

ENGE 5714 Course Notes 2021

Dr. Katz

2021-03-14

Contents

Preamble	5
1 Week 1: Introductions	7
2 Week 2: Intro stats, Data & Distributions, Intro R & RStudio	9
2.1 First steps in R	9
2.2 Getting your R environment set up	10
2.3 Reading in data	11
2.4 Exploring the data	13
2.5 Plotting data	16
2.6 Some brief stats	20
3 Week 3: Data Cleaning, Organizing, Describing, and Communicating	25
3.1 Visualizing your data	25
3.2 Joining two datasets	33
3.3 Discrete Predictor, Continuous Outcome	38
3.4 Continuous predictor and continuous outcome	40
3.5 Mutating Variables	46
3.6 Filtering and Selecting	49
3.7 Grouping and Summarizing	53

4	Week 4: Assumptions and Correlations	57
4.1	Assumptions	57
4.2	Correlation	58
4.3	Another worked example for cleaning and prelim analysis	59
5	Week 5: Simple Regression	69
5.1	General Modeling Philosophy	69
5.2	Data generation demo - one set sample size	72
5.3	Data generation demo - one set sample size;	82
5.4	Data generation with three different sample sizes	91
6	Week 6: Regression II	103
6.1	Explore the child aggression data set	103
6.2	Multiple regression	109
7	Week 7: Logistic Regression	135
7.1	Round 1 - No systematic variation in outcomes	135
7.2	Round 2 - Systematic variation in outcomes as a function of discipline	139
7.3	Round 3 - Systematic variation in outcomes as a function of discipline and gpa	143
7.4	Round 4 - GPA and persistence vary by discipline	148
8	Week 8: Comparing two means (t-tests)	157
8.1	Demo 1 - Comparing salary data for chemical engineering and environmental engineering	157
8.2	Demo 2 - Student SAT scores	162

Preamble

This book will be a living document of notes for ENGE 5714 - Quantitative Research Methods in Engineering Education. I will try to keep it updated and post alerts about updates in our course Slack workspace.

Chapter 1

Week 1: Introductions

In week one we just reviewed some of the materials from the fall semester. By the end we discussed R and RStudio, but this first week was primarily about getting to know each other and the outline of the course.

Chapter 2

Week 2: Intro stats, Data & Distributions, Intro R & RStudio

This week, we discuss some very basic ideas related to statistics, data, and working in R.

2.1 First steps in R

We can create a new variable by assigning it a value with the `<-` operator. Let's create a vector of numbers 1 to 10 with the `seq()` function and then a separate vector that takes each of the `x` values, multiplies it by 2, and adds 3.

```
x <- seq(1:10)
y <- 2* x + 3
```

Just to make sure everything worked as expected, we can then just type `x` and `y` and R will print their values. We could also look in the “environment” window to see whether those variables (and their expected values) were actually created.

```
x
```

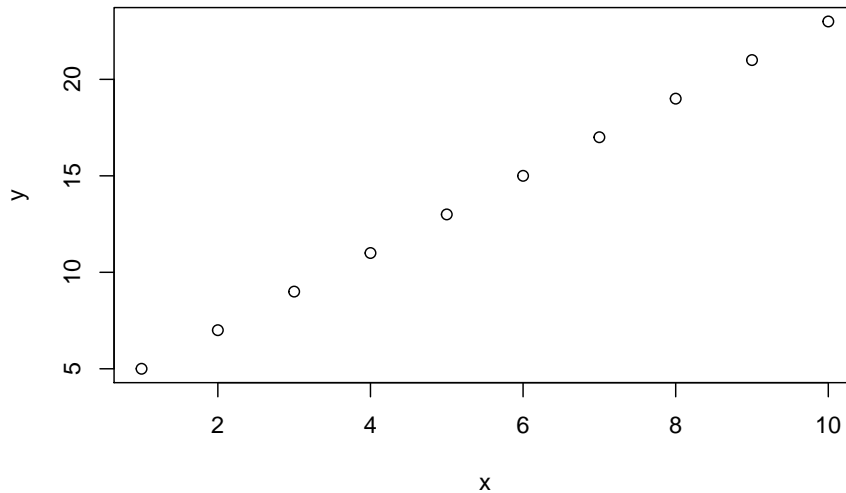
```
## [1] 1 2 3 4 5 6 7 8 9 10
```

```
y
```

```
## [1] 5 7 9 11 13 15 17 19 21 23
```

So far, so good. If we want to quickly visualize this, we could create a simple scatter plot with the `plot()` command (note: we will come back to plotting data much more in week 3).

```
plot(x, y)
```



2.2 Getting your R environment set up

One of the first things you will have in any script or .rmd file is a section to load all the libraries that you use in that script.

You can install a library by using the `install.packages()` function, for example:

```
install.packages("tidyverse"), install.packages("janitor"), and  
install.packages("psych")
```

with this installed, you can then load the package using the `library()` function

```
library(tidyverse)
```

```
## -- Attaching packages --- tidyverse 1.3.0 --
```

```
## v ggplot2 3.3.2      v purrr  0.3.4
## v tibble  3.0.3      v dplyr  1.0.0
## v tidyr   1.1.0      v stringr 1.4.0
## v readr   1.3.1      v forcats 0.5.0
```

```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```
library(janitor)
```

```
## Warning: package 'janitor' was built under R version 4.0.3
```

```
##
## Attaching package: 'janitor'
```

```
## The following objects are masked from 'package:stats':
##
##      chisq.test, fisher.test
```

```
library(psych)
```

```
##
## Attaching package: 'psych'
```

```
## The following objects are masked from 'package:ggplot2':
##
##      %+%, alpha
```

2.3 Reading in data

A good first step when working in R is to check which directory you are working in with the `getwd()` function. You should get a directory in response.

```
getwd()
```

```
## [1] "C:/Users/akatz4/Desktop/test_course_note"
```

You can also check which files are in that directory with `list.files()`.

```
list.files()
```

```
## [1] "_book"
## [2] "_bookdown.yml"
## [3] "_bookdown_files"
## [4] "_output.yml"
## [5] "01-Week_01.Rmd"
## [6] "02-Week_02.Rmd"
## [7] "03-Week_03.Rmd"
## [8] "04-Week_04.Rmd"
## [9] "05-Week_05.Rmd"
## [10] "06-Week_06.Rmd"
## [11] "07-Week_07.Rmd"
## [12] "08-Week_08.Rmd"
## [13] "book.bib"
## [14] "ChildAggression.dat"
## [15] "docs"
## [16] "ENGE_5714_2021_pre_survey.csv"
## [17] "Free Reduced Lunch by Schools and Grade Structures 2008-2017_final.csv"
## [18] "index.Rmd"
## [19] "packages.bib"
## [20] "preamble.tex"
## [21] "README.md"
## [22] "RExam.dat"
## [23] "seniorsurvey.csv"
## [24] "student_happiness.csv"
## [25] "style.css"
## [26] "survey_student_info.csv"
## [27] "test_course_note.Rproj"
## [28] "test_course_notes.Rmd"
## [29] "test_course_notes_files"
```

If you notice that the file you are looking for is not there, then you can use `setwd()` to change your working directory

```
setwd("./Week 2/")
```

After that, make sure you have switched to the correct working directory `getwd()` and then `list.files()`.

Assuming you have directed yourself to the correct place, you can now read in the file(s) that you want to be working with. There are a *lot* of ways to do this. Since we will be spending a lot of time in class working with .csv files, we will focus on using the `read_csv()` function from the `readr` package (part of the tidyverse collection of packages). This function will read in the .csv file and store the data as a tibble (a tidyverse version of a data frame, which we can think of as a collection of observations stored in rows with values for variables for each observation stored in columns).

```
prior_survey <- read_csv("ENGE_5714_2021_pre_survey.csv")
```

```
## Parsed with column specification:
## cols(
##   .default = col_character(),
##   student_id = col_double()
## )
```

```
## See spec(...) for full column specifications.
```

2.4 Exploring the data

Now that we have loaded in the data, let's take a look at the csv. If we just run a line with the name of the tibble - i.e., `prior_survey` then we should receive a printout that shows the first several rows of that tibble and a listing of all the columns, along with the data types (i.e., double for numeric values, character for strings, etc) of each column.

```
prior_survey
```

```
## # A tibble: 24 x 49
##   student_id `I have taken a~`I am intereste~`I know what a ~
##         <dbl> <chr>          <chr>          <chr>
## 1           1 Somewhat disagr~ Somewhat agree Strongly disagr~
## 2           2 Strongly disagr~ Neither agree n~ Somewhat agree
## 3           3 Strongly disagr~ Somewhat agree Somewhat agree
## 4           4 Somewhat disagr~ Strongly agree Strongly disagr~
## 5           5 Somewhat agree Strongly agree Somewhat agree
## 6           6 Somewhat disagr~ Somewhat agree Somewhat disagr~
## 7           7 Strongly disagr~ Somewhat agree Strongly disagr~
## 8           8 Somewhat agree Somewhat agree Somewhat agree
## 9           9 Strongly disagr~ Strongly agree Somewhat agree
## 10          10 Neither agree n~ Strongly agree Somewhat agree
## # ... with 14 more rows, and 45 more variables: `I know what a type II error`
```

```
## # is` <chr>, `I know what a (statistical) confidence level is` <chr>, `I know
## # what a p value is` <chr>, `I know what p-hacking means` <chr>, `I know what
## # statistical power means` <chr>, `I have heard of frequentist statistics
## # before` <chr>, `I have heard of Bayesian statistics before` <chr>, `I have
## # heard the term "parametric statistics" before` <chr>, `I have heard the
## # term "non-parametric statistics" before` <chr>, `I know what a histogram
## # is.` <chr>, `I know what a probability distribution is.` <chr>, `I know
## # what a random variable is.` <chr>, `I know what a probability distribution
## # function is.` <chr>, `I know what a cumulative distribution function
## # is.` <chr>, `I know what the expectation of a random variable is.` <chr>,
## # `I know how to calculate the variance of a random variable.` <chr>, `I know
## # what a z score is.` <chr>, `I know how to calculate the correlation between
## # two variables.` <chr>, `I know how to interpret the correlation coefficient
## # between two variables` <chr>, `I have heard of linear regression` <chr>, `I
## # know how to run a linear regression (in some software...or by hand, if I'm
## # feeling wild).` <chr>, `I know how to interpret a linear
## # regression.` <chr>, `I have heard of multiple regression` <chr>, `I know
## # how to perform a multiple regression` <chr>, `I know how to interpret a
## # multiple regression` <chr>, `I have heard of logistic regression.` <chr>,
## # `I understand when to use a logistic regression.` <chr>, `I know how to
## # interpret the results of a logistic regression` <chr>, `I have heard of
## # t-tests` <chr>, `I have performed a t-test before` <chr>, `I know how to
## # interpret the results of a t-test` <chr>, `I have heard of Analysis of
## # Variance.` <chr>, `I understand when to run an Analysis of Variance
## # (ANOVA)` <chr>, `I know how to interpret the results from an ANOVA` <chr>,
## # `I have heard of a chi-square test` <chr>, `I have used a chi-square test
## # before` <chr>, `I know how to interpret the results of a chi-square
## # test` <chr>, `I have heard of cluster analysis before` <chr>, `I have used
## # cluster analysis before` <chr>, `I know how to interpret the results of a
## # cluster analysis` <chr>, `I have heard of factor analysis (either
## # exploratory or confirmatory)` <chr>, `I have used factor analysis (either
## # exploratory or confirmatory)` <chr>, `I know how to interpret the results
## # of a factor analysis (either exploratory or confirmatory)` <chr>, `I
## # already have R and Rstudio downloaded to my computer.` <chr>, `I have used
## # R before` <chr>
```

When we do this, we see that there are a bunch of columns that have spaces in their names. This is okay (in the sense that R can handle this), but it can be a little frustrating to work with. Let's try cleaning the column names with `clean_names()` from the `janitor` package. This function will replace the spaces in the column names with underscores and make everything lower case. So, a column name like "I have take a statistics course before" will be changed to "i_have_taken_a_statistics_course_before".

```
prior_survey <- prior_survey %>% clean_names() # from janitor package
```

Look at the data in `prior_survey` again and see if anything looks different (hint: it should).

```
prior_survey
```

```
## # A tibble: 24 x 49
##   student_id i_have_taken_a_ i_am_interested i_know_what_a_t~
##         <dbl> <chr>          <chr>          <chr>
## 1           1 Somewhat disagr~ Somewhat agree  Strongly disagr~
## 2           2 Strongly disagr~ Neither agree n~ Somewhat agree
## 3           3 Strongly disagr~ Somewhat agree  Somewhat agree
## 4           4 Somewhat disagr~ Strongly agree  Strongly disagr~
## 5           5 Somewhat agree  Strongly agree  Somewhat agree
## 6           6 Somewhat disagr~ Somewhat agree  Somewhat disagr~
## 7           7 Strongly disagr~ Somewhat agree  Strongly disagr~
## 8           8 Somewhat agree  Somewhat agree  Somewhat agree
## 9           9 Strongly disagr~ Strongly agree  Somewhat agree
## 10          10 Neither agree n~ Strongly agree  Somewhat agree
## # ... with 14 more rows, and 45 more variables:
## #   i_know_what_a_type_ii_error_is <chr>,
## #   i_know_what_a_statistical_confidence_level_is <chr>,
## #   i_know_what_a_p_value_is <chr>, i_know_what_p_hacking_means <chr>,
## #   i_know_what_statistical_power_means <chr>,
## #   i_have_heard_of_frequentist_statistics_before <chr>,
## #   i_have_heard_of_bayesian_statistics_before <chr>,
## #   i_have_heard_the_term_parametric_statistics_before <chr>,
## #   i_have_heard_the_term_non_parametric_statistics_before <chr>,
## #   i_know_what_a_histogram_is <chr>,
## #   i_know_what_a_probability_distribution_is <chr>,
## #   i_know_what_a_random_variable_is <chr>,
## #   i_know_what_a_probability_distribution_function_is <chr>,
## #   i_know_what_a_cumulative_distribution_function_is <chr>,
## #   i_know_what_the_expectation_of_a_random_variable_is <chr>,
## #   i_know_how_to_calculate_the_variance_of_a_random_variable <chr>,
## #   i_know_what_a_z_score_is <chr>,
## #   i_know_how_to_calculate_the_correlation_between_two_variables <chr>,
## #   i_know_how_to_interpret_the_correlation_coefficient_between_two_variables <chr>,
## #   i_have_heard_of_linear_regression <chr>,
## #   i_know_how_to_run_a_linear_regression_in_some_software_or_by_hand_if_im_feeling_wild <chr>,
## #   i_know_how_to_interpret_a_linear_regression <chr>,
## #   i_have_heard_of_multiple_regression <chr>,
## #   i_know_how_to_perform_a_multiple_regression <chr>,
```

```
## # i_know_how_to_interpret_a_multiple_regression <chr>,
## # i_have_heard_of_logistic_regression <chr>,
## # i_understand_when_to_use_a_logistic_regression <chr>,
## # i_know_how_to_interpret_the_results_of_a_logistic_regression <chr>,
## # i_have_heard_of_t_tests <chr>, i_have_performed_a_t_test_before <chr>,
## # i_know_how_to_interpret_the_results_of_a_t_test <chr>,
## # i_have_heard_of_analysis_of_variance <chr>,
## # i_understand_when_to_run_an_analysis_of_variance_anova <chr>,
## # i_know_how_to_interpret_the_results_from_an_anova <chr>,
## # i_have_heard_of_a_chi_square_test <chr>,
## # i_have_used_a_chi_square_test_before <chr>,
## # i_know_how_to_interpret_the_results_of_a_chi_square_test <chr>,
## # i_have_heard_of_cluster_analysis_before <chr>,
## # i_have_used_cluster_analysis_before <chr>,
## # i_know_how_to_interpret_the_results_of_a_cluster_analysis <chr>,
## # i_have_heard_of_factor_analysis_either_exploratory_or_confirmatory <chr>,
## # i_have_used_factor_analysis_either_exploratory_or_confirmatory <chr>,
## # i_know_how_to_interpret_the_results_of_a_factor_analysis_either_exploratory_or_c
## # i_already_have_r_and_rstudio_downloaded_to_my_computer <chr>,
## # i_have_used_r_before <chr>
```

