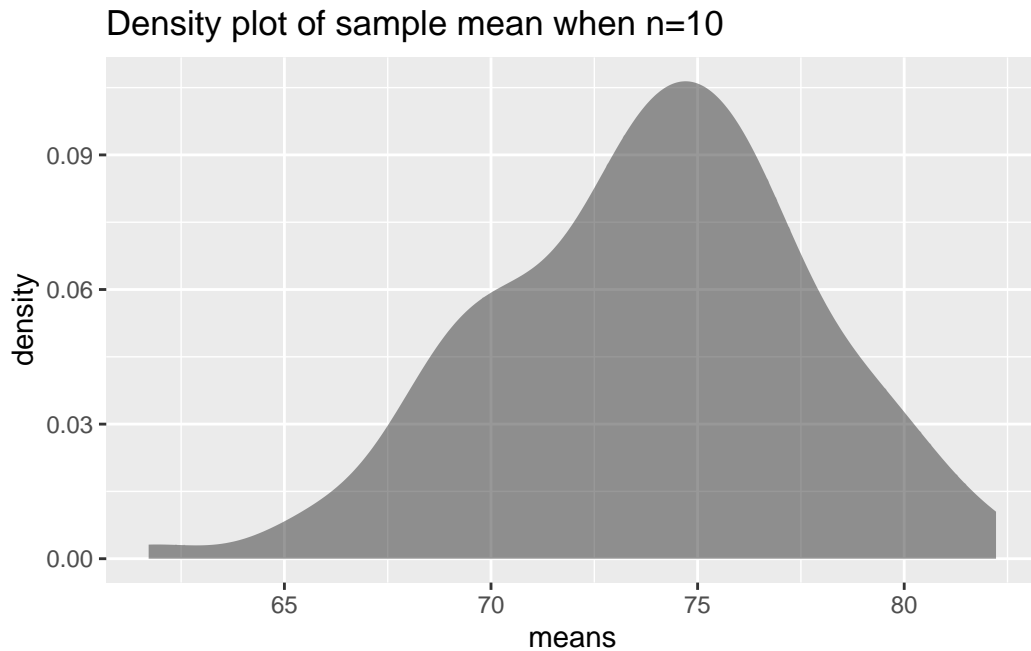


Figure 7.22: Density Plot of Sample Means When $n = 10$

This distribution is a sampling distribution. That is all a sampling distribution is. It is a distribution created from statistics.

Notice the distribution does look a great deal like the distribution of the original random variable. Notice the mean of the sample means $\mu_{\bar{x}} = 73.8$ bpm which is almost the same of as the mean of the population. The standard deviation of the sample means, $\sigma_{\bar{x}} = 4.35$ bpm is about $\frac{1}{3}$ of the population standard deviation.

What does this distribution look like if instead of repeating the experiment 10 times you repeat it 50 times instead?

This density plot of the sampling distribution is displayed in Figure ??

```
Trials <- do(100) * { NHANES |>
  sample_n(size = 50) |>
  df_stats( ~Pulse, means = mean) }
gf_density( ~means, data = Trials, title="Sample means when n=50")|>
  gf_lims(x=c(68,79))
df_stats(~means, data=Trials, mean, sd)
```

	response	mean	sd
1	means	73.66942	1.790344

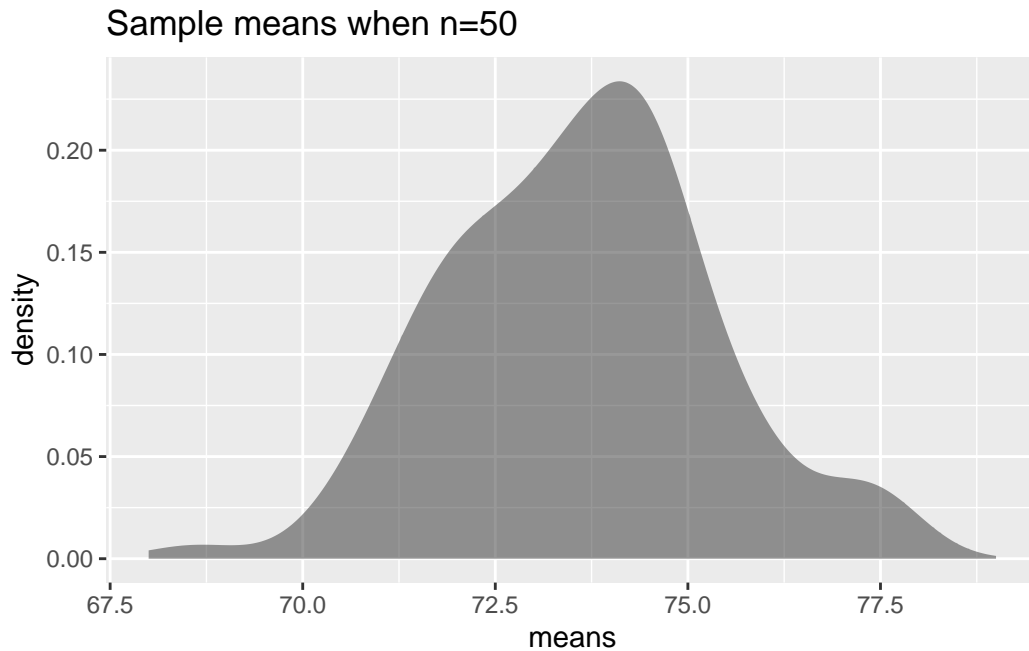


Figure 7.23: Density Plot of Sample Means When $n = 50$

Notice this density plot of the sample mean looks approximately symmetrical and could almost be called normal. Notice, the mean of the sample means is 73.6 bpm which is approximately what the population mean is. The standard deviation of the sample means is 1.77 bpm which is around $\frac{1}{7}$ of the population standard deviation. What if you keep increasing n ? What will the sampling distribution of the sample mean look like? In other words, what does the sampling distribution of \bar{x} look like as n gets even larger?

This depends on how the original distribution is distributed. In Example: Sampling Distribution, the random variable was approximately normally distributed. When n was 10, the distribution of the mean looked approximately normal. What if the original distribution wasn't normal? How big would n have to be? Consider a different variable in the NHANES data frame that isn't normally distributed such as age when a participant started to smoke cigarettes (SmokeAge). The density plot for the large sample is in Figure ???. The mean for the large sample is 17.8 years and the standard deviation is 5.3 years, so $\mu = 17.8$ years and $\sigma = 5.3$ years approximately.

```
gf_density(~SmokeAge, data=NHANES, title = "Density Plot of Age when Person Started Smoking")
df_stats(~SmokeAge, data=NHANES, mean, sd)
```

```
response    mean    sd
1 SmokeAge 17.82662 5.32666
```

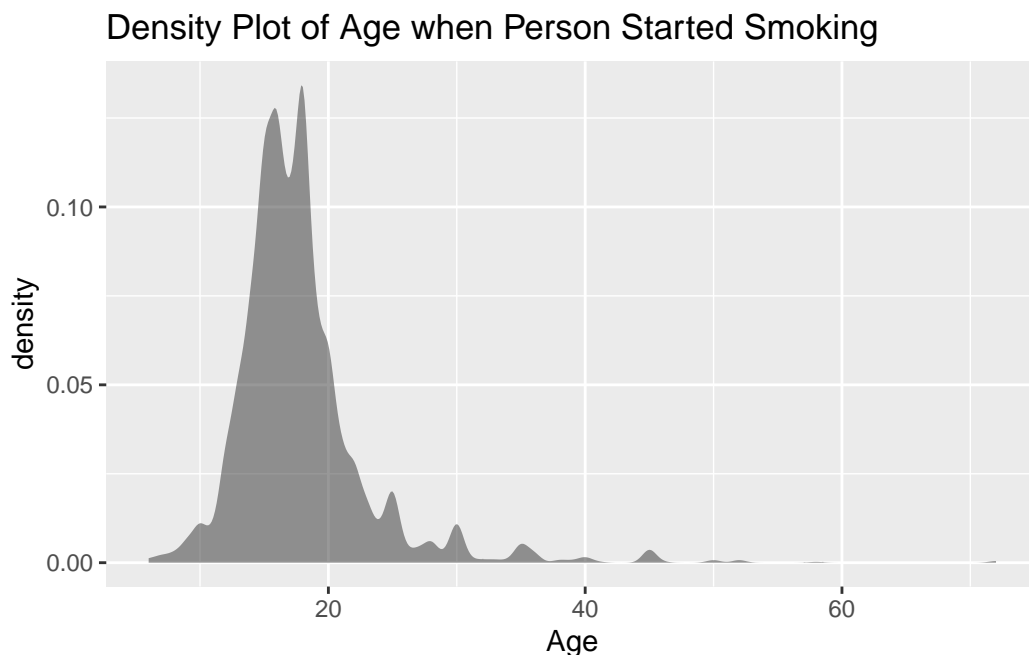


Figure 7.24: Density Plot of Age When Person Started Smoking.

Now take 100 samples of size 50 individuals from the NHANES Data Frame. Then graph a density plot of `SmokeAge`, the age when someone started to smoke. Notice the the sampling distribution of the sample means looks fairly normally distributed even though the original random variable was not normally distributed. The mean of the sample mean, $\mu_{\bar{x}} = 17.8$ years and the standard deviation of the sample mean, $\sigma_{\bar{x}} = 1.56$ years. The mean of the sample mean is the same as the mean of the population, but the standard deviation of the sample mean is much less than the standard deviation of the original data.

	response	mean	sd
1	means	17.93834	1.366059

One question is, why is the mean of the sample means the same as the mean of the population? Suppose you have a random variable that has a population mean, μ , and a population standard deviation, σ . If a sample of size n is taken, then the sample mean, has a mean $\mu_{\bar{x}} = \mu$ and standard deviation of $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$. The standard deviation of the sample mean is lower because by taking the mean you are averaging out the extreme values, which makes the distribution of the sample mean less spread out.

You now know the center and the variability of \bar{x} . You also want to know the shape of the distribution of \bar{x} . You hope it is normal, since you know how to find probabilities using the normal curve. The following theorem tells you the requirement to have \bar{x} be normally distributed.

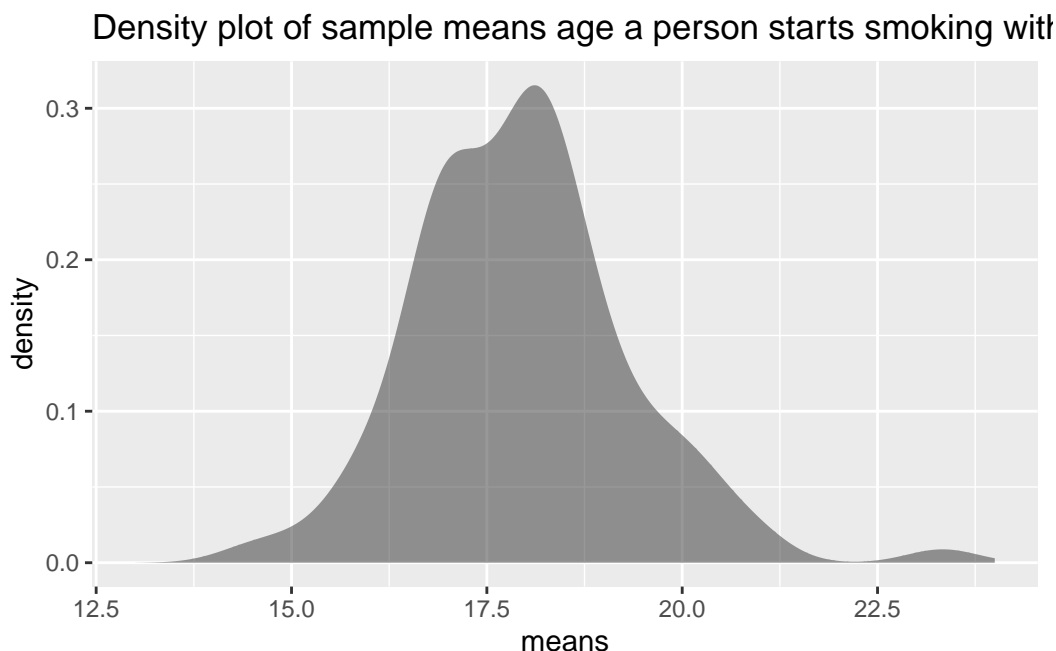


Figure 7.25: Density Plot of Age When Person Started Smoking when sample size is 50.

7.4.2 Central Limit Theorem

Suppose a random variable is from any distribution. If a sample of size n is taken, the the sample mean, \bar{x} , becomes normally distributed as n increases.

What this says is that no matter what x looks like, \bar{x} would look normal if n is large enough. Now, what size of n is large enough? That depends on how x is distributed in the first place. If the original random variable is normally distributed, then n just needs to be 2 or more data points. If the original random variable is somewhat mound shaped and symmetrical, then n needs to be greater than or equal to 30. Sometimes the sample size can be smaller, but this is a good general rule to use. The sample size may have to be much larger if the original random variable is really skewed one way or another.

Now that you know when the sample mean will look like a normal distribution, then you can find the probability related to the sample mean. Remember that the mean of the sample mean is just the mean of the original data ($\mu_{\bar{x}} = \mu$), but the standard deviation of the sample mean, $\sigma_{\bar{x}}$, also known as the standard error of the mean, is actually $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$. Make sure you use this in all calculations. If you are using the z -score, the formula when working with \bar{x} is $z = \frac{\bar{x} - \mu_{\bar{x}}}{\sigma_{\bar{x}}} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$. To use rStudio to calculate probabilities use $P(\bar{x} < a) = \text{pnorm}(a, \mu_{\bar{x}}, \sigma_{\bar{x}}, \text{lower.tail} = \text{TRUE})$ $P(\bar{x} > a) = \text{pnorm}(a, \mu_{\bar{x}}, \sigma_{\bar{x}}, \text{lower.tail} = \text{FALSE})$.

7.4.3 Example: Finding Probabilities for Sample Means

The birth weight of boy babies of European descent who were delivered at 40 weeks is normally distributed with a mean of 3687.6 g with a standard deviation of 410.5 g (Janssen, Thiessen, Klein, Whitfield, MacNab & Cullis-Kuhl, 2007). Suppose there were nine European descent boy babies born on a given day and the mean birth weight is calculated.

- State the random variable.
- What is the mean of the sample mean?
- What is the standard deviation of the sample mean?
- What distribution is the sample mean distributed as?
- Find the probability that the mean weight of the nine boy babies born was less than 3500.4 g.
- Find the probability that the mean weight of the nine babies born was less than 3452.5 g.

7.4.3.1 Solution

- State the random variable.

x = birth weight of boy babies (Note: the random variable is something you measure, and it is not the mean birth weight. Mean weight is calculated.)

- What is the mean of the sample mean?

$$\mu_{\bar{x}} = \mu = 3687.4g$$

- What is the standard deviation of the sample mean?

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{410.5g}{\sqrt{9}} = 136.8g$$

- What distribution is the sample mean distributed as?

Since the original random variable is distributed normally, then the sample mean is distributed normally.

- Find the probability that the mean weight of the nine boy babies born was less than 3500.4 g.

To find $P(\bar{x} < 3500.4) = 0.086$. use the rStudio command

```
pnorm(3500.4, 3687.6, 410.5/sqrt(9), lower.tail = TRUE)
```

```
[1] 0.08564231
```

There is an 8.6% chance that the mean birth weight of the nine boy babies born would be less than 3500.4 g. Since this is more than 5%, this is not unusual.

- f. Find the probability that the mean weight of the nine babies born was less than 3452.5 g.

You are looking for the $P(\bar{x} < 3452.5)$.

To find in rStudio, $P(\bar{x} < 3452.5) = 0.043$ use the command

```
pnorm(3452.5, 3687.4, 410.5/sqrt(9), lower.tail = TRUE)
```

```
[1] 0.04301819
```

There is a 4.3% chance that the mean birth weight of the nine boy babies born would be less than 3452.5 g. Since this is less than 5% this would be an unusual event. If it actually happened, then you may think there is something unusual about this sample. Maybe some of the nine babies were born as multiples, which brings the mean weight down, or some or all of the babies were not of European descent (in fact the mean weight of South Asian boy babies is 3452.5 g), or some were born before 40 weeks, or the babies were born at high altitudes.

7.4.4 Example: Finding Probabilities for Sample Means

For Americans that smoke, the average age that they started smoking is 17.8 years, with a standard deviation of approximately 1.56 years from the NHANES data. This random variable is not normally distributed, though it is somewhat mound shaped.

- State the random variable.
- Suppose a sample of 35 smoking American's is taken. Find the probability that the mean age that these 35 smoking Americans started to smoke is more than 21 years.

7.4.4.1 Solution

- State the random variable.

x = age that smoking Americans started to smoke

- Suppose a sample of 35 smoking American's is taken. Find the probability that the mean age that these 35 smoking Americans started to smoke is more than 21 years.

Even though the original random variable is not normally distributed, the sample size is over 30, by the central limit theorem the sample mean will be normally distributed. The mean of the sample mean is $\mu_{\bar{x}} = 17.8$. The standard deviation of the sample mean is $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{1.56}{\sqrt{35}}$. You have all the information you need to use the normal command using rStudio. Without the central limit theorem, you couldn't use the normal command, and you would not be able to answer this question.

The probability that the mean age that 35 smoking Americans start to smoke is more than 21 years, is the mathematical statement $P(\bar{x} > 21)$

To find $P(\bar{x} > 21) = 3.42 \times 10^{-34}$ using r studio, use the command:

```
pnorm(21, 17.8, 1.56/sqrt(35), lower.tail=FALSE)
```

```
[1] 3.422499e-34
```

The probability of a sample mean of 35 smoking Americans being more than 21 years when they smoked for the first time is very small. This is extremely unlikely to happen. If it does, it may make you wonder about the sample. Could the population mean have increased from the 17.8 years as was stated? Could the sample not have been random, and instead have been a group of smoking Americans who had started to smoke much later? These questions, and more, are ones that you would want to ask as a researcher

7.4.5 Homework for Sampling Distribution and the Central Limit Theorem Section

1. A random variable is not normally distributed, but it is mound shaped. It has a mean of 14 and a standard deviation of 3.
 - a. If you take a sample of size 10, can you say what the shape of the sampling distribution for the sample mean is? Why?
 - b. For a sample of size 10, state the mean of the sample mean and the standard deviation of the sample mean.
 - c. If you take a sample of size 35, can you say what the shape of the distribution of the sample mean is? Why?
 - d. For a sample of size 35, state the mean of the sample mean and the standard deviation of the sample mean.
2. A random variable is normally distributed. It has a mean of 245 and a standard deviation of 21.
 - a. If you take a sample of size 10, can you say what the shape of the distribution for the sample mean is? Why?

- b. For a sample of size 10, state the mean of the sample mean and the standard deviation of the sample mean.
 - c. For a sample of size 10, find the probability that the sample mean is more than 241.
 - d. If you take a sample of size 35, can you say what the shape of the distribution of the sample mean is? Why?
 - e. For a sample of size 35, state the mean of the sample mean and the standard deviation of the sample mean.
 - f. For a sample of size 35, find the probability that the sample mean is more than 241.
 - g. Compare your answers in part c and f. Why is one smaller than the other?
3. The mean starting salary for nurses is \\$67,694 nationally (“Staff nurse -,” 2013). The standard deviation is approximately \$10,333. The starting salary is not normally distributed but it is mound shaped. A sample of 42 starting salaries for nurses is taken.
 - a. State the random variable.
 - b. What is the mean of the sample mean?
 - c. What is the standard deviation of the sample mean?
 - d. What is the shape of the sampling distribution of the sample mean? Why?
 - e. Find the probability that the sample mean is more than \$75,000.
 - f. Find the probability that the sample mean is less than \$60,000.
 - g. If you did find a sample mean of more than \$75,000 would you find that unusual? What could you conclude?
 4. According to the WHO MONICA Project the mean blood pressure for people in China is 128 mmHg with a standard deviation of 23 mmHg (Kuulasmaa, Hense & Tolonen, 1998). Blood pressure is normally distributed.
 - a. State the random variable.
 - b. Suppose a sample of size 15 is taken. State the shape of the distribution of the sample mean.
 - c. Suppose a sample of size 15 is taken. State the mean of the sample mean.
 - d. Suppose a sample of size 15 is taken. State the standard deviation of the sample mean.
 - e. Suppose a sample of size 15 is taken. Find the probability that the sample mean blood pressure is more than 135 mmHg.
 - f. Would it be unusual to find a sample mean of 15 people in China of more than 135 mmHg? Why or why not?
 - g. If you did find a sample mean for 15 people in China to be more than 135 mmHg, what might you conclude?
 5. The size of fish is very important to commercial fishing. A study conducted in 2012 found the length of Atlantic cod caught in nets in Karlskrona to have a mean of 49.9 cm and a standard deviation of 3.74 cm (Ovegard, Berndt & Lunneryd, 2012). The length of fish is normally distributed. A sample of 15 fish is taken.

- a. State the random variable.
 - b. Find the mean of the sample mean.
 - c. Find the standard deviation of the sample mean
 - d. What is the shape of the distribution of the sample mean? Why?
 - e. Find the probability that the sample mean length of the Atlantic cod is less than 52 cm.
 - f. Find the probability that the sample mean length of the Atlantic cod is more than 74 cm.
 - g. If you found sample mean length for Atlantic cod to be more than 74 cm, what could you conclude?
6. The mean cholesterol levels of women age 45-59 in Ghana, Nigeria, and Seychelles is 5.1 mmol/l and the standard deviation is 1.0 mmol/l (Lawes, Hoorn, Law & Rodgers, 2004). Assume that cholesterol levels are normally distributed.
- a. State the random variable.
 - b. Find the probability that a woman age 45-59 in Ghana has a cholesterol level above 6.2 mmol/l (considered a high level).
 - c. Suppose doctors decide to test the woman's cholesterol level again and average the two values. Find the probability that this woman's mean cholesterol level for the two tests is above 6.2 mmol/l.
 - d. Suppose doctors being very conservative decide to test the woman's cholesterol level a third time and average the three values. Find the probability that this woman's mean cholesterol level for the three tests is above 6.2 mmol/l.
 - e. If the sample mean cholesterol level for this woman after three tests is above 6.2 mmol/l, what could you conclude?
7. In the United States, males between the ages of 40 and 49 eat on average 103.1 g of fat every day with a standard deviation of 4.32 g ("What we eat," 2012). The amount of fat a person eats is not normally distributed but it is relatively mound shaped.
- a. State the random variable.
 - b. Find the probability that a sample mean amount of daily fat intake for 35 men age 40-59 in the U.S. is more than 100 g.
 - c. Find the probability that a sample mean amount of daily fat intake for 35 men age 40-59 in the U.S. is less than 93 g.
 - d. If you found a sample mean amount of daily fat intake for 35 men age 40-59 in the U.S. less than 93 g, what would you conclude?
8. A dishwasher has a mean life of 12 years with an estimated standard deviation of 1.25 years ("Appliance life expectancy," 2013). The life of a dishwasher is normally distributed. Suppose you are a manufacturer and you take a sample of 10 dishwashers that you made.
- a. State the random variable.
 - b. Find the mean of the sample mean.

- c. Find the standard deviation of the sample mean.
- d. What is the shape of the sampling distribution of the sample mean? Why?
- e. Find the probability that the sample mean of the dishwashers is less than 6 years.
- f. If you found the sample mean life of the 10 dishwashers to be less than 6 years, would you think that you have a problem with the manufacturing process? Why or why not?

8 One Sample Inference

Now that you have all this information about descriptive statistics and probabilities, it is time to start inferential statistics. There are two branches of inferential statistics: hypothesis testing and confidence intervals.

Hypothesis Testing: making a decision about a parameter(s) based on a statistic(s).

Confidence Interval: estimating a parameter(s) based on a statistic(s).

This chapter will describe hypothesis testing, but as was stated in Chapter 1, the American Statistical Association (ASA) is suggesting not discussing statistical significance and p-values. So this chapter is mostly for background to understand previously published studies.

8.1 Basics of Hypothesis Testing

To understand the process of a hypothesis tests, you need to first have an understanding of what a hypothesis is, which is an educated guess about a parameter. Once you have the hypothesis, you collect data and use the data to make a determination to see if there is enough evidence to show that the hypothesis is true. However, in hypothesis testing you actually assume something else is true, and then you look at your data to see how likely it is to get an event that your data demonstrates with that assumption. If the event is very unusual, then you might think that your assumption is actually false. If you are able to say this assumption is false, then your hypothesis must be true. This is known as a proof by contradiction. You assume the opposite of your hypothesis is true and show that it can't be true. If this happens, then your hypothesis must be true. All hypothesis tests go through the same process. Once you have the process down, then the concept is much easier. It is easier to see the process by looking at an example. Concepts that are needed will be detailed in this example.

8.1.1 Example: Basics of Hypothesis Testing

Suppose a manufacturer of the XJ35 battery claims the mean life of the battery is 500 days with a standard deviation of 25 days. You are the buyer of this battery and you think this claim is incorrect. You would like to test your belief because without a good reason you can't get out of your contract.

8.1.1.1 Solution

What do you do?

Well first, you should know what you are trying to measure. Define the random variable.

Let x = life of a XJ35 battery

Now you are not just trying to find different x values. You are trying to find what the true mean is. Since you are trying to find it, it must be unknown. You don't think it is 500 days. If you did, you wouldn't be doing any testing. The true mean, μ , is unknown. That means you should define that too.

Let μ = mean life of a XJ35 battery

Now what?

You may want to collect a sample. What kind of sample?

You could ask the manufacturers to give you batteries, but there is a chance that there could be some bias in the batteries they pick. To reduce the chance of bias, it is best to take a random sample.

How big should the sample be?

A sample of size 30 or more means that you can use the central limit theorem. Pick a sample of size 50.

Table ?? contains the data for the sample you collected:

```
Battery<- read.csv( "https://krkozak.github.io/MAT160/battery.csv")
knitr::kable(head(Battery))
```

Table 8.1: Data on Battery Life

life
491
485
503
492
482
490

Now what should you do? Looking at the data set, you see some of the times are above 500 and some are below. But looking at all of the numbers is too difficult. It might be helpful to calculate the mean for this sample.


```
df_stats(~life, data=Battery, mean)
```

```
response mean
1      life  490
```

The sample mean is 491.42 days. Looking at the sample mean, one might think that you are right. However, the standard deviation and the sample size also plays a role, so maybe you are wrong.

Before going any farther, it is time to formalize a few definitions.