One other function that we will see more in the future is the `table()` function, which will create a table with the counts of the values for a variable. For example, if we wanted to quickly know how students answered the “I have taken a quantitative research methods course before” question, we can run the following:

```
table(prior_survey$i_have_taken_a_quantitative_research_methods_course_before)
```

```
##
## Neither agree nor disagree          Somewhat agree
##                2                    5
##          Somewhat disagree          Strongly agree
##                5                    2
##          Strongly disagree
##                10
```

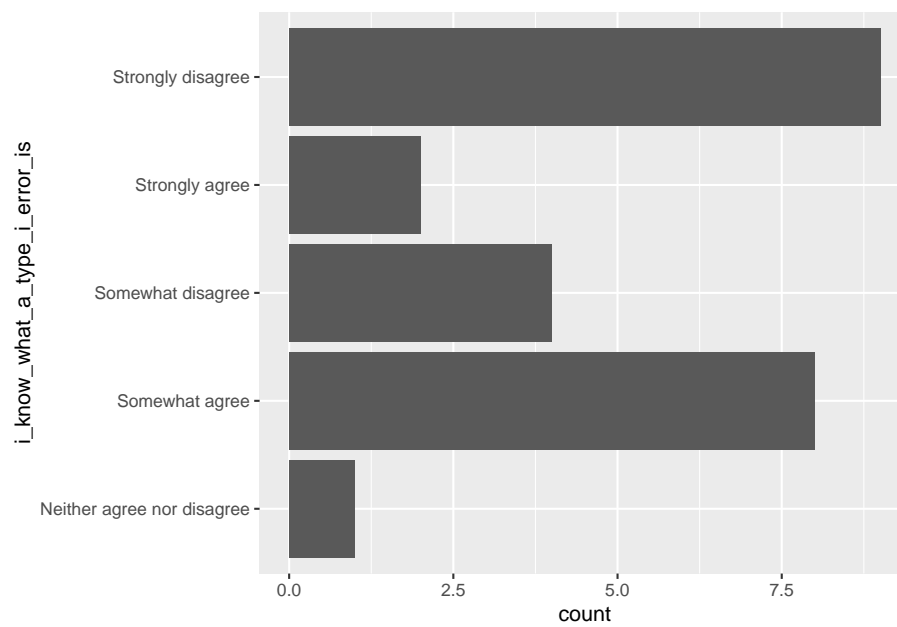
2.5 Plotting data

We will discuss plotting more next week, but here is a brief preview of what’s to come...

There are multiple ways to plot data. Focusing on using `ggplot`, here are two.

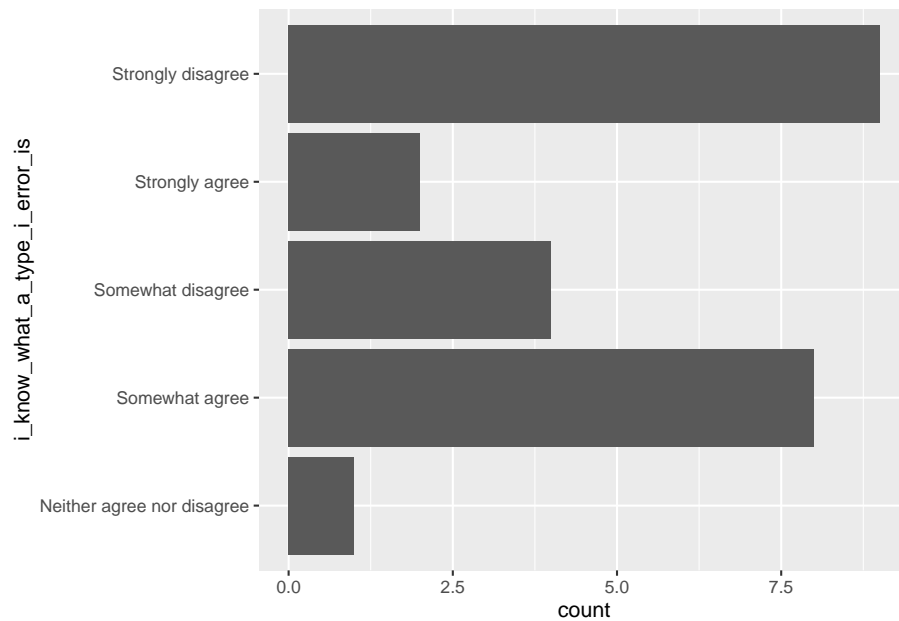
The first way passes the `prior_survey` dataframe explicitly to `ggplot`


```
ggplot(data = prior_survey, mapping = aes(x = i_know_what_a_type_i_error_is)) +  
  geom_bar() +  
  coord_flip()
```



The second way does this implicitly, using the pipe operator. Note that the results should be the same.

```
prior_survey %>%  
  ggplot(mapping = aes(x = i_know_what_a_type_i_error_is)) +  
  geom_bar() +  
  coord_flip()
```



If we wanted to get extra fancy, we could first convert the data from a wide format to a long format and then start plotting all the items together.

Converting to long format would produce something like this:

```
prior_survey %>%
  gather(key = "survey_item", value = "survey_response")
```

```
## # A tibble: 1,176 x 2
##   survey_item survey_response
##   <chr>      <chr>
## 1 student_id 1
## 2 student_id 2
## 3 student_id 3
## 4 student_id 4
## 5 student_id 5
## 6 student_id 6
## 7 student_id 7
## 8 student_id 8
## 9 student_id 9
## 10 student_id 10
## # ... with 1,166 more rows
```

Then we can combine that with the `group_by()` and `summarize()` functions and plot the results.

```
prior_survey %>%
  gather(key = "survey_item", value = "survey_response") %>%
  group_by(survey_item, survey_response) %>%
  summarize(n = n()) %>%
  ggplot(mapping = aes(x = survey_response, y = survey_item, fill = n)) +
  geom_tile()

## `summarise()` regrouping output by 'survey_item' (override with `.groups` argument)
```



This plot is okay for giving a general sense of what is going on in these plots but there are a bunch of other ways to go about doing this.

First, maybe we want to rename the response categories to a numerical scale. We can accomplish this with a `mutate()` and `case_when()`.

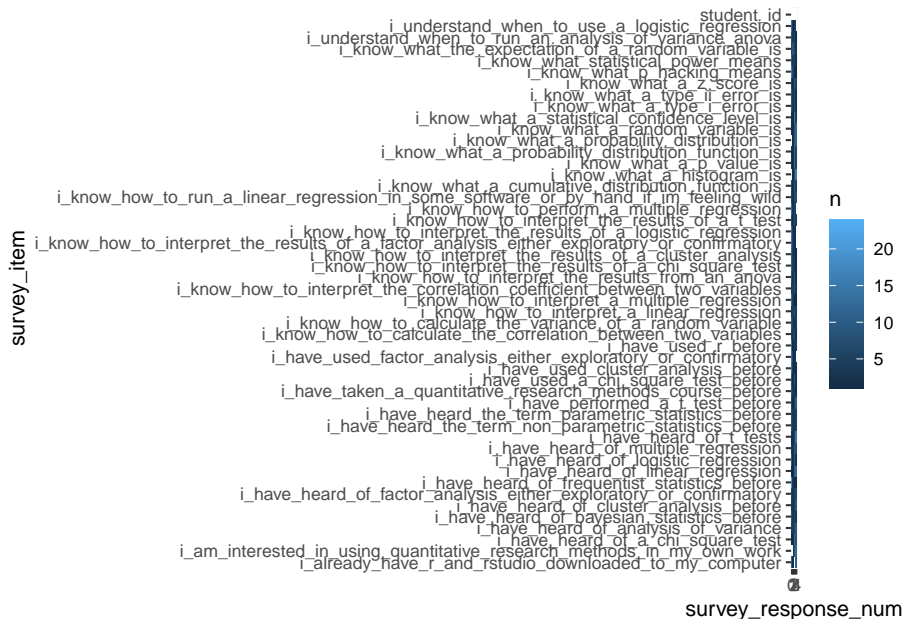
```
prior_survey <- prior_survey %>%
  gather(key = "survey_item", value = "survey_response") %>%
  mutate(survey_response_num = case_when(survey_response == "Strongly disagree" ~ 0,
                                         survey_response == "Somewhat disagree" ~ 1,
                                         survey_response == "Neither agree nor disagree" ~ 2,
                                         survey_response == "Somewhat agree" ~ 3,
                                         survey_response == "Strongly agree" ~ 4,
                                         ))
```

Then we plot the same data but with the numerical scale along the x-axis.

```
prior_survey %>%
  group_by(survey_item, survey_response_num) %>%
  summarize(n = n()) %>%
  ggplot(mapping = aes(x = survey_response_num, y = survey_item, fill = n)) +
  geom_tile()
```

```
## `summarise()` regrouping output by 'survey_item' (override with `.groups` argument)
```

```
## Warning: Removed 3 rows containing missing values (geom_tile).
```

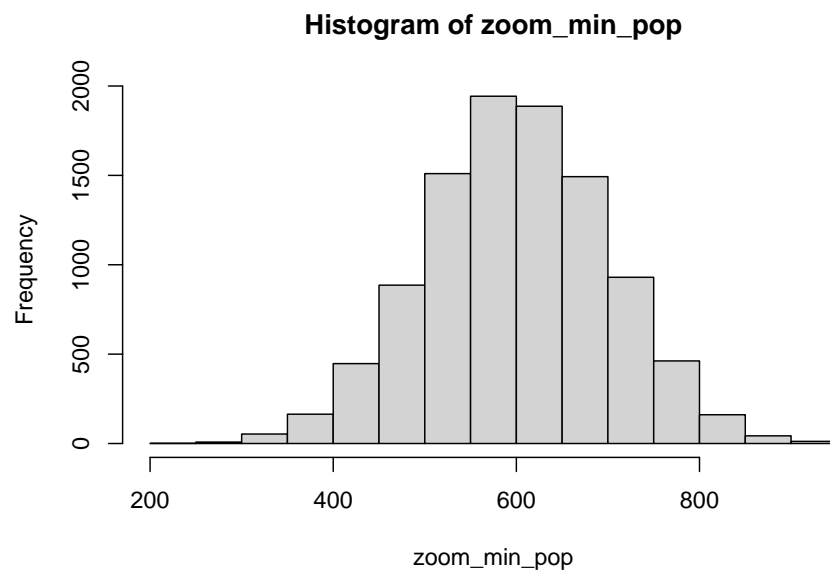


2.6 Some brief stats

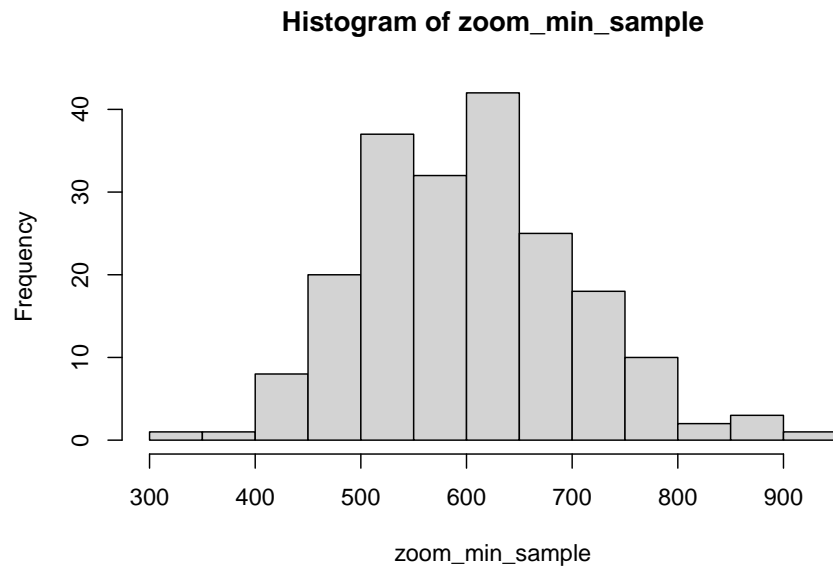
In this week's reading, there was also discussion about standard errors and the central limit theorem. These are fairly important theoretical concepts to grasp. To some extent they deal with the scenario where you go out and repeatedly sample from a population and calculate a statistic from each of those samples. The distributions *of that statistic* is what we will call the sampling distribution (as opposed to the sample distribution, which would more accurately describe the distribution of the data that we get in any one sample that we draw from the population).

2.6.1 Central Limit Theorem and Standard Error Demo

```
pop_students <- 10000  
zoom_min_pop <- rnorm(n = pop_students, mean = 600, sd = 100)  
hist(zoom_min_pop)
```



```
zoom_min_sample <- sample(x = zoom_min_pop,  
                           size = 200,  
                           replace = FALSE)  
hist(zoom_min_sample)
```



```
mean(zoom_min_sample)
```

```
## [1] 603.3624
```

```
sd(zoom_min_sample)
```

```
## [1] 102.4502
```

As a brief aside, let's review the idea of a loop

```
num_reps <- 100
```

```
data_vec <- rep(NA, num_reps) # this creates an empty vector of size num_reps with NA
```

```
# this loops through the vector starting at position 1 and ending at the final position
```

```
for (i in 1:num_reps){
  data_vec[i] <- i
}
```

Okay, so that's how we create an empty vector and how we loop through its different entries. For this demo, we will also need to remember how to generate random numbers from a norm distribution with a specified mean and standard deviation.

```
rmnorm(n = 10, mean = 5, sd = 2) # n is the number of random numbers we draw from this normal dist
```

```
## [1] 4.099294 3.411509 6.232366 5.549321 5.418897 4.410174 5.330722 6.588594
## [9] 7.336503 3.075376
```

Okay, so that's not bad. Now, that command will produce a vector with 10 random numbers. We can calculate the mean and standard deviation of those 10 numbers (which should be close to the values that we specified in `rmnorm()` with the `mean()` and `sd()` functions.

```
mean(rmnorm(n = 10, mean = 5, sd = 2))
```

```
## [1] 5.288278
```

```
sd(rmnorm(n = 10, mean = 5, sd = 2))
```

```
## [1] 2.544623
```

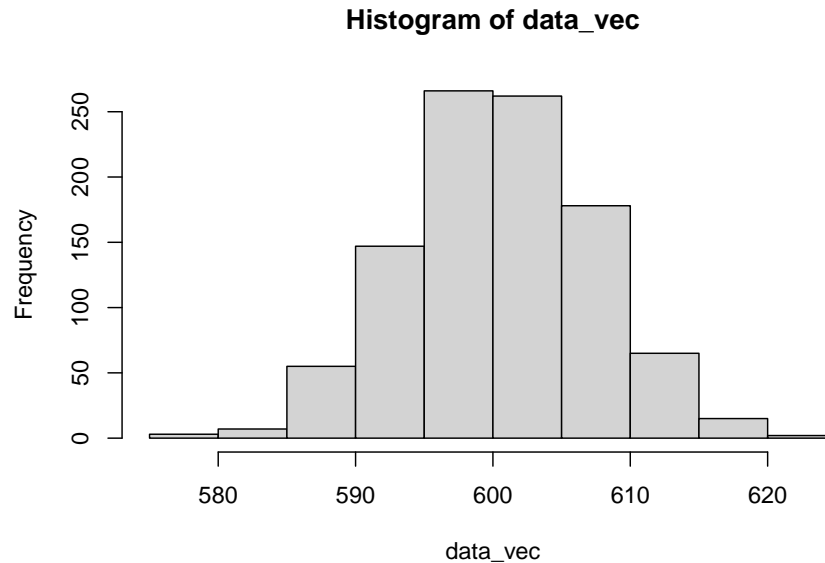
Next, let's act as if we are drawing a certain sample of size `samp_size` of data points for `num_reps` number of times. Keep in mind that, in practice, when we are collecting data ourselves in our own research, `num_reps` will almost always be 1. We are just demonstrating the underlying assumptions for how we can calculate some of the statistics that we use.

```
num_reps <- 1000 # specify how many times to take a sample
samp_size <- 200 # specify the size of each sample
data_vec <- rep(NA, num_reps) # create an empty vector of size num_reps with NA in each entry.
for (i in 1:num_reps){
  data_vec[i] <- mean(rmnorm(n = samp_size, mean = 600, sd = 100)) # store the mean of each of the
}
```

With this, we have a vector `data_vec` of size `num_reps` with the mean of each of our samples that we drew. This vector contains our sampling distribution of our sample means. **NOTE:** The standard deviation of this sampling mean is what we are calling our *standard error*.

We can plot a histogram of this sampling distribution and calculate the standard deviation of the sampling mean.

```
hist(data_vec)
```



```
sd(data_vec)
```

```
## [1] 6.876998
```

On your own, try copying this code and changing the `num_reps` and `sample_size` variables to larger and smaller values. Focus on how the x-axis values in your histogram change when you change the `num_reps` and `samp_size` variables.

Hint: CLT will explain the normal distribution of the sampling mean (the shape you see in the histogram) while the Weak Law of Large Numbers will explain the concentration around the true mean as `samp_size` increases (i.e., when we draw a larger sample size from the population, our sample mean gets closer to the population mean).

```
## Quick note on the rep() function: notice what happens when you specify "each" vs "t.
rep(c(1, 2), times = 5)
```

```
## [1] 1 2 1 2 1 2 1 2 1 2
```

```
rep(c(1, 2), each = 5)
```

```
## [1] 1 1 1 1 1 2 2 2 2 2
```


Chapter 3

Week 3: Data Cleaning, Organizing, Describing, and Communicating

This week we focus on different steps you will often take when you first start working with your data. These tend to fall under the umbrella of “data processing” and often need to happen before you can start doing any kind of analysis.

3.1 Visualizing your data

Once your data have been cleaned, you are ready to start visualizing what you are working with. There is a huge range of what you can do with these plots. That’s great! On the other hand, it can quickly start to feel overwhelming. To help get this under control and make it more manageable, it is convenient to think about the *types* of data that you have. In particular, are your variables nominal, ordinal, interval, or ratio variables?

3.1.1 One continuous variable (either predictor or outcome variable)

When you have one continuous variable, a standard option is to plot a histogram. These are plots that show the frequency of each of the values that the variable takes. Oftentimes it is helpful to create bins of values so that any number that falls in the 0-4 range counts in one bin, numbers from 5-9 are in a second bin, and so on.

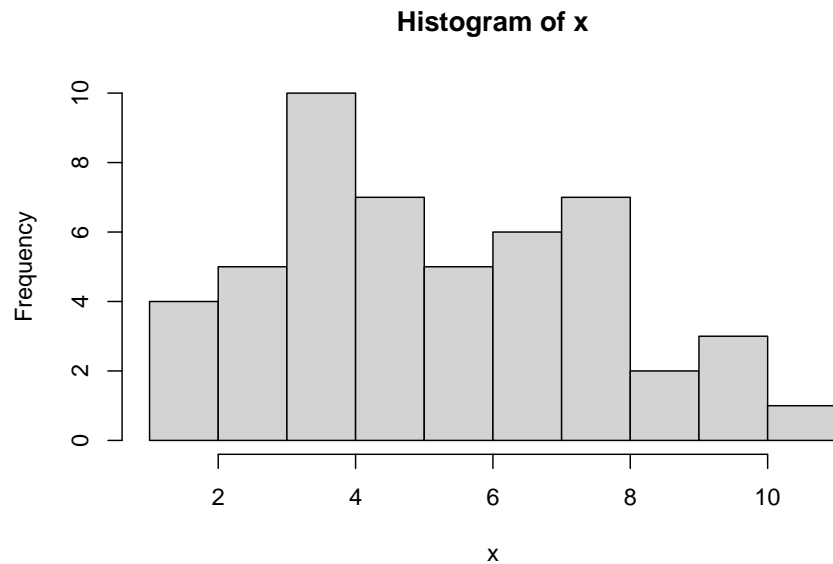
For this example, we will start by generating some data using `rnorm()`, which generates a random number (or in our case, `num` numbers) from a normal distribution with mean `mu` and standard deviation `stdev`.

```
num <- 50
mu <- 5
stdev <- 2

x <- rnorm(n = num, mean = mu, sd = stdev)
```

With these data generated, we can then quickly plot the histogram with `hist()`. This will use base R graphics.

```
hist(x)
```

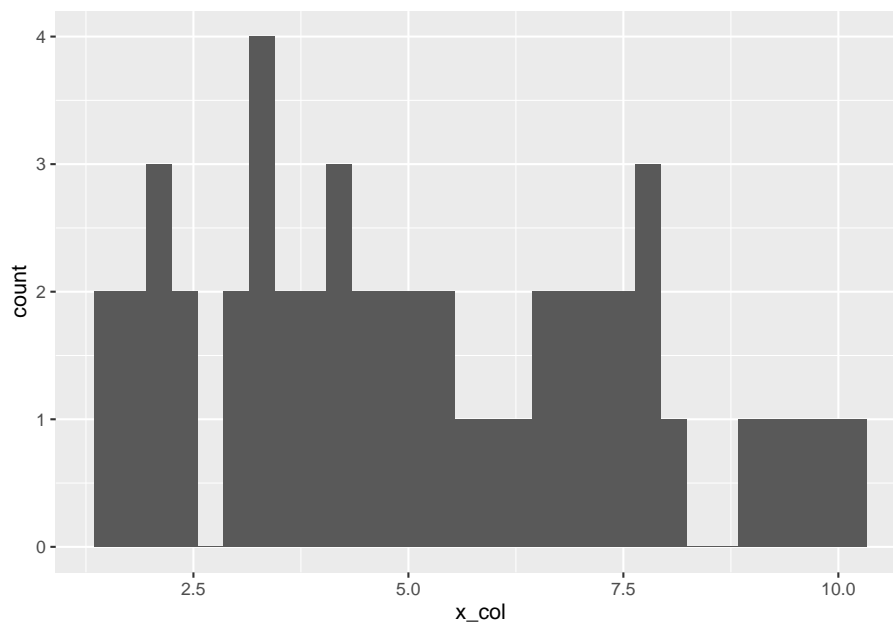


You can also do this using `ggplot` rather than base R graphics.

```
x_df <- tibble(x_col = x)

ggplot(data = x_df, mapping = aes(x = x_col)) +
  geom_histogram()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

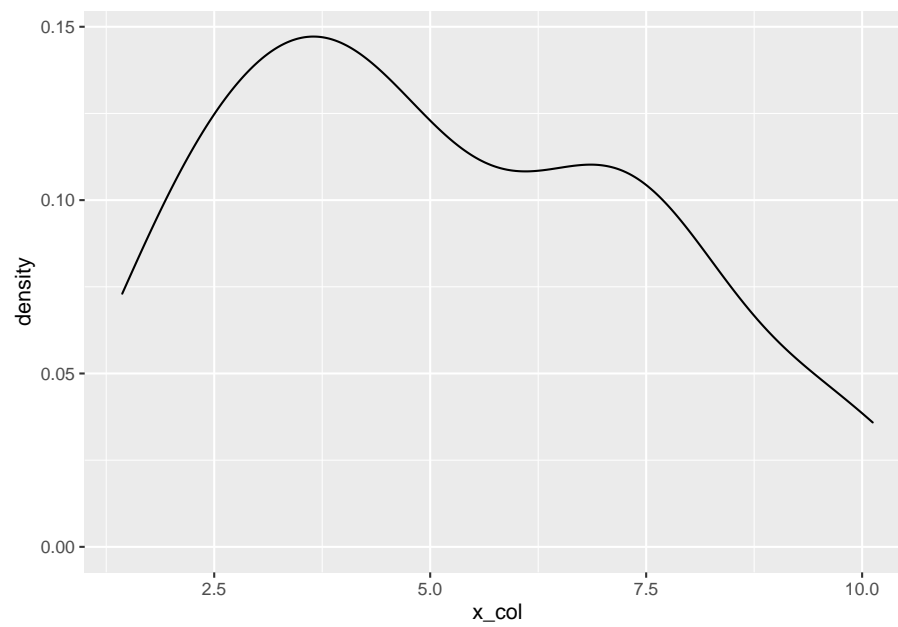


The histogram is a very standard plot, and you should consider it a go-to option in your toolkit. Alternatively, you can use `geom_density()` instead of `geom_histogram()` to get a smooth graph rather than one with discrete bins. We will use the same data that we generated before.

We will write this two ways to demonstrate how the pipe `%>%` operator works.

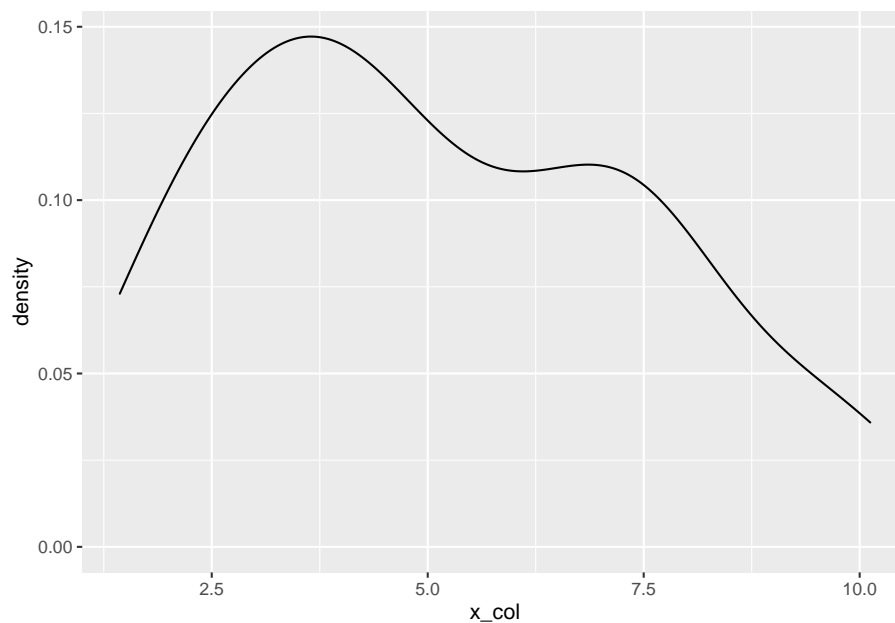
First way:

```
ggplot(data = x_df, mapping = aes(x = x_col)) +  
  geom_density()
```



Second way:

```
x_df %>%  
  ggplot(mapping = aes(x = x_col)) +  
  geom_density()
```



Just for fun, look at what happens to the the plot if you increase the sample size

First, we will generate the data with a sample size of 5,000 rather than 50.

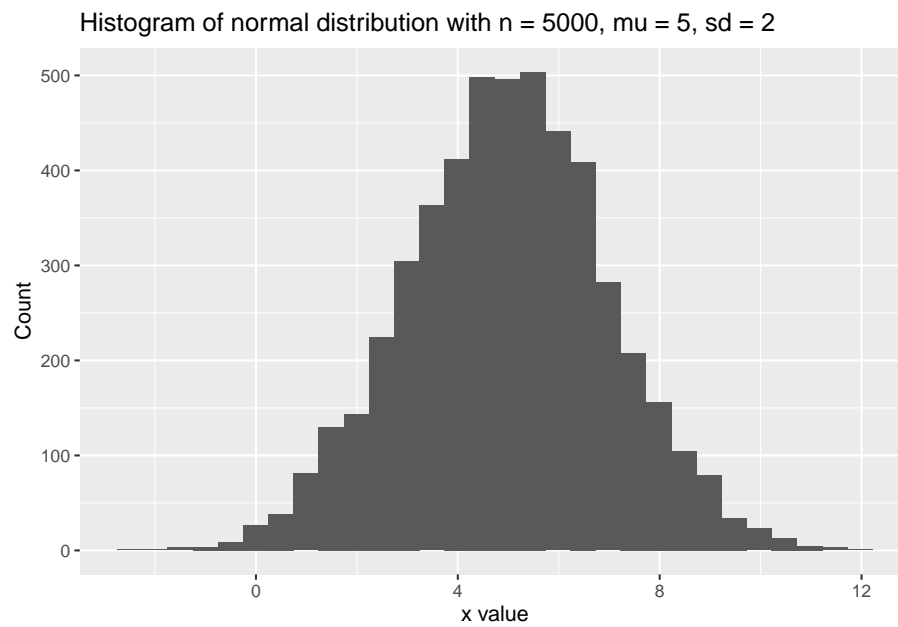
```
num <- 5000
mu <- 5
stdev <- 2

x <- rnorm(n = num, mean = mu, sd = stdev)
x_df <- tibble(x_col = x)
```

Then we will plot the histogram

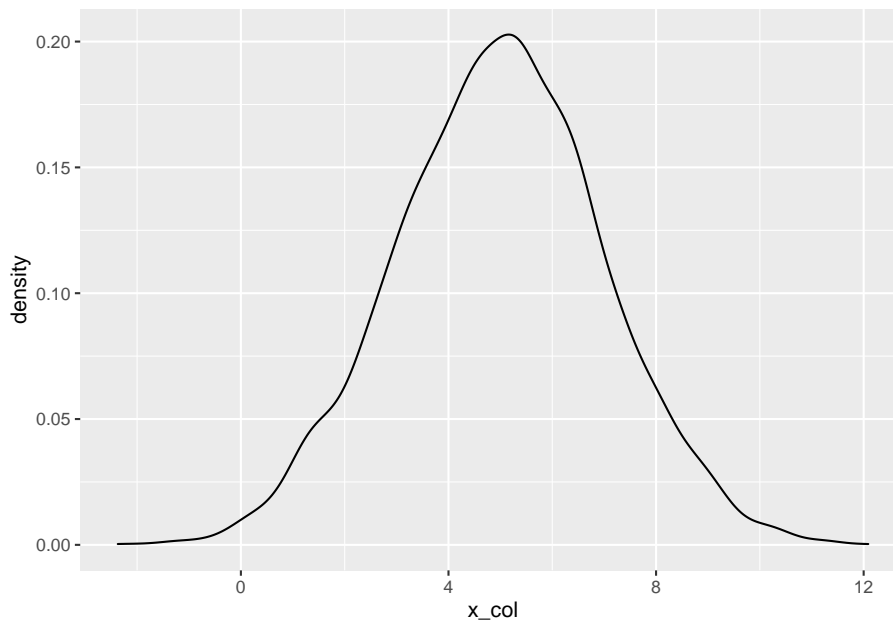
```
x_df %>%
  ggplot(aes(x = x_col)) +
  geom_histogram() +
  labs(x = "x value",
       y = "Count",
       title = "Histogram of normal distribution with n = 5000, mu = 5, sd = 2")
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



And, finally, we can make the density plot instead of the histogram, if that's our jam.

```
ggplot(data = x_df, mapping = aes(x = x_col)) +  
  geom_density()
```



3.1.2 One Discrete Variable (either predictor or outcome)

What if instead of a continuous (i.e., interval or ratio) variable we have a discrete variable such as a nominal (e.g., major, university) or ordinal (e.g., Likert scale item, level of education) variable? For that we can use something like `geom_bar()` or `geom_col()` to plot the counts of observations within each of those categories.