You have a guess that the mean life of a battery is not 500 days. This is opposed to what the manufacturer claims. There really are two hypotheses, which are just guesses here — the one that the manufacturer claims and the one that you believe. It is helpful to have names for them.

Null Hypothesis: historical value, claim, or product specification. The symbol used is H_o .

Alternate Hypothesis: what you want to prove. This is what you want to accept as true when you reject the null hypothesis. There are two symbols that are commonly used for the alternative hypothesis: H_a or H_1 . The symbol H_a will be used in this book.

In general, the hypotheses look something like this:

$$H_o : \mu = \mu_o$$

$$H_a : \mu \neq \mu_o$$

where μ_o just represents the value that the claim says the population mean is actually equal to.

Also, H_a can be less than, greater than, or not equal to, though not equal to is more common these days.

For this problem:

$H_o : \mu = 500$ days, since the manufacturer says the mean life of a battery is 500 days.

$H_a : \mu \neq 500$ days, since you believe that the mean life of the battery is not 500 days.

Now back to the mean. You have a sample mean of 491.42 days. Is this different enough to believe that you are right and the manufacturer is wrong? How different does it have to be?

If you calculated a sample mean of 235 or 690, you would definitely believe the population mean is not 500. But even if you had a sample mean of 435 or 575 you would probably believe that the true mean was not 500. What about 475? or 535? Or 483? or 514? There is some point where you would stop being so sure that the population mean is not 500. That point

separates the values of where you are sure or pretty sure that the mean is not 500 from the area where you are not so sure. How do you find that point?

Well it depends on how much error you want to make. Of course you don't want to make any errors, but unfortunately that is unavoidable in statistics. You need to figure out how much error you made with your sample. Take the sample mean, and find the probability of getting another sample mean less than it, assuming for the moment that the manufacturer is right. The idea behind this is that you want to know what is the chance that you could have come up with your sample mean even if the population mean really is 500 days.

Chances are probabilities. So you want to find the probability that the sample mean of 491.42 is unusual given that the population mean is really 500 days. To compute this probability, you need to know how the sample mean is distributed. Since the sample size is at least 30, then you know the sample mean is approximately normally distributed. Now, you want to find the z -value. The z -value is $z = \frac{491.42 - 500}{\frac{25}{\sqrt{50}}} = -2.43$.

This is more than 2 standard deviations below the mean, so that seems that the sample mean is unusual. It might be helpful to find the probability though. Since you are saying that the sample mean is different from 500 days, then you are asking if it is greater than or less than. This means that you are in the tails of the normal curve. So the probability you want to find is the probability being more than 2.43 or less than -2.43 . This is $P(-2.43 < z) + P(z > 2.43) = 0.015$

```
pnorm(-2.43, 0, 1, lower.tail=TRUE)+pnorm(2.43, 0, 1, lower.tail=FALSE)
```

```
[1] 0.01509882
```

So the probability of being in the tails is 0.015. This probability is known as a p-value for probability-value. This is unusual, so it is unlikely to get a sample mean of 491.42 if the population mean is 500 days.

So it appears the assumption that the population mean is 500 days is wrong, and you can reject the manufacturer's claim.

But how do you quantify really small? Is 5% or 10% or 15% really small? How do you decide?

Before you answer that question, a couple more definitions are needed.

Test statistic: $z = \frac{\bar{x} - \mu_o}{\frac{\sigma}{\sqrt{n}}}$ since it is calculated as part of the testing of the hypothesis

p - value: probability that the test statistic will take on more extreme values than the observed test statistic, given that the null hypothesis is true. It is the probability that was calculated above.

Now, how small is small enough? To answer that, you really want to know the types of errors you can make.

There are actually only two errors that can be made. The first error is if you say that is false, when in fact it is true. This means you reject when was true. The second error is if you say that is true, when in fact it is false. This means you fail to reject when is false. The following table organizes this for you:

8.1.2 Type of errors:

Table 8.2: Type of errors

	Ho true	Ho false
Reject Ho	Type I error	no error
Fail to reject Ho	no error	Type II error

Thus

Type I Error is rejecting H_o when H_o is true, and

Type II Error is failing to reject H_o when is H_o false.

Since these are the errors, then one can define the probabilities attached to each error.

$$\alpha = P(\text{type I error}) = P(\text{rejecting } H_o \text{ given it is true})$$

$$\beta = P(\text{type II error}) = P(\text{failing to reject } H_o \text{ given it is false})$$

α is also called the **level of significance**.

Another common concept that is used is Power = $1 - \beta$

Now there is a relationship between α and β . They are not complements of each other. How are they related?

If α increases that means the chances of making a type I error will increase. It is more likely that a type I error will occur. It makes sense that you are less likely to make type II errors, only because you will be rejecting more often. You will be failing to reject less, and therefore, the chance of making a type II error will decrease. Thus, as α increases, β will decrease, and vice versa. That makes them seem like complements, but they aren't complements. What gives? Consider one more factor -- sample size.

Consider if you have a larger sample that is representative of the population, then it makes sense that you have more accuracy then with a smaller sample. Think of it this way, which would you trust more, a sample mean of 490 if you had a sample size of 35 or sample size of 350 (assuming a representative sample)? Of course the 350 because there are more data points

and so more accuracy. If you are more accurate, then there is less chance that you will make any error. By increasing the sample size of a representative sample, you decrease both α and β .

Summary of all of this:

1. For a certain sample size, α increases, β decreases.
2. For a certain level of significance, α , if n increases, β decreases.

Now how do you find α and β ? Well α is actually chosen. There are only two values that are usually picked for α : 0.01 and 0.05. It is very difficult to find β , so usually it isn't found. If you want to make sure it is small you take as large of a sample as you can afford provided it is a representative sample. This is one use of the Power. You want to be small and the Power of the test is large. The Power word sounds good.

Which pick of α do you pick? Well that depends on what you are working on. Remember in this example you are the buyer who is trying to get out of a contract to buy these batteries. If you create a type I error, you said that the batteries are bad when they aren't, most likely the manufacturer will sue you. You want to avoid this. You might pick α to be 0.01. This way you have a small chance of making a type I error. Of course this means you have more of a chance of making a type II error. No big deal right? What if the batteries are used in pacemakers and you tell the person that their pacemaker's batteries are good for 500 days when they actually last less, that might be bad. If you make a type II error, you say that the batteries do last 500 days when they last less, then you have the possibility of killing someone. You certainly do not want to do this. In this case you might want to pick α as 0.05. If both errors are equally bad, then pick α as 0.05.

The above discussion is why the choice of depends on what you are researching. As the researcher, you are the one that needs to decide what level to use based on your analysis of the consequences of making each error is.

If a type I error is really bad, then pick $\alpha = 0.01$.

If a type II error is really bad, then pick $\alpha = 0.05$

If neither error is bad, or both are equally bad, then pick $\alpha = 0.05$

Usually α is picked to be 0.05 in most cases.

The main thing is to always pick the α before you collect the data and start the test.

The above discussion was long, but it is really important information. If you don't know what the errors of the test are about, then there really is no point in making conclusions with the tests. Make sure you understand what the two errors are and what the probabilities are for them.

Now it is time to go back to the example and put this all together. This is the basic structure of testing a hypothesis, usually called a hypothesis test. Since this one has a test statistic

involving z , it is also called a z -test. And since there is only one sample, it is usually called a one-sample z -test.

8.1.3 Example: Battery Example Revisited.

Steps of a hypothesis test:

1. State the random variable and the parameter in words
2. State the null and alternative hypothesis and the level of significance
3. State and check the conditions for a hypothesis test
4. Find the sample statistic, test statistic, and p-value
5. Conclusion:
6. Interpretation:

8.1.3.1 Solution

1. State the random variable and the parameter in words

x = life of battery

μ = mean life of a XJ35 battery

2. State the null and alternative hypothesis and the level of significance

$$H_o : \mu = 500$$

$$H_a : \mu \neq 500$$

$$\alpha = 0.05 \text{ (from above discussion about consequences)}$$

3. State and check the conditions for a hypothesis test

Every hypothesis has some conditions that be met to make sure that the results of the hypothesis are valid. The conditions are different for each test. This test has the following conditions.

- a. A random sample of size n is taken.

This occurred in this example, since it was stated that a random sample of 50 battery lives were taken.

- b. The population standard deviation is known.

This is true, since it was given in the problem.

- c. The sample size is at least 30 or the population of the random variable is normally distributed.

The sample size was 30, so this condition is met.

4. Find the sample statistic, test statistic, and p-value

The test statistic depends on how many samples there are, what parameter you are testing, and conditions that need to be checked. In this case, there is one sample and you are testing the mean. The conditions were checked above.

Sample statistic:

```
df_stats(~life, data=Battery, mean)
```

```
response mean
1    life  490
```

Test statistic: The z-value is $z = \frac{491.42 - 400}{\frac{25}{\sqrt{n}}} = -2.43$.

p-value: $P(-2.43 < z) + P(z > 2.43) = 0.015$

5. Conclusion:

Now what? Well, this p-value is 0.015. This is a lot smaller than the amount of error you would accept in the problem $\alpha = 0.05$. That means that finding a sample mean less than 490 days is unusual to happen if it is true. This should make you think that it is not true. You should reject H_o .

In fact, in general:

Reject H_o if the p-value $< \alpha$

Fail to reject H_o if the p-value $\geq \alpha$.

6. Interpretation:

Since you rejected H_o , what does this mean in the real world? That is what goes in the interpretation. Since you rejected the claim by the manufacturer that the mean life of the batteries is 500 days, then you now can believe that your hypothesis was correct. In other words, there is enough evidence to support that the mean life of the battery is less than 500 days.

Now that you know that the batteries last less than 500 days, should you cancel the contract? Statistically, there is evidence that the batteries do not last as long as the manufacturer says they should. However, based on this sample there are only ten days less on average than the

batteries last. There may not be practical significance in this case. Ten days do not seem like a large difference. In reality, if the batteries are used in pacemakers, then you would probably tell the patient to have the batteries replaced every year. You have a large buffer whether the batteries last 490 days or 500 days. It seems that it might not be worth it to break the contract over ten days. What if the 10 days was practically significant? Are there any other things you should consider? You might look at the business relationship with the manufacturer. You might also look at how much it would cost to find a new manufacturer. These are also questions to consider before making any changes. What this discussion should show you is that just because a hypothesis has statistical significance does not mean it has practical significance. The hypothesis test is just one part of a research process. There are other pieces that you need to consider.

That's it. That is what a hypothesis test looks like. All hypothesis tests are done with the same six steps. Those general six steps are outlined below.

8.1.4 Steps for hypothesis test

1. State the random variable and the parameter in words. This is where you are defining what the unknowns are in this problem.

x = random variable

μ = mean of random variable, if the parameter of interest is the mean. There are other parameters you can test, and you would use the appropriate symbol for that parameter.

2. State the null and alternative hypotheses and the level of significance

$H_o : \mu = \mu_o$, where μ_o is the known mean

$H_a : \mu \neq \mu_o$, You can replace \neq with $<$ or $>$ but usually you use \neq

Also, state your level here.

3. State and check the conditions for a hypothesis test

Each hypothesis test has its own conditions. They will be stated when the different hypothesis tests are discussed.

4. Find the sample statistic, test statistic, and p-value

This depends on what parameter you are working with, how many samples, and the conditions of the test. Technology will be used to find the sample statistic, test statistic, and p-value.

5. Conclusion

This is where you write reject H_o or fail to reject H_o . The rule is: if the p-value $< \alpha$, then reject H_o . If the p-value $\geq \alpha$, then fail to reject H_o .

6. Interpretation

This is where you interpret in real world terms the conclusion to the test. The conclusion for a hypothesis test is that you either have enough evidence to support H_a , or you do not have enough evidence to support H_a .

Sorry, one more concept about the conclusion and interpretation. First, the conclusion is that you reject or you fail to reject H_0 . Why was it said like this? It is because you never accept the null hypothesis. If you wanted to accept the null hypothesis, then why do the test in the first place? In the interpretation, you either have enough evidence to support H_a , or you do not have enough evidence to support H_a . You wouldn't want to go to all this work and then find out you wanted to accept the claim. Why go through the trouble? You always want to have enough evidence to support the alternative hypothesis. Sometimes you can do that and sometimes you can't. If you don't have enough evidence to support H_a , it doesn't mean you support the null hypothesis; it just means you can't support the alternative hypothesis. Here is an example to demonstrate this.

8.1.5 Example: Conclusions in Hypothesis Tests

In the U.S. court system a jury trial could be set up as a hypothesis test. To really help you see how this works, let's use OJ Simpson as an example. In the court system, a person is presumed innocent until he/she is proven guilty, and this is your null hypothesis. OJ Simpson was a football player in the 1970s. In 1994 his ex-wife and her friend were killed. OJ Simpson was accused of the crime, and in 1995 the case was tried. The prosecutors wanted to prove OJ was guilty of killing his wife and her friend, and that is the alternative hypothesis. In this case, a verdict of not guilty was given. That does not mean that he is innocent of this crime. It means there was not enough evidence to prove he was guilty. Many people believe that OJ was guilty of this crime, but the jury did not feel that the evidence presented was enough to show there was guilt. The verdict in a jury trial is always guilty or not guilty!

The same is true in a hypothesis test. There is either enough or not enough evidence to support the alternative hypothesis. It is not that you proved the null hypothesis true.

When identifying hypothesis, it is important to state your random variable and the appropriate parameter you want to make a decision about. If you count something, then the random variable is the number of whatever you counted. The parameter is the proportion of what you counted. If the random variable is something you measured, then the parameter is the mean of what you measured. (Note: there are other parameters you can calculate, and some analysis of those will be presented in later chapters.)

8.1.6 Example: Stating Hypotheses

Identify the hypotheses necessary to test the following statements:

- a. The average salary of a teacher is different from \\$30,000.
- b. The proportion of students who like math is not 10%.
- c. The average age of students in this class differs from 21.

8.1.6.1 Solution

- a. The average salary of a teacher is different from \\$30,000.

x = salary of teacher

μ = mean salary of teacher

The guess is that $\mu \neq 30000$ and that is the alternative hypothesis.

The null hypothesis has the same parameter and number with an equal sign.

$$H_o : \mu = 30000 \quad H_a : \mu \neq 30000$$

- b. The proportion of students who like math is not 10%.

x = number of students who like math

p = proportion of students who like math

The guess is that p is not 0.10 and that is the alternative hypothesis. $H_a : p \neq 0.10$ and the null hypothesis would be $H_o : p = 0.10$

- c. The average age of students in this class differs from 21.

x = age of students in this class

μ = mean age of students in this class

The guess is that $\mu \neq 21$ and that is the alternative hypothesis. $H_a : \mu \neq 21$ and the null hypothesis would be $H_o : \mu = 21$

8.1.7 Example: Stating Type I and II Errors and Picking Level of Significance

- a. The plant-breeding department at a major university developed a new hybrid raspberry plant called YumYum Berry. Based on research data, the claim is made that from the time shoots are planted 90 days on average are required to obtain the first berry with a standard deviation of 9.2 days. A corporation that is interested in marketing the product tests 60 shoots by planting them and recording the number of days before each plant produces its first berry. The sample mean is 92.3 days. The corporation wants to know if the mean number of days is more than the 90 days claimed. State the type I and type II errors in terms of this problem, consequences of each error, and state which level of significance to use.

- b. A concern was raised in Australia that the percentage of deaths of Aboriginal prisoners was higher than the percent of deaths of non-indigenous prisoners, which is 0.27%. State the type I and type II errors in terms of this problem, consequences of each error, and state which level of significance to use.

8.1.7.1 Solution

- a. The plant-breeding department at a major university developed a new hybrid raspberry plant called YumYum Berry. Based on research data, the claim is made that from the time shoots are planted 90 days on average are required to obtain the first berry with a standard deviation of 9.2 days. A corporation that is interested in marketing the product tests 60 shoots by planting them and recording the number of days before each plant produces its first berry. The sample mean is 92.3 days. The corporation wants to know if the mean number of days is more than the 90 days claimed. State the type I and type II errors in terms of this problem, consequences of each error, and state which level of significance to use.

x = time to first berry for YumYum Berry plant

μ = mean time to first berry for YumYum Berry plant

Type I Error: If the corporation does a type I error, then they will say that the plants take longer to produce than 90 days when they don't. They probably will not want to market the plants if they think they will take longer. They will not market them even though in reality the plants do produce in 90 days. They may have loss of future earnings, but that is all.

Type II error: The corporation do not say that the plants take longer than 90 days to produce when they do take longer. Most likely they will market the plants. The plants will take longer, and so customers might get upset and then the company would get a bad reputation. This would be really bad for the company.

Level of significance: It appears that the corporation would not want to make a type II error. Pick a 5% level of significance, $\alpha = 0.05$.

- b. A concern was raised in Australia that the percentage of deaths of Aboriginal prisoners was higher than the percent of deaths of non-indigenous prisoners, which is 0.27%. State the type I and type II errors in terms of this problem, consequences of each error, and state which level of significance to use.

x = number of Aboriginal prisoners who have died

p = proportion of Aboriginal prisoners who have died

Type I error: Rejecting that the proportion of Aboriginal prisoners who died was 0.27%, when in fact it was 0.27%. This would mean you would say there is a problem when there isn't one.

You could anger the Aboriginal community, and spend time and energy researching something that isn't a problem.

Type II error: Failing to reject that the proportion of Aboriginal prisoners who died was 0.27%, when in fact it is higher than 0.27%. This would mean that you wouldn't think there was a problem with Aboriginal prisoners dying when there really is a problem. You risk causing deaths when there could be a way to avoid them.

Level of significance: It appears that both errors may be issues in this case. You wouldn't want to anger the Aboriginal community when there isn't an issue, and you wouldn't want people to die when there may be a way to stop it. It may be best to pick a 5% level of significance, $\alpha = 0.05$.

Hint -- hypothesis testing is really easy if you follow the same recipe every time. The only differences in the various problems are the conditions of the test and the test statistic you calculate so you can find the p-value. Do the same steps, in the same order, with the same words, every time and these problems become very easy.

8.1.8 Homework for Basics of Hypothesis Testing Section

For the problems in this section, a question is being asked. This is to help you understand what the hypotheses are. You are not to run any hypothesis tests nor come up with any conclusions in this section.

1. The Arizona Republic/Morrison/Cronkite News poll published on Monday, October 20, 2016, found 390 of the registered voters surveyed favor Proposition 205, which would legalize marijuana for adults. The statewide telephone poll surveyed 779 registered voters between Oct. 10 and Oct. 15. (Sanchez, 2016) Fifty-five percent of Colorado residents supported the legalization of marijuana. Does the data provide evidence that the percentage of Arizona residents who support legalization of marijuana is different from the proportion of Colorado residents who support it? State the random variable, population parameter, and hypotheses.
2. According to the February 2008 Federal Trade Commission report on consumer fraud and identity theft, 23% of all complaints in 2007 were for identity theft. In that year, Alaska had 321 complaints of identity theft out of 1,432 consumer complaints ("Consumer fraud and," 2008). Does this data provide enough evidence to show that Alaska had a different proportion of identity theft than 23%? State the random variable, population parameter, and hypotheses.
3. The Kyoto Protocol was signed in 1997, and required countries to start reducing their carbon emissions. The protocol became enforceable in February 2005. In 2004, the mean CO₂ emission was 4.87 metric tons per capita. Is there enough evidence to show that the mean CO₂ emission is different in 2010 than in 2004? State the random variable, population parameter, and hypotheses.

4. The FDA regulates that fish that is consumed is allowed to contain 1.0 mg/kg of mercury. In Florida, bass fish were collected in 53 different lakes to measure the amount of mercury in the fish. Do the data provide enough evidence to show that the fish in Florida lakes has a different amount of mercury than the allowable amount? State the random variable, population parameter, and hypotheses.
5. The Arizona Republic/Morrison/Cronkite News poll published on Monday, October 20, 2016, found 390 of the registered voters surveyed favor Proposition 205, which would legalize marijuana for adults. The statewide telephone poll surveyed 779 registered voters between Oct. 10 and Oct. 15. (Sanchez, 2016) Fifty-five percent of Colorado residents supported the legalization of marijuana. Does the data provide evidence that the percentage of Arizona residents who support legalization of marijuana is different from the proportion of Colorado residents who support it. State the type I and type II errors in this case, consequences of each error type for this situation from the perspective of the manufacturer, and the appropriate alpha level to use. State why you picked this alpha level.
6. According to the February 2008 Federal Trade Commission report on consumer fraud and identity theft, 23% of all complaints in 2007 were for identity theft. In that year, Alaska had 321 complaints of identity theft out of 1,432 consumer complaints (\“Consumer fraud and,\” 2008). Does this data provide enough evidence to show that Alaska had a different proportion of identity theft than 23%? State the type I and type II errors in this case, consequences of each error type for this situation from the perspective of the state of Alaska, and the appropriate alpha level to use. State why you picked this alpha level.
7. The Kyoto Protocol was signed in 1997, and required countries to start reducing their carbon emissions. The protocol became enforceable in February 2005. In 2004, the mean CO₂ emission was 4.87 metric tons per capita. Is there enough evidence to show that the mean CO₂ emission is lower in 2010 than in 2004? State the type I and type II errors in this case, consequences of each error type for this situation from the perspective of the agency overseeing the protocol, and the appropriate alpha level to use. State why you picked this alpha level.
8. The FDA regulates that fish that is consumed is allowed to contain 1.0 mg/kg of mercury. In Florida, bass fish were collected in 53 different lakes to measure the amount of mercury in the fish. Do the data provide enough evidence to show that the fish in Florida lakes has different amount of mercury than the allowable amount? State the type I and type II errors in this case, consequences of each error type for this situation from the perspective of the FDA, and the appropriate alpha level to use. State why you picked this alpha level.

8.2 One-Sample Proportion Test

There are many different parameters that you can test. There is a test for the mean, such as was introduced with the t -test. There is also a test for the population proportion, p . This is where you might be curious if the proportion of students who smoke at your school is lower than the proportion in your area. Or you could question if the proportion of accidents caused by teenage drivers who do not have a drivers' education class is more than the national proportion.

To test a population proportion, there are a few things that need to be defined first. Usually, Greek letters are used for parameters and Latin letters for statistics. When talking about proportions, it makes sense to use p for proportion. The Greek letter for p is π , but that is too confusing to use. Instead, it is best to use p for the population proportion. That means that a different symbol is needed for the sample proportion. The convention is to use, \hat{p} , known as p-hat. This way you know that p is the population proportion, and that \hat{p} is the sample proportion related to it.

Now proportion tests are about looking for the percentage of individuals who have a particular attribute. You are really looking for the number of successes that happen. Thus, a proportion test involves a binomial distribution.

8.2.1 Hypothesis Test for One Population Proportion (1-Prop Test)

1. State the random variable and the parameter in words.

x = number of successes

p = proportion of successes

2. State the null and alternative hypotheses and the level of significance

$H_o : p = p_o$, where p_o is the known proportion

$H_a : p \neq p_o$, you can also use $<$ or $>$, but \neq is the more common one to use.

Also, state your α level here.

3. State and check the conditions for a hypothesis test

- a. State: A simple random sample of size n is taken. Check: describe how the sample was collected
- b. State: The conditions for the binomial experiment are satisfied. Check: Show all four properties are true.

- c. State: The sampling distribution of \hat{p} is normally distributed. Check: you need to show that $p * n \geq 5$ and $q * n \geq 5$, where $q = 1 - p$. If this requirement is true, then the sampling distribution of \hat{p} is well approximated by a normal curve.

4. Find the sample statistic, test statistic, and p-value

This will be computed on r Studio using the command

```
prop.test(r, n, p=what_Ho_says)
```

where r =observed number of successes and n = number of trials.

5. Conclusion

This is where you write reject or fail to reject H_o . The rule is: if the p-value $< \alpha$, then reject H_o . If the p-value $\geq \alpha$, then fail to reject H_o .

6. Interpretation

This is where you interpret in real world terms the conclusion to the test. The conclusion for a hypothesis test is that you either have enough evidence to support H_a , or you do not have enough evidence to support H_a .

8.2.2 Example: Hypothesis Test for One Proportion

A concern was raised in Australia that the percentage of deaths of Aboriginal prisoners was different than the percent of deaths of non-Aboriginal prisoners, which is 0.27%. A sample of six years (1990-1995) of data was collected, and it was found that out of 14,495 Aboriginal prisoners, 51 died (\“Indigenous deaths in,\” 1996). Do the data provide enough evidence to show that the proportion of deaths of Aboriginal prisoners is different from 0.27%?

8.2.2.1 Solution

1. State the random variable and the parameter in words.

x = number of Aboriginal prisoners who die

p = proportion of Aboriginal prisoners who die

2. State the null and alternative hypotheses and the level of significance

$$H_o : p = 0.0027$$

$$H_a : p \neq 0.0027$$

From Example: Stating Type I and II Errors and Picking Level of Significance part b, the argument was made to pick 5% for the level of significance. So $\alpha = 0.05$

3. State and check the conditions for a hypothesis test
 - a. A simple random sample of 14,495 Aboriginal prisoners was taken. Check: The sample was not a random sample, since it was data from six years. It is the numbers for all prisoners in these six years, but the six years were not picked at random. Unless there was something special about the six years that were chosen, the sample is probably a representative sample. This condition is probably met.
 - b. The properties of a binomial experiment are met. There are 14,495 prisoners in this case. Check: The prisoners are all Aboriginals, so you are not mixing Aboriginal with non-Aboriginal prisoners. There are only two outcomes, either the prisoner dies or doesn't. The chance that one prisoner dies over another may not be constant, but if you consider all prisoners the same, then it may be close to the same probability. Thus the conditions for the binomial distribution are satisfied
 - c. The sampling distribution of \hat{p} can be approximated with a normal distributed. Check: In this case $p = 0.0027$ and $n = 14,495$. $n * p = 39.1365 \geq 5$ and $n * q = 14455.86 \geq 5$. So, the sampling distribution for \hat{p} is normally distributed.
4. Find the sample statistic, test statistic, and p-value

Use the following command in rStudio:

```
prop.test(51, 14495, p=0.0027)
```

1-sample proportions test with continuity correction

```
data: 51 out of 14495
X-squared = 3.3084, df = 1, p-value = 0.06893
alternative hypothesis: true p is not equal to 0.0027
95 percent confidence interval:
 0.002647440 0.004661881
sample estimates:
      p
0.003518455
```

Sample Proportion: $\hat{p} = 0.0035$

Test Statistic: $\chi^2 = 3.3085$

p-value: $p - value = 0.06893$

5. Conclusion

Since the $p - value \geq 0.05$, then fail to reject H_o .

6. Interpretation

There is not enough evidence to support that the proportion of deaths of Aboriginal prisoners is different from non-Aboriginal prisoners.

8.2.3 Example: Hypothesis Test for One Proportion

A researcher who is studying the effects of income levels on breastfeeding of infants hypothesizes that countries with a low income level have a different rate of infant breastfeeding than higher income countries. It is known that in Germany, considered a high-income country by the World Bank, 22% of all babies are breastfeed. In Tajikistan, considered a low-income country by the World Bank, researchers found that in a random sample of 500 new mothers that 125 were breastfeeding their infant. At the 5% level of significance, does this show that low-income countries have a different incident of breastfeeding?

8.2.3.1 Solution

1. State you random variable and the parameter in words.

x = number of woman who breastfeed in a low-income country

p = proportion of woman who breastfeed in a low-income country

2. State the null and alternative hypotheses and the level of significance

$$H_o : p = 0.22$$

$$H_a : p \neq 0.22$$

$$\alpha = 0.05$$

3. State and check the conditions for a hypothesis test

- a. A simple random sample of 500 breastfeeding habits of woman in a low-income country was taken. Check: This was stated in the problem.
- b. The properties of a Binomial Experiment have been met. Check: There were 500 women in the study. The women are considered identical, though they probably have some differences. There are only two outcomes, either the woman breastfeeds or she doesn't. The probability of a woman breastfeeding is probably not the same for each woman, but it is probably not very different for each woman. The conditions for the binomial distribution are satisfied

- c. The sampling distribution of \hat{p} can be approximated with a normal distributed. Check: In this case, $n = 500$ and $p = 0.22$. $n * p = 110 \geq 5$ and $n * q = 390 \geq 5$, so the sampling distribution of \hat{p} is well approximated by a normal curve.

4. Find the sample statistic, test statistic, and p-value

On r studio, use the following command

```
prop_test(125, 500, p=0.22)
```

1-sample proportions test with continuity correction

```
data: 125 out of 500
X-squared = 2.4505, df = 1, p-value = 0.1175
alternative hypothesis: true p is not equal to 0.22
95 percent confidence interval:
 0.2131062 0.2908059
sample estimates:
      p
0.25
```

Sample Statistic: $\hat{p} = 0.25$

test Statistic: $\chi^2 = 2.4505$

p-value: $p - value = 0.1175$

5. Conclusion

Since the p-value is more than 0.05, you fail to reject H_o .

6. Interpretation

There is not enough evidence to support that the proportion of women who breastfeed in low-income countries is different from the proportion of women in high-income countries who breastfeed.

Notice, the conclusion is that there wasn't enough evidence to support H_a . The conclusion was not that you support H_o . There are many reasons why you can't say that H_o is true. It could be that the countries you chose were not very representative of what truly happens. If you instead looked at all high-income countries and compared them to low-income countries, you might have different results. It could also be that the sample you collected in the low-income country was not representative. It could also be that income level is not an indication of breastfeeding habits. It could be that the sample that was taken didn't show evidence but

another sample would show evidence. There could be other factors involved. This is why you can't say that you support H_o . There are too many other factors that could be the reason that you failed to reject H_o .

8.2.4 Homework for One-Sample Proportion Test Section

In each problem show all steps of the hypothesis test. If some of the conditions are not met, note that the results of the test may not be correct and then continue the process of the hypothesis test.

1. The Arizona Republic/Morrison/Cronkite News poll published on Monday, October 20, 2016, found 390 of the registered voters surveyed favor Proposition 205, which would legalize marijuana for adults. The statewide telephone poll surveyed 779 registered voters between Oct. 10 and Oct. 15. (Sanchez, 2016) Fifty-five percent of Colorado residents supported the legalization of marijuana. Does the data provide evidence that the percentage of Arizona residents who support legalization of marijuana is different from the proportion of Colorado residents who support it. Test at the 1% level.
2. In July of 1997, Australians were asked if they thought unemployment would increase, and 47% thought that it would increase. In November of 1997, they were asked again. At that time 284 out of 631 said that they thought unemployment would increase (\“Morgan Gallup poll,\” 2013). At the 5% level, is there enough evidence to show that the proportion of Australians in November 1997 who believe unemployment would increase is different from the proportion who felt it would increase in July 1997?
3. According to the February 2008 Federal Trade Commission report on consumer fraud and identity theft, 23% of all complaints in 2007 were for identity theft. In that year, Arkansas had 1,601 complaints of identity theft out of 3,482 consumer complaints (\“Consumer fraud and,\” 2008). Does this data provide enough evidence to show that Arkansas had a different percentage of identity theft than 23%? Test at the 5% level.
4. According to the February 2008 Federal Trade Commission report on consumer fraud and identity theft, 23% of all complaints in 2007 were for identity theft. In that year, Alaska had 321 complaints of identity theft out of 1,432 consumer complaints (\“Consumer fraud and,\” 2008). Does this data provide enough evidence to show that Alaska had a different proportion of identity theft than 23%? Test at the 5% level.
5. In 2001, the Gallup poll found that 81% of American adults believed that there was a conspiracy in the death of President Kennedy. In 2013, the Gallup poll asked 1,039 American adults if they believe there was a conspiracy in the assassination, and found that 634 believe there was a conspiracy (\“Gallup news service,\” 2013). Do the data show that the proportion of Americans who believe in this conspiracy has changed? Test at the 1% level.

6. In 2008, there were 507 children in Arizona out of 32,601 who were diagnosed with Autism Spectrum Disorder (ASD) (“Autism and developmental,” 2008). Nationally 1 in 88 children are diagnosed with ASD (“CDC features -,” 2013). Is there sufficient data to show that the incident of ASD is different in Arizona than nationally? Test at the 1% level.

8.3 One-Sample Test for the Mean

It is time to go back to look at the test for the mean that was introduced in section 7.1 called the z -test. In the example, you knew what the population standard deviation, σ , was. What if you don't know σ ?

If you don't know σ , then you don't know the sampling distribution of the mean. Can it be found another way? The answer is of course, yes. One way is to use a method called resampling. The following example explains how resampling is performed.

8.3.1 Example: Resampling

A random sample of 10 body mass index (BMI) were taken from the NHANES Data frame. The mean BMI of Americans is 27.2 kg/m^2 . Is there evidence that Americans have a different BMI from people in Australia. Test at the 5% level.

8.3.1.1 Solution

The standard deviation of BMI is not known for Americans. To answer this questions, first look at the sample from NHANES Table ??.

```
sample_NHANES_10<-  
  NHANES |>  
  slice_sample(n=10)  
knitr::kable(head(sample_NHANES_10))
```

Table 8.3: Sample of size 10 from NHA

[illegible]

The mean BMI from this sample is

```
df_stats(~BMI, data=sample NHANES 10, mean)
```

	response	mean
1	BMI	26.751

The sample mean for Americans is different from the mean BMI for Australians, but could it just be by chance. Suppose you take another sample of size 10, but you only have these 10 BMIs to work with. So how could you do this. One way is to assume that the sample you took is representative of the entire population, and so you create a population by copying this sample over and over again. So you could have over 1000 copies of this sample of 10 BMIs. Then take a sample of size 10 from this created population. When doing this, you could conceivably choose the same number several times that was in the original sample and not choose some of the numbers that were in the original sample. Instead of physically creating this new population, you could just take samples from your original sample but with replacement. This means that you randomly pick the first number, record it, and then put it back that value back before collecting the next number. This kind a sampling is called **randomization sampling**. A sample using randomization could be Table ??.

```
knitr::kable(resample(sample_NHANES_10))
```

Table 8.4: Resample from NHANES S

[illegible]

Notice that some of the unit of observations are repeated. That is what happens when you resample. Now one resampling isn't enough. So you want to resample many times so you can create a resampling distribution Figure ??

```
# mutate NHANES to subtract 27.2 (Australia's BMI) from US BMI measurements
mutate_NHANES <- NHANES |>
  mutate(NewBMI=BMI-27.2)
# Generate the single sample
Single_sample<- mutate_NHANES |>
  sample_n(size = 10)
```

```

#Calculate the mean age of the single sample
Single_sample_mean <-
  Single_sample |>
  df_stats( ~ NewBMI, means = mean)
#Take 200 resamples from the single sample
Trials_resample <-
  do(200) * { Single_sample |>
    resample() |>
    df_stats( ~ NewBMI, means = mean) }
# Plot the resample distribution of means
gf_density( ~ means, data = Trials_resample, bins = 10) |>
  gf_lims(x = c(-5, 10)) |>
  gf_labs(title = "Resampling Distribution") |>
  gf_vline(data = Single_sample_mean, xintercept = ~ means, color="red")
df_stats( ~ means, data = Trials_resample, mean, sd)

```

```

      response      mean      sd
1      means -3.56641 2.338371

```

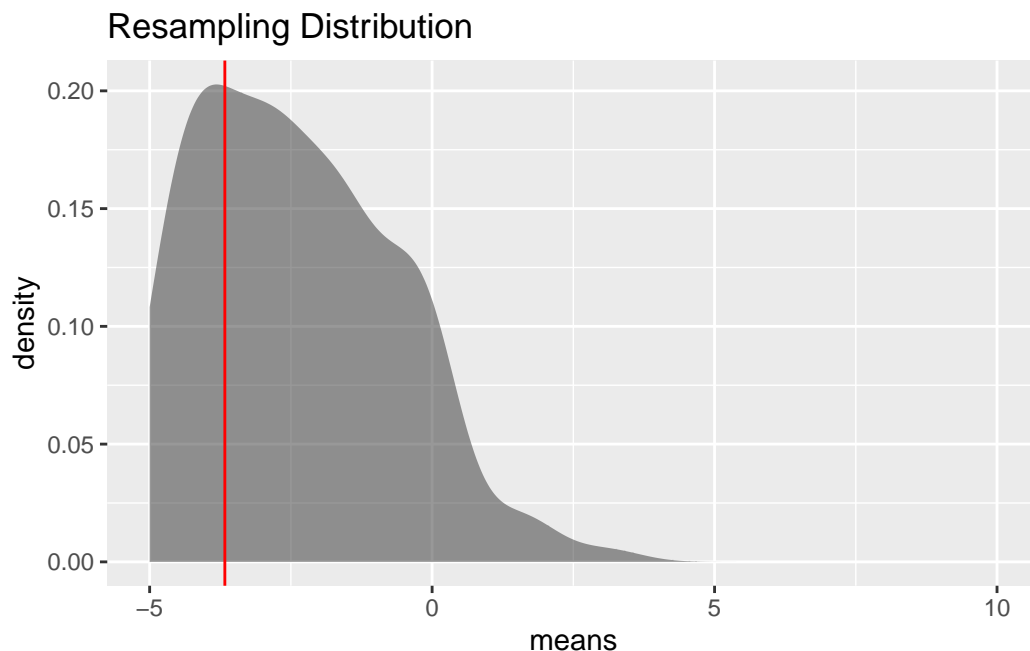


Figure 8.1: Resampling distribution of mean BMI with sample size 10

Notice the sample mean from the resampling is very close to 0, so that means that the US BMI are not that different from the Australian BMI. There doesn't seem to be enough evidence to

show that the US BMI is different from the Australian BMI. One note, the sample size used here was 10 so you could see the sample, but really the sample size should be more than 100 for this method to be valid.

So this is one way to answer the question about if there is evidence to show a population mean is different from a value. This is actually the method that Ronald Fisher developed when he created all the foundation work that he did in statistics in the early 1900s. However, at the time, computers didn't exist, so taking 100 resampling samples was not possible at that time. So other methods had to be developed that could be computed during that time. One method was developed by William (W.S) Gossett, a Chemist who worked for Guinness as their head brewer. Gossett developed a distribution called the Student's T-distribution. His process was to use the sample standard deviation, s , as an approximation of σ . This means the test statistic is now $t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$. This new test statistic is actually distributed as a Student's t-distribution, developed by W.S. Gossett. There are some conditions that must be made for this formula to be a Student's t-distribution. These are outlined in the following theorem. Note: the t-distribution is called the Student's t-distribution because that is the name he published under because he couldn't publish under his own name due to his employer not wanting him to publish under his own name. His employer by the way was Guinness and they didn't want competitors knowing they had a chemist/statistician working for them. It is not called the Student's t-distribution because it is only used by students.

Theorem: If the following conditions are met

- a. A random sample of size n is taken.
- b. The distribution of the random variable is normal.

Then the distribution of is a Student's t-distribution with $n - 1$ degrees of freedom.

Explanation of degrees of freedom: Recall the formula for sample standard deviation is $\sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$. Notice the denominator is $n - 1$. This is the same as the degrees of freedom. This is no accident. The reason the denominator and the degrees of freedom are both comes from how the standard deviation is calculated. First you take each data value and subtract \bar{x} . If you add up all of these new values, you will get 0. This must happen. Since it must happen, the first $n - 1$ data values you have "freedom of choice", but the n th data value, you have no freedom to choose. Hence, you have $n - 1$ degrees of freedom. Another way to think about it is that if you five people and five chairs, the first four people have a choice of where they are sitting, but the last person does not. They have no freedom of where to sit. Only $n - 1$ people have freedom of choice.

The Student's t-distribution is bell-shape that is more spread out than the normal distribution. There are many t -distributions, one for each different degree of freedom.

Figure ?? is of the normal distribution and the Student's t-distribution for $df = 1$, $df = 3$, $df=8$, $df=30$.

Comparison of t Distributions

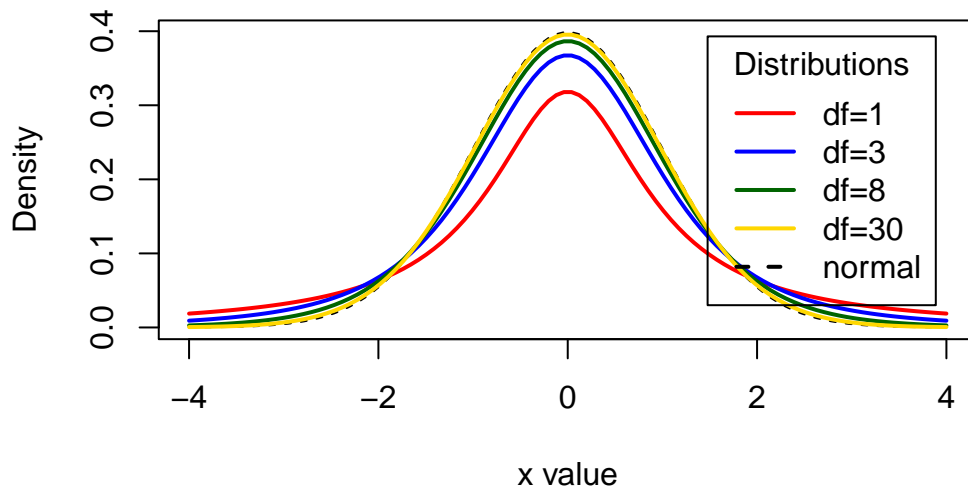


Figure 8.2: Typical Student t-Distributions

As the degrees of freedom increases, the student's t-distribution looks more like the normal distribution.

To find probabilities for the t-distribution, again technology can do this for you. There are many technologies out there that you can use.

8.3.2 Hypothesis Test for One Population Mean (t-Test)

1. State the random variable and the parameter in words.

x = random variable

μ = mean of random variable

2. State the null and alternative hypotheses and the level of significance

$H_o : \mu = \mu_o$, where μ_o is the known mean

$H_a : \mu \neq \mu_o$, you can also use $<$ or $>$, but \neq is the more modern one to use.

Also, state your α level here.

3. State and check the conditions for a hypothesis test

- a. State: A random sample of size n is taken. Check: Describe the process taken to collect the sample.

- b. State: The population of the random variable is normally distributed. Check: examine density graph and normal quantile plot. Note: The t-test is fairly robust to the condition if the sample size is large. This means that if this condition isn't met, but your sample size is quite large, then the results of the t-test are valid.

4. Find the sample statistic, test statistic, and p-value

On rStudio, the command is

```
t.test(~variable, data=data_frame, mu=what_Ho_says)
```

5. Conclusion

This is where you write reject or fail to reject H_o . The rule is: if the p-value $< \alpha$, then reject H_o . If the p-value $\geq \alpha$, then fail to reject H_o .

6. Interpretation

This is where you interpret in real world terms the conclusion to the test. The conclusion for a hypothesis test is that you either have enough evidence to support H_a , or you do not have enough evidence to support H_a .

Note: if the conditions behind this test are not valid, then the conclusions you make from the test are not valid. If you do not have a random sample, that is your fault. Make sure the sample you take is as random as you can make it following sampling techniques from chapter 1. If the population of the random variable is not normal, then take a larger sample. If you cannot afford to do that, or if it is not logistically possible, then you do different tests called non-parametric tests or you can try resampling. The advantage of resampling is that you don't need to know the underlying distribution of the random variable.

8.3.3 Example: Test of the Mean Using One Sample T-test

A random sample of 50 body mass index (BMI) were taken from the NHANES Data frame Table ???. The mean BMI of Australians is 27.2 kg/m^2 . Is there evidence that Americans have a different BMI from people in Australia. Test at the 5% level.

```
sample_NHANES_50<- sample_n(NHANES, size=50)
knitr::kable(head(sample_NHANES_50))
```

Table 8.5: BMI of Americans

[illegible]

8.3.3.1 Solution

1. State the random variable and the parameter in words.

x = BMI of an American

μ = mean BMI of Americans

2. State the null and alternative hypotheses and the level of significance

$$H_o : \mu = 27.2$$

$$H_a : \mu \neq 27.2$$

level of significance $\alpha = 0.05$

3. State and check the conditions for a hypothesis test

- A random sample of 50 BMI levels was taken. Check: A random sample was taken from the NHANES data frame using `r Studio`
- The population of BMI levels is normally distributed. Check:

(ref:sample-NHANES-50-density-cap) Density Plot of BMI from NHANES sample

```
gf_density(~BMI, data=sample_NHANES_50, title="Body Mass Index", xlab="Body Mass Index")  
gf_qq(~BMI, data=sample_NHANES_50, title="Body Mass Index", xlab="Body Mass Index")
```

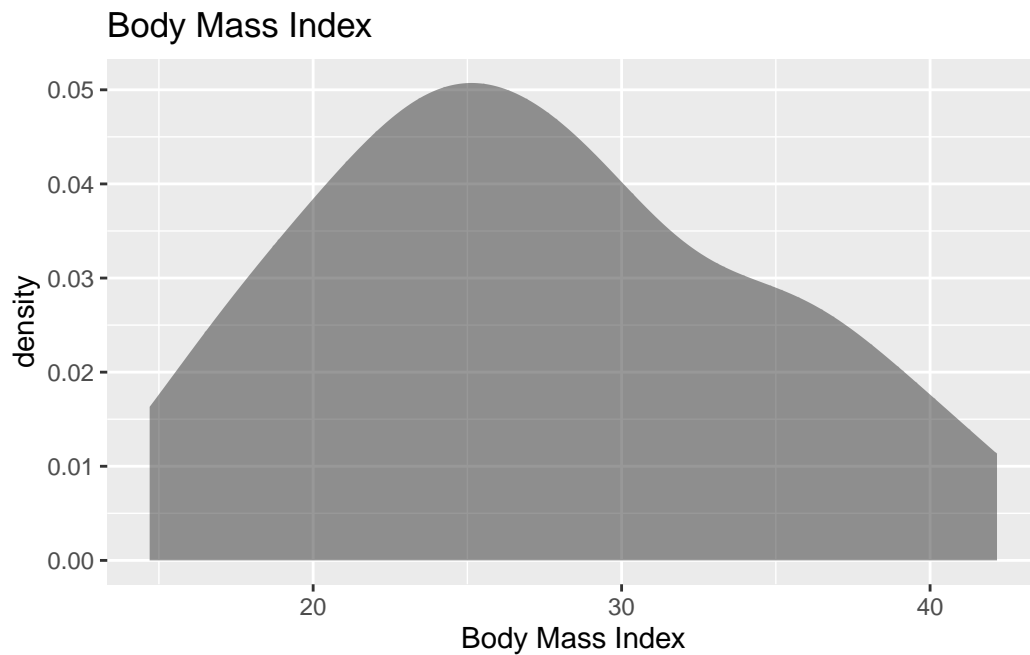


Figure 8.3: Density Plot of BMI from NHANES sample

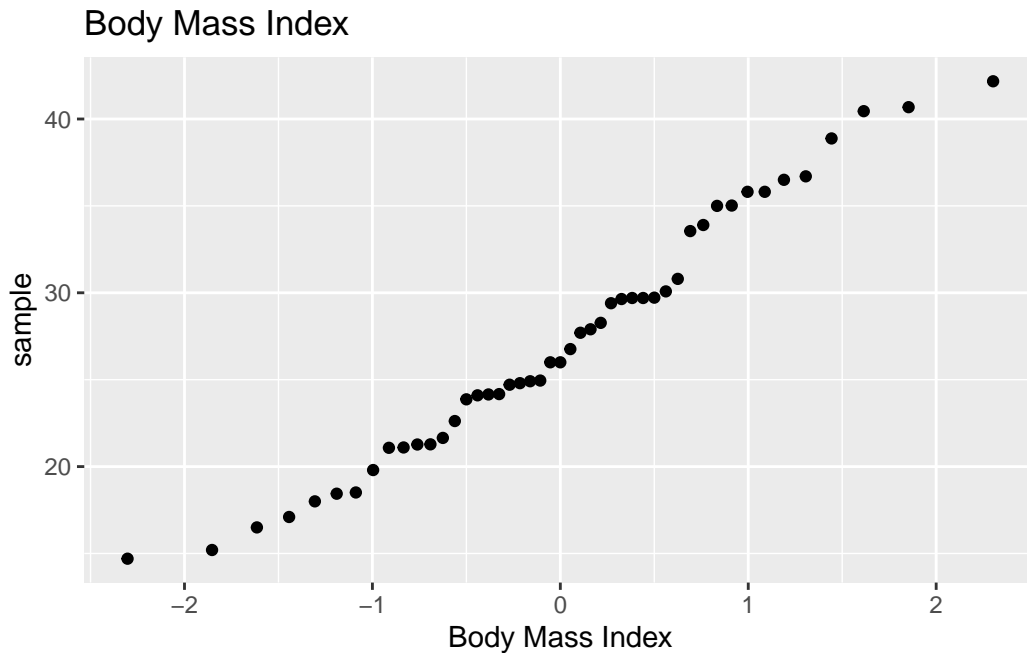


Figure 8.4: Density Plot of BMI from NHANES sample

The density plot looks somewhat skewed right and the normal quantile plot looks somewhat linear. However, there doesn't seem to be strong evidence that the sample comes from a population that is normally distributed. However, since the sample is moderate to large, the t -test is robust to this condition not being met. So the results of the test are probably valid.

4. Find the sample statistic, test statistic, and p -value

On rStudio, the command would be

```
t.test(~BMI, data= sample_NHANES_50, mu=27.2)
```

One Sample t -test

```
data: BMI
t = 0.013202, df = 46, p-value = 0.9895
alternative hypothesis: true mean is not equal to 27.2
95 percent confidence interval:
 25.10528 29.32238
sample estimates:
mean of x
 27.21383
```

The test statistic is the t in the output, the sample statistic is the mean of x in the output, and the p -value is the p -value is the output.

5. Conclusion

Since the p -value is not less than 5%, then fail to reject H_o .

6. Interpretation

There is not enough evidence to support that Americans have a different BMI from Australians.

Note: this is the same conclusion that was found when using resampling. So the two method could give similar conclusions.

8.3.4 Example: Test of the Mean Using One Sample T-test

In 2011, the average life expectancy for a woman in Europe was 79.8 years. The data in Table ?? are the life expectancies for all people in European countries (\“WHO life expectancy,” 2013). The Table ?? filtered the data frame for just males and just year 2000. The year 2000 was randomly chosen as the year to use. Do the data indicate that men’s life expectancy is different from women’s? Test at the 1% level.

```
Expectancy<-read.csv( "https://krkozak.github.io/MAT160/Life_expectancy_Europe.csv")
knitr::kable(head(Expectancy))
```

Table 8.6: Life Expectancies for European Countries

	year	WHO_region	country	sex	expect
	1990	Europe	Albania	Male	67
	1990	Europe	Albania	Female	71
	1990	Europe	Albania	Both sexes	69
	2000	Europe	Albania	Male	68
	2000	Europe	Albania	Female	73
	2000	Europe	Albania	Both sexes	71

```
Expectancy_male<-
  Expectancy |>
  filter(sex=="Male", year=="2000")
knitr::kable(head(Expectancy_male))
```

Table 8.7: Life Expectancies of males in European Countries in 2000

year	WHO_region	country	sex	expect
2000	Europe	Albania	Male	68
2000	Europe	Andorra	Male	76
2000	Europe	Armenia	Male	68
2000	Europe	Austria	Male	75
2000	Europe	Azerbaijan	Male	64
2000	Europe	Belarus	Male	63

Code book for data frame Expectancy

Description This data extract has been generated by the Global Health Observatory of the World Health Organization. The data was extracted on 2013-09-19 13:10:20.0.

This data frame contains the following columns:

year: year for life expectancies

WHO_region: World Health Organizations designation for the location of the country

country: country where the epectancies are from

sex: sex of the group that expectancies are calculated for

expect: average life expectancies of the different groups of the different countries.

Source http://apps.who.int/gho/athena/data/download.xml?format=xml&target=GHO/WHOSIS_000001&pro

References World Health Organization (WHO).

8.3.4.1 Solution

1. State the random variable and the parameter in words.

x = life expectancy for a European man

μ = mean life expectancy for European men

2. State the null and alternative hypotheses and the level of significance

$$H_o : \mu = 79.8$$

$$H_a : \mu \neq 79.8$$

$$\alpha = 0.01$$

3. State and check the conditions for a hypothesis test

- a. State: A random sample of 53 life expectancies of European men in 2000 was taken.

Check: The data is actually all of the life expectancies for every country that is considered part of Europe by the World Health Organization in the year 2000. Since the year 2000 was picked at random, then the sample is a random sample.

- b. State: The distribution of life expectancies of European men in 2000 is normally distributed.