To demonstrate this, we first need some data to work with. We will use the pre-semester, prior knowledge survey that everyone took. I have combined this year's results with last year's results in order to increase the sample size. After reading in the data, I will also use the `clean_names()` function from the `janitor` package.

```
## load in the data
survey_df <- read_csv("ENGE_5714_2021_pre_survey.csv")
```

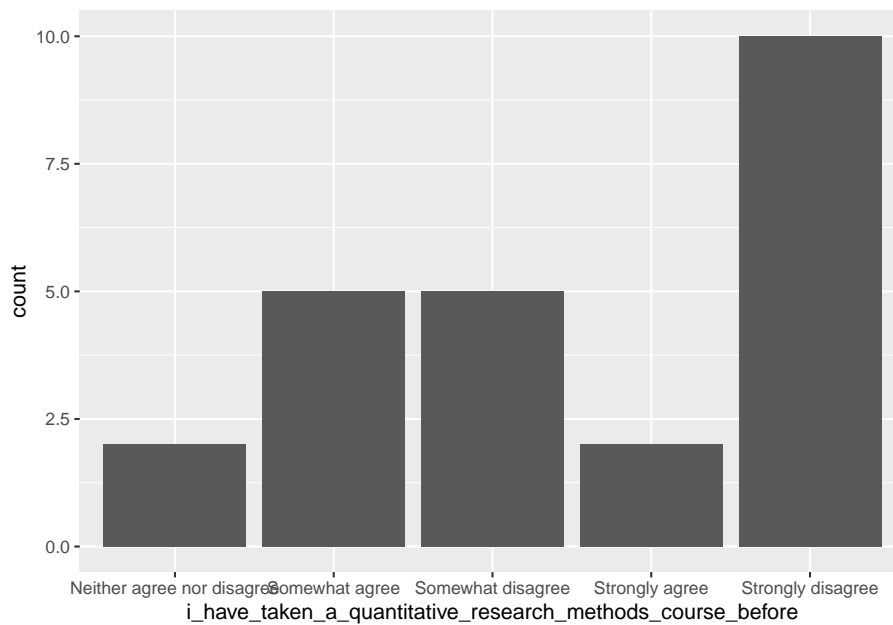
```
## Parsed with column specification:
## cols(
##   .default = col_character(),
##   student_id = col_double()
## )
```

```
## See spec(...) for full column specifications.
```

```
survey_df <- survey_df %>% clean_names()
```

Next, we can go ahead and make a bar plot with the following code:

```
survey_df %>%
  ggplot(aes(x = i_have_taken_a_quantitative_research_methods_course_before)) +
  geom_bar()
```



Notice that the ordering is not quite what we would want. It is alphabetical. Let's try to fix this.

Here is one way: we first specify the levels of that variable (i.e., the different values that it could take) and store that in the variable `q_levels`. Then, we pass that to the `factor()` function, which will tell R that we want whichever variable is passed to `factor()` two things. First, it will say that we want to make that variable a factor variable with `levels = ...`. Second, we set `ordered = TRUE` to tell R that there is a specific ordering to that variable. This way, whenever there is something like a plot that we make, the ordering will persist in the labeling and R will not show the labels in alphabetical order.

Here is an example of that in action:

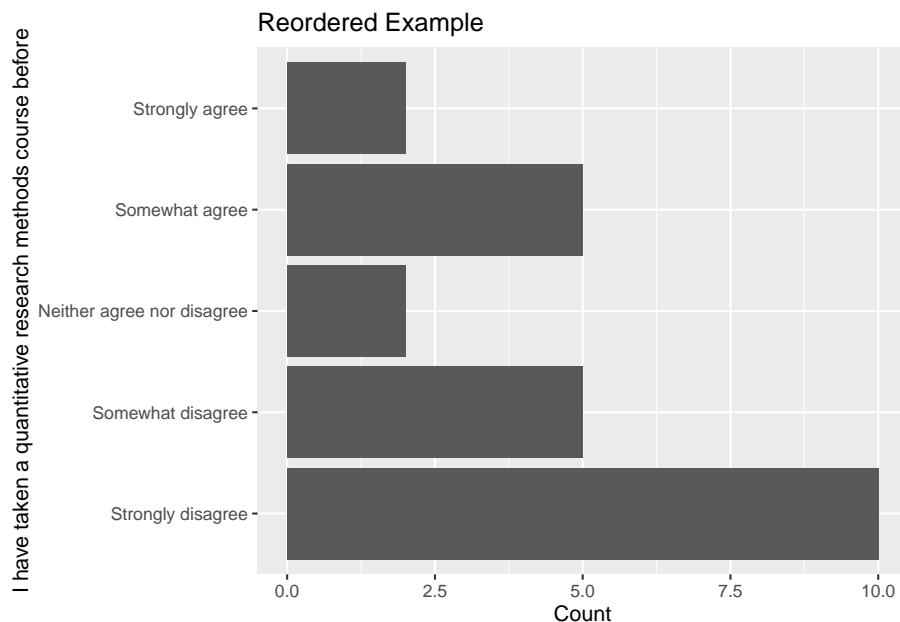
```
q_levels <- c("Strongly disagree", "Somewhat disagree", "Neither agree nor disagree",
              "Somewhat agree", "Strongly agree")
```



```
survey_df$i_have_taken_a_quantitative_research_methods_course_before <- factor(survey_df$i_have_t
                                                                    levels = q_levels,
                                                                    ordered = TRUE)
```

Now try plotting these data. We will also add in a `coord_flip()` to plot the categories along the y-axis. This is a common move to avoid text from the different levels overlapping with each other. Finally, we will also change the x, y, and title labels with `labs()`.

```
survey_df %>%
  ggplot(aes(x = i_have_taken_a_quantitative_research_methods_course_before)) +
  geom_bar() +
  coord_flip() +
  labs(x = "I have taken a quantitative research methods course before",
       y = "Count",
       title = "Reordered Example")
```



3.2 Joining two datasets

Let's imagine that we have a separate dataset that has information about the students who completed the pre-course prior knowledge survey.

First, we will load in that dataset

```
survey_info_df <- read_csv("survey_student_info.csv")
```

```
## Parsed with column specification:
## cols(
##   student_id = col_double(),
##   standing = col_character(),
##   college = col_character(),
##   required = col_character()
## )
```

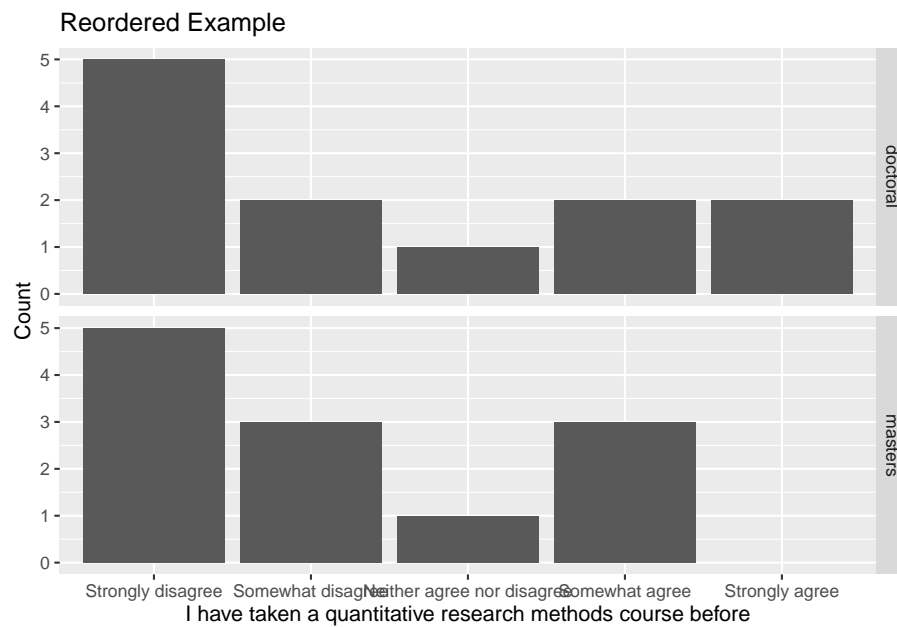
Next, let's join the two datasets based on the student id column, which is in each of the two dataframes.

```
survey_df <- survey_df %>% inner_join(survey_info_df, by = "student_id")
```

Now we should have both datasets joined into one and saved as `survey_df`.

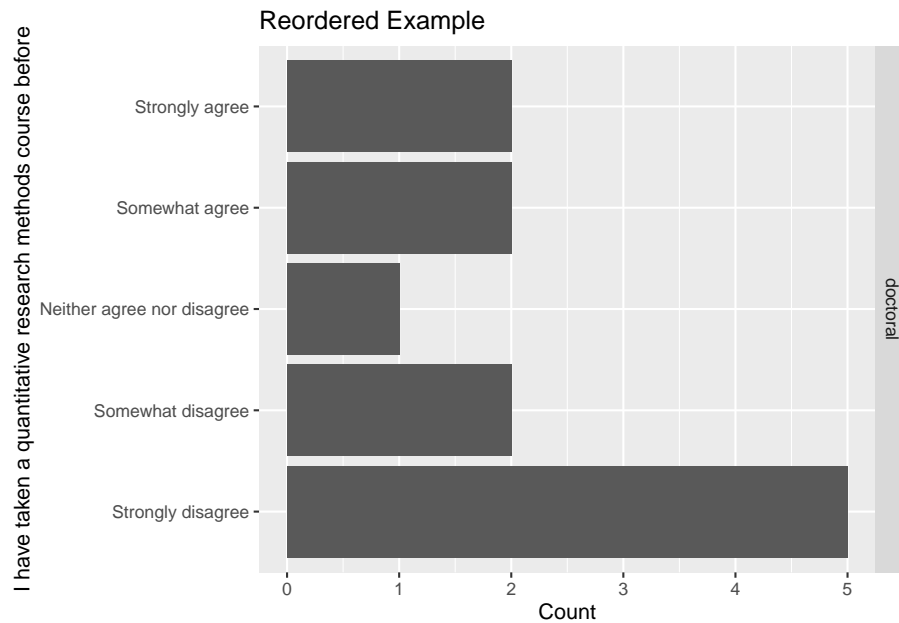
With this, we can make some nicer plots and do something like use `facet_grid()` to look at students who are masters and doctoral students, for example.

```
survey_df %>%
  ggplot(aes(x = i_have_taken_a_quantitative_research_methods_course_before)) +
  geom_bar() +
  facet_grid(standing ~.) +
  labs(x = "I have taken a quantitative research methods course before",
       y = "Count",
       title = "Reordered Example")
```



The x axis looks a little crowded. What if we try `coord_flip()`?

```
survey_df %>%
  filter(standing == "doctoral") %>%
  ggplot(aes(x = i_have_taken_a_quantitative_research_methods_course_before)) +
  geom_bar() +
  coord_flip() +
  facet_grid(standing ~.) +
  labs(x = "I have taken a quantitative research methods course before",
       y = "Count",
       title = "Reordered Example")
```



That looks much better.

A quick note on filters

If you want to look at only a subset of your data, you will want to use the `filter()` function. The general idea is that you can look at observations (rows) that match a certain criteria. For example, you may want to only look at students from a certain region or year or major. In our case, with the prior knowledge survey, let's say we only want to look at student who have to take the course (i.e., there is a “yes” for them for the `required` variable). We can do that with the first line. The second line just stores the result as a new dataframe called `filtered_df`.

```
survey_df %>% filter(required == "yes")
```

```
## # A tibble: 12 x 52
##   student_id i_have_taken_a_ i_am_interested i_know_what_a_t
##   <dbl> <ord> <chr> <chr>
## 1         1 Somewhat disagr~ Somewhat agree Strongly disagr~
## 2         2 Strongly disagr~ Neither agree n~ Somewhat agree
## 3         4 Somewhat disagr~ Strongly agree Strongly disagr~
## 4         8 Somewhat agree Somewhat agree Somewhat agree
## 5         9 Strongly disagr~ Strongly agree Somewhat agree
## 6        11 Strongly disagr~ Strongly agree Strongly disagr~
```

```

## 7      16 Strongly agree   Strongly agree   Somewhat agree
## 8      17 Strongly disagr~ Strongly agree   Strongly disagr~
## 9      18 Somewhat disagr~ Somewhat agree   Somewhat disagr~
## 10     20 Strongly disagr~ Neither agree n~ Neither agree n~
## 11     22 Strongly disagr~ Strongly agree   Strongly disagr~
## 12     23 Somewhat agree   Strongly agree   Somewhat agree
## # ... with 48 more variables: i_know_what_a_type_ii_error_is <chr>,
## #   i_know_what_a_statistical_confidence_level_is <chr>,
## #   i_know_what_a_p_value_is <chr>, i_know_what_p_hacking_means <chr>,
## #   i_know_what_statistical_power_means <chr>,
## #   i_have_heard_of_frequentist_statistics_before <chr>,
## #   i_have_heard_of_bayesian_statistics_before <chr>,
## #   i_have_heard_the_term_parametric_statistics_before <chr>,
## #   i_have_heard_the_term_non_parametric_statistics_before <chr>,
## #   i_know_what_a_histogram_is <chr>,
## #   i_know_what_a_probability_distribution_is <chr>,
## #   i_know_what_a_random_variable_is <chr>,
## #   i_know_what_a_probability_distribution_function_is <chr>,
## #   i_know_what_a_cumulative_distribution_function_is <chr>,
## #   i_know_what_the_expectation_of_a_random_variable_is <chr>,
## #   i_know_how_to_calculate_the_variance_of_a_random_variable <chr>,
## #   i_know_what_a_z_score_is <chr>,
## #   i_know_how_to_calculate_the_correlation_between_two_variables <chr>,
## #   i_know_how_to_interpret_the_correlation_coefficient_between_two_variables <chr>,
## #   i_have_heard_of_linear_regression <chr>,
## #   i_know_how_to_run_a_linear_regression_in_some_software_or_by_hand_if_im_feeling_wild <chr>,
## #   i_know_how_to_interpret_a_linear_regression <chr>,
## #   i_have_heard_of_multiple_regression <chr>,
## #   i_know_how_to_perform_a_multiple_regression <chr>,
## #   i_know_how_to_interpret_a_multiple_regression <chr>,
## #   i_have_heard_of_logistic_regression <chr>,
## #   i_understand_when_to_use_a_logistic_regression <chr>,
## #   i_know_how_to_interpret_the_results_of_a_logistic_regression <chr>,
## #   i_have_heard_of_t_tests <chr>, i_have_performed_a_t_test_before <chr>,
## #   i_know_how_to_interpret_the_results_of_a_t_test <chr>,
## #   i_have_heard_of_analysis_of_variance <chr>,
## #   i_understand_when_to_run_an_analysis_of_variance_anova <chr>,
## #   i_know_how_to_interpret_the_results_from_an_anova <chr>,
## #   i_have_heard_of_a_chi_square_test <chr>,
## #   i_have_used_a_chi_square_test_before <chr>,
## #   i_know_how_to_interpret_the_results_of_a_chi_square_test <chr>,
## #   i_have_heard_of_cluster_analysis_before <chr>,
## #   i_have_used_cluster_analysis_before <chr>,
## #   i_know_how_to_interpret_the_results_of_a_cluster_analysis <chr>,
## #   i_have_heard_of_factor_analysis_either_exploratory_or_confirmatory <chr>,
## #   i_have_used_factor_analysis_either_exploratory_or_confirmatory <chr>,

```

```
## #   i_know_how_to_interpret_the_results_of_a_factor_analysis_either_exploratory_or_
## #   i_already_have_r_and_rstudio_downloaded_to_my_computer <chr>,
## #   i_have_used_r_before <chr>, standing <chr>, college <chr>, required <chr>
```

```
filtered_df <- survey_df %>% filter(required == "yes")
```

A little more about plotting

We are going to shift gears again and look at a few different kinds of plots. The main thing to remember here is that you want to think about whether the variables you have are nominal, ordinal, or continuous (that includes interval and ratio).