Check:

```
gf_density(~expect, data=Expectancy_male, title="Life Expectancies of Males in Europe in 2000")
```

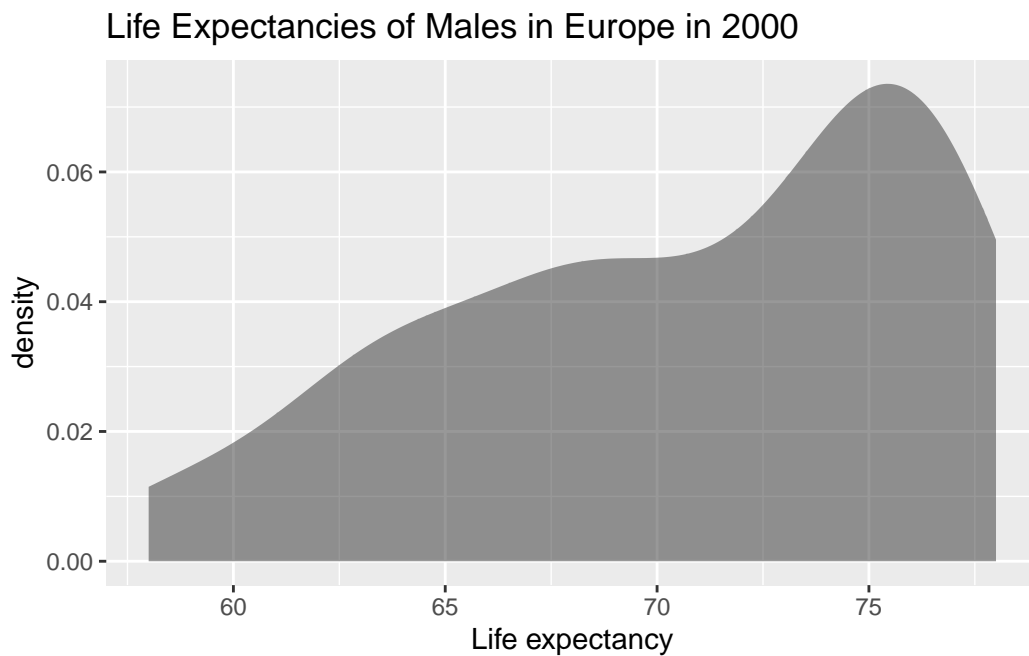


Figure 8.5: Density Plot of Life Expectancy of Males in Europe in 2000

```
gf_qq(~expect, data=Expectancy_male, title="Life Expectancies of Males in Europe in 2000")
```

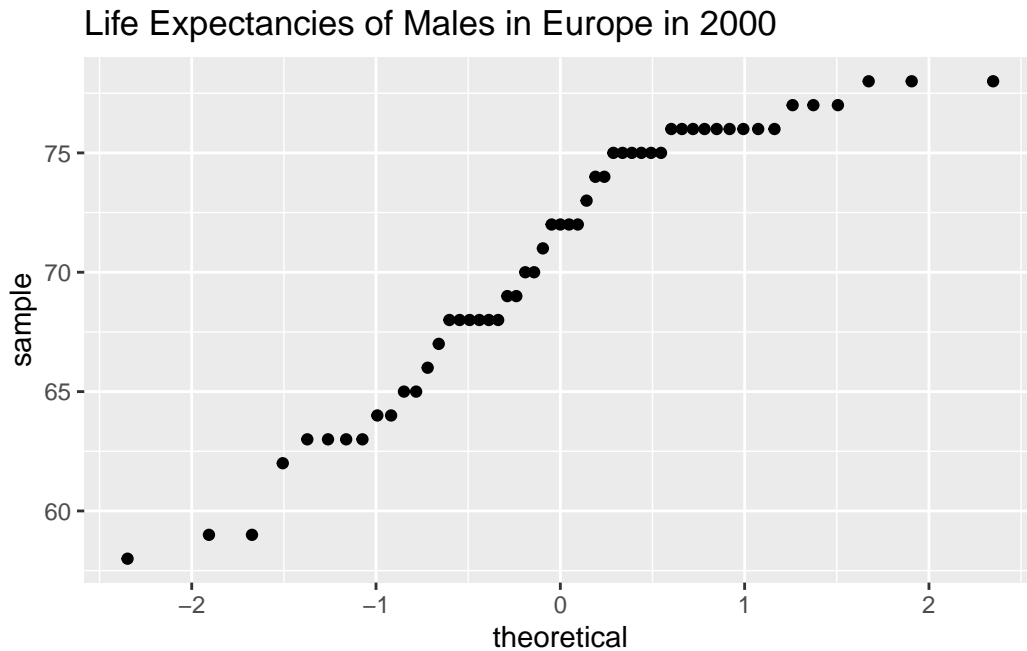


Figure 8.6: Quantile Plot of Life Expectancy of Males in Europe in 2000

This sample does not appear to come from a population that is normally distributed. This sample is moderate to large, so it is good that the t-test is robust.

4. Find the sample statistic, test statistic, and p -value

On rStudio, the command is

```
t.test(~expect, data=Expectancy_male, mu=79.8)
```

One Sample t-test

```
data: expect
t = -11.733, df = 52, p-value = 3.145e-16
alternative hypothesis: true mean is not equal to 79.8
95 percent confidence interval:
 69.11930 72.23919
sample estimates:
mean of x
 70.67925
```

Sample statistic is 70.68 years, test statistic is $t = -11.733$, and $p\text{-value} = 3.14 \times 10^{-16}$.

5. Conclusion

Since the p-value is less than 1%, then reject H_0 .

6. Interpretation

There is enough evidence to support that the mean life expectancy for European men is different than the mean life expectancy for European women of 79.8 years.

Note: if you want to conduct a hypothesis test with $H_a : \mu > \mu_o$, then the rStudio command would be

```
t.test(~variable, data=Data_Frame, mu=number  $H_0$  equals, alternative="greater")
```

If you want to conduct a hypothesis test with $H_a : \mu < \mu_o$, then the r Studio command would be

```
t.test(~variable, data=Data Frame, mu=number  $H_0$  equals, alternative="less")
```

8.3.5 Homework for One-Sample Test for the Mean Section

In each problem show all steps of the hypothesis test. If some of the conditions are not met, note that the results of the test may not be correct and then continue the process of the hypothesis test.

1. The Kyoto Protocol was signed in 1997, and required countries to start reducing their carbon emissions. The protocol became enforceable in February 2005. In 2004, the mean CO2 emission was 4.87 metric tons per capita. The Table ?? contains a random sample of CO2 emissions in 2010 (CO2 emissions (metric tons per capita), 2018). Is there enough evidence to show that the mean CO2 emission is different in 2010 than in 2004? Test at the 1% level.

```
Emission <- read.csv("https://krkozak.github.io/MAT160/CO2_emission.csv")
knitr::kable(head(Emission))
```

Table 8.8: CO2 Emissions (in metric tons per capita) in 2010

[illegible]

[illegible]

Description Carbon dioxide emissions are those stemming from the burning of fossil fuels and the manufacture of cement. They include carbon dioxide produced during consumption of solid, liquid, and gas fuels and gas flaring.

country: country around the world

Source CO2 emissions (metric tons per capita). (n.d.). Retrieved July 18, 2019, from <https://data.worldbank.org/indicator/EN.ATM.CO2E.PC>

References Carbon Dioxide Information Analysis Center, Environmental Sciences Division, Oak Ridge National Laboratory, Tennessee, United States.

- The amount of sugar in a Krispy Kream glazed donut is 10 g. Many people feel that cereal is a healthier alternative for children over glazed donuts. The Table ?? contains the amount of sugar in a sample of cereal (breakfast cereal, 2019). Is there enough evidence to show that the mean amount of sugar in children's cereal is different than in a glazed donut? Test at the 5% level.

Table 8.9: Nutrition Amounts in Cereal

name	manf	age	type	calories	protein	fat	sodium	fiber	carbs	sugars	shelf	potassium	vitamin	weight	serving
100%_Bran	Nabisco	adult	cold	70	4	1	130	10.0	5.0	6	3	280	25	1	0.33
100%_Natural	Quaker	adult	cold	120	3	5	15	2.0	8.0	8	3	135	0	1	-
															1.00
All-Bran	Kelloggs	adult	cold	70	4	1	260	9.0	7.0	5	3	320	25	1	0.33
All-Bran with Extra Fiber	Kelloggs	adult	cold	50	4	0	140	14.0	8.0	0	3	330	25	1	0.50

Table 8.9: Nutrition Amounts in Cereal

name	manf	age	type	calories	protein	fat	sodium	fiber	carb	sugar	shelf	potassium	vit	weight	serving
Almond_Delishious	Post	Adult	Cold	110	2	2	200	1.0	14.0	8	3	-1	25	1	0.75
Apple_Cinnamon	General Mills	Children	Cold	110	2	2	180	1.5	10.5	10	1	70	25	1	0.75

Code book for data frame Sugar

Description Nutritional information about cereals.

This data frame contains the following columns:

name: the cereal brand

manf: manufacturer

age: whether the cereal is geared towards children or adults

type: whether the cereal is considered a hot or cold cereal

calories: the number of calories in the cereal (number)

protein: the amount of protein in a serving of the cereal (g)

fat: the amount of fat a serving of the cereal (g)

sodium: the amount of sodium in a serving of the cereal (mg)

fiber: the amount of fiber in a serving of the cereal (g)

carb: the amount of complex carbohydrates in a serving of the cereal (g)

sugars: the amount of sugar in a serving of the cereal (g)

display shelf: what shelf the cereal is on counting from the floor

potassium: the amount of potassium in a serving of the cereal (mg)

vit: the amount of vitamins and minerals in a serving of the cereal (0, 25, or 100)

weight: weight in ounces of one serving

serving: cups per serving

Source (n.d.). Retrieved July 18, 2019, from <https://www.idvbook.com/teaching-aid/datasets/the-breakfast-cereal-data-set/> The Best Kids' Cereal. (n.d.). Retrieved July 18, 2019, from <https://www.ranker.com/list/best-kids-cereal/ranker-food>

References Interactive Data Visualization Foundations, Techniques, Applications (Matthew Ward | Georges Grinstein | Daniel Keim)

A new data frame Table ?? will need to be created of just cereal for children. To create that use the following command in rStudio

```
Sugar_children<-
  Sugar%>%
  filter(age=="child")
knitr::kable(head(Sugar_children))
```

Table 8.10: Nutrition Amounts in Children's Cereal

name	manf	age	type	calorie	protein	fat	sodium	fiber	carb	sugar	shelf	potassium	vitamin	weight	serving
Apple_Cinnamon	General Mills	child	cold	110	2	2	180	1.5	10.5	10	1	70	25	1	0.75
Apple_Jacks	Kelloggs	child	cold	110	2	0	125	1.0	11.0	14	2	30	25	1	1.00
Bran_Cheerios	Ralston	child	cold	90	2	1	200	4.0	15.0	6	1	125	25	1	0.67
Cap'n_Crunch	Quaker	child	cold	120	1	2	220	0.0	12.0	12	2	35	25	1	0.75
Cheerios	General Mills	child	cold	110	6	2	290	2.0	17.0	1	1	105	25	1	1.25
Cinnamon_Tostitos	General Mills	child	cold	120	1	3	210	0.0	13.0	9	2	45	25	1	0.75

- The FDA regulates that fish that is consumed is allowed to contain 1.0 mg/kg of mercury. In Florida, bass fish were collected in 53 different lakes to measure the health of the lakes. The data frame of measurements from Florida lakes is in Table ?? (NISER 081107 ID Data, 2019). Do the data provide enough evidence to show that the fish in Florida lakes has different amounts of mercury than the allowable amount? Test at the 10% level.

```
Mercury<- read.csv( "https://krkozak.github.io/MAT160/mercury.csv")
knitr::kable(head(Mercury))
```

Table 8.11: Health of Florida lake Fish

ID	lake	alkalinity	ph	calcium	chlorophyll	mercury	no.samples	min	max	X3_yr_standard	age	data
1	Alligator	5.9	6.1	3.0	0.7	1.23	5	0.85	1.43	1.53	1	
2	Annie	3.5	5.1	1.9	3.2	1.33	7	0.92	1.90	1.33	0	
3	Apopka	116.0	9.1	44.1	128.3	0.04	6	0.04	0.06	0.04	0	
4	Blue_Cypress	39.4	6.9	16.4	3.5	0.44	12	0.13	0.84	0.44	0	
5	Brick	2.5	4.6	2.9	1.8	1.20	12	0.69	1.50	1.33	1	
6	Bryant	19.6	7.3	4.5	44.1	0.27	14	0.04	0.48	0.25	1	

Code book for data frame Mercury

Description Largemouth bass were studied in 53 different Florida lakes to examine the factors that influence the level of mercury contamination. Water samples were collected from the

surface of the middle of each lake in August 1990 and then again in March 1991. The pH level, the amount of chlorophyll, calcium, and alkalinity were measured in each sample. The average of the August and March values were used in the analysis. Next, a sample of fish was taken from each lake with sample sizes ranging from 4 to 44 fish. The age of each fish and mercury concentration in the muscle tissue was measured. (Note: Since fish absorb mercury over time, older fish will tend to have higher concentrations). Thus, to make a fair comparison of the fish in different lakes, the investigators used a regression estimate of the expected mercury concentration in a three year old fish as the standardized value for each lake. Finally, in 10 of the 53 lakes, the age of the individual fish could not be determined and the average mercury concentration of the sampled fish was used instead of the standardized value. (Reference: Lange, Royals, & Connor. (1993))

This data frame contains the following columns:

ID: ID number

Lake: Name of lake

alkalinity: Alkalinity (mg/L as Calcium Carbonate)

pH: pH

calcium: calcium (mg/l)

chlorophyll: chlorophyll (mg/l)

mercury: Average mercury concentration (parts per million) in the muscle tissue of the fish sampled from that lake

no.samples: How many fish were sampled from the lake

min: Minimum mercury concentration among the sampled fish

max: Maximum mercury concentration among the sampled fish

X3_yr_Standard_mercury: Regression estimate of the mercury concentration in a 3 year old fish from the lake (or = Avg Mercury when age data was not available)

age_data: Indicator of the availability of age data on fish sampled

Source Lange TL, Royals HE, Connor LL (1993) Influence of water chemistry on mercury concentration in largemouth bass from Florida lakes. Trans Am Fish Soc 122:74-84. Michael K. Saiki, Darell G. Slotton, Thomas W. May, Shaun M. Ayers, and Charles N. Alpers (2000) Summary of Total Mercury Concentrations in Fillets of Selected Sport Fishes Collected during 2000–2003 from Lake Natoma, Sacramento County, California (Raw data is included in appendix), U.S. Geological Survey Data Series 103, 1-21. NISER 081107 ID Data. (n.d.). Retrieved July 18, 2019, from http://wiki.stat.ucla.edu/socr/index.php/NISER_081107_ID_Data

References NISER 081107 ID Data

- The data frame Table ?? contains various variables about a person including their pulse rates before the subject exercised and after the subject ran in place for one minute. The mean pulse rate after running for 1 minute of females who do not drink is 97 beats per minute. Do the data show that the mean pulse rate of females who do drink alcohol is higher than the mean pulse rate of females who do not drink? Test at the 5% level.

Code book for data frame **Pulse** is below Table ??.

Create a data frame Table ?? that contains only females who drink alcohol. Then test the pulse after for woman who do drink alcohol to the known value for females who do not drink alcohol. To create a new data frame with just females who drink alcohol use the following command, where the new name is Females:

```
Females<- Pulse%>% filter(gender=="female", alcohol=="yes")
knitr::kable(head(Females))
```

Table 8.12: Pulse Rates Before and After Exercise of Females who do drink Alcohol

height	weight	age	gender	smokes	alcohol	exercise	ran	pulse_before	pulse_after	year
165	60	19	female	yes	yes	low	ran	88	120	98
163	47	23	female	yes	yes	low	ran	71	125	98
173	57	18	female	no	yes	moderate	sat	86	88	93
179	58	19	female	no	yes	moderate	ran	82	150	93
167	62	18	female	no	yes	high	ran	96	176	93
173	64	18	female	no	yes	low	sat	90	88	93

- The economic dynamism is an index of productive growth in dollars. Economic data for many countries are in Table ?? (SOCR Data 2008 World CountriesRankings, 2019). Countries that are considered high-income have a mean economic dynamism of 60.29.

```
Economics <- read.csv( "https://krkozak.github.io/MAT160/Economics_country.csv")
knitr::kable(head(Economics))
```

Table 8.13: Economic Data for Countries

Id	incGroup	key	name	popGroup	region	key2	ED	Edu	HI	QOL	PE	OA	Relig
0	Low	al	Albania	Small	Southern_Europe	popS	34.0862	81.0164	71.0244	67.9240	58.6742	57	39
1	Middle	dz	Algeria	Medium	North_Africa	popM	25.8057	74.8027	66.1951	60.9347	32.6054	85	95
2	Middle	ar	Argentina	Medium	South_America	popM	7.4511	69.8825	78.2683	68.1559	68.6647	66	66
3	High	au	Australia	Medium	Australia	popM	71.4888	91.4802	95.1707	90.5729	90.9629	4	65
4	High	at	Austria	Small	Central_Europe	popS	53.9431	90.4578	90.3415	87.5630	91.2073	18	20

Table 8.13: Economic Data for Countries

Id	incGroup	key	name	popGroup	region	key2	ED	Edu	HI	QOL	PE	OA	Relig
5	Low	az	Azerbaijan	Small	central_Asia	popS	53.6457	68.9880	58.9512	68.9572	40.0390	69	50

Code book for data frame Economics

Description These data represent commonly accepted measures for ranking Countries on variety of factors which affect the country's internal and external international perception of the country's rank relative the to rest of the World.

This data frame contains the following columns:

id: Unique country identifier

incGroup: Income group: Low: GNI per capita < \ \$3,946, Middle: \ \$3,946 < GNI per capita < \ \$12,195, High: GNI per capita > \ \$12,196

key: unique 2-letter country code

name: Country Name

popGroup: Population Group: Small: Population < 20 million, Medium: 20 million < Population < 50 million, Large: Population > 50 million

region: Relative geographic position of the Country

key2: Country Group Classification Label: world: All countries, g7: G7, g20: G20, latin: Latin America & Caribbean, eu: European Union, centasia: Europe & Central Asia, pacasia: East Asia & Pacific, asean: Asean, sasias: South Asia, mideast: Middle East & North Africa, africa: Sub-Saharan Africa, bric: Brazil, Russia, India and China (BRIC)

ED: Economic Dynamism: Index of Productive growth in dollars (GDP/capita at PPP, Avg of GDP/capita growth rate over last ten years, GDP/capita growth rate over next ten years, Economic Dynamism: Manufacturing percent of GDP, Services percent of GDP percent (100=best, 0=worst).

Edu: Education/Literacy Rate (percent of population able to read and write at a specified age)

HI: Health Index: The average number of years a person lives in full health, taking into account years lived in less than full health

QOL: Quality of Life: Population percent living on < \ \$2/day

PE: Political Environment: Freedom house rating of political participation (qualitative assessment of voter participation/turn-out for national elections, citizens engagement with politics)

OA: Overall country ranking taking all measures into account.

Relig: Religiosity of the Country as a percent (%) of the population.

Source SOCR Data 2008 World CountriesRankings. (n.d.). Retrieved July 19, 2019, from http://wiki.stat.ucla.edu/socr/index.php/SOCR_Data_2008_World_CountriesRankings#SOCR_Data_-_Ranking_of_the_top_100_Countries_based_on_Political.2C_Economic.2C_Health.2C_and_Quality-of-Life_Factors

References SOCR Data 2008 World CountriesRankings, Amazon Web-Services World's Best Countries.

Create a data frame that contains only middle income countries. Do the data show that the mean economic dynamism of middle-income countries is less than the mean for high-income countries? Test at the 5% level. To create a new data frame Table ?? with just middle income countries use the following command, where the new name is Middle_economics:

```
Middle_economics<-
  Economics |>
  filter(incGroup=="Middle")
knitr::kable(head(Middle_economics))
```

Table 8.14: Economic Data for Middle income Countries

Id	incGroup	key	name	popGroup	region	key2	ED	Edu	HI	QOL	PE	OA	Relig
1	Middle	dz	Algeria	Medium	North_Africa	popM	25.805	774.802	766.195	160.934	732.605	485	95
2	Middle	ar	Argentina	Medium	South_America	popM	37.451	169.882	578.268	368.155	968.664	746	66
7	Middle	by	Belarus	Small	central_Asia	popS	51.915	086.615	566.195	174.146	734.050	156	34
10	Middle	bw	Botswana	Small	Africa	popS	43.695	273.460	834.804	950.087	572.683	380	80
11	Middle	br	Brazil	Large	South_America	popL	47.850	671.373	571.024	462.423	867.413	148	87
12	Middle	bg	Bulgaria	Small	Southern_Europe	popS	43.717	882.227	775.853	773.119	773.168	638	50

6. In 1999, the average percentage of women who received prenatal care per country is 80.1%. Table ?? contains the percentage of woman receiving prenatal care in a sample of countries over several years. (births per woman), 2019). Do the data show that the average percentage of women receiving prenatal care in 2009 (p2009) is different than in 1999? Test at the 5% level.

Code book for Data frame Fert_prenatal is below Table ??.

7. Maintaining your balance may get harder as you grow older. A study was conducted to see how steady the elderly is on their feet. They had the subjects stand on a force platform and have them react to a noise. The force platform then measured how much they swayed forward and backward, and the data is in Table ?? (Maintaining Balance

while Concentrating, 2019). Do the data show that the elderly sway more than the mean forward sway of younger people, which is 18.125 mm? Test at the 5% level. Follow the filtering methods in other homework problems to create a data frame for only Elderly.

```
Sway <- read.csv( "https://krkozak.github.io/MAT160/sway.csv")
knitr::kable(head(Sway))
```

Table 8.15: Sway (in mm) of Elderly Subjects

age	fbsway	sidesway
Elderly	19	14
Elderly	30	41
Elderly	20	18
Elderly	19	11
Elderly	29	16
Elderly	25	24

Code book for data frame Sway

Description How difficult is it to maintain your balance while concentrating? It is more difficult when you are older? Nine elderly (6 men and 3 women) and eight young men were subjects in an experiment. Each subject stood barefoot on a “force platform” and was asked to maintain a stable upright position and to react as quickly as possible to an unpredictable noise by pressing a hand held button. The noise came randomly and the subject concentrated on reacting as quickly as possible. The platform automatically measured how much each subject swayed in millimeters in both the forward/backward and the side-to-side directions.

This data frame contains the following columns:

Age: Elderly or Young

FBSway: Sway in forward/backward direction

SideSwayy: Sway in side to side direction

Source Maintaining Balance while Concentrating. (n.d.). Retrieved July 19, 2019, from <http://www.statsci.org/data/general/balaconc.html>

References Teasdale, N., Bard, C., La Rue, J., and Fleury, M. (1993). On the cognitive penetrability of posture control. *Experimental Aging Research* 19, 1-13. The data was obtained from the DASL Data and Story Line online database.

9 Estimation

In hypothesis tests, the purpose was to make a decision about a parameter, in terms of it being greater than, less than, or not equal to a value. But what if you want to actually know what the parameter is. You need to do estimation. There are two types of estimation -- point estimator and confidence interval. The American Statistical Association (ASA) is recommending that confidence intervals are the process that should be followed when analyzing data.

9.1 Basics of Confidence Intervals

A point estimator is just the statistic that you have calculated previously. As an example, when you wanted to estimate the population mean, μ , the point estimator is the sample mean, \bar{x} . To estimate the population proportion, p , you use the sample proportion, \hat{p} . In general, if you want to estimate any population parameter, we will call it θ , you use the sample statistic, $\hat{\theta}$.

Point estimators are really easy to find, but they have some drawbacks. First, if you have a large sample size, then the estimate is better. But with a point estimator, you don't know what the sample size is. Also, you don't know how accurate the estimate is. Both of these problems are solved with a confidence interval.

Confidence interval: This is where you have an interval surrounding your parameter, and the interval has a chance of being a true statement. In general, a confidence interval looks like: $\hat{\theta} \pm E$, where $\hat{\theta}$ is the point estimator and E is the margin of error term that is added and subtracted from the point estimator. Thus making an interval.

9.1.1 Interpreting a confidence interval:

The statistical interpretation is that the confidence interval has a probability $C = (1 - \alpha)$ (where α is the complement of the confidence level) of containing the population parameter. As an example, if you have a 95% confidence interval of $0.65 < p < 0.73$, then you would say, "you are 95% confident that the interval 0.65 to 0.73 contains the true population proportion." This means that if you have 100 intervals, 95 of them will contain the true proportion, and 5 will not. The wrong interpretation is that there is a 95% confidence that the true value of p will fall between 0.65 and 0.73. The reason that this interpretation is wrong is that the true

value is fixed out there somewhere. You are trying to capture it with this interval. So this is the chance that your interval captures it, and not that the true value falls in the interval.

There is also a real world interpretation that depends on the situation. It is where you are telling people what numbers you found the parameter to lie between. So your real world is where you tell what values your parameter is between. There is no probability attached to this statement. That probability is in the statistical interpretation.

The common probabilities used for confidence intervals are 90%, 95%, and 99%. These are known as the confidence level. The confidence level and the alpha level are related. If you are conducting a hypothesis test with $H_a : \mu \neq \mu_o$, then the confidence level is $C = 1 - \alpha$. This is because the α is both tails and the confidence level is area between the two tails. As an example, for a hypothesis test $H_a : \mu \neq \mu_o$ with α equal to 0.05, the confidence level would be 0.95 or 95%. If you have a hypothesis test with $H_a : \mu < \mu_o$, then your α is only one tail of the curve. Because of symmetry the other tail is also α . You have 2α with both tails. So the confidence level, which is the area between the two tails, is $C = 1 - 2\alpha$.

9.1.2 Example: Stating the Statistical and Real World Interpretations for a Confidence Interval

- Suppose you have a 95% confidence interval for the mean age a woman gets married in 2013 is $26 < \mu < 28$. State the statistical and real world interpretations of this statement.
- Suppose a 99% confidence interval for the proportion of Americans who have tried marijuana as of 2013 is $0.35 < p < 0.41$. State the statistical and real world interpretations of this statement.

9.1.2.1 Solution

- Suppose you have a 95% confidence interval for the mean age a woman gets married in 2013 is $26 < \mu < 28$. State the statistical and real world interpretations of this statement.

Statistical Interpretation: You are 95% confident that the interval contains the mean age in 2013 that a woman gets married.

Real World Interpretation: The mean age that a woman married in 2013 is between 26 and 28 years of age.

- Suppose a 99% confidence interval for the proportion of Americans who have tried marijuana as of 2013 is $0.35 < p < 0.41$. State the statistical and real world interpretations of this statement.

Statistical Interpretation: You are 99% confident that the interval contains the proportion of Americans who have tried marijuana as of 2013.

Real World Interpretation: The proportion of Americans who have tried marijuana as of 2013 is between 0.35 and 0.41.

One last thing to know about confidence is how the sample size and confidence level affect how wide the interval is. The following discussion demonstrates what happens to the width of the interval as you get more confident.

Think about shooting an arrow into the target. Suppose you are really good at that and that you have a 90% chance of hitting the bull's eye. Now the bull's eye is very small. Since you hit the bull's eye approximately 90% of the time, then you probably hit inside the next ring out 95% of the time. You have a better chance of doing this, but the circle is bigger. You probably have a 99% chance of hitting the target, but that is a much bigger circle to hit. You can see, as your confidence in hitting the target increases, the circle you hit gets bigger. The same is true for confidence intervals. This is demonstrated in Image \#8.1.1.

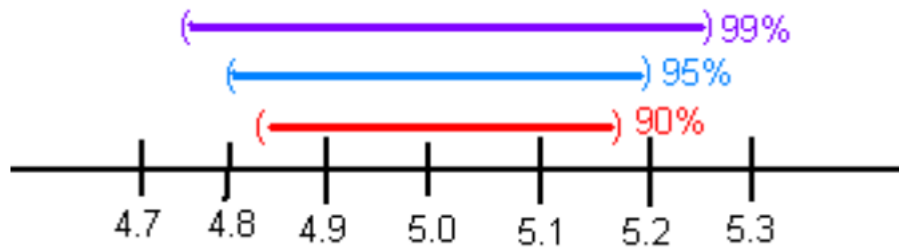


Figure 9.1: Image \#8.1.1 Confidence Level Effect

The higher level of confidence makes a wider interval. There's a trade off between width and confidence level. You can be really confident about your answer but your answer will not be very precise. Or you can have a precise answer (small margin of error) but not be very confident about your answer.

Now look at how the sample size affects the size of the interval. Suppose Image \#8.1.2 represents confidence intervals calculated on a 95% interval. A larger sample size from a representative sample makes the width of the interval narrower. This makes sense. Large samples are closer to the true population so the point estimate is pretty close to the true value.

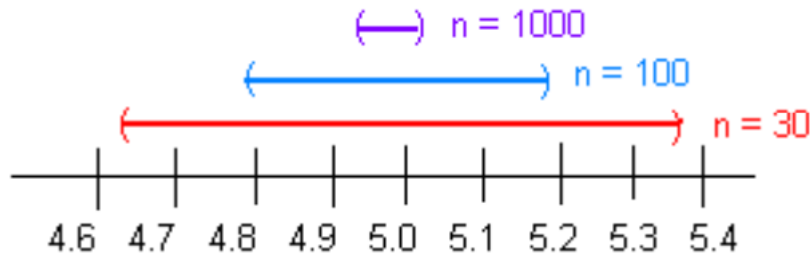


Figure 9.2: Image #8.1.2 Effect of Sample size

Now you know everything you need to know about confidence intervals except for the actual formula. The formula depends on which parameter you are trying to estimate. With different situations you will be given the confidence interval for that parameter.

9.1.3 Homework for Basics of Confidence Intervals Section

1. Suppose you compute a confidence interval with a sample size of 25. What will happen to the confidence interval if the sample size increases to 50?
2. Suppose you compute a 95% confidence interval. What will happen to the confidence interval if you increase the confidence level to 99%?
3. Suppose you compute a 95% confidence interval. What will happen to the confidence interval if you decrease the confidence level to 90%?
4. Suppose you compute a confidence interval with a sample size of 100. What will happen to the confidence interval if the sample size decreases to 80?
5. A 95% confidence interval is $6353km < \mu < 6384km$, where μ is the mean diameter of the Earth. State the statistical interpretation.
6. A 95% confidence interval is $6353km < \mu < 6384km$, where μ is the mean diameter of the Earth. State the real world interpretation.
7. In 2013, Gallup conducted a poll and found a 95% confidence interval of $0.52 < p < 0.60$, where p is the proportion of Americans who believe it is the government's responsibility for health care. Give the real world interpretation.
8. In 2013, Gallup conducted a poll and found a 95% confidence interval of $0.52 < p < 0.60$, where p is the proportion of Americans who believe it is the government's responsibility for health care. Give the statistical interpretation.

9.2 One-Sample Interval for the Proportion

Suppose you want to estimate the population proportion, p . As an example you may be curious what proportion of students at your school smoke. Or you could wonder what is the proportion of accidents caused by teenage drivers who do not have a drivers' education class.

9.2.1 Confidence Interval for One Population Proportion (1-Prop Interval)

1. State the random variable and the parameter in words.

x = number of successes

p = proportion of successes

2. State and check the conditions for the confidence interval
 - a. State: A simple random sample of size n is taken. Check: describe how sample was taken.
 - b. State: The condition for the binomial distribution are satisfied. Check: argue that each condition has been met.
 - c. State: The sampling distribution of \hat{p} can be approximated by a normal distributed. check: To determine the sampling distribution of \hat{p} is normally distributed, you need to show that $n * \hat{p} \geq 5$ and $n * \hat{q} \geq 5$ where $\hat{q} = 1 - \hat{p}$. If this requirement is true, then the sampling distribution of \hat{p} is well approximated by a normal curve. (In reality this is not really true, since the correct condition deals with p . However, in a confidence interval you do not know p , so you must use \hat{p} .)
3. Find the sample statistic and the confidence interval

This will be conducted using rStudio. The command is

```
prop.test(r, n, conf.level=C) #type C as a decimal
```

4. Statistical Interpretation: In general this looks like, “you are C% confident that $\hat{p} \pm E$ contains the true proportion.”
5. Real World Interpretation: This is where you state what interval contains the true proportion.

9.2.2 Example: Confidence Interval for the Population Proportion

A concern was raised in Australia that the percentage of deaths of Aboriginal prisoners was higher than the percent of deaths of non-Aboriginal prisoners, which is 0.27%. A sample of six years (1990-1995) of data was collected, and it was found that out of 14,495 Aboriginal prisoners, 51 died (“Indigenous deaths in,” 1996). Find a 95% confidence interval for the proportion of Aboriginal prisoners who died.

9.2.2.1 Solution

1. State the random variable and the parameter in words.

x = number of Aboriginal prisoners who die

p = proportion of Aboriginal prisoners who die

2. State and check the conditions for the confidence interval

- a. State: A simple random sample of 14,495 Aboriginal prisoners was taken. Check: The sample was not a random sample, since it was data from six years. It is the numbers for all prisoners in these six years, but the six years were not picked at random. Unless there was something special about the six years that were chosen, the sample is probably a representative sample. This condition is probably met.
- b. State: The properties of the binomial experiment have been met. Check: There are 14,495 prisoners in this case. The prisoners are all Aboriginals, so you are not mixing Aboriginal with non-Aboriginal prisoners. There are only two outcomes, either the prisoner dies or doesn't. The chance that one prisoner dies over another may not be constant, but if you consider all prisoners the same, then it may be close to the same probability. Thus the properties of the binomial experiment are satisfied
- c. State: The sampling distribution of \hat{p} can be approximated with a normal distribution. Check: $\hat{p} * n = \frac{51}{14495} * 14495 = 51 \geq 5$ and $\hat{q} * n = \frac{14495-51}{14495} * 14495 = 14444 \geq 5$. The sampling distribution of \hat{p} can be approximated with a normal distribution.

3. Find the sample statistic and the confidence interval

The command in r Studio for a confidence interval for a proportion is

```
prop.test(51,14495, conf.level = 0.95)
```

```

1-sample proportions test with continuity correction

data:  51 out of 14495
X-squared = 14290, df = 1, p-value < 2.2e-16
alternative hypothesis: true p is not equal to 0.5
95 percent confidence interval:
 0.002647440 0.004661881
sample estimates:
              p
0.003518455

```

the 95% confidence level is $0.002647440 < p < 0.004661881$.

4. Statistical Interpretation: You are 95% confident that the interval $0.0026 < p < 0.0047$ contains the proportion of Aboriginal prisoners who have died in prison.
5. Real World Interpretation: The proportion of Aboriginal prisoners who died in prison is between 0.26% and 0.47%.

9.2.3 Example: Confidence Interval for the Population Proportion

A researcher who is studying the effects of income levels on breastfeeding of infants hypothesizes that countries with a low income level have a different rate of infant breastfeeding than higher income countries. It is known that in Germany, considered a high-income country by the World Bank, 22% of all babies are breastfeed. In Tajikistan, considered a low-income country by the World Bank, researchers found that in a random sample of 500 new mothers that 125 were breastfeeding their infant. Find a 90% confidence interval of the proportion of mothers in low-income countries who breastfeed their infants?

9.2.3.1 Solution

1. State you random variable and the parameter in words.

x = number of woman who breastfeed in a low-income country

p = proportion of woman who breastfeed in a low-income country

2. State and check the conditions for the confidence interval
 - a. State: A simple random sample of 500 breastfeeding habits of woman in a low-income country was taken. Check: This was stated in the problem.

- b. State: The properties of a Binomial Experiment have been met. Check: There were 500 women in the study. The women are considered identical, though they probably have some differences. There are only two outcomes, either the woman breastfeeds or she doesn't. The probability of a woman breastfeeding is probably not the same for each woman, but it is probably not very different for each woman. The conditions for the binomial distribution are satisfied
- c. State: The sampling distribution of \hat{p} can be approximated with a normal distributed. Check: $n * \hat{p} = 500 * \frac{125}{500} = 125 \geq 5$ and $n * \hat{q} = 500 * \frac{500-125}{500} = 375 \geq 5$, so the sampling distribution of \hat{p} is well approximated by a normal distribution.
4. Find the sample statistic and confidence interval

On rstudio, use the following command

```
prop.test(125, 500, conf.level = .90)
```

```
1-sample proportions test with continuity correction

data: 125 out of 500
X-squared = 124, df = 1, p-value < 2.2e-16
alternative hypothesis: true p is not equal to 0.5
90 percent confidence interval:
 0.2185980 0.2841772
sample estimates:
      p 
0.25
```

90% confidence interval for p is $0.2185980 < p < 0.2841772$.

4. Statistical Interpretation: You are 90% confident that $0.2185980 < p < 0.2841772$ contains the proportion of women in low-income countries who breastfeed their infants.
5. Real World Interpretation: The proportion of women in low-income countries who breastfeed their infants is between 0.219 and 0.284.

9.2.4 Homework for One-Sample Interval for the Proportion Section

In each problem show all steps of the confidence interval. If some of the conditions are not met, note that the results of the interval may not be correct and then continue the process of the confidence interval.

1. The Arizona Republic/Morrison/Cronkite News poll published on Monday, October 20, 2016, found 390 of the registered voters surveyed favor Proposition 205, which would legalize marijuana for adults. The statewide telephone poll surveyed 779 registered voters between Oct. 10 and Oct. 15. (Sanchez, 2016) Find a 99% confidence interval for the proportion of Arizona's who supported legalizing marijuana for adults.
2. In November of 1997, Australians were asked if they thought unemployment would increase. At that time 284 out of 631 said that they thought unemployment would increase (\“Morgan gallup poll,\” 2013). Estimate the proportion of Australians in November 1997 who believed unemployment would increase using a 95% confidence interval?
3. According to the February 2008 Federal Trade Commission report on consumer fraud and identity theft, Arkansas had 1,601 complaints of identity theft out of 3,482 consumer complaints (\“Consumer fraud and,\” 2008). Calculate a 90% confidence interval for the proportion of identity theft in Arkansas.
4. According to the February 2008 Federal Trade Commission report on consumer fraud and identity theft, Alaska had 321 complaints of identity theft out of 1,432 consumer complaints (\“Consumer fraud and,\” 2008). Calculate a 90% confidence interval for the proportion of identity theft in Alaska.
5. In 2013, the Gallup poll asked 1,039 American adults if they believe there was a conspiracy in the assassination of President Kennedy, and found that 634 believe there was a conspiracy (\“Gallup news service,\” 2013). Estimate the proportion of American's who believe in this conspiracy using a 98% confidence interval.
6. In 2008, there were 507 children in Arizona out of 32,601 who were diagnosed with Autism Spectrum Disorder (ASD) (\“Autism and developmental,\” 2008). Find the proportion of ASD in Arizona with a confidence level of 99%.

9.3 One-Sample Interval for the Mean

Suppose you want to estimate the mean height of Americans, or you want to estimate the mean salary of college graduates. A confidence interval for the mean would be the way to estimate these means.

9.3.1 Confidence Interval for One Population Mean (t-Interval)

1. State the random variable and the parameter in words.

x = random variable

μ = mean of random variable

9.3.2.1 Solution

1. State the random variable and the parameter in words.

x = BMI of an American

μ = mean BMI of Americans

2. State and check the conditions for the confidence interval
 - a. A random sample of 50 BMI levels was taken. Check: A random sample was taken from the NHANES data frame using r Studio
 - b. The population of BMI levels is normally distributed. Check:

```
gf_density(~BMI, data=sample_NHANES_50, title="BMI of an American", xlab="Body Mass Index")
```

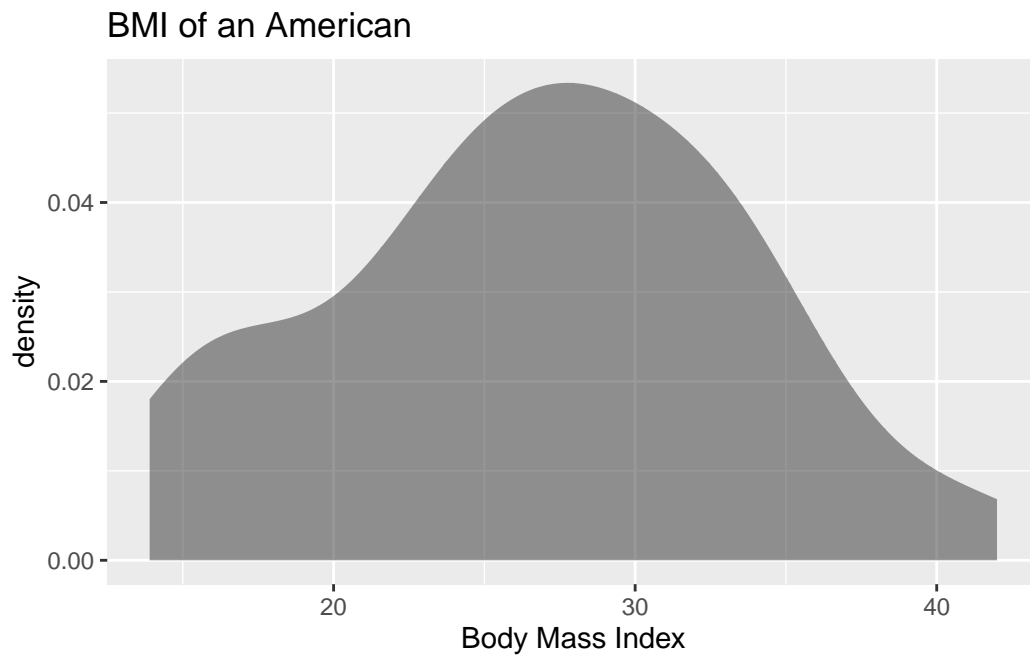


Figure 9.3: Density Plot of BMI from NHANES sample

```
gf_qq(~BMI, data=sample_NHANES_50, title="BMI of an American")
```

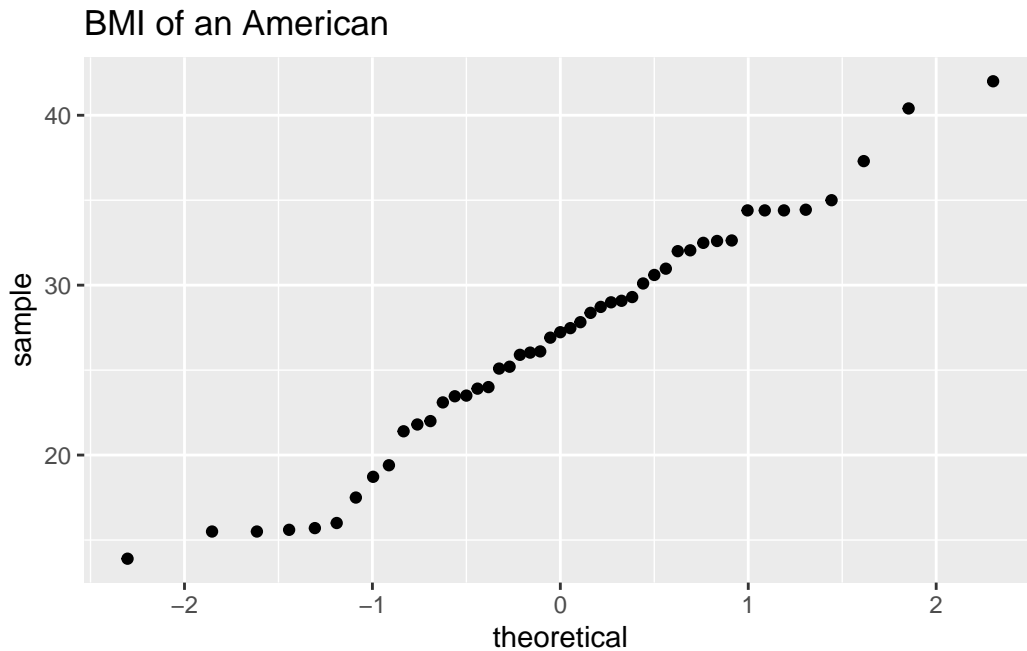


Figure 9.4: Normal quantile Plot of BMI from NHANES sample

The density plot looks somewhat skewed right and the normal quantile plot looks somewhat linear. There doesn't seem to be strong evidence that the sample comes from a population that is normally distributed. However, since the sample is moderate to large, the t-test is robust to this condition not being met. So the results of the test are probably valid.

4. Find the sample statistic and confidence interval

On r Studio, the command would be

```
t.test(~BMI, data= sample_NHANES_50, conf.level=0.95)
```

One Sample t-test

```
data: BMI
t = 26.696, df = 46, p-value < 2.2e-16
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 24.76710 28.80652
sample estimates:
mean of x
 26.78681
```

The sample statistic is the mean of x in the output, and confidence interval is under the words 95 percent confidence interval.

4. Statistical Interpretation: You are 95% confident that $24.87190 < \mu < 28.71422$ contains the mean BMI of Americans.
5. Real World Interpretation: The mean BMI of Americans is between 24.87 and 28.71 kg/m^2 .

Notice that in example the chapter 7, you were asked if the mean BMI of Americans was different from Australians' mean BMI of 27.2 kg/m^2 . The interval that Example: Confidence Interval for the Population Mean calculated does contain the value of 27.2. So you can't say that Americans' mean BMI and Australians' mean BMI are different. This means that you can just use confidence intervals and not conduct hypothesis tests at all if you prefer.

Note: When creating this book, the random samples may change. So the answers may be different from what is said in the interpretations. This shows sampling variability, so it was not adjusted to show that this could happen.

9.3.3 Example: Confidence Interval for the Population Mean

The data in Table ?? are the life expectancies for all people in European countries (“WHO life expectancy,” 2013). The data in Table ?? filtered the data frame for just males and just year 2000. The year 2000 was randomly chosen as the year to use. Estimate the mean life expectancy for a man in Europe at the 99% level.

Code book for data frame Expectancy is below Table ??.

9.3.3.1 Solution

1. State the random variable and the parameter in words.

x = life expectancy for a European man

μ = mean life expectancy for European men

2. State and check the conditions for the confidence interval

- a. State: A random sample of 53 life expectancies of European men in 2000 was taken.

Check: The data is actually all of the life expectancies for every country that is considered part of Europe by the World Health Organization in the year 2000. Since the year 2000 was picked at random, then the sample is a random sample.

- b. State: The distribution of life expectancies of European men in 2000 is normally distributed.

Check:

```
gf_density(~expect, data=Expectancy_male, title="Life Expectancy of a male", xlab="Life Expectancy of a Male")
```



Figure 9.5: Density Plot of Life Expectancy of Males in Europe in 2000

```
gf_qq(~expect, data=Expectancy_male, title="Male Life Expectancy")
```

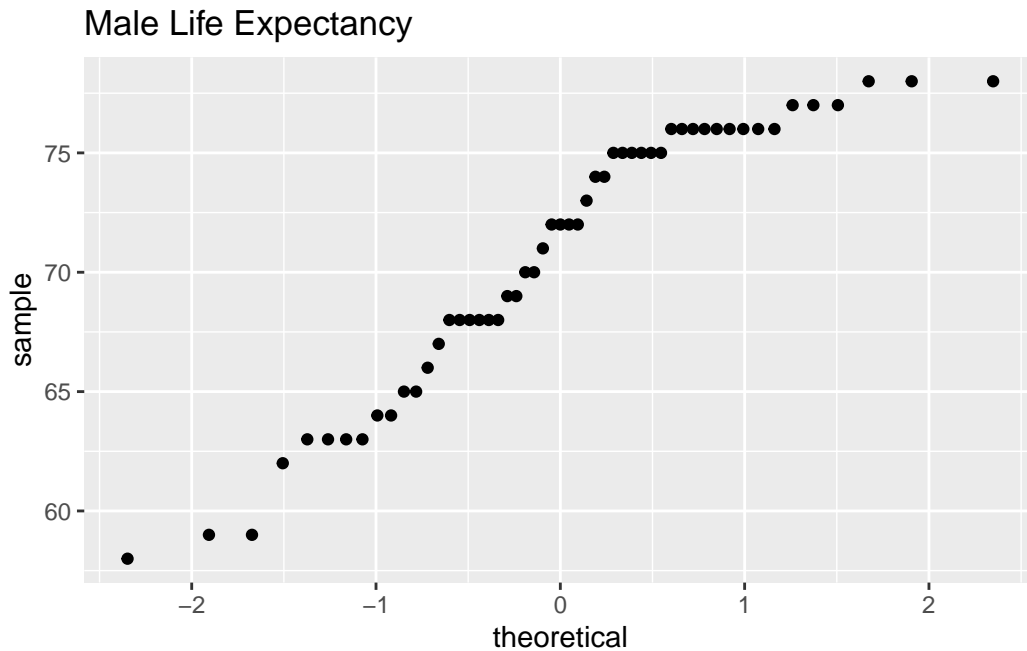


Figure 9.6: Quantile Plot of Life Expectancy of Males in Europe in 2000

This sample does not appear to come from a population that is normally distributed. This sample is moderate to large, so it is good that the t-test is robust.

3. Find the sample statistic and confidence interval

On rStudio, the command would be

```
t.test(~expect, data=Expectancy_male, conf.level=0.99)
```

One Sample t-test

```
data: expect
t = 90.919, df = 52, p-value < 2.2e-16
alternative hypothesis: true mean is not equal to 0
99 percent confidence interval:
 68.60071 72.75778
sample estimates:
mean of x
 70.67925
```

Sample statistic is 70.68 years, and the confidence interval is $68.60071 < \mu < 72.75778$.

4. Statistical Interpretation: You are 99% confident that $68.60071 < \mu < 72.75778$ contains the mean life expectancy of European men.
5. Real World Interpretation: The mean life expectancy of European men is between 68.60 and 72.76 years.

9.3.4 Homework for One-Sample Interval for the Mean Section

In each problem show all steps of the confidence interval. If some of the conditions are not met, note that the results of the interval may not be correct and then continue the process of the confidence interval.

1. The Kyoto Protocol was signed in 1997, and required countries to start reducing their carbon emissions. The protocol became enforceable in February 2005. Table ?? contains a random sample of CO₂ emissions in 2010 (CO₂ emissions (metric tons per capita), 2018). Find a 99% confidence interval for the mean CO₂ emissions in 2010.

Code book for data frame Emission is below Table ??.

2. The amount of sugar in a Krispy Kream glazed donut is 10 g. Many people feel that cereal is a healthier alternative for children over glazed donuts. Table ?? contains the amount of sugar in a sample of cereal that is geared towards children (breakfast cereal, 2019). Estimate the mean amount of sugar in children's cereal at the 95% confidence level.

Code book for data frame Sugar is below Table ??.

A new data frame will need to be created of just cereal for children. It is Table ??.

3. The FDA regulates that fish that is consumed is allowed to contain 1.0 mg/kg of mercury. In Florida, bass fish were collected in 53 different lakes to measure the health of the lakes. The data frame of measurements from Florida lakes is in Table ?? (NISER 081107 ID Data, 2019). Calculate with 90% confidence the mean amount of mercury in fish in Florida lakes. Is there too much mercury in the fish in Florida?

Code book for data frame Mercury is below Table ??.

4. The data frame Table ?? contains various variables about a person including their pulse rates before the subject exercised and after the subject ran in place for one minute. Estimate the mean pulse rate before exercise of females who do drink alcohol with a 95% level of confidence?

Code book for data frame Pulse below Table ??.

A new data frame with just females who drink alcohol is Table ?? from chapter 7.

5. The economic dynamism is an index of productive growth in dollars. Economic data for many countries are in Table ?? (SOCR Data 2008 World CountriesRankings, 2019).

Code book for data frame Economics is below Table ??.

A data frame that contains only middle income countries was created in chapter 7 and is Table ?. Find a 95% confidence interval for the mean economic dynamism for middle income countries.

6. Table ? contains the percentage of woman receiving prenatal care in a sample of countries over several years. (births per woman), 2019). Estimate the average percentage of women receiving prenatal care in 2009 (p2009) with a 95% confidence interval?

Code book for Data frame Fert_prenatal is below Table ?.

7. Maintaining your balance may get harder as you grow older. A study was conducted to see how steady the elderly is on their feet. They had the subjects stand on a force platform and have them react to a noise. The force platform then measured how much they swayed forward and backward, and the data is in Table ? (Maintaining Balance while Concentrating, 2019). Find the mean forward/backward sway of elderly person? Use a 95% confidence level. Follow the filtering methods in other homework problems to create a data frame for only Elderly.

Code book for data frame Sway is below Table ?.