3.3 Discrete Predictor, Continuous Outcome

So far we have looked at plots for one variable, but of course we want to have ways to plot multiple variables simultaneously. We will start with the scenario where we want to plot a continuous variable against a discrete variable. This can arise when you want to plot something like an assessment score and you think it may differ across groups in some way (maybe you intentionally introduced a difference by exposing the two groups to different interventions, for example).

In these scenarios, a boxplot is a very standard way to go.

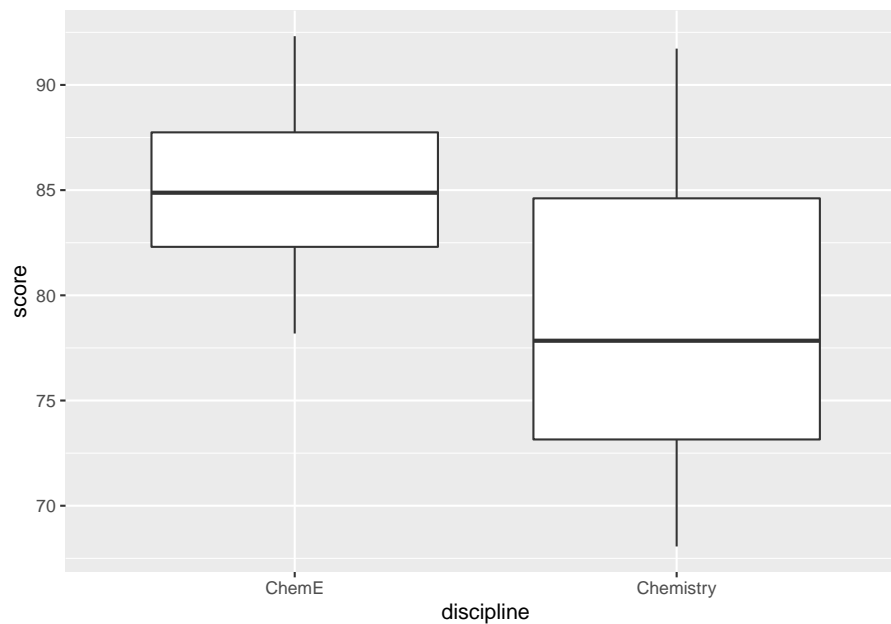
To demonstrate this, we will simulate a situation in which we want to look at differences on an assessment. We are specifically interested in differences between chemistry and chemical engineering students. Let's go ahead and create the data by creating two groups of 20 students each. The chemical engineering students will have scores generated from a normal distribution with $\mu = 85$ and $\sigma = 4$ (i.e., a mean of 85 and a standard deviation of 4). We will say the chemistry students have scores from a normal distribution with $\mu = 78$ and $\sigma = 6$. This about what these distributions might look like in your head.

```
group_size <- 20
chem_e_scores <- rnorm(n = group_size, mean = 85, sd = 4)
chem_scores <- rnorm(n = group_size, mean = 78, sd = 6)

data_df <- tibble(
  discipline = rep(c("ChemE", "Chemistry"), each = group_size),
  score = c(chem_e_scores, chem_scores)
)
```

With these data, we can then create a boxplot using `geom_boxplot()`

```
data_df %>%  
  ggplot(aes(x = discipline, y = score)) +  
  geom_boxplot()
```

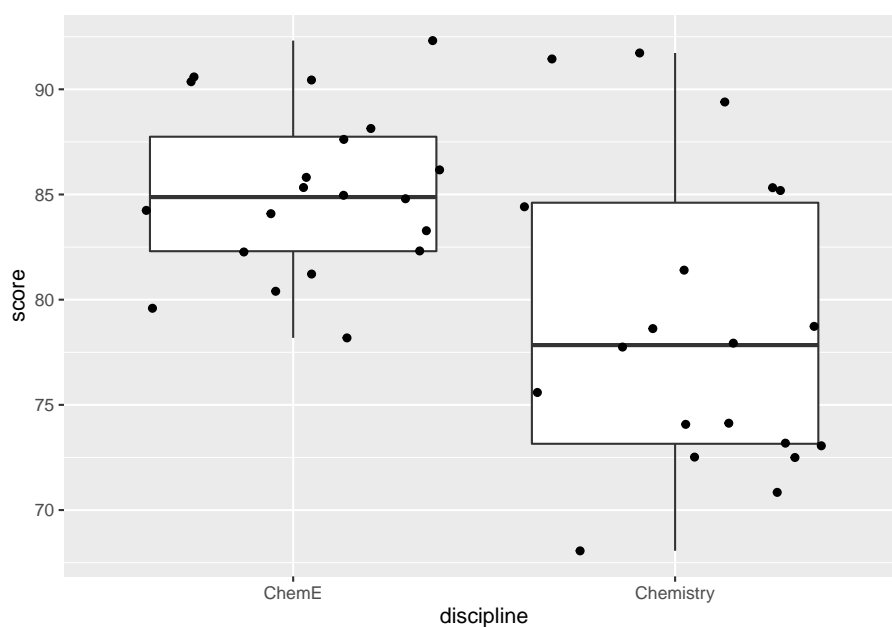


You can make a few modifications to possibly make this plot easier to read.

The first is to put the discrete category on the y axis instead of the x axis.

The second is to use `geom_jitter()` in addition to `geom_boxplot()` to show the individual points in each group.

```
data_df %>%  
  ggplot(aes(y = score, x = discipline)) +  
  geom_boxplot() +  
  geom_jitter()
```



3.4 Continuous predictor and continuous outcome

First, let's re-do a lot of the steps in this week's script for reading in data and transforming it a little

```
mydata <- read_csv("Free Reduced Lunch by Schools and Grade Structures 2008-2017_final
```

```
## Parsed with column specification:
## cols(
##   .default = col_character(),
##   div_num = col_double()
## )
```

```
## See spec(...) for full column specifications.
```

Check the structure of the data (this output is a bit long).

```
str(mydata)
```

```
## tibble [2,101 x 137] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
```



```

## $ sch_id      : chr [1:2101] "001-0070" "001-0080" "001-0530" "001-0540" ...
## $ div_num     : num [1:2101] 1 1 1 1 1 1 1 1 1 ...
## $ div_name    : chr [1:2101] "Accomack County" "Accomack County" "Accomack County" "Accomack County" ...
## $ school_num  : chr [1:2101] "0070<U+00A0>" "0080<U+00A0>" "0530<U+00A0>" "0540<U+00A0>" ...
## $ school_name : chr [1:2101] "NANDUA HIGH" "CHINCOTEAGUE ELEM" "TANGIER COMBINED" "ARCADIA ...
## $ school_name2 : chr [1:2101] NA NA NA NA ...
## $ type0809    : chr [1:2101] "SCH-HIGH" "SCH-ELEM" "SCH-COMB" "SCH-HIGH" ...
## $ lowgrade_2008 : chr [1:2101] "9" "PK" "KG" "9" ...
## $ higrade_2008 : chr [1:2101] "12" "5" "12" "12" ...
## $ totalFT_2008 : chr [1:2101] "731" "263" "80" "638" ...
## $ total_2008   : chr [1:2101] "731" "263" "80" "638" ...
## $ snp_0809     : chr [1:2101] "659" "257" "80" "622" ...
## $ free_elig_0809 : chr [1:2101] "306" "95" "38" "289" ...
## $ free_per_0809 : chr [1:2101] "46.43%" "36.96%" "47.50%" "46.46%" ...
## $ red_elig_0809 : chr [1:2101] "64" "8" "0" "56" ...
## $ red_per_0809  : chr [1:2101] "9.71%" "3.11%" "0.00%" "9.00%" ...
## $ totalFRL_0809 : chr [1:2101] "370" "103" "38" "345" ...
## $ totalper_0809 : chr [1:2101] "56.15%" "40.08%" "47.50%" "55.47%" ...
## $ type0910     : chr [1:2101] "SCH-HIGH" "SCH-ELEM" "SCH-COMB" "SCH-HIGH" ...
## $ lowgrade_2009 : chr [1:2101] "9" "PK" "KG" "9" ...
## $ higrade_2009 : chr [1:2101] "12" "5" "12" "12" ...
## $ totalFT_2009 : chr [1:2101] "654" "266" "78" "634" ...
## $ total_2009   : chr [1:2101] "654" "266" "78" "634" ...
## $ snp_0910     : chr [1:2101] "655" "266" "78" "635" ...
## $ free_elig_0910 : chr [1:2101] "290" "99" "36" "286" ...
## $ free_per_0910 : chr [1:2101] "44.27%" "37.22%" "46.15%" "45.04%" ...
## $ red_elig_0910 : chr [1:2101] "37" "14" "0" "66" ...
## $ red_per_0910  : chr [1:2101] "5.65%" "5.26%" "0.00%" "10.39%" ...
## $ totalFRL_0910 : chr [1:2101] "327" "113" "36" "352" ...
## $ totalper_0910 : chr [1:2101] "49.92%" "42.48%" "46.15%" "55.43%" ...
## $ type1011     : chr [1:2101] "SCH-HIGH" "SCH-ELEM" "SCH-COMB" "SCH-HIGH" ...
## $ lowgrade_2010 : chr [1:2101] "9" "PK" "KG" "9" ...
## $ higrade_2010 : chr [1:2101] "12" "5" "12" "12" ...
## $ totalFT_2010 : chr [1:2101] "603" "268" "74" "614" ...
## $ total_2010   : chr [1:2101] "603" "268" "74" "614" ...
## $ snp_1011     : chr [1:2101] "603" "277" "74" "606" ...
## $ free_elig_1011 : chr [1:2101] "285" "108" "32" "308" ...
## $ free_per_1011 : chr [1:2101] "47.26%" "38.99%" "43.24%" "50.83%" ...
## $ red_elig_1011 : chr [1:2101] "46" "8" "0" "50" ...
## $ red_per_1011  : chr [1:2101] "7.63%" "2.89%" "0.00%" "8.25%" ...
## $ totalFRL_1011 : chr [1:2101] "331" "116" "32" "358" ...
## $ totalper_1011 : chr [1:2101] "54.89%" "41.88%" "43.24%" "59.08%" ...
## $ type1112     : chr [1:2101] "SCH-HIGH" "SCH-ELEM" "SCH-COMB" "SCH-HIGH" ...
## $ lowgrade_2011 : chr [1:2101] "9" "PK" "KG" "9" ...
## $ higrade_2011 : chr [1:2101] "12" "5" "12" "12" ...
## $ totalFT_2011 : chr [1:2101] "593" "276" "73" "605" ...

```

```

## $ total_2011      : chr [1:2101] "593" "276" "73" "605" ...
## $ snp_1112        : chr [1:2101] "593" "281" "73" "611" ...
## $ free_elig_1112: chr [1:2101] "289" "116" "31" "318" ...
## $ free_per_1112  : chr [1:2101] "48.74%" "41.28%" "42.47%" "52.05%" ...
## $ red_elig_1112  : chr [1:2101] "50" "14" "0" "44" ...
## $ red_per_1112   : chr [1:2101] "8.43%" "4.98%" "0.00%" "7.20%" ...
## $ totalFRL_1112  : chr [1:2101] "339" "130" "31" "362" ...
## $ totalper_1112  : chr [1:2101] "57.17%" "46.26%" "42.47%" "59.25%" ...
## $ type1213       : chr [1:2101] "SCH-HIGH" "SCH-ELEM" "SCH-COMB" "SCH-HIGH" ...
## $ lowgrade_2012  : chr [1:2101] "9" "PK" "KG" "9" ...
## $ higrade_2012   : chr [1:2101] "12" "5" "12" "12" ...
## $ totalFT_2012   : chr [1:2101] "637" "258" "68" "579" ...
## $ total_2012     : chr [1:2101] "637" "258" "68" "579" ...
## $ snp_1213       : chr [1:2101] "633" "259" "68" "579" ...
## $ free_elig_1213: chr [1:2101] "324" "117" "21" "348" ...
## $ free_per_1213  : chr [1:2101] "51.18%" "45.17%" "30.88%" "60.10%" ...
## $ red_elig_1213  : chr [1:2101] "42" "20" "5" "33" ...
## $ red_per_1213   : chr [1:2101] "6.64%" "7.72%" "7.35%" "5.70%" ...
## $ totalFRL_1213  : chr [1:2101] "366" "137" "26" "381" ...
## $ totalper_1213  : chr [1:2101] "57.82%" "52.90%" "38.24%" "65.80%" ...
## $ type1314       : chr [1:2101] "SCH-HIGH" "SCH-ELEM" "SCH-COMB" "SCH-HIGH" ...
## $ lowgrade_2013  : chr [1:2101] "9" "PK" "KG" "9" ...
## $ higrade_2013   : chr [1:2101] "12" "5" "12" "12" ...
## $ totalFT_2013   : chr [1:2101] "670" "238" "66" "582" ...
## $ total_2013     : chr [1:2101] "670" "238" "66" "582" ...
## $ snp_1314       : chr [1:2101] "668" "239" "56" "589" ...
## $ free_elig_1314: chr [1:2101] "346" "102" "12" "347" ...
## $ free_per_1314  : chr [1:2101] "51.80%" "42.68%" "21.43%" "58.91%" ...
## $ red_elig_1314  : chr [1:2101] "44" "19" "4" "54" ...
## $ red_per_1314   : chr [1:2101] "6.59%" "7.95%" "7.14%" "9.17%" ...
## $ totalFRL_1314  : chr [1:2101] "390" "121" "16" "401" ...
## $ totalper_1314  : chr [1:2101] "58.38%" "50.63%" "28.57%" "68.08%" ...
## $ type1415       : chr [1:2101] NA NA NA NA ...
## $ lowgrade_2014  : chr [1:2101] "9" "PK" "KG" "9" ...
## $ higrade_2014   : chr [1:2101] "12" "5" "12" "12" ...
## $ totalFT_2014   : chr [1:2101] "685" "251" "65" "581" ...
## $ total_2014     : chr [1:2101] "685" "251" "65" "581" ...
## $ snp_1415       : chr [1:2101] "672" "239" "61" "586" ...
## $ free_elig_1415: chr [1:2101] "361" "93" "14" "351" ...
## $ free_per_1415  : chr [1:2101] "53.72%" "38.91%" "22.95%" "59.90%" ...
## $ red_elig_1415  : chr [1:2101] "40" "17" "4" "40" ...
## $ red_per_1415   : chr [1:2101] "5.95%" "7.11%" "6.56%" "6.83%" ...
## $ totalFRL_1415  : chr [1:2101] "401" "110" "18" "391" ...
## $ totalper_1415  : chr [1:2101] "59.67%" "46.03%" "29.51%" "66.72%" ...
## $ CEP_1516       : chr [1:2101] "#NULL!" "#NULL!" "#NULL!" "#NULL!" ...
## $ type1516       : chr [1:2101] "SCH-HIGH" "SCH-ELEM" "SCH-COMB" "SCH-HIGH" ...