10 Two Sample Inference

Chapter 7 discussed methods of hypothesis testing about one-population parameters. Chapter 8 discussed methods of estimating population parameters from one sample using confidence intervals. This chapter will look at methods of confidence intervals and hypothesis testing for two populations. Since there are two populations, there are two random variables, two means or proportions, and two samples (though with paired samples you usually consider there to be one sample with pairs collected). Examples of where you would do this are:

Testing and estimating the difference in testosterone levels of men before and after they had children (Gettler, McDade, Feranil & Kuzawa, 2011).

Testing the claim that a diet works by looking at the weight before and after subjects are on the diet.

Estimating the difference in proportion of those who approve of President Obama in the age group 18 to 26 year old and the 55 and over age group.

All of these are examples of hypothesis tests or confidence intervals for two populations. The methods to conduct these hypothesis tests and confidence intervals will be explored in this chapter. As a reminder, all hypothesis tests are the same process. The only thing that changes is the formula that you use and the conditions. Confidence intervals are also the same process, except that the formula is different.

10.1 Two Proportions

There are times you want to test a claim about two population proportions or construct a confidence interval estimate of the difference between two population proportions. As with all other hypothesis tests and confidence intervals, the process is the same though the formulas and conditions are different.

10.1.1 Hypothesis Test for Two Population Proportion (2-Prop Test)

1. State the random variables and the parameters in words.

x_1 = number of successes from group 1

x_2 = number of successes from group 2

p_1 = proportion of successes in group 1

p_2 = proportion of successes in group 2

2. State the null and alternative hypotheses and the level of significance

$H_o : p_1 = p_2$

$H_a : p_1 \neq p_2$. the \neq can be replaced with $<$ or $>$ depending on the question.

Also, state your α level here.

3. State and check the conditions for a hypothesis test

- a. State: A simple random sample of size n_1 is taken from population 1, and a simple random sample of size n_2 is taken from population 2. Check: describe how each sample was collected.
- b. State: The samples are independent. Check: describe why the two samples are independent.
- c. State: The properties for the binomial distribution are satisfied for both populations. Check: describe how each population meets all the properties.
- d. State: The sampling distribution of \hat{p}_1 can be approximated as a normal distribution. Check: To determine the sampling distribution of \hat{p}_1 , you need to show that $p_1 * n_1 \geq 5$ and $q_1 * n_1 \geq 5$ where $q_1 = 1 - p_1$. If this requirement is true, then the sampling distribution of \hat{p}_1 is well approximated by a normal curve. State: The sampling distribution of \hat{p}_2 can be approximated as a normal distribution. Check: To determine the sampling distribution of \hat{p}_2 , you need to show that $p_2 * n_2 \geq 5$ and $q_2 * n_2 \geq 5$ where $q_2 = 1 - p_2$. If this requirement is true, then the sampling distribution of \hat{p}_2 is well approximated by a normal curve. However, if you do not know p_1 and p_2 , you will need to use \hat{p}_1 and \hat{p}_2 instead. This is not perfect, but it is the best you can do.

4. Find the sample statistics, test statistic, and p-value

On rStudio, use the command

```
prop.test(c(x1,x2), c(n1, n2))
```

5. Conclusion

This is where you write reject or fail to reject H_o . The rule is: if the p-value $< \alpha$, then reject H_o . If the p-value $\geq \alpha$, then fail to reject H_o .

6. Interpretation

This is where you interpret in real world terms the conclusion to the test. The conclusion for a hypothesis test is that you either have enough evidence to support H_a , or you do not have enough evidence to support H_a .

10.1.2 Confidence Interval for the Difference Between Two Population Proportion (2-Prop Interval)

The confidence interval for the difference in proportions has the same random variables and proportions and the same conditions as the hypothesis test for two proportions. If you have already completed the hypothesis test, then you do not need to state them again. If you haven't completed the hypothesis test, then state the random variables and proportions and state and check the conditions before completing the confidence interval step.

1. Find the sample statistics and the confidence interval

The confidence interval estimate of the difference is found using the following command in R Studio:

`prop.test(c(x1,x2), c(n1, n2), conf.level=C)` Type C as a decimal

2. Statistical Interpretation: In general this looks like, "You are C% confident that the confidence interval contains the true difference in proportions."
3. Real World Interpretation: This is where you state how much more (or less) the first proportion is from the second proportion.

10.1.3 Example: Hypothesis Test for Two Population Proportions

Do husbands cheat on their wives in a different proportion from the proportion of wives cheat on their husbands ("Statistics brain," 2013)? Suppose you take a group of 1000 randomly selected husbands and find that 231 had cheated on their wives. Suppose in a group of 1200 randomly selected wives, 176 cheated on their husbands. Do the data show that the proportion of husbands who cheat on their wives is different from the proportion of wives who cheat on their husbands. Test at the 5% level.

10.1.3.1 Solution

1. State the random variables and the parameters in words.

x_1 = number of husbands who cheat on his wife

x_2 = number of wives who cheat on her husband

p_1 = proportion of husbands who cheat on his wife

p_2 = proportion of wives who cheat on her husband

2. State the null and alternative hypotheses and the level of significance

$$H_o : p_1 = p_2$$

$$H_a : p_1 \neq p_2$$

level of significance is $\alpha = 0.05$

3. State and check the conditions for a hypothesis test

- a. State: A simple random sample of 1000 responses about cheating from husbands is taken. Check: This was stated in the problem. State: A simple random sample of 1200 responses about cheating from wives is taken. Check: This was stated in the problem.
- b. State: The samples are independent. Check: The samples are independent. This is true since the samples involved different genders.
- c. State: The properties of the binomial distribution are satisfied in both populations. Check: This is true since there are only two responses, there are a fixed number of trials, the probability of a success is the same, and the trials are independent.
- d. State: The sampling distributions of \hat{p}_1 and \hat{p}_2 can be approximated with a normal distribution. Check: $n_1 * p_1$, $n_2 * p_2$, $n_1 * q_1$, and $n_2 * q_2$ are all greater than or equal to 5. So both sampling distributions of \hat{p}_1 and \hat{p}_2 can be approximated with a normal distribution.

4. Find the sample statistics, test statistic, and p-value

On r use the command:

```
prop.test(c(231,176), c(1000, 1200))
```

2-sample test for equality of proportions with continuity correction

```
data:  c out of c231 out of 1000176 out of 1200
X-squared = 25.173, df = 1, p-value = 5.241e-07
alternative hypothesis: two.sided
95 percent confidence interval:
 0.05050705 0.11815962
```

```
sample estimates:
  prop 1    prop 2
0.2310000 0.1466667
```

5. Conclusion

Reject H_o , since the p-value is less than 5%.

6. Interpretation

This is enough evidence to support that the proportion of husbands having affairs is different from the proportion of wives having affairs.

10.1.4 Example: Confidence Interval for Two Population Proportions

What is the difference in proportion that husbands cheat on their wives than wives cheat on the husbands (“Statistics brain,” 2013)? Suppose you take a group of 1000 randomly selected husbands and find that 231 had cheated on their wives. Suppose in a group of 1200 randomly selected wives, 176 cheated on their husbands. Estimate the difference in the proportion of husbands and wives who cheat on their spouses using a 95% confidence level.

10.1.4.1 Solution

1. State the random variables and the parameters in words.

These were stated in Example: Hypothesis Test for Two Population Proportions.

2. State and check the conditions for the confidence interval

The conditions were stated and checked in Example: Hypothesis Test for Two Population Proportions.

3. Find the sample statistics and the confidence interval

On r use the command:

```
prop.test(c(231,176), c(1000, 1200), conf.level = .95)
```

2-sample test for equality of proportions with continuity correction

```
data:  c out of c231 out of 1000176 out of 1200
X-squared = 25.173, df = 1, p-value = 5.241e-07
```

```

alternative hypothesis: two.sided
95 percent confidence interval:
 0.05050705 0.11815962
sample estimates:
  prop 1      prop 2
0.2310000 0.1466667

```

4. Statistical Interpretation: You are 95% confident that $0.05050705 < p_1 - p_2 < 0.11815962$ contains the true difference in proportions.
5. Real World Interpretation: The proportion of husbands who cheat on their wives is anywhere from 5.05% to 11.82% higher than the proportion of wives who cheat on their husband.

10.1.5 Homework for Two Proportions Section

In each problem show all steps of the hypothesis test or confidence interval. If some of the conditions are not met, note that the results of the test or interval may not be correct and then continue the process of the hypothesis test or confidence interval.

1. Many high school students take the AP tests in different subject areas. In 2007, of the 144,796 students who took the biology exam 84,199 of them were female. In that same year, of the 211,693 students who took the calculus AB exam 102,598 of them were female ("AP exam scores," 2013). Is there enough evidence to show that the proportion of female students taking the biology exam is different than the proportion of female students taking the calculus AB exam? Test at the 5% level.
2. Many high school students take the AP tests in different subject areas. In 2007, of the 144,796 students who took the biology exam 84,199 of them were female. In that same year, of the 211,693 students who took the calculus AB exam 102,598 of them were female ("AP exam scores," 2013). Estimate the difference in the proportion of female students taking the biology exam and female students taking the calculus AB exam using a 90% confidence level.
3. Many high school students take the AP tests in different subject areas. In 2007, of the 211,693 students who took the calculus AB exam 102,598 of them were female and 109,095 of them were male ("AP exam scores," 2013). Is there enough evidence to show that the proportion of female students taking the calculus AB exam is different from the proportion of male students taking the calculus AB exam? Test at the 5% level.
4. Many high school students take the AP tests in different subject areas. In 2007, of the 211,693 students who took the calculus AB exam 102,598 of them were female and 109,095 of them were male ("AP exam scores," 2013). Estimate using a 90% level the

difference in proportion of female students taking the calculus AB exam versus male students taking the calculus AB exam.

5. Are there more children diagnosed with Autism Spectrum Disorder (ASD) in states that have larger urban areas over states that are mostly rural? In the state of Pennsylvania, a fairly urban state, there are 245 eight year old diagnosed with ASD out of 18,440 eight year old evaluated. In the state of Utah, a fairly rural state, there are 45 eight year old diagnosed with ASD out of 2,123 eight year old evaluated (“Autism and developmental,” 2008). Is there enough evidence to show that the proportion of children diagnosed with ASD in Pennsylvania is different than the proportion in Utah? Test at the 1% level.
6. Are there more children diagnosed with Autism Spectrum Disorder (ASD) in states that have larger urban areas over states that are mostly rural? In the state of Pennsylvania, a fairly urban state, there are 245 eight year old diagnosed with ASD out of 18,440 eight year old evaluated. In the state of Utah, a fairly rural state, there are 45 eight year old diagnosed with ASD out of 2,123 eight year old evaluated (“Autism and developmental,” 2008). Estimate the difference in proportion of children diagnosed with ASD between Pennsylvania and Utah. Use a 98% confidence level.
7. A child dying from an accidental poisoning is a terrible incident. Is it more likely that a male child will get into poison than a female child? To find this out, data was collected that showed that out of 1830 children between the ages one and four who pass away from poisoning, 1031 were males and 799 were females (Flanagan, Rooney & Griffiths, 2005). Do the data show that there is different proportion of male children dying of poisoning than female children? Test at the 1% level.
8. A child dying from an accidental poisoning is a terrible incident. Is it more likely that a male child will get into poison than a female child? To find this out, data was collected that showed that out of 1830 children between the ages one and four who pass away from poisoning, 1031 were males and 799 were females (Flanagan, Rooney & Griffiths, 2005). Compute a 99% confidence interval for the difference in proportions of poisoning deaths of male and female children ages one to four.

10.2 Paired Samples for Two Means

Are two populations the same? Is the average height of men taller than the average height of women? Is the mean weight less after a diet than before?

You can compare populations by comparing their means. You take a sample from each population and compare the statistics.

Anytime you compare two populations you need to know if the samples are independent or dependent. The formulas you use are different for different types of samples.

If how you choose one sample has no effect on the way you choose the other sample, the two samples are **independent**. The way to think about it is that in independent samples, the observations from one sample are overall different from the observations from the other sample. This will mean that sample one has no affect on sample two. The sample values from one sample are not related or paired with values from the other sample.

If you choose the samples so that a measurement in one sample is paired with a measurement from the other sample, the samples are **dependent** or **matched** or **paired**. (Often a before and after situation.) You want to make sure there is a meaning for pairing data values from one sample with a specific data value from the other sample. One way to think about it is that in dependent samples, the observations from one sample are the same observations from the other sample, though there can be other reasons to pair values. This makes the sample values from each sample paired.

In tidy data, remember each row is a unit of observation, and each column is a variable. In paired samples, you would have two variables that you are working with. In independent samples, you would have a variable that distinguishes an observation from another observation. As an example, in the Pulse data frame, consider the variables `pulse_before` and `pulse_after`. Since they are measured off the same observation, then comparing the two variables would be a paired samples analysis. However, consider the `pulse_after` and whether a person smokes would be comparing the variable `pulse_after` against the variable `smokes` to see if smoking effects a person's pulse rate after exercise. In this case, the observations would be different based on smoking yes or smoking no. Consider the variable `smoking` to be the factor that one is interested in seeing how it effects pulse rate in the data frame Table ??.

10.2.1 Example: Independent or Dependent Samples

Determine if the following are dependent or independent samples.

- Randomly choose 5 men and 6 women and compare their heights
- Choose 10 men and weigh them. Give them a new diet drug and later weigh them again.
- Take 10 people and measure the strength of their dominant arm and their non-dominant arm.

10.2.1.1 Solution

- Randomly choose 5 men and 6 women and compare their heights

Independent, since there is no reason that one value belongs to another. The units of observations are not the same for both samples. The units of observations are definitely different. A way to think about this is that the knowledge that a man is chosen in one sample does not give any information about any of the woman chosen in the other sample.

- b. Choose 10 men and weigh them. Give them a new diet drug and later weigh them again.

Dependent, since each person's before weight can be matched with their after weight. The units of observations are the same for both samples. A way to think about this is that the knowledge that a person weighs 400 pounds at the beginning will tell you something about their weight after the diet drug.

- c. Take 10 people and measure the strength of their dominant arm and their non-dominant arm.

Dependent, since you can match the two arm strengths. The units of observations are the same for both samples. So the knowledge of one person's dominant arm strength will tell you something about the strength of their non-dominant arm.

To analyze data when there are matched or paired samples, called dependent samples, you conduct a paired t-test. Since the samples are matched, you can find the difference between the values of the two random variables.

10.2.2 Hypothesis Test for Two Sample Paired t-Test

1. State the random variables and the parameters in words.

x_1 = random variable 1

x_2 = random variable 2

μ_1 = mean of random variable 1

μ_2 = mean of random variable 2

2. State the null and alternative hypotheses and the level of significance

The hypotheses would be

$$H_o : \mu_1 = \mu_2 \text{ or } H_o : \mu_1 - \mu_2 = 0$$

$$H_a : \mu_1 \neq \mu_2 \text{ or } H_a : \mu_1 - \mu_2 \neq 0$$

However, since you are finding the differences, then you can actually think of $\mu_1 - \mu_2 = \mu_d$.

So the hypotheses could become

$$H_o : \mu_d = 0$$

$$H_a : \mu_d \neq 0 \text{ Remember, you can replace } \neq \text{ with } < \text{ or } >.$$

Also, state your α level here.

3. State and check the conditions for the hypothesis test

- a. State: A random sample of n pairs is taken. Check: state how the sample was collected.
 - b. Check: The population of the difference between random variables is normally distributed. Check: In this case the population you are interested in has to do with the differences that you find. It does not matter if each random variable is normally distributed. It is only important if the differences you find are normally distributed. Just as before, the t-test is fairly robust to the condition if the sample size is large. This means that if this condition isn't met, but your sample size is quite large, then the results of the t-test are valid.
4. Find the sample statistic, test statistic, and p-value

Realize that a paired test is a one sample t-test on the difference between two variables. So you are running a one-sample t-test on a new variable known as the difference variable. You need to create this difference variable by creating a new data frame. This is done on rStudio by doing the following command (The following shows how to create the variable difference for pulse_after-pulse_before on the data frame Pulse. Change the variables used and data frame used to your data frame and variables):

```
Pulse<-
  Pulse |>
  mutate(difference=pulse_after-pulse_before)
knitr::kable(head(Pulse))
```

Table 10.1: Pulse Data frame with Difference Column Added

height	weight	age	gender	smokes	alcohol	exercise	ran	pulse_before	pulse_after	year	difference
170	68	22	male	yes	yes	moderate	sat	70	71	93	1
182	75	26	male	yes	yes	moderate	sat	80	76	93	-4
180	85	19	male	yes	yes	moderate	ran	68	125	95	57
182	85	20	male	yes	yes	low	sat	70	68	95	-2
167	70	22	male	yes	yes	low	sat	92	84	96	-8
178	86	21	male	yes	yes	low	sat	76	80	98	4

Notice rStudio added a new variable called difference to the data frame Table ??. Now to conduct a paired t-test use the rStudio command

```
t.test(~difference_variable, data=Data_Frame)
```

Note: if the H_a is $<$, then the command becomes

```
t.test(~difference_variable, data=Data_Frame, alternative="less")
```

Similarly for $>$ put alternative="greater"

5. Conclusion

This is where you write reject H_o or fail to reject H_o . The rule is: if the p-value $< \alpha$, then reject H_o . If the p-value $\geq \alpha$, then fail to reject H_o .

6. Interpretation

This is where you interpret in real world terms the conclusion to the test. The conclusion for a hypothesis test is that you either have enough evidence to support H_a , or you do not have enough evidence to support H_a .

10.2.3 Confidence Interval for Difference in Means from Paired Samples (t-Interval)

The confidence interval for the difference in means has the same random variables and means and the same conditions as the hypothesis test for two paired samples. If you have already completed the hypothesis test, then you do not need to state them again. If you haven't completed the hypothesis test, then state the random variables and means, and state and check the conditions before completing the confidence interval step.

1. Find the sample statistic and confidence interval. Again, you will need to create a new data frame with a difference variable. Then on rStudio the command is
`t.test(~difference_variable, data=Data_Frame, conf.level=C)` Type C as a decimal
2. Statistical Interpretation: In general this looks like, "You are C% confident that the statement contains the true mean difference."
3. Real World Interpretation: This is where you state what interval contains the true mean difference.

10.2.4 Example: Hypothesis Test for Paired Samples

Is the pulse rate after exercise different from the pulse rate before exercise for a woman who drinks alcohol? Use the data frame Table ???. Test at the 5% level.

Code book for data frame Pulse below Table ??.

10.2.4.1 Solution

1. State the random variables and the parameters in words.

x_1 = pulse of a smoking woman who drinks alcohol after exercise

x_2 = pulse of a smoking woman who drinks alcohol before exercise

μ_1 = mean pulse of a smoking woman who drinks alcohol after exercise

μ_2 = mean pulse of a smoking woman who drinks alcohol before exercise

2. State the null and alternative hypotheses and the level of significance

$$H_o : \mu_1 = \mu_2$$

$$H_a : \mu_1 \neq \mu_2$$

level of significance, $\alpha = 0.05$

3. State and check the conditions for the hypothesis test

- a. State: A random sample of 110 pairs of pulse rates after and before exercise was taken. Check: The data frame says that the data was collected from students in classes for several years. Though this was not a random sample, it is probably a representative sample.
- b. State: The population of the difference in after and before pulse rates is normally distributed. Check: To see if this is true, look at the density plot and the normal quantile plot for the difference between after and before. This variable must be created before the density plot and normal quantile plot can be created. The data frame Table ?? is females who drink alcohol.

```
Pulse_female<-
  Pulse |>
  filter(gender=="female", alcohol=="yes")
knitr::kable(head(Pulse_female))
```

Table 10.2: Pulse Rates Before and After Exercise of Females who do drink Alcohol with Difference

height	weight	age	gender	smokes	alcohol	exercise	ran	pulse_before	pulse_after	year	difference
165	60	19	female	yes	yes	low	ran	88	120	98	32
163	47	23	female	yes	yes	low	ran	71	125	98	54
173	57	18	female	no	yes	moderate	sat	86	88	93	2
179	58	19	female	no	yes	moderate	ran	82	150	93	68

Table 10.2: Pulse Rates Before and After Exercise of Females who do drink Alcohol with Difference

height	weight	age	gender	smokes	alcohol	exercise	ran	pulse_before	pulse_after	year	difference
167	62	18	female	no	yes	high	ran	96	176	93	80
173	64	18	female	no	yes	low	sat	90	88	93	-2

Now mutate Table ?? data frame to include a difference variable.

```
Pulse_female<-
  Pulse_female |>
  mutate(difference=pulse_after-pulse_before)
knitr::kable(head(Pulse_female))
```

Table 10.3: Pulse Rates Before and After Exercise of Females who do drink Alcohol with Difference

height	weight	age	gender	smokes	alcohol	exercise	ran	pulse_before	pulse_after	year	difference
165	60	19	female	yes	yes	low	ran	88	120	98	32
163	47	23	female	yes	yes	low	ran	71	125	98	54
173	57	18	female	no	yes	moderate	sat	86	88	93	2
179	58	19	female	no	yes	moderate	ran	82	150	93	68
167	62	18	female	no	yes	high	ran	96	176	93	80
173	64	18	female	no	yes	low	sat	90	88	93	-2

Using Table ?? create a density plot and normal quantile plot on the difference variable.

```
gf_density(~difference, data=Pulse_female, title = "Difference in Pulse Rates for Females who drink Alcohol")
```

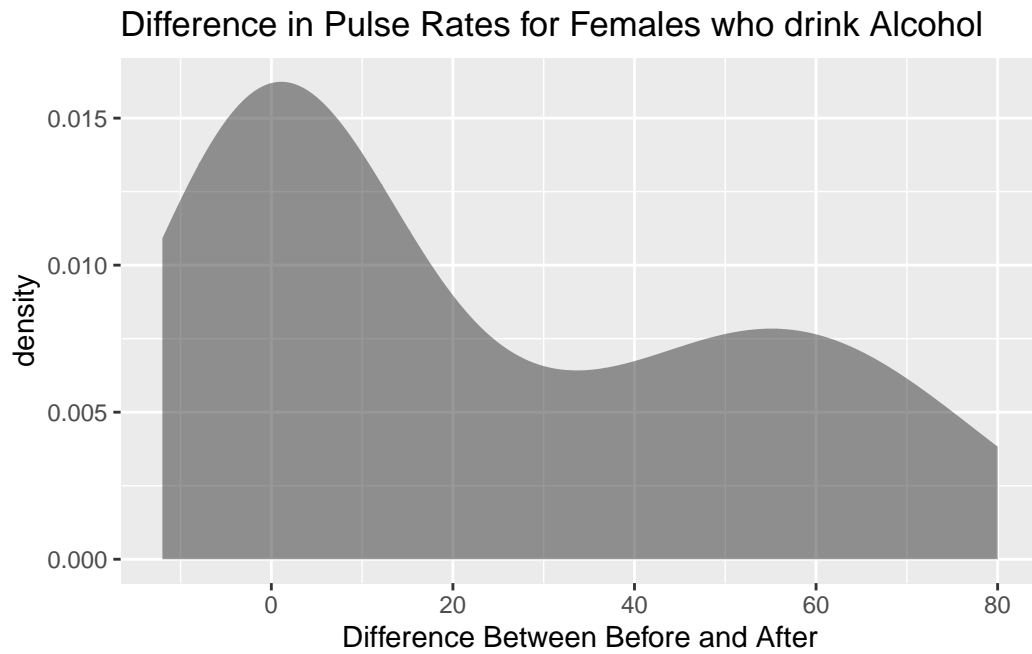


Figure 10.1: Density plot of differences in pulse rates

```
gf_qq(~difference, data=Pulse_female, title = "Difference in Pulse Rates for Females who drink Alcohol")
```

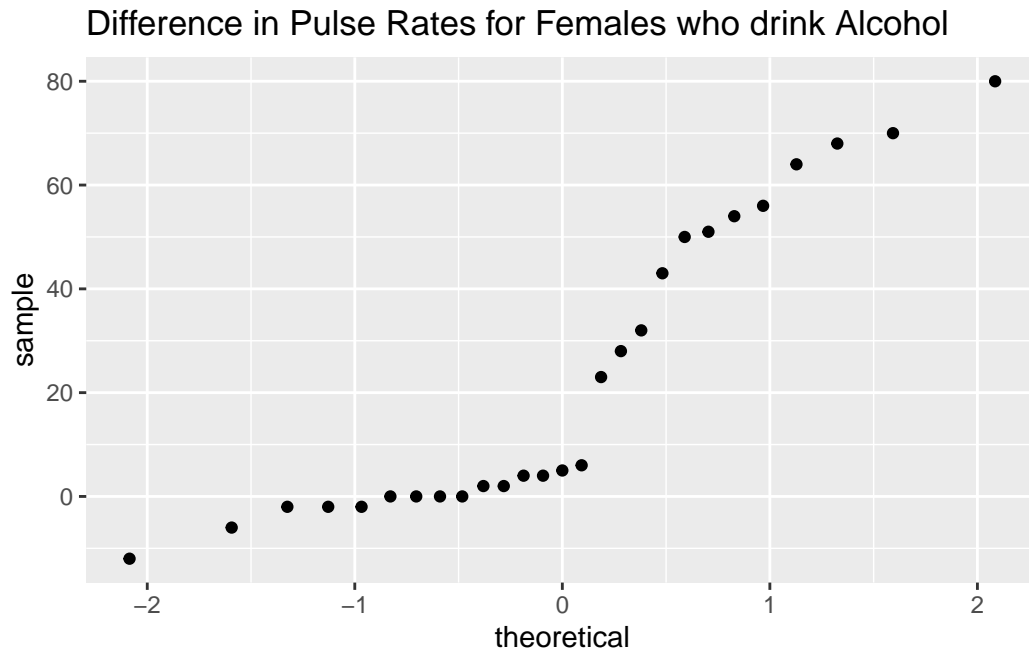



Figure 10.2: Normal Quantile Plot of Differences in Pulse Rates

The density plot is not symmetrical and the normal quantile plot on the differences is not linear. So you cannot assume that the distribution of the difference in pulse rates is normal. It is good that the t-test is robust if there is a large sample. The sample is of size 110, so that should be adequate to assume the conclusion is valid.

4. Find the sample statistic, test statistic, and p-value On r Studio, use the command:

```
t.test(~difference, data=Pulse_female)
```

One Sample t-test

```
data: difference
t = 4.1353, df = 26, p-value = 0.0003283
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 11.51152 34.26625
sample estimates:
mean of x
 22.88889
```

5. Conclusion

Since the p-value < 0.05 , reject H_o .