```

```

## $ lowgrade_2015 : chr [1:2101] "9" "PK" "KG" "9" ...
## $ higrade_2015  : chr [1:2101] "12" "5" "12" "12" ...
## $ totalFT_2015  : chr [1:2101] "737" "259" "65" "621" ...
## $ total_2015    : chr [1:2101] "737" "259" "65" "621" ...
## $ snp_1516      : chr [1:2101] "728" "268" "67" "608" ...
## $ free_elig_1516: chr [1:2101] "362" "109" "12" "339" ...
## $ free_per_1516 : chr [1:2101] "49.73%" "40.67%" "17.91%" "55.76%" ...
## [list output truncated]
## - attr(*, "spec")=
## .. cols(
## ..   sch_id = col_character(),
## ..   div_num = col_double(),
## ..   div_name = col_character(),
## ..   school_num = col_character(),
## ..   school_name = col_character(),
## ..   school_name2 = col_character(),
## ..   type0809 = col_character(),
## ..   lowgrade_2008 = col_character(),
## ..   higrade_2008 = col_character(),
## ..   totalFT_2008 = col_character(),
## ..   total_2008 = col_character(),
## ..   snp_0809 = col_character(),
## ..   free_elig_0809 = col_character(),
## ..   free_per_0809 = col_character(),
## ..   red_elig_0809 = col_character(),
## ..   red_per_0809 = col_character(),
## ..   totalFRL_0809 = col_character(),
## ..   totalper_0809 = col_character(),
## ..   type0910 = col_character(),
## ..   lowgrade_2009 = col_character(),
## ..   higrade_2009 = col_character(),
## ..   totalFT_2009 = col_character(),
## ..   total_2009 = col_character(),
## ..   snp_0910 = col_character(),
## ..   free_elig_0910 = col_character(),
## ..   free_per_0910 = col_character(),
## ..   red_elig_0910 = col_character(),
## ..   red_per_0910 = col_character(),
## ..   totalFRL_09010 = col_character(),
## ..   totalper_0910 = col_character(),
## ..   type1011 = col_character(),
## ..   lowgrade_2010 = col_character(),
## ..   higrade_2010 = col_character(),
## ..   totalFT_2010 = col_character(),
## ..   total_2010 = col_character(),
## ..   snp_1011 = col_character(),

```

```

## .. free_elig_1011 = col_character(),
## .. free_per_1011 = col_character(),
## .. red_elig_1011 = col_character(),
## .. red_per_1011 = col_character(),
## .. totalFRL_1011 = col_character(),
## .. totalper_1011 = col_character(),
## .. type1112 = col_character(),
## .. lowgrade_2011 = col_character(),
## .. higrade_2011 = col_character(),
## .. totalFT_2011 = col_character(),
## .. total_2011 = col_character(),
## .. snp_1112 = col_character(),
## .. free_elig_1112 = col_character(),
## .. free_per_1112 = col_character(),
## .. red_elig_1112 = col_character(),
## .. red_per_1112 = col_character(),
## .. totalFRL_1112 = col_character(),
## .. totalper_1112 = col_character(),
## .. type1213 = col_character(),
## .. lowgrade_2012 = col_character(),
## .. higrade_2012 = col_character(),
## .. totalFT_2012 = col_character(),
## .. total_2012 = col_character(),
## .. snp_1213 = col_character(),
## .. free_elig_1213 = col_character(),
## .. free_per_1213 = col_character(),
## .. red_elig_1213 = col_character(),
## .. red_per_1213 = col_character(),
## .. totalFRL_1213 = col_character(),
## .. totalper_1213 = col_character(),
## .. type1314 = col_character(),
## .. lowgrade_2013 = col_character(),
## .. higrade_2013 = col_character(),
## .. totalFT_2013 = col_character(),
## .. total_2013 = col_character(),
## .. snp_1314 = col_character(),
## .. free_elig_1314 = col_character(),
## .. free_per_1314 = col_character(),
## .. red_elig_1314 = col_character(),
## .. red_per_1314 = col_character(),
## .. totalFRL_1314 = col_character(),
## .. totalper_1314 = col_character(),
## .. type1415 = col_character(),
## .. lowgrade_2014 = col_character(),
## .. higrade_2014 = col_character(),
## .. totalFT_2014 = col_character(),

```

```

## .. total_2014 = col_character(),
## .. snp_1415 = col_character(),
## .. free_elig_1415 = col_character(),
## .. free_per_1415 = col_character(),
## .. red_elig_1415 = col_character(),
## .. red_per_1415 = col_character(),
## .. totalFRL_1415 = col_character(),
## .. totalper_1415 = col_character(),
## .. CEP_1516 = col_character(),
## .. type1516 = col_character(),
## .. lowgrade_2015 = col_character(),
## .. higrade_2015 = col_character(),
## .. totalFT_2015 = col_character(),
## .. total_2015 = col_character(),
## .. snp_1516 = col_character(),
## .. free_elig_1516 = col_character(),
## .. free_per_1516 = col_character(),
## .. red_elig_1516 = col_character(),
## .. red_Per_1516 = col_character(),
## .. totalFRL_1516 = col_character(),
## .. totalper_1516 = col_character(),
## .. CEP_1617 = col_character(),
## .. type1617 = col_character(),
## .. lowgrade_2016 = col_character(),
## .. higrade_2016 = col_character(),
## .. totalFT_2016 = col_character(),
## .. total_2016 = col_character(),
## .. snp_2016 = col_character(),
## .. free_elig_1617 = col_character(),
## .. free_per_1617 = col_character(),
## .. red_elig_1617 = col_character(),
## .. red_per_1617 = col_character(),
## .. totalFRL_1617 = col_character(),
## .. totalper_1617 = col_character(),
## .. CEP_1718 = col_character(),
## .. type1718 = col_character(),
## .. lowgrade_2017 = col_character(),
## .. higrade_2017 = col_character(),
## .. totalFT_2017 = col_character(),
## .. total_2017 = col_character(),
## .. snp_1718 = col_character(),
## .. free_elig_1718 = col_character(),
## .. free_per_1718 = col_character(),
## .. red_elig_1718 = col_character(),
## .. red_per_1718 = col_character(),
## .. totalFRL_1718 = col_character(),

```

```
## .. totalper_1718 = col_character(),
## .. stable = col_character(),
## .. new = col_character(),
## .. closed = col_character(),
## .. close_yr = col_character(),
## .. reuseid = col_character(),
## .. gradechg = col_character(),
## .. gradechg_yr = col_character(),
## .. grchgyr_2 = col_character()
## .. )
```

Or just check the structure of one specific variable.

```
str(mydata$total_2017)
```

```
## chr [1:2101] "742" "236" "60" "624" "286" "485" "583" "550" "600" "514" ...
```

NOTE: When you have a lot of variables, running this `str()` function is not a great idea - the output is a little too cumbersome

3.5 Mutating Variables

Note that almost all of the data reads in as a “character” data type which are just strings. This can create issues.

We know that many of the columns are actually storing numbers or “numeric” values as R refers to them. We need to fix this.

Let’s tell R that these columns (at least the two we are going to use) are numeric.

We are going to see two interchangeable ways to do this.

First, we use the `$` operator which lets me specify a specific column within my data frame in combination with the `as.numeric()` function

```
mydata$total_2017<-as.numeric(mydata$total_2017)
mydata$totalFRL_1718<-as.numeric(mydata$totalFRL_1718)
```

Some columns have a percent symbol, which you will need to remove before coercing to numeric data type

```
mydata <- mydata %>%
  mutate(totalper_0809 = str_remove(totalper_0809, "%"))
```

Then we can change the column from character to numeric

```
mydata$totalper_0809 <- as.numeric(mydata$totalper_0809)
```

```
## Warning: NAs introduced by coercion
```

Check to make sure it converted the column type correctly using `str()`.

```
str(mydata$totalper_0809)
```

```
##  num [1:2101] 56.1 40.1 47.5 55.5 33.4 ...
```

Second, alternatively, we can do this for a whole set of variables at once. We just need to specify a matching criteria.

```
newdf <- mydata %>%
  mutate_at(vars(starts_with("total")), as.numeric)
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```



```
## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion
```

```
newdf <- newdf %>%
  mutate_at(vars(starts_with("totalFRL")), as.numeric)
```

Check whether the old and new variables are stored differently (old as a character, new as a numeric variable)

```
str(mydata$total_2008)
```

```
## chr [1:2101] "731" "263" "80" "638" "333" "536" "610" "490" "585" "450" ...
```

```
str(newdf$total_2008)
```

```
## num [1:2101] 731 263 80 638 333 536 610 490 585 450 ...
```

3.6 Filtering and Selecting

A basic operation we do a lot is to filter the data so that we are working with a subset of all that we have.

We can do this with the `filter()` function, part of the `dplyr` package (in the tidyverse collection of packages).

Let's say we want to look at the schools with `div_num` values less than 50.

```
newdf %>% filter(div_num < 50)
```

```
## # A tibble: 800 x 137
```

```
##   sch_id div_num div_name school_num school_name school_name2 type0809
##   <chr>    <dbl> <chr>    <chr>      <chr>      <chr>      <chr>
## 1 001-0~      1 Accomac~ "0070\xa0" NANDUA HIGH <NA>      SCH-HIGH
## 2 001-0~      1 Accomac~ "0080\xa0" CHINCOTEAG~ <NA>      SCH-ELEM
## 3 001-0~      1 Accomac~ "0530\xa0" TANGIER CO~ <NA>      SCH-COMB
## 4 001-0~      1 Accomac~ "0540\xa0" ARCADIA HI~ <NA>      SCH-HIGH
## 5 001-0~      1 Accomac~ "0580\xa0" CHINCOTEAG~ <NA>      SCH-COMB
## 6 001-0~      1 Accomac~ "0590\xa0" PUNGOTEAGU~ <NA>      SCH-ELEM
## 7 001-0~      1 Accomac~ "0600\xa0" KEGOTANK E~ <NA>      SCH-ELEM
## 8 001-0~      1 Accomac~ "0701\xa0" ACCAWMACKE~ <NA>      SCH-ELEM
## 9 001-0~      1 Accomac~ "0702\xa0" METOMPKIN ~ <NA>      SCH-ELEM
## 10 001-0~      1 Accomac~ "0703\xa0" NANDUA MID~ <NA>      SCH-MID
```

```
## # ... with 790 more rows, and 130 more variables: lowgrade_2008 <chr>,
## #   higrade_2008 <chr>, totalFT_2008 <dbl>, total_2008 <dbl>, snp_0809 <chr>,
## #   free_elig_0809 <chr>, free_per_0809 <chr>, red_elig_0809 <chr>,
## #   red_per_0809 <chr>, totalFRL_0809 <dbl>, totalper_0809 <dbl>,
## #   type0910 <chr>, lowgrade_2009 <chr>, higrade_2009 <chr>,
## #   totalFT_2009 <dbl>, total_2009 <dbl>, snp_0910 <chr>, free_elig_0910 <chr>,
## #   free_per_0910 <chr>, red_elig_0910 <chr>, red_per_0910 <chr>,
## #   totalFRL_0910 <dbl>, totalper_0910 <dbl>, type1011 <chr>,
## #   lowgrade_2010 <chr>, higrade_2010 <chr>, totalFT_2010 <dbl>,
## #   total_2010 <dbl>, snp_1011 <chr>, free_elig_1011 <chr>,
## #   free_per_1011 <chr>, red_elig_1011 <chr>, red_per_1011 <chr>,
## #   totalFRL_1011 <dbl>, totalper_1011 <dbl>, type1112 <chr>,
## #   lowgrade_2011 <chr>, higrade_2011 <chr>, totalFT_2011 <dbl>,
## #   total_2011 <dbl>, snp_1112 <chr>, free_elig_1112 <chr>,
## #   free_per_1112 <chr>, red_elig_1112 <chr>, red_per_1112 <chr>,
## #   totalFRL_1112 <dbl>, totalper_1112 <dbl>, type1213 <chr>,
## #   lowgrade_2012 <chr>, higrade_2012 <chr>, totalFT_2012 <dbl>,
## #   total_2012 <dbl>, snp_1213 <chr>, free_elig_1213 <chr>,
## #   free_per_1213 <chr>, red_elig_1213 <chr>, red_per_1213 <chr>,
## #   totalFRL_1213 <dbl>, totalper_1213 <dbl>, type1314 <chr>,
## #   lowgrade_2013 <chr>, higrade_2013 <chr>, totalFT_2013 <dbl>,
## #   total_2013 <dbl>, snp_1314 <chr>, free_elig_1314 <chr>,
## #   free_per_1314 <chr>, red_elig_1314 <chr>, red_per_1314 <chr>,
## #   totalFRL_1314 <dbl>, totalper_1314 <dbl>, type1415 <chr>,
## #   lowgrade_2014 <chr>, higrade_2014 <chr>, totalFT_2014 <dbl>,
## #   total_2014 <dbl>, snp_1415 <chr>, free_elig_1415 <chr>,
## #   free_per_1415 <chr>, red_elig_1415 <chr>, red_per_1415 <chr>,
## #   totalFRL_1415 <dbl>, totalper_1415 <dbl>, CEP_1516 <chr>, type1516 <chr>,
## #   lowgrade_2015 <chr>, higrade_2015 <chr>, totalFT_2015 <dbl>,
## #   total_2015 <dbl>, snp_1516 <chr>, free_elig_1516 <chr>,
## #   free_per_1516 <chr>, red_elig_1516 <chr>, red_per_1516 <chr>,
## #   totalFRL_1516 <dbl>, totalper_1516 <dbl>, CEP_1617 <chr>, type1617 <chr>,
## #   lowgrade_2016 <chr>, higrade_2016 <chr>, ...
```

Or, if we want to look at schools where the highest grade in 2008 was grade five, we can try:

```
newdf %>% filter(higrade_2008 == "5") # this returns a subsetting dataframe with 878 rows
```

```
## # A tibble: 878 x 137
##   sch_id div_num div_name school_num school_name school_name2 type0809
##   <chr>    <dbl> <chr>    <chr>    <chr>    <chr>    <chr>
## 1 001-0~      1 Accomac~ "0080\xa0" CHINCOTEAG~ <NA>    SCH-ELEM
## 2 001-0~      1 Accomac~ "0590\xa0" PUNGOTEAGU~ <NA>    SCH-ELEM
## 3 001-0~      1 Accomac~ "0600\xa0" KEGOTANK E~ <NA>    SCH-ELEM
```

```
## 4 001-0~      1 Accomac~ "0701\xa0" ACCAWMACKE~ <NA>          SCH-ELEM
## 5 001-0~      1 Accomac~ "0702\xa0" METOMPKIN ~ <NA>          SCH-ELEM
## 6 002-0~      2 Albemar~ "0010\xa0" HOLLYMEAD ~ <NA>          SCH-ELEM
## 7 002-0~      2 Albemar~ "0030\xa0" SCOTTSVILL~ <NA>          SCH-ELEM
## 8 002-0~      2 Albemar~ "0040\xa0" MARY CARR ~ <NA>          SCH-ELEM
## 9 002-0~      2 Albemar~ "0100\xa0" BROADUS WO~ <NA>          SCH-ELEM
## 10 002-0~     2 Albemar~ "0150\xa0" PAUL H CAL~ <NA>          SCH-ELEM
## # ... with 868 more rows, and 130 more variables: lowgrade_2008 <chr>,
## #   higrade_2008 <chr>, totalFT_2008 <dbl>, total_2008 <dbl>, snp_0809 <chr>,
## #   free_elig_0809 <chr>, free_per_0809 <chr>, red_elig_0809 <chr>,
## #   red_per_0809 <chr>, totalFRL_0809 <dbl>, totalper_0809 <dbl>,
## #   type0910 <chr>, lowgrade_2009 <chr>, higrade_2009 <chr>,
## #   totalFT_2009 <dbl>, total_2009 <dbl>, snp_0910 <chr>, free_elig_0910 <chr>,
## #   free_per_0910 <chr>, red_elig_0910 <chr>, red_per_0910 <chr>,
## #   totalFRL_0910 <dbl>, totalper_0910 <dbl>, type1011 <chr>,
## #   lowgrade_2010 <chr>, higrade_2010 <chr>, totalFT_2010 <dbl>,
## #   total_2010 <dbl>, snp_1011 <chr>, free_elig_1011 <chr>,
## #   free_per_1011 <chr>, red_elig_1011 <chr>, red_per_1011 <chr>,
## #   totalFRL_1011 <dbl>, totalper_1011 <dbl>, type1112 <chr>,
## #   lowgrade_2011 <chr>, higrade_2011 <chr>, totalFT_2011 <dbl>,
## #   total_2011 <dbl>, snp_1112 <chr>, free_elig_1112 <chr>,
## #   free_per_1112 <chr>, red_elig_1112 <chr>, red_per_1112 <chr>,
## #   totalFRL_1112 <dbl>, totalper_1112 <dbl>, type1213 <chr>,
## #   lowgrade_2012 <chr>, higrade_2012 <chr>, totalFT_2012 <dbl>,
## #   total_2012 <dbl>, snp_1213 <chr>, free_elig_1213 <chr>,
## #   free_per_1213 <chr>, red_elig_1213 <chr>, red_per_1213 <chr>,
## #   totalFRL_1213 <dbl>, totalper_1213 <dbl>, type1314 <chr>,
## #   lowgrade_2013 <chr>, higrade_2013 <chr>, totalFT_2013 <dbl>,
## #   total_2013 <dbl>, snp_1314 <chr>, free_elig_1314 <chr>,
## #   free_per_1314 <chr>, red_elig_1314 <chr>, red_per_1314 <chr>,
## #   totalFRL_1314 <dbl>, totalper_1314 <dbl>, type1415 <chr>,
## #   lowgrade_2014 <chr>, higrade_2014 <chr>, totalFT_2014 <dbl>,
## #   total_2014 <dbl>, snp_1415 <chr>, free_elig_1415 <chr>,
## #   free_per_1415 <chr>, red_elig_1415 <chr>, red_per_1415 <chr>,
## #   totalFRL_1415 <dbl>, totalper_1415 <dbl>, CEP_1516 <chr>, type1516 <chr>,
## #   lowgrade_2015 <chr>, higrade_2015 <chr>, totalFT_2015 <dbl>,
## #   total_2015 <dbl>, snp_1516 <chr>, free_elig_1516 <chr>,
## #   free_per_1516 <chr>, red_elig_1516 <chr>, red_per_1516 <chr>,
## #   totalFRL_1516 <dbl>, totalper_1516 <dbl>, CEP_1617 <chr>, type1617 <chr>,
## #   lowgrade_2016 <chr>, higrade_2016 <chr>, ...
```

Note that we had to set it equal to the character value “5” rather than the numeric value 5. Why?

If we wanted to filter on numeric values instead, we would want to do something like this:

```

newdf %>%
  mutate(higrade_2008 = as.numeric(higrade_2008)) %>%
  filter(higrade_2008 == 5) # again, this returns a subsetting dataframe with 878 rows

## Warning in mask$eval_all_mutate(dots[[i]]): NAs introduced by coercion

## # A tibble: 878 x 137
##   sch_id div_num div_name school_num school_name school_name2 type0809
##   <chr>    <dbl> <chr>    <chr>    <chr>    <chr>    <chr>
## 1 001-0~      1 Accomac~ "0080\xa0" CHINCOTEAG~ <NA>    SCH-ELEM
## 2 001-0~      1 Accomac~ "0590\xa0" PUNGOTEAGU~ <NA>    SCH-ELEM
## 3 001-0~      1 Accomac~ "0600\xa0" KEGOTANK E~ <NA>    SCH-ELEM
## 4 001-0~      1 Accomac~ "0701\xa0" ACCAWMACKE~ <NA>    SCH-ELEM
## 5 001-0~      1 Accomac~ "0702\xa0" METOMPKIN ~ <NA>    SCH-ELEM
## 6 002-0~      2 Albemar~ "0010\xa0" HOLLYMEAD ~ <NA>    SCH-ELEM
## 7 002-0~      2 Albemar~ "0030\xa0" SCOTTSVILL~ <NA>    SCH-ELEM
## 8 002-0~      2 Albemar~ "0040\xa0" MARY CARR ~ <NA>    SCH-ELEM
## 9 002-0~      2 Albemar~ "0100\xa0" BROADUS WO~ <NA>    SCH-ELEM
## 10 002-0~      2 Albemar~ "0150\xa0" PAUL H CAL~ <NA>    SCH-ELEM
## # ... with 868 more rows, and 130 more variables: lowgrade_2008 <chr>,
## #   higrade_2008 <dbl>, totalFT_2008 <dbl>, total_2008 <dbl>, snp_0809 <chr>,
## #   free_elig_0809 <chr>, free_per_0809 <chr>, red_elig_0809 <chr>,
## #   red_per_0809 <chr>, totalFRL_0809 <dbl>, totalper_0809 <dbl>,
## #   type0910 <chr>, lowgrade_2009 <chr>, higrade_2009 <chr>,
## #   totalFT_2009 <dbl>, total_2009 <dbl>, snp_0910 <chr>, free_elig_0910 <chr>,
## #   free_per_0910 <chr>, red_elig_0910 <chr>, red_per_0910 <chr>,
## #   totalFRL_0910 <dbl>, totalper_0910 <dbl>, type1011 <chr>,
## #   lowgrade_2010 <chr>, higrade_2010 <chr>, totalFT_2010 <dbl>,
## #   total_2010 <dbl>, snp_1011 <chr>, free_elig_1011 <chr>,
## #   free_per_1011 <chr>, red_elig_1011 <chr>, red_per_1011 <chr>,
## #   totalFRL_1011 <dbl>, totalper_1011 <dbl>, type1112 <chr>,
## #   lowgrade_2011 <chr>, higrade_2011 <chr>, totalFT_2011 <dbl>,
## #   total_2011 <dbl>, snp_1112 <chr>, free_elig_1112 <chr>,
## #   free_per_1112 <chr>, red_elig_1112 <chr>, red_per_1112 <chr>,
## #   totalFRL_1112 <dbl>, totalper_1112 <dbl>, type1213 <chr>,
## #   lowgrade_2012 <chr>, higrade_2012 <chr>, totalFT_2012 <dbl>,
## #   total_2012 <dbl>, snp_1213 <chr>, free_elig_1213 <chr>,
## #   free_per_1213 <chr>, red_elig_1213 <chr>, red_per_1213 <chr>,
## #   totalFRL_1213 <dbl>, totalper_1213 <dbl>, type1314 <chr>,
## #   lowgrade_2013 <chr>, higrade_2013 <chr>, totalFT_2013 <dbl>,
## #   total_2013 <dbl>, snp_1314 <chr>, free_elig_1314 <chr>,
## #   free_per_1314 <chr>, red_elig_1314 <chr>, red_per_1314 <chr>,
## #   totalFRL_1314 <dbl>, totalper_1314 <dbl>, type1415 <chr>,
## #   lowgrade_2014 <chr>, higrade_2014 <chr>, totalFT_2014 <dbl>,
## #   total_2014 <dbl>, snp_1415 <chr>, free_elig_1415 <chr>,

```

```
## #   free_per_1415 <chr>, red_elig_1415 <chr>, red_per_1415 <chr>,
## #   totalFRL_1415 <dbl>, totalper_1415 <dbl>, CEP_1516 <chr>, type1516 <chr>,
## #   lowgrade_2015 <chr>, higrade_2015 <chr>, totalFT_2015 <dbl>,
## #   total_2015 <dbl>, snp_1516 <chr>, free_elig_1516 <chr>,
## #   free_per_1516 <chr>, red_elig_1516 <chr>, red_Per_1516 <chr>,
## #   totalFRL_1516 <dbl>, totalper_1516 <dbl>, CEP_1617 <chr>, type1617 <chr>,
## #   lowgrade_2016 <chr>, higrade_2016 <chr>, ...
```

3.7 Grouping and Summarizing

Let's shift gears to a different combination of operations...

Let's go ahead and try using tidyverse to narrow to what we want. Imagine we want to see the county level aggregate numbers for FRL in the 2017-2018 school year.

We will start out with our entire data frame and then use pipes (the %>% operator) to work from there. The final result will be stored in our new data frame that we are creating, called `county_level_aggregate`.

First, `select` will pick columns Next, `group_by` and `summarize` work together to get us our aggregate totals.

```
county_level_aggregate <- newdf %>%
  select(div_name, total_2017, totalFRL_1718) %>%
  group_by(div_name) %>%
  summarize(totalstudents = sum(total_2017),
            totalFRL = sum(totalFRL_1718))
```

```
## `summarise()` ungrouping output (override with `.groups` argument)
```

Now, we can compute percentages if we like and we can specify a new column by referring to. One that doesn't exist yet but will after we run this code. We will do this two interchangeable ways.

First, the old school way:

```
county_level_aggregate$percent_FRL <- county_level_aggregate$totalFRL / county_level_aggregate$totalstudents
```

Second, the tidyverse way:

```
county_level_aggregate <- county_level_aggregate %>%
  mutate(percent_frl = totalFRL / totalstudents * 100)
```

Just for fun, let's see how this could have been incorporated into our summarize call

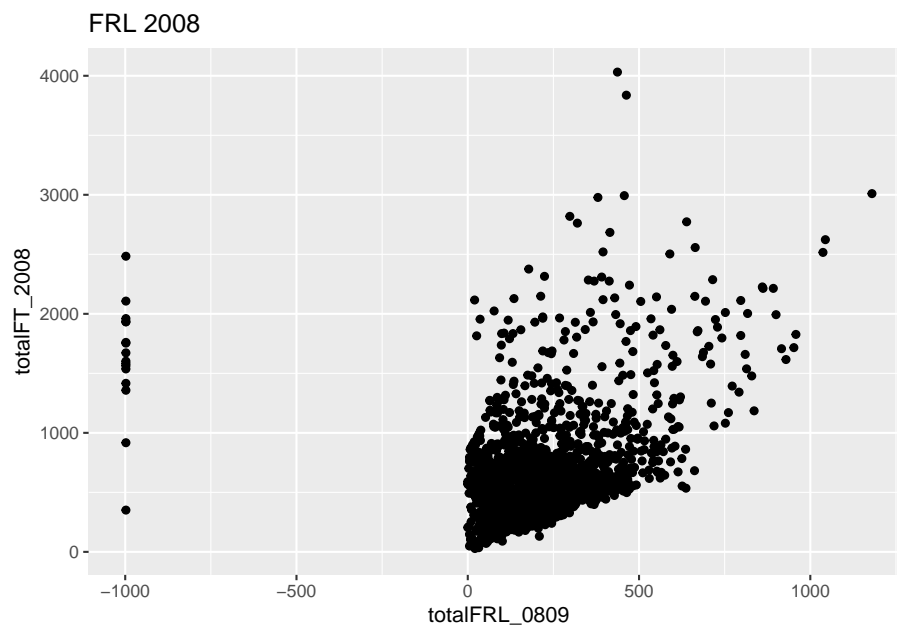
```
county_level_percents <- newdf %>%
  select(div_name, total_2017, totalFRL_1718) %>%
  group_by(div_name) %>%
  summarize(percentFRL=sum(totalFRL_1718)/sum(total_2017) * 100)

## `summarise()` ungrouping output (override with `.groups` argument)
```

Something is going to look weird with this plot

```
newdf %>%
  ggplot(aes(totalFRL_0809, totalFT_2008)) +
  geom_point() +
  labs(title = "FRL 2008", x = "totalFRL_0809")

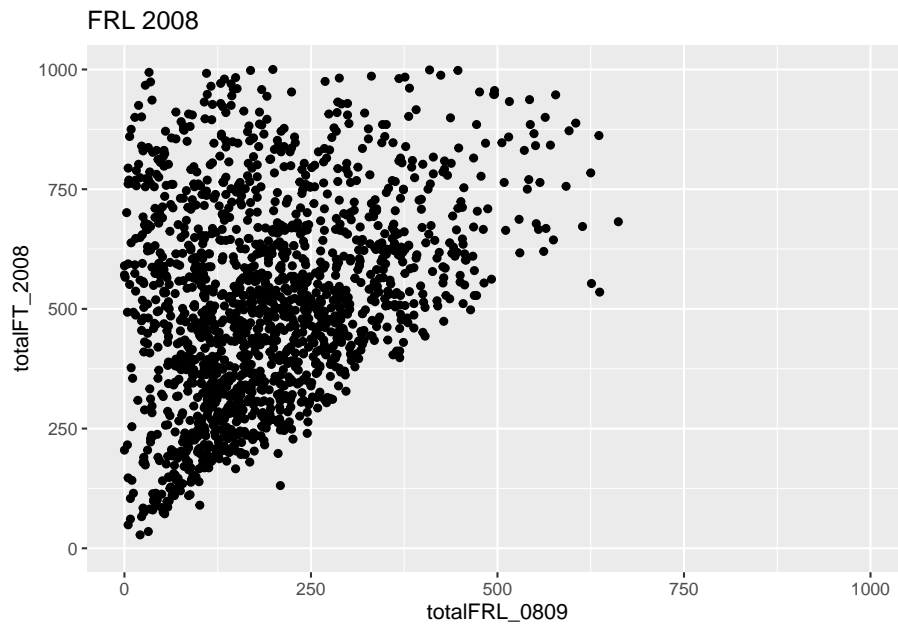
## Warning: Removed 236 rows containing missing values (geom_point).
```



Let's see if we can fix it

```
newdf %>%  
  filter(!is.na(totalFRL_0809)) %>%  
  ggplot(aes(totalFRL_0809, totalFT_2008)) +  
  geom_point() +  
  labs(title = "FRL 2008",  
        x = "totalFRL_0809") +  
  xlim(0, 1000) +  
  ylim(0, 1000)
```