6. Interpretation

There is enough evidence to support that there is a difference in pulse rate before and after exercise of females who smoke.

10.2.5 Example: Hypothesis Test for Paired Samples

The New Zealand Air Force purchased a batch of flight helmets. They then found out that the helmets didn't fit. In order to make sure that they order the correct size helmets, they measured the head size of recruits. To save money, they wanted to use cardboard calipers, but were not sure if they will be accurate enough. So they took 18 recruits and measured their heads with the cardboard calipers and also with metal calipers. The data frame is in Table ?? (Helmet Sizes for New Zealand Airforce, 2019). Do the data provide enough evidence to show that there is a difference in measurements between the cardboard and metal calipers? Use a 5% level of significance.

```
Helmet<-read.csv( "https://krkozak.github.io/MAT160/helmet.csv")
knitr::kable(head(Helmet))
```

Table 10.4: Helmet Head Measurements

Cardboard	Metal
146	145
151	153
163	161
152	151
151	145
151	150

Code book for data frame **Helmet**

Description After purchasing a batch of flight helmets that did not fit the heads of many pilots, the NZ Airforce decided to measure the head sizes of all recruits. Before this was carried out, information was collected to determine the feasibility of using cheap cardboard calipers to make the measurements, instead of metal ones which were expensive and uncomfortable. The data lists the head diameters of 18 recruits measured once using cardboard calipers and again using metal calipers. One question is whether there is any systematic difference between the two sets of calipers. One might also ask whether there is more variability in the cardboard calipers measurement than that of the metal calipers.

This data frame contains the following columns:

Cardboard: measurement using cardboard calipers (cm)

Metal: measurement using metal calipers (cm)

Source Helmet Sizes for New Zealand Airforce. (n.d.). Retrieved July 20, 2019, from <http://www.statsci.org/data/oz/nzhelmet.html>

References Data courtesy of Dr Stephen Legg. Seber and Lee (1998). Page 545.

10.2.5.1 Solution

1. State the random variables and the parameters in words.

x_1 = head measurement of recruit using cardboard caliper

x_2 = head measurement of recruit using metal caliper

μ_1 = mean head measurement of recruit using cardboard caliper

μ_2 = mean head measurement of recruit using metal caliper

2. State the null and alternative hypotheses and the level of significance

$$H_o : \mu_1 = \mu_2$$

$$H_a : \mu_1 \neq \mu_2$$

level of significance, $\alpha = 0.05$

3. State and check the conditions for the hypothesis test

- a. State: A random sample of 18 pairs of head measures of recruits with cardboard and metal caliper was taken. Check: This was not stated, but probably could be safely assumed.
- b. State: The population of the difference in head measurements between cardboard and metal calipers is normally distributed. Check: First create the difference variable, then the density plot and normal quantile plot.

```
Helmet<-  
  Helmet |>  
  mutate(difference=Cardboard-Metal)  
knitr::kable(head(Helmet))
```

Table 10.5: Helmet Head Measurements

Cardboard	Metal	difference
146	145	1
151	153	-2
163	161	2
152	151	1
151	145	6
151	150	1

```
gf_density(~difference, data=Helmet, title="Differences in Head Measurements", xlab="Differences in Head Measurements")
```

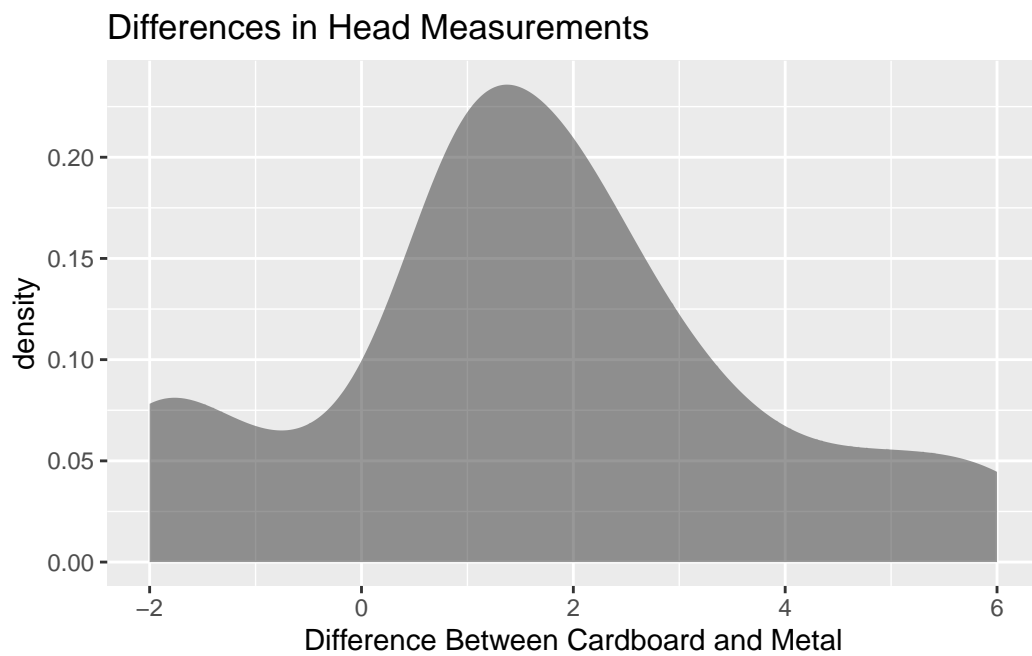


Figure 10.3: Density plot of differences in head measurements

```
gf_qq(~difference, data=Helmet, title="Differences in Head Measurements")
```

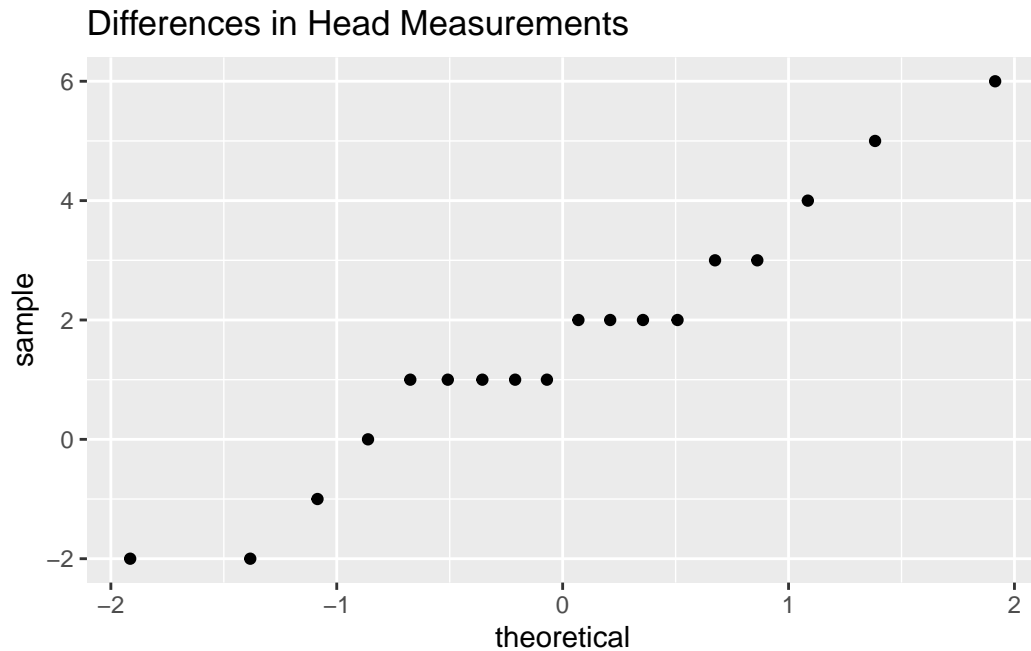


Figure 10.4: Normal Quantile Plot of Differences in Head Measurements

This density plot Figure ?? looks somewhat bell shaped. The normal quantile plot Figure ?? on the differences looks somewhat linear. So you can assume that the distribution of the difference in weights is normal.

4. Find the sample statistic, test statistic, and p-value

Using rStudio the command is

```
t.test(~difference, data=Helmet)
```

One Sample t-test

```
data: difference
t = 3.1854, df = 17, p-value = 0.005415
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 0.5440163 2.6782060
sample estimates:
mean of x
 1.611111
```

The sample statistic is 1.6111, the test statistic is 3.1854, and the p-value is 0.005415.

5. Conclusion

Since the p-value < 0.05 , reject H_o .

6. Interpretation

There is enough evidence to support that the mean head measurements using the cardboard calipers are not the same as when using the metal calipers. So it looks like the New Zealand Air Force shouldn't use the cardboard calipers.

10.2.6 Example: Confidence Interval for Paired Samples

The New Zealand Air Force purchased a batch of flight helmets. They then found out that the helmets didn't fit. In order to make sure that they order the correct size helmets, they measured the head size of recruits. To save money, they wanted to use cardboard calipers, but were not sure if they will be accurate enough. So they took 18 recruits and measured their heads with the cardboard calipers and also with metal calipers. The data frame is in Table ?? (Helmet Sizes for New Zealand Airforce, 2019). Estimate the difference in measurements between the cardboard and metal calipers using a 95% confidence interval.

10.2.6.1 Solution

1. State the random variables and the parameters in words.

These were stated in Example: Hypothesis Test for Paired Samples.

2. State and check the conditions for the confidence interval

The conditions were stated and checked in Example: Hypothesis Test for Paired Samples.

3. Find the sample statistic and confidence interval

Using the data frame Table ?? the rStudio the command is

```
t.test(~difference, data=Helmet, conf.level=0.95)
```

One Sample t-test

```
data: difference  
t = 3.1854, df = 17, p-value = 0.005415
```

alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
0.5440163 2.6782060
sample estimates:
mean of x
1.611111

4. Statistical Interpretation: You are 95% confidence that $0.5440163 < \mu_1 - \mu_2 < 2.6782060$ contains the true mean difference in head measurement between using the cardboard and metal calibers.
5. Real World Interpretation: The mean head measurement using the cardboard calibers is anywhere from 0.54 cm to 2.68 cm more than the head measurement using the metal calibers.

Examples 9.2.6 and 9.2.7 use the same data set, but one is conducting a hypothesis test and the other is conducting a confidence interval. Notice that the hypothesis test's conclusion was to reject and say that there was a difference in the means, and the confidence interval does not contain the number 0. If the confidence interval did contain the number 0, then that would mean that the two means could be the same. Since the interval did not contain 0, then you could say that the means are different just as in the hypothesis test. This means that the hypothesis test and the confidence interval can produce the same interpretation. Do be careful though, you can run a hypothesis test with a particular significance level and a confidence interval with a confidence level that is not compatible with your significance level. This will mean that the conclusion from the confidence interval would not be the same as with a hypothesis test. So if you want to estimate the mean difference, then conduct a confidence interval. If you want to show that the means are different, then conduct a hypothesis test. As a reminder, the American Statistical Association (ASA) suggests not conducting hypothesis tests and just create confidence intervals.

10.2.7 Homework for Paired Samples for Two Means Section

In each problem show all steps of the hypothesis test or confidence interval. If some of the conditions are not met, note that the results of the test or interval may not be correct and then continue the process of the hypothesis test or confidence interval.

1. The cholesterol level of patients who had heart attacks was measured multiple times after the heart attack. The researchers want to see if the cholesterol level of patients who have heart attacks changes as the time since their heart attack increases. The data is in Table ???. Do the data show that the mean cholesterol level of patients that have had a heart attack changes as the time increases since their heart attack? Use day2 and day4 variables to answer the question. Test at the 1% level.

Code book for Data Frame Cholesterol is below Table ??.

2. The cholesterol level of patients who had heart attacks was measured multiple times after the heart attack. The researchers want to see if the cholesterol level of patients who have heart attacks changes as the time since their heart attack increases. The data is in Table ??. Calculate a 98% confidence interval for the mean difference in cholesterol levels from day two to day four.
3. All Fresh Seafood is a wholesale fish company based on the east coast of the U.S. Catalina Offshore Products is a wholesale fish company based on the west coast of the U.S. Table ?? contains prices from both companies for specific fish types (\“Seafood online,\” 2013) (\“Buy sushi grade,\” 2013). Do the data provide enough evidence to show that fish cost different from west coast fish wholesaler and east coast wholesaler? Test at the 5% level.

```
Price <- read.csv( "https://krkozak.github.io/MAT160/price.csv")
knitr::kable(head(Price))
```

Table 10.6: Wholesale Prices of Fish in Dollars

	fish	east	west
Cod		19.99	17.99
Tilapi		6.00	13.99
Farmed Salmon		19.99	22.99
Organic Salmon		24.99	24.99
Grouper Fillet		29.99	19.99
Tuna		28.99	31.99

Code book for data frame Price

Description Price of fish was collected from two websites. One for Catalina Offshore Products (west coast) and the other for All Fresh Seafood (east coast) in 2013.

This data frame contains the following columns:

fish: type of fish for sale

east: price of fish from east coast supplier (\$)

west: price of fish from west coast supplier (\$)

Source Seafood online. (2013, November 20). Retrieved from <http://www.allfreshseafood.com/>

Buy sushi grade fish online. (2013, November 20). Retrieved from <http://www.catalinaop.com/>

References Websites of Catalina Offshore Products and All Fresh Seafood

4. All Fresh Seafood is a wholesale fish company based on the east coast of the U.S. Catalina Offshore Products is a wholesale fish company based on the west coast of the U.S. Table ?? contains prices from both companies for specific fish types (\“Seafood online,\” 2013) (\“Buy sushi grade,\” 2013). Find a 95% confidence interval for the mean difference in wholesale price between the east coast and west coast suppliers.
5. The British Department of Transportation studied to see if people avoid driving or shopping, or have more accidents on Friday the 13th. They collected data from different locations (Friday the 13th, 2019). The data for each location on the two different dates is in Table ?. Do the data show that on average different number of people are engaged in activities on Friday the 13th? Test at the 5% level.

```
Traffic <- read.csv( "https://krkozak.github.io/MAT160/traffic.csv")
knitr::kable(head(Traffic))
```

Table 10.7: Traffic Count

	source	year	month	X6th	X13th	location
	traffic	1990,	July	139246	138548	7 to 8
	traffic	1990,	July	134012	132908	9 to 10
	traffic	1991,	September	137055	136018	7 to 8
	traffic	1991,	September	133732	131843	9 to 10
	traffic	1991,	December	123552	121641	7 to 8
	traffic	1991,	December	121139	118723	9 to 10

Code book for data frame Traffic

Description This file consists of three separate data sets, all of which address the issues of how superstitions regarding Friday the 13th affect human behavior, and whether Friday the 13th is an unlucky day. Scanlon, et al. collected data on traffic and shopping patterns and accident frequency for Fridays the 6th and 13th between October of 1989 and November of 1992.

For the first data set, the researchers obtained information from the British Department of Transport regarding the traffic flows between junctions 7 to 8 and junctions 9 to 10 of the M25 motorway. They collected the numbers of shoppers in nine different supermarkets in southeast England for the second data set. The third data set contains numbers of emergency admissions to hospitals due to transport accidents.

We present the three data sets in a combined format, with the variable “Data set” as an identifier that may be used to separate them.

This data frame contains the following columns:

source: which data set the data were obtained from

year: which year the data was collected from

Month: the month that the Friday was in

x6th: Number of cars passing through junction (traffic data set), shoppers for each supermarket (shopping data set), or admissions due to transport accidents (accident data set) on Friday the 6th

x13th: Number of cars passing through junction (traffic data set), shoppers for each supermarket (shopping data set), or admissions due to transport accidents (accident data set) on Friday the 13th

location: Motorway junction (traffic data set), supermarket location (shopping data set) or hospital (accident data set) to which the data correspond

Source (n.d.). Retrieved from <https://www3.nd.edu/~busiforc/handouts/Data and Stories/test/Friday The Thirteenth/Friday The Thirteenth Data.html>

References Scanlon, T.J., Luben, R.N., Scanlon, F.L., Singleton, N. (1993), "Is Friday the 13th Bad For Your Health?," BMJ, 307, 1584-1586.

6. The British Department of Transportation studied to see if people avoid driving or shopping, or have more accidents on Friday the 13th. They collected data from different locations (Friday the 13th, 2019). The data for each location on the two different dates is in Table ?? . Do the data show that on average different number of people are engaged in activities on Friday the 13th? Estimate the mean difference in activity count between the 6th and the 13th using a 95% level.
7. To determine if Reiki is an effective method for treating pain, a pilot study was carried out where a certified second-degree Reiki therapist provided treatment on volunteers. Pain was measured using a visual analogue scale (VAS) and a likert scale immediately before and after the Reiki treatment (Olson & Hanson, 1997). The data is in Table ?? . Do the data show that Reiki treatment reduces pain? Test at the 5% level.

Code book for data frame Reiki is below Table ??.

8. To determine if Reiki is an effective method for treating pain, a pilot study was carried out where a certified second-degree Reiki therapist provided treatment on volunteers. Pain was measured using a visual analogue scale (VAS) and a likert scale immediately before and after the Reiki treatment (Olson & Hanson, 1997). The data is in Table ?? . Compute a 90% confidence level for the mean difference in VAS score from before and after Reiki treatment.

9. The female labor force participation rates (FLFPR) of women in countries from 1990 to 2018 are in table 9.2.8.5 (Labor force participation rate, female (% of female population ages 15+) (modeled ILO estimate), 2019). Do the data show that the mean female labor force participation rate in 1990 is different from that in the 2018 using a 5% level of significance?

```
Labor <- read.csv( "https://krkozak.github.io/MAT160/labor.csv")
knitr::kable(head(Labor))
```

Table 10.8: Female Labor Force Participation Rates

Country	1990	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018
Arab World High income & Caribbean	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN
Africa South Asian	43.43	43.18	43.00	42.80	42.60	42.40	42.20	42.00	41.80	41.60	41.40	41.20	41.00	40.80	40.60	40.40	40.20	40.00	39.80	39.60	39.40	39.20	39.00	38.80	38.60	38.40	38.20	38.00	37.80
Angola Sub-Saharan Africa in-come	74.50	73.80	73.10	72.40	71.70	71.00	70.30	69.60	68.90	68.20	67.50	66.80	66.10	65.40	64.70	64.00	63.30	62.60	61.90	61.20	60.50	59.80	59.10	58.40	57.70	57.00	56.30	55.60	54.90
Albania Europe & middle income	58.56	58.50	58.40	58.30	58.20	58.10	58.00	57.90	57.80	57.70	57.60	57.50	57.40	57.30	57.20	57.10	57.00	56.90	56.80	56.70	56.60	56.50	56.40	56.30	56.20	56.10	56.00	55.90	55.80
Andorra High & in-come	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN	AN
Arab World	19.18	19.09	19.01	18.93	18.85	18.77	18.69	18.61	18.53	18.45	18.37	18.29	18.21	18.13	18.05	17.97	17.89	17.81	17.73	17.65	17.57	17.49	17.41	17.33	17.25	17.17	17.09	17.01	16.93

Code book for data frame Labor

Description Labor force participation rate, female (% of female population ages 15+)

This data frame contains the following columns:

Country Name: The name of a country around the world

Country Code: The 3 letter country code

Region: The location of the country in the world

IncomeGroup: The World Bank's income classification

y1990-y2018: Labor force participation rate, female (% of female population ages 15+) for the years 1990–2018

Source Labor force participation rate, female (% of female population ages 15) (modeled ILO estimate). (n.d.). Retrieved July 20, 2019, from <https://data.worldbank.org/indicator/SL.TLF.CACT.FE.ZS>

References International Labour Organization, ILOSTAT database. Data retrieved in April 2019.

10. The female labor force participation rates (FLFPR) of women in countries from 1990 to 2018 are in Table ?? (Labor force participation rate, female (% of female population ages 15+) (modeled ILO estimate), 2019). Estimate the mean difference in the female labor force participation rate in 1990 to 2018 using a 95% confidence level?
11. Is the pulse rate after exercise different from the pulse rate before exercise for a man who drinks alcohol but doesn't smoke? Use the data frame Pulse Table ?. Test at the 5% level.

Code book for data frame Pulse is below Table ??.

12. Table ?? contains pulse rates Compute a 95% confidence interval for the mean difference in pulse rates from before and after exercise for males who drink but do not smoke.

10.3 Independent Samples for Two Means

This section will look at how to analyze when two samples are collected that are independent. As with all other hypothesis tests and confidence intervals, the process is the same though the formulas and conditions are different.

10.3.1 Hypothesis Test for the Difference in Means from Two Independent Samples

1. State the random variables and the parameters in words.

x_1 = random variable 1

x_2 = random variable 2

μ_1 = mean of random variable 1

μ_2 = mean of random variable 2

2. State the null and alternative hypotheses and the level of significance

The hypotheses would be

$$H_o : \mu_1 = \mu_2$$

$$H_a : \mu_1 \neq \mu_2, \text{ the } \neq \text{ can be replaced with } < \text{ or } >$$

Also, state your α level here.

3. State and check the conditions for the hypothesis test
 - a. State: A random sample of size n_1 is taken from population 1. A random sample of size n_2 is taken from population 2. Check: describe how both samples are collected. Note: the samples do not need to be the same size, but the test is more robust if they are.
 - b. State: The two samples are independent. Check: describe why the samples are independent of each other.
 - c. State: Population 1 is normally distributed. Population 2 is normally distributed. Check: draw the density graph and normal quantile plot for both samples and discuss if they meet the criteria. Just as before, the t-test is fairly robust to the condition if the sample size is large. This means that if this condition isn't met, but your sample sizes are quite large, then the results of the t-test are valid.
 - d. State: The population variances are unknown and not assumed to be equal. The old condition is that the variances are equal. However, this condition is no longer a condition that most statisticians use. This is because it isn't really realistic to assume that the variances are equal. So just assume the condition of the variances being unknown and not assumed to be equal is true, and it will not be checked.
4. Find the sample statistic, test statistic, and p-value

The command using r is

```
t.test(variable~factor, data=Data_Frame)
```

Note: if the H_a is $<$, then the command becomes

```
t.test(variable~factor, data=Data_Frame, alternative="less")
```

Similarly for $>$ put `alternative="greater"`

5. Conclusion

This is where you write reject or fail to reject H_0 . The rule is: if the p-value $< \alpha$, then reject H_0 . If the p-value $\geq \alpha$, then fail to reject H_0 .

6. Interpretation

This is where you interpret in real world terms the conclusion to the test. The conclusion for a hypothesis test is that you either have enough evidence to support H_a , or you do not have enough evidence to support H_a .

10.3.2 Confidence Interval for the Difference in Means from Two Independent Samples

The confidence interval for the difference in means has the same random variables and means and the same conditions as the hypothesis test for independent samples. If you have already completed the hypothesis test, then you do not need to state them again. If you haven't completed the hypothesis test, then state the random variables and means and state and check the conditions before completing the confidence interval step.

Find the sample statistic and confidence interval

On r Studio, the command is

```
t.test(variable~factor, data=Data_Frame, conf.level=C) type C as a decimal
```

2. Statistical Interpretation: In general this looks like, "You are C% confident that the interval contains the true mean difference."
3. Real World Interpretation: This is where you state what interval contains the true difference in means, though often you state how much more (or less) the first mean is from the second mean.

10.3.3 Example: Hypothesis Test for Two Means

The cholesterol level of people vary for many reasons. The question is do people with diabetes have different cholesterol levels from people who do not have diabetes? Use the NHANES data frame. Test at the 5% level.

```
names(NHANES) #displays the names of the variables in a data frame
```

[1] "ID"	"SurveyYr"	"Gender"	"Age"
[5] "AgeDecade"	"AgeMonths"	"Race1"	"Race3"
[9] "Education"	"MaritalStatus"	"HHIncome"	"HHIncomeMid"
[13] "Poverty"	"HomeRooms"	"HomeOwn"	"Work"
[17] "Weight"	"Length"	"HeadCirc"	"Height"
[21] "BMI"	"BMICatUnder20yrs"	"BMI_WHO"	"Pulse"
[25] "BPSysAve"	"BPDiaAve"	"BPSys1"	"BPDia1"
[29] "BPSys2"	"BPDia2"	"BPSys3"	"BPDia3"
[33] "Testosterone"	"DirectChol"	"TotChol"	"UrineVol1"
[37] "UrineFlow1"	"UrineVol2"	"UrineFlow2"	"Diabetes"
[41] "DiabetesAge"	"HealthGen"	"DaysPhysHlthBad"	"DaysMentHlthBad"
[45] "LittleInterest"	"Depressed"	"nPregnancies"	"nBabies"
[49] "Age1stBaby"	"SleepHrsNight"	"SleepTrouble"	"PhysActive"
[53] "PhysActiveDays"	"TVHrsDay"	"CompHrsDay"	"TVHrsDayChild"
[57] "CompHrsDayChild"	"Alcohol12PlusYr"	"AlcoholDay"	"AlcoholYear"
[61] "SmokeNow"	"Smoke100"	"Smoke100n"	"SmokeAge"
[65] "Marijuana"	"AgeFirstMarij"	"RegularMarij"	"AgeRegMarij"
[69] "HardDrugs"	"SexEver"	"SexAge"	"SexNumPartnLife"
[73] "SexNumPartYear"	"SameSex"	"SexOrientation"	"PregnantNow"

Code book for data frame NHANES type `help("NHANES")` in the `r` Console.

10.3.3.1 Solution

1. State the random variables and the parameters in words.

x_1 = Cholesterol level of people with diabetes

x_2 = Cholesterol level of people without diabetes

μ_1 = mean cholesterol level of people with diabetes

μ_2 = mean cholesterol level of people without diabetes

2. State the null and alternative hypotheses and the level of significance

The hypotheses would be

$$H_o : \mu_1 = \mu_2$$

$$H_a : \mu_1 \neq \mu_2$$

level of significance, $\alpha = 0.05$

3. State and check the conditions for the hypothesis test

a. State: A random sample of cholesterol levels of people with diabetes is taken. A random sample of cholesterol levels of people without diabetes is taken.

Check: The NHANES data frame uses cluster sampling which incorporates random sampling, so the sample is probably representative. This condition has been met.

b. State: The two samples are independent.

Check: This is because either they were dealing with people who have diabetes or not.

c. State: Population of all cholesterol levels of people who have diabetes is normally distributed. Population of all cholesterol levels of people without diabetes is normally distributed.

Check:

```
NHANES_no_NA<-  
  NHANES |>  
  drop_na(Diabetes)  
gf_density(~TotChol|Diabetes, data=NHANES_no_NA, title = "Cholesterol of a person with and w
```

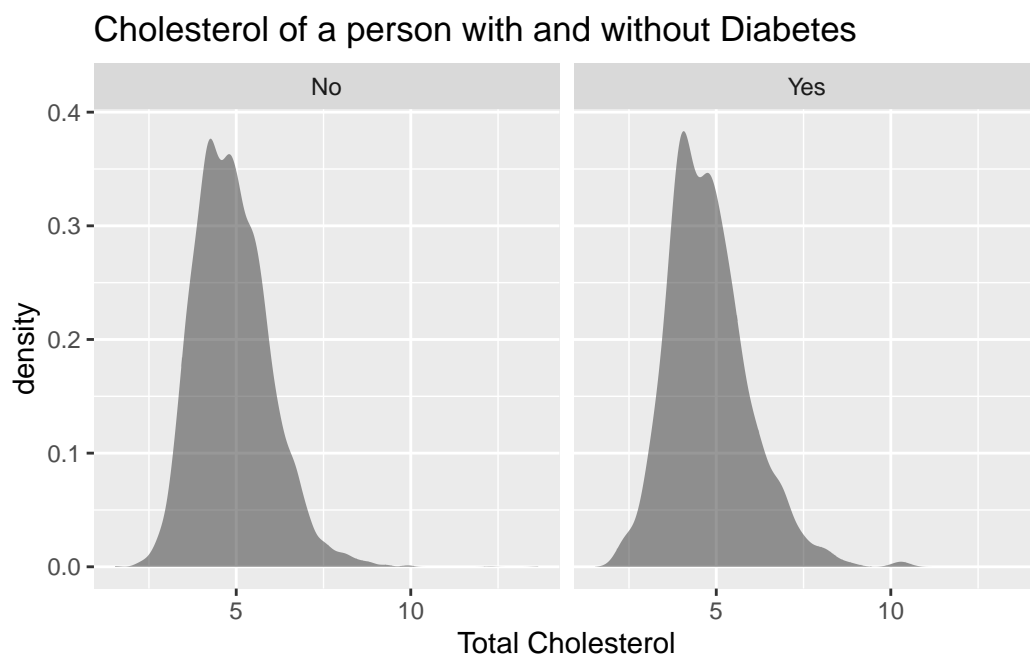



Figure 10.5: Density Plot of Cholesterol of a person with and without Diabetes

Both the yes group and the no group look somewhat bell shaped.

```
gf_qq(~TotChol|Diabetes, data=NHANES_no_NA, title = "Cholesterol of a person with and without Diabetes")
```

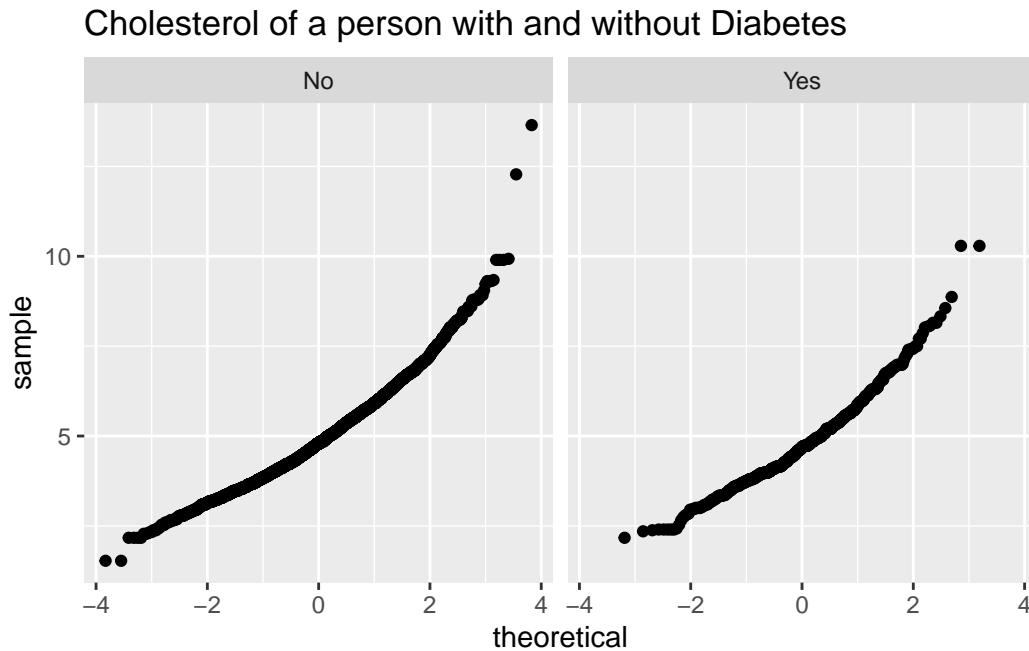


Figure 10.6: quantile Plot of Cholesterol of a person with and without Diabetes

Both the yes group and the no group look somewhat linear.

The population of all cholesterol levels of people who have diabetes is probably normally distributed. The population of all cholesterol levels of people who do not have diabetes is probably normally distributed.

4. Find the sample statistic, test statistic, and p-value

The variable is cholesterol (TotChol) and separating based on if a person has diabetes or not. So the factor is Diabetes. Using r Studio the command would be

```
t.test(TotChol~Diabetes, data=NHANES)
```

Welch Two Sample t-test

```
data: TotChol by Diabetes
```

```
t = 2.4286, df = 809.7, p-value = 0.01537
```

```
alternative hypothesis: true difference in means between group No and group Yes is not equal
```

```
95 percent confidence interval:
```

```
0.02105115 0.19851114
```

```
sample estimates:
```

mean in group No	mean in group Yes
4.887936	4.778155

5. Conclusion

Reject H_o since the p-value $< \alpha$.

6. Interpretation

There is enough evidence to support that people who have diabetes have different cholesterol levels on average from people who do not have diabetes.

10.3.4 Example: Confidence Interval in Two Samples

The cholesterol level of people vary for many reasons. The question is how different is the cholesterol levels of people with diabetes from people who do not have diabetes? Use the NHANES data frame. Compute a 95% confidence interval.

10.3.4.1 Solution

1. State the random variables and the parameters in words.

These were stated in Example: Hypothesis Test for Two Means.

2. State and check the conditions for the hypothesis test

The conditions were stated and checked in Example: Hypothesis Test for Two Means.

3. Find the sample statistic and confidence interval

The variable is cholesterol (TotChol) and separating based on if a person has diabetes or not. So the factor is Diabetes. Using rStudio the command would be

```
t.test(TotChol~Diabetes, data=NHANES, conf.level=0.95)
```

Welch Two Sample t-test

data: TotChol by Diabetes

t = 2.4286, df = 809.7, p-value = 0.01537

alternative hypothesis: true difference in means between group No and group Yes is not equal

95 percent confidence interval:

0.02105115 0.19851114

sample estimates:

mean in group No	mean in group Yes
4.887936	4.778155

4. Statistical Interpretation: You are 95% confident that the interval $0.02105115 < \mu_1 - \mu_2 < 0.19851114$ contains the true difference in means.
5. Real World Interpretation: The mean cholesterol level for people with diabetes is anywhere from 0.021 mmol/L to 0.199 mmol/L more than the mean cholesterol level for people without diabetes.

10.3.5 Example: Hypothesis Test for Two Means

The amount of sodium in beef and poultry hot dogs was measured. (\“SOCR 012708 id,\” 2013). The data is in Table ???. Is there enough evidence to show that beef has different amounts of sodium on average than poultry hot dogs? Use a 5% level of significance.

```
Hotdog<-read.csv( "https://krkozak.github.io/MAT160/hotdog_beef_poultry.csv")
knitr::kable(head(Hotdog))
```

Table 10.9: Hot dog Data

type	calories	sodium
Beef	186	495
Beef	181	477
Beef	176	425
Beef	149	322
Beef	184	482
Beef	190	587

Code book for data frame Hot dog

Description Results of a laboratory analysis of calories and sodium content of major hot dog brands. Researchers for Consumer Reports analyzed three types of hot dog: beef, poultry, and meat (mostly pork and beef, but up to 15% poultry meat). The meat was left off this data frame so a two-sample t-test could be performed.

This data frame contains the following columns:

type: Type of hot dog (beef or poultry)

calories: Calories per hot dog

sodium: Milligrams of sodium per hot dog

Source SOCR 012708 id data hotdogs. (2013, November 13). Retrieved from <http://wiki.stat.ucla.edu/socr/index>

References SOCR Home page: <http://www.socr.ucla.edu>

10.3.5.1 Solution

1. State the random variables and the parameters in words.

x_1 = sodium level in beef hot dogs

x_2 = sodium level in poultry hot dogs

μ_1 = mean sodium level in beef hot dogs

μ_2 = mean sodium level in poultry hot dogs

2. State the null and alternative hypotheses and the level of significance

The hypotheses would be

$$H_o : \mu_1 = \mu_2$$

$$H_o : \mu_1 \neq \mu_2$$

level of significance: $\alpha = 0.05$

3. State and check the conditions for the hypothesis test

- a. State: A random sample of 20 sodium levels in beef hot dogs is taken. A random sample of 20 sodium levels in poultry hot dogs.

Check: The code does not state if either sample was randomly selected, but since Consumer Reports performed the test, it is safe to assume the samples were both random.

- b. State: The two samples are independent.

Check: These are different types of hot dogs so this is true.

- c. State; Population of all sodium levels in beef hot dogs is normally distributed. Population of all sodium levels in poultry hot dogs is normally distributed.

Check:

```
gf_density(~sodium|type, data=Hotdog, title="Sodium amount in Hot Dogs facettted by Type of Meat")
```

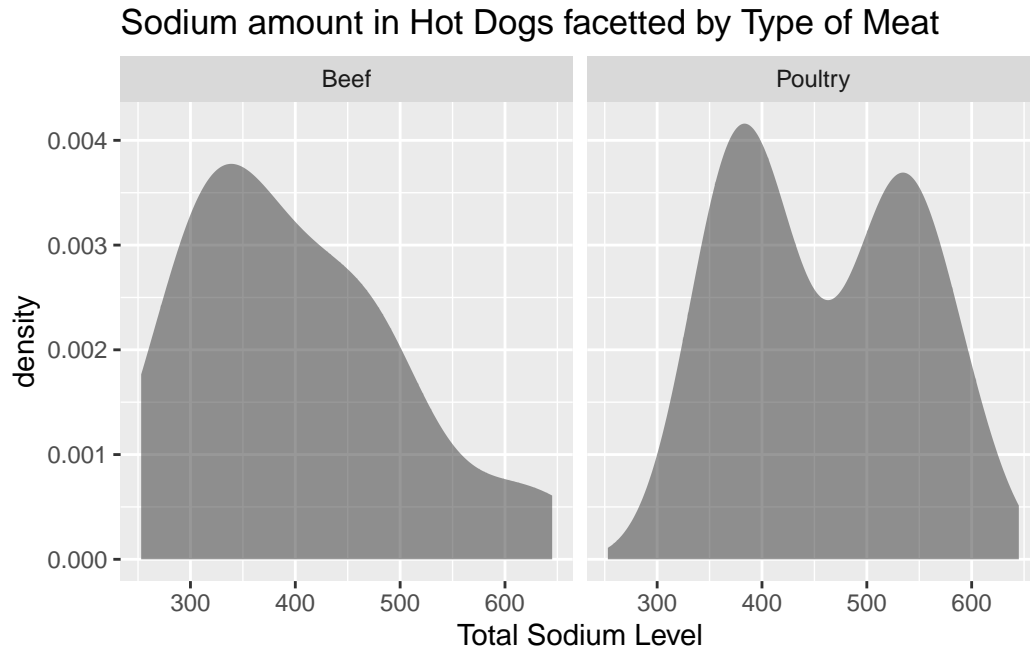


Figure 10.7: Density Plot of Sodium Amount in Hot Dogs facettted by Type of Meat

The density plot for beef hot dogs looks somewhat bell shaped, but the density plot for poultry hot dogs does not look bell shaped.

```
gf_qq(~sodium|type, data=Hotdog, title="Sodium amoount in Hot Dogs facettted by Type of Meat")
```

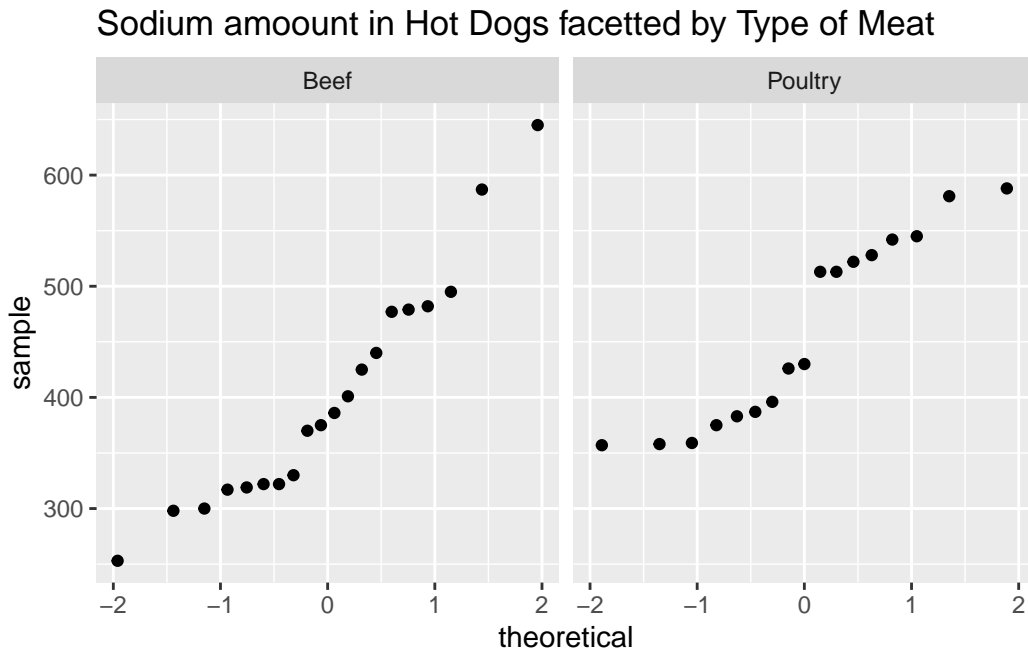


Figure 10.8: Quantile Plot of Sodium Amount in Hot Dogs facetted by Type of Meat

The normal quantile plot Figure ?? for the sodium level in beef hot dogs looks somewhat linear. The normal quantile plot Figure ?? for the sodium level in poultry hot dogs does not look linear. The population of all sodium levels in beef hot dogs may be normally distributed, but the population of all sodium levels in poultry hot dogs is probably not normally distributed. The sample size is not very large so the results of the test may not be valid. A larger sample would be a good idea.

4. Find the sample statistic, test statistic, and p-value

Using rStudio the variable is sodium levels (sodium) in different types of hot dogs. So the factor is type. The command is

```
t.test(sodium~type, data=Hotdog)
```

Welch Two Sample t-test

data: sodium by type

t = -1.8798, df = 34.983, p-value = 0.06848

alternative hypothesis: true difference in means between group Beef and group Poultry is not

95 percent confidence interval:

-120.325706 4.625706

sample estimates:

mean in group Beef	mean in group Poultry
401.15	459.00

5. Conclusion: Fail to reject H_o since the p-value $\geq \alpha$.

6. Interpretation

This is not enough evidence to support that beef hot dogs' sodium level is different from poultry hot dogs. (Though do realize that the population conditions is not valid, so this interpretation may be invalid.)

10.3.6 Example: Confidence Interval for Two Independent Samples

The amount of sodium in beef and poultry hot dogs was measured. ("SOCR 012708 id," 2013). The data is in Table ???. Find a 95% confidence interval for the mean difference in sodium levels between beef and poultry hot dogs.

10.3.6.1 Solution

1. State the random variables and the parameters in words.

These were stated in Example: Hypothesis Test for Two Means.

2. State and check the conditions for the hypothesis test

The conditions were stated and checked in Example: Hypothesis Test for Two Means.

3. Find the sample statistic and confidence interval Using r Studio the variable is sodium levels (sodium) in different types of hot dogs. So the factor is type. The command is

```
t.test(sodium~type, data=Hotdog, conf.level=0.95)
```

Welch Two Sample t-test

data: sodium by type

t = -1.8798, df = 34.983, p-value = 0.06848

alternative hypothesis: true difference in means between group Beef and group Poultry is not
95 percent confidence interval:

-120.325706 4.625706

sample estimates:

mean in group Beef	mean in group Poultry
401.15	459.00

4. Statistical Interpretation: You are 95% confident that the interval $-120.325706 < \mu_1 - \mu_2 < 4.625706$ contains the true difference in mean sodium level between beef and poultry hot dogs.
5. Real World Interpretation: The mean sodium level of beef hot dogs is anywhere from 120.33 mg less than the mean sodium level of poultry hot dogs to 4.63 mg more. (The negative sign on the lower limit implies that the first mean is less than the second mean. The positive sign on the upper limit implies that the first mean is greater than the second mean.)

Do realize that the population conditions is not valid, so this interpretation may be invalid.

10.3.7 Homework for Independent Samples for Two Means Section

In each problem show all steps of the hypothesis test or confidence interval. If some of the conditions are not met, note that the results of the test or interval may not be correct and then continue the process of the hypothesis test or confidence interval.

1. The NHANES data contains many variables. One variable is the income of households derived from the middle income of different income categories. The variable is called HHIIncomeMid. Is there enough evidence to show that the mean income of males is different from the mean income of females? Test at the 1% level.

`names(NHANES)`

[1] "ID"	"SurveyYr"	"Gender"	"Age"
[5] "AgeDecade"	"AgeMonths"	"Race1"	"Race3"
[9] "Education"	"MaritalStatus"	"HHIncome"	"HHIncomeMid"
[13] "Poverty"	"HomeRooms"	"HomeOwn"	"Work"
[17] "Weight"	"Length"	"HeadCirc"	"Height"
[21] "BMI"	"BMICatUnder20yrs"	"BMI_WHO"	"Pulse"
[25] "BPSysAve"	"BPDiaAve"	"BPSys1"	"BPDia1"
[29] "BPSys2"	"BPDia2"	"BPSys3"	"BPDia3"
[33] "Testosterone"	"DirectChol"	"TotChol"	"UrineVol1"
[37] "UrineFlow1"	"UrineVol2"	"UrineFlow2"	"Diabetes"
[41] "DiabetesAge"	"HealthGen"	"DaysPhysHlthBad"	"DaysMentHlthBad"
[45] "LittleInterest"	"Depressed"	"nPregnancies"	"nBabies"
[49] "Age1stBaby"	"SleepHrsNight"	"SleepTrouble"	"PhysActive"
[53] "PhysActiveDays"	"TVHrsDay"	"CompHrsDay"	"TVHrsDayChild"
[57] "CompHrsDayChild"	"Alcohol12PlusYr"	"AlcoholDay"	"AlcoholYear"
[61] "SmokeNow"	"Smoke100"	"Smoke100n"	"SmokeAge"

[65]	"Marijuana"	"AgeFirstMarij"	"RegularMarij"	"AgeRegMarij"
[69]	"HardDrugs"	"SexEver"	"SexAge"	"SexNumPartnLife"
[73]	"SexNumPartYear"	"SameSex"	"SexOrientation"	"PregnantNow"

2. The NHANES data contains many variables. One variable is the income of households derived from the middle income of different income categories. The variable is called HHIIncomeMid. Estimate with 95% confidence the mean difference in incomes between males and females in the U.S.
3. A study was conducted that measured the total brain volume (TBV) of patients that had schizophrenia and patients that do not have schizophrenia. Table ?? contains the TBV of the all patients ("SOCR data oct2009,\" 2013). Is there enough evidence to show that the patients with schizophrenia have a different TBV on average than a patient without schizophrenia? Test at the 10% level.

```
Brain <- read.csv( "https://krkozak.github.io/MAT160/brain.csv")
knitr::kable(head(Brain))
```

Table 10.10: Total Brain Volume of Patients

type	volume
n	1663407
n	1583940
n	1299470
n	1535137
n	1431890
n	1578698

Code book for data frame Brain

Description A study to measure the total brain volume (TBV) (in) of patients that had schizophrenia and patients that do not have schizophrenia.

This data frame contains the following columns:

type: whether the patient had schizophrenia (s) or did not have schizophrenia (n)

volume: the total brain volume of a patient.(mm^3)

Source SOCR data Oct2009 id ni. (2013, November 16). Retrieved from <http://wiki.stat.ucla.edu/socr/index.php>

References "SOCR data nips," 2013

4. A study was conducted that measured the total brain volume (TBV) of patients that had schizophrenia and patients that do not have schizophrenia. Table ?? contains the TBV of the all patients (“SOCR data oct2009,” 2013). Is there enough evidence to show that the patients with schizophrenia have a different TBV on average than a patient without schizophrenia? Test at the 10% level. Compute a 90% confidence interval for the difference in TBV of patients with Schizophrenia and patients without Schizophrenia.
5. The lengths (in kilometers) of rivers on the South Island of New Zealand and what body of water they flow into are listed in Table ?? (Lee, 1994). Do the data provide enough evidence to show on average that the rivers that travel to the Pacific Ocean are different length than the rivers that travel to the Tasman Sea? Use a 5% level of significance.

Code book for data frame Length below Table ??.

6. The lengths (in kilometers) of rivers on the South Island of New Zealand and what body of water they flow into are listed in Table ?? (Lee, 1994). Estimate the difference in mean lengths of rivers between rivers in New Zealand that travel to the Pacific Ocean and ones that travel to the Tasman Sea. Use a 95% confidence level.
 7. A vitamin K shot is given to infants soon after birth. Nurses at Northbay Healthcare were involved in a study to see if how they handle the infants could reduce the pain the infants feel (\“SOCR data nips,\” 2013). The data frame is in Table ??.
- Is there enough evidence to show that infants cried a different amount on average when they are held by their mothers than if held using conventional methods? Test at the 5% level.

10.3.7.1 Table: Crying Time of Infants Given Shots Using New Methods

```
Crying<- read.csv( "https://krkozak.github.io/MAT160/crying.csv")
knitr::kable(head(Crying))
```

Table 10.11: Crying Time of Infants Given Shots Using New Methods

method	crying
convent	63
convent	0
convent	2
convent	46
convent	33
convent	33

Code book for data frame Crying

Description Nurses at Northbay Healthcare were involved in a study to see if how they handle the infants could reduce the pain the infants feel. One of the measurements taken was how long, in seconds, the infant cried after being given the shot. A random sample was taken from the group that was given the shot using conventional methods, and a random sample was taken from the group that was given the shot where the mother held the infant prior to and during the shot.

This data frame contains the following columns:

method: whether the infant was given the conventional method (convent) or the new method (new) prior to being given the vitamin K shot.

crying: how long the infant cried after given a vitamin K shot. (seconds)

Source SOCR data nips infantvitK shotdata. (2013, November 16). Retrieved from http://wiki.stat.ucla.edu/socr/index.php/SOCR_Data_NIPS_InfantVitK_ShotData

References \“SOCR data nips,\” 2013

8. A vitamin K shot is given to infants soon after birth. Nurses at Northbay Healthcare were involved in a study to see if how they handle the infants could reduce the pain the infants feel (\“SOCR data nips,\” 2013). The data frame is in Table ???. Calculate a 95% confidence interval for the mean difference in mean crying time after being given a vitamin K shot between infants held using conventional methods and infants held by their mothers.

10.4 Which Analysis Should You Conduct?

One of the most important concept that you need to understand is deciding which analysis you should conduct for a particular situation. To help you to figure out the analysis to conduct, there are a series of questions you should ask yourself.

1. Does the problem deal with mean or proportion?

Sometimes the problem states explicitly the words mean or proportion, but other times you have to figure it out based on the information you are given. If you counted number of individuals that responded in the affirmative to a question, then you are dealing with proportion. If you measured something, then you are dealing with mean.

2. Does the problem have one or two samples?

So look to see if one group was measured or if two groups were measured. You need to decide if the problem describes collecting data from one group or from two groups, or if you are comparing two different groups.

3. If you have two samples, then you need to determine if the samples are independent or dependent.

If the individuals are different for both samples, then most likely the samples are independent. If you can't tell, then determine if a data value from the first sample influences the data value in the second sample. In other words, can you pair data values together so you can find the difference, and that difference has meaning. If the answer is yes, then the samples are paired. Otherwise, the samples are independent.

4. Does the situation involve a hypothesis test or a confidence interval?

If the problem talks about “do the data show”, “is there evidence of”, “test to see”, then you are doing a hypothesis test. If the problem talks about “find the value”, “estimate the” or “find the interval”, then you are doing a confidence interval.

So if you have a situation that has two samples, independent samples, involving the mean, and is a hypothesis test, then you have a two-sample independent t-test. Now you look up the conditions and the technology process for doing this test. Every hypothesis test involves the same six steps, and you just have to use the correct conditions and calculations. Every confidence interval has the same five steps, and again you just need to use the correct conditions and calculations. So this is why it is so important to figure out what analysis you should conduct.

11 Regression

The previous chapter looked at comparing populations to see if there is a difference between the two. That involved two random variables that are similar measures. This chapter will look at two random variables that do not need to be similar measures, and see if there is a relationship between the two or more variables. To do this, you look at regression, which finds the linear relationship, and correlation, which measures the strength of a linear relationship.

Please note: there are many other types of relationships besides linear that can be found for the data. This book will only explore linear, but realize that there are other relationships that can be used to describe data.

11.1 Regression

When comparing different variables, two questions come to mind: “Is there a relationship between two variables?” and “How strong is that relationship?” These questions can be answered using **regression** and **correlation**. Regression answers whether there is a relationship (again this book will explore linear only) and correlation answers how strong the linear relationship is. The variable that are used to explain the change in the other variable is called the **explanatory variable** while the variable that is changing is called the **response variable**. Other variables that help to explain the changes are known as **covariates**. To introduce the concepts of regression and correlation it is easier to look at a set of data.

11.1.1 Example: Determining if there is a Relationship

Is there a relationship between the alcohol content and the number of calories in 12-ounce beer? To determine if there is one, we explore a sample of 227 beers’ alcohol content and their calories (Find Out How Many Calories in Beer?, 2019). Table ?? shows the first five rows of the dataset.

Table 11.1: Alcohol and Calorie Content in Beer

beer	brewery	location	alcohol	calories	carbs
American Amber Lager	Straub Brewery	domestic	0.04	136	10.5
American Lager	Straub Brewery	domestic	0.04	132	10.5