```
## Warning: Removed 324 rows containing missing values (geom_point).
```



Chapter 4

Week 4: Assumptions and Correlations

This week we will be discussing Chapters 5 and 6 from DSUR. These notes will pull out some of the important pieces from each chapter.

4.1 Assumptions

These assumptions that we are making are helpful when determining whether we should be using parametric vs non-parametric statistical tests. What does “parametric” mean here? It means that the data are from a parameterized distribution (i.e., parameters characterize the distribution that the data come from). An example of a parameterized distribution that we have already seen is the normal distribution. The two parameters for the normal distribution are μ for the mean and σ for the standard deviation. We have seen this altogether with this kind of notation to denote that x_i is from a normal distribution:

$$x_i \sim \mathcal{N}(\mu, \sigma^2)$$

4.1.1 Normally distributed data

This assumption is about the normality of the sampling distribution. The big idea here is that we tend to operate under the belief that if our *sampled* data are normally distributed then the underlying *sampling distribution* is also normally distributed. Also, keep in mind that this becomes less of a concern as our sample size increases (thank, Central Limit Theorem!).

There are several tests for normality that we will discuss, which include either (a) calculations or (b) visual examination. We will discuss both.

4.1.1.1 Visual check of normality

You can accomplish this with a histogram (e.g., `hist()` or `geom_histogram()`) or a q-q plot `qqplot()` (which stands for quantile-quantile).

4.1.2 Homogeneity of variance

Here, you want to know whether the variance of a variable is the same across different groups. For example, if you are looking at test scores in chemistry and chemical engineering students, you want to know if the variances (spread) of the test scores in the chemistry group and the chemical engineering group are close to each other.

4.1.3 Interval data

This might be a little redundant given that we want normally distributed data, but you want at least interval data (ratio data are also fine, but in practice very few things we work with actually qualify as ratio variables). If you have ordinal or nominal variables, you might be in trouble with this assumption...

4.1.4 Independence

This assumption is about the observations not being related to each other or affecting each other in some way. In practice, this can also be a little tricky. For example, if you are sampling students from different classrooms, depending on the variables you are measuring, you might actually have reason to believe that students in one classroom are more related to each other than students in a different classroom. In practice, you can handle this with a multi-level model (aka hierarchical model), but that is beyond the scope of this class.

4.2 Correlation

4.2.1 Covariance

First, start with the observation that variance is calculated with: $Variance(s^2) = \frac{\sum(x_i - \bar{x})^2}{N-1} = \frac{\sum(x_i - \bar{x})(x_i - \bar{x})}{N-1}$

But now let's say that we want to know how, for each observation we have, how does the value of x vary with the value of y on average. For example, when the value of x increases, does the value of y also increase? This could happen when x represents the number of hours of sleep you get each night and y is your average grade on an exam you take the next day. The opposite could arise

when x increases but we expect y to go down. An example of this might be when x is the number of hilarious jokes that a teacher tells in class and y is the number of students who fall asleep in class. As the number of jokes increases, we might expect/hope that it keeps students' attention and keeps them from dozing. This generally process of considering how one variable changes when another variable changes is where the notion of covariance comes in.

In practice, what we really want to know is: when x_i is above its average value in a sample (\bar{x}), how does y_i change? Does it also tend to be above the sample average for y (\bar{y})? This is expressed in the general formula for covariance:

$$\text{cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

While covariance can be a helpful value to work with in many settings, for most of what we do in this class, we will be using correlation coefficients instead of covariance. This is because covariance is an unnormalized value, which can make comparisons across different ranges of values difficult.

4.2.2 Correlation coefficient

In order to standardize the covariance to a value that is easier to work with across ranges of values, we use the correlation coefficient. There are several versions of this, depending on the type of data you are working with. The most basic version is the Pearson correlation coefficient. It is calculated by dividing the covariance by the standard deviations of your two variables of interest:

$$r = \frac{\text{cov}_{xy}}{s_x s_y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{(N-1)s_x s_y}$$

This is a *bivariate* correlation coefficient because it is looking at the correlation between *two variables*. There are also partial correlation coefficients, which look at the correlation between two variables while controlling for other variables.

We can calculate the correlation between two variables using the `cor()` or `cor.test()` functions, which are part of base R.

4.3 Another worked example for cleaning and prelim analysis

This script takes an incomplete subset of senior data from a .csv file, cleans it, computes factor scores, and prepares it for analysis.

If you have not already done so, make sure that you have run `library(tidyverse)` and `library(psych)` since we will be using functions from both of those packages.

4.3.1 Loading in data

First, as usual, load in your data. We will use the file `seniorsurvey.csv` for this demo.

`file_path <- "YOUR PATH HERE"` `setwd(file_path)` use this command to change the working directory to the folder where you have your file `list.files()` run this to make sure that your file is in your current working directory

```
seniorSurvey_df <- read_csv("seniorsurvey.csv") # replace text in the parentheses with
```

```
## Parsed with column specification:
## cols(
##   .default = col_double()
## )
```

```
## See spec(...) for full column specifications.
```

4.3.2 Data prep and cleaning

After loading, it is always nice to just see how things loaded in. Functions like `str()` and `describe()` from the `psych` package are nice for this. For example, if we use `describe()`, we can see the following (we deleted some variables):

```
psych::describe(seniorSurvey_df)
```

	vars	n	mean	sd	median
## What is your PRIMARY MAJOR?	1	1849	31.79	20.43	29
## Internship, field experience, co-op, or practicum	2	1121	1.00	0.00	1
## ParticipateServiceL	3	489	1.00	0.00	1
## ParticipateCService	4	1296	1.00	0.00	1
## ParticipateStudyAbroadSemester	5	142	1.00	0.00	1
## SJ1	6	1733	2.32	1.03	2
## SJ2	7	1732	2.08	0.96	2
## SJ3	8	1731	2.77	0.88	3
## SJ4	9	1726	2.27	1.01	2
## SJ5	10	1728	3.27	0.92	3
## SJ6	11	1719	3.50	0.83	4
## SJ7	12	1719	4.01	0.79	4
## SJ8	13	1719	4.15	0.83	4
## DA1	14	1719	2.23	0.93	2
## DA2	15	1719	2.86	0.95	3
## DA3	16	1720	1.97	0.81	2

4.3. ANOTHER WORKED EXAMPLE FOR CLEANING AND PRELIM ANALYSIS61

## DA4	17	1721	4.20	0.72	4
## DA5	18	1721	4.06	0.81	4
## LocalRole	19	1453	3.49	0.94	4
## LocalFinance	20	1453	3.22	0.90	3
## LocalTime	21	1453	3.58	0.89	4
## GlobalRole	22	1446	3.57	0.99	4
## GlobalFinance	23	1449	3.19	1.00	3
## GlobalTime	24	1449	3.42	0.98	3
## Your gender?	25	1678	1.49	0.50	1
##		trimmed	mad	min	max range
## What is your PRIMARY MAJOR?		31.25	28.17	1	70 69
## Internship, field experience, co-op, or practicum		1.00	0.00	1	1 0
## ParticipateServiceL		1.00	0.00	1	1 0
## ParticipateCService		1.00	0.00	1	1 0
## ParticipateStudyAbroadSemester		1.00	0.00	1	1 0
## SJ1		2.25	1.48	1	5 4
## SJ2		1.98	1.48	1	5 4
## SJ3		2.79	1.48	1	5 4
## SJ4		2.20	1.48	1	5 4
## SJ5		3.27	1.48	1	5 4
## SJ6		3.54	1.48	1	5 4
## SJ7		4.07	0.00	1	5 4
## SJ8		4.23	1.48	1	5 4
## DA1		2.15	1.48	1	5 4
## DA2		2.89	1.48	1	5 4
## DA3		1.90	0.00	1	5 4
## DA4		4.27	1.48	1	5 4
## DA5		4.11	1.48	1	5 4
## LocalRole		3.51	1.48	1	5 4
## LocalFinance		3.23	1.48	1	5 4
## LocalTime		3.63	1.48	1	5 4
## GlobalRole		3.62	1.48	1	5 4
## GlobalFinance		3.18	1.48	1	5 4
## GlobalTime		3.44	1.48	1	5 4
## Your gender?		1.49	0.00	1	2 1
##		skew	kurtosis	se	
## What is your PRIMARY MAJOR?		0.18	-1.33	0.48	
## Internship, field experience, co-op, or practicum		NaN	NaN	0.00	
## ParticipateServiceL		NaN	NaN	0.00	
## ParticipateCService		NaN	NaN	0.00	
## ParticipateStudyAbroadSemester		NaN	NaN	0.00	
## SJ1		0.46	-0.51	0.02	
## SJ2		0.64	-0.16	0.02	
## SJ3		-0.03	-0.14	0.02	
## SJ4		0.52	-0.40	0.02	
## SJ5		-0.19	-0.03	0.02	

## SJ6	-0.40	0.15	0.02
## SJ7	-0.75	1.05	0.02
## SJ8	-0.86	0.69	0.02
## DA1	0.57	-0.11	0.02
## DA2	-0.09	-0.41	0.02
## DA3	0.70	0.31	0.02
## DA4	-0.66	0.54	0.02
## DA5	-0.59	0.16	0.02
## LocalRole	-0.44	-0.33	0.02
## LocalFinance	-0.19	-0.27	0.02
## LocalTime	-0.59	0.22	0.02
## GlobalRole	-0.47	-0.30	0.03
## GlobalFinance	-0.11	-0.33	0.03
## GlobalTime	-0.35	-0.17	0.03
## Your gender?	0.03	-2.00	0.01

Upon examining this, we can notice a few things: Primary Major variable is all messed up. We won't fix it here, but basically there is a numeric code needed (e.g., 13 = underwater basket weaving)

Columns 3 and 5 have lots of missing values (note the small N's) – this means that this was asked via checkbox so (1) is true and missing is not missing but False

SJ1-8 and DA1-5 all look essentially ok – about the same N (some survey fatigue or skips) but all values in range (1-5)

Now, we know that SJ and DA are scales from the literature and we want to compute scale scores for those. Typically for attitude scales like these we just report means across the items. So, we will use the “psych” package to use a built in function to help us with this. If you have not used psych yet, be sure it is installed using the command `install.packages(“psych”)` – you need only do this once and then in subsequent uses you only need `library(psych)` to tell R to look in that package for the functions you will be using.

```
library(psych)
```

Subset out only the SJ and DA items in their own dataframe and then use tools in the psych package to compute scale means

The first method to do this - use numbering of the columns:

```
seniorSurveyScales_df <- seniorSurvey_df[6:18]
```

A second method to do this - use `select()` from `dplyr`

4.3. ANOTHER WORKED EXAMPLE FOR CLEANING AND PRELIM ANALYSIS63

```
seniorSurveyScales_df <- seniorSurvey_df %>% select(SJ1:DA5)
```

Use the `make.keys()` function from `psych` package to key-in how the scales are built (mapping items to scales, use - for reverse scored items)

```
my_keys <- make.keys(seniorSurveyScales_df, list(SJCa=c(-1,-2,-3,-4),SJCh=c(5,6,7),DA=c(-9,-10,-11,-12,-13,-14,-15,-16,-17,-18,-19,-20,-21,-22,-23,-24,-25,-26,-27,-28,-29,-30,-31,-32,-33,-34,-35,-36,-37,-38,-39,-40,-41,-42,-43,-44,-45,-46,-47,-48,-49,-50,-51,-52,-53,-54,-55,-56,-57,-58,-59,-60,-61,-62,-63,-64,-65,-66,-67,-68,-69,-70,-71,-72,-73,-74,-75,-76,-77,-78,-79,-80,-81,-82,-83,-84,-85,-86,-87,-88,-89,-90,-91,-92,-93,-94,-95,-96,-97,-98,-99,-100,-101,-102,-103,-104,-105,-106,-107,-108,-109,-110,-111,-112,-113,-114,-115,-116,-117,-118,-119,-120,-121,-122,-123,-124,-125,-126,-127,-128,-129,-130,-131,-132,-133,-134,-135,-136,-137,-138,-139,-140,-141,-142,-143,-144,-145,-146,-147,-148,-149,-150,-151,-152,-153,-154,-155,-156,-157,-158,-159,-160,-161,-162,-163,-164,-165,-166,-167,-168,-169,-170,-171,-172,-173,-174,-175,-176,-177,-178,-179,-180,-181,-182,-183,-184,-185,-186,-187,-188,-189,-190,-191,-192,-193,-194,-195,-196,-197,-198,-199,-200,-201,-202,-203,-204,-205,-206,-207,-208,-209,-210,-211,-212,-213,-214,-215,-216,-217,-218,-219,-220,-221,-222,-223,-224,-225,-226,-227,-228,-229,-230,-231,-232,-233,-234,-235,-236,-237,-238,-239,-240,-241,-242,-243,-244,-245,-246,-247,-248,-249,-250,-251,-252,-253,-254,-255,-256,-257,-258,-259,-260,-261,-262,-263,-264,-265,-266,-267,-268,-269,-270,-271,-272,-273,-274,-275,-276,-277,-278,-279,-280,-281,-282,-283,-284,-285,-286,-287,-288,-289,-290,-291,-292,-293,-294,-295,-296,-297,-298,-299,-300,-301,-302,-303,-304,-305,-306,-307,-308,-309,-310,-311,-312,-313,-314,-315,-316,-317,-318,-319,-320,-321,-322,-323,-324,-325,-326,-327,-328,-329,-330,-331,-332,-333,-334,-335,-336,-337,-338,-339,-340,-341,-342,-343,-344,-345,-346,-347,-348,-349,-350,-351,-352,-353,-354,-355,-356,-357,-358,-359,-360,-361,-362,-363,-364,-365,-366,-367,-368,-369,-370,-371,-372,-373,-374,-375,-376,-377,-378,-379,-380,-381,-382,-383,-384,-385,-386,-387,-388,-389,-390,-391,-392,-393,-394,-395,-396,-397,-398,-399,-400,-401,-402,-403,-404,-405,-406,-407,-408,-409,-410,-411,-412,-413,-414,-415,-416,-417,-418,-419,-420,-421,-422,-423,-424,-425,-426,-427,-428,-429,-430,-431,-432,-433,-434,-435,-436,-437,-438,-439,-440,-441,-442,-443,-444,-445,-446,-447,-448,-449,-450,-451,-452,-453,-454,-455,-456,-457,-458,-459,-460,-461,-462,-463,-464,-465,-466,-467,-468,-469,-470,-471,-472,-473,-474,-475,-476,-477,-478,-479,-480,-481,-482,-483,-484,-485,-486,-487,-488,-489,-490,-491,-492,-493,-494,-495,-496,-497,-498,-499,-500,-501,-502,-503,-504,-505,-506,-507,-508,-509,-510,-511,-512,-513,-514,-515,-516,-517,-518,-519,-520,-521,-522,-523,-524,-525,-526,-527,-528,-529,-530,-531,-532,-533,-534,-535,-536,-537,-538,-539,-540,-541,-542,-543,-544,-545,-546,-547,-548,-549,-550,-551,-552,-553,-554,-555,-556,-557,-558,-559,-560,-561,-562,-563,-564,-565,-566,-567,-568,-569,-570,-571,-572,-573,-574,-575,-576,-577,-578,-579,-580,-581,-582,-583,-584,-585,-586,-587,-588,-589,-590,-591,-592,-593,-594,-595,-596,-597,-598,-599,-600,-601,-602,-603,-604,-605,-606,-607,-608,-609,-610,-611,-612,-613,-614,-615,-616,-617,-618,-619,-620,-621,-622,-623,-624,-625,-626,-627,-628,-629,-630,-631,-632,-633,-634,-635,-636,-637,-638,-639,-640,-641,-642,-643,-644,-645,-646,-647,-648,-649,-650,-651,-652,-653,-654,-655,-656,-657,-658,-659,-660,-661,-662,-663,-664,-665,-666,-667,-668,-669,-670,-671,-672,-673,-674,-675,-676,-677,-678,-679,-680,-681,-682,-683,-684,-685,-686,-687,-688,-689,-690,-691,-692,-693,-694,-695,-696,-697,-698,-699,-700,-701,-702,-703,-704,-705,-706,-707,-708,-709,-710,-711,-712,-713,-714,-715,-716,-717,-718,-719,-720,-721,-722,-723,-724,-725,-726,-727,-728,-729,-730,-731,-732,-733,-734,-735,-736,-737,-738,-739,-740,-741,-742,-743,-744,-745,-746,-747,-748,-749,-750,-751,-752,-753,-754,-755,-756,-757,-758,-759,-760,-761,-762,-763,-764,-765,-766,-767,-768,-769,-770,-771,-772,-773,-774,-775,-776,-777,-778,-779,-780,-781,-782,-783,-784,-785,-786,-787,-788,-789,-790,-791,-792,-793,-794,-795,-796,-797,-798,-799,-800,-801,-802,-803,-804,-805,-806,-807,-808,-809,-810,-811,-812,-813,-814,-815,-816,-817,-818,-819,-820,-821,-822,-823,-824,-825,-826,-827,-828,-829,-830,-831,-832,-833,-834,-835,-836,-837,-838,-839,-840,-841,-842,-843,-844,-845,-846,-847,-848,-849,-850,-851,-852,-853,-854,-855,-856,-857,-858,-859,-860,-861,-862,-863,-864,-865,-866,-867,-868,-869,-870,-871,-872,-873,-874,-875,-876,-877,-878,-879,-880,-881,-882,-883,-884,-885,-886,-887,-888,-889,-890,-891,-892,-893,-894,-895,-896,-897,-898,-899,-900,-901,-902,-903,-904,-905,-906,-907,-908,-909,-910,-911,-912,-913,-914,-915,-916,-917,-918,-919,-920,-921,-922,-923,-924,-925,-926,-927,-928,-929,-930,-931,-932,-933,-934,-935,-936,-937,-938,-939,-940,-941,-942,-943,-944,-945,-946,-947,-948,-949,-950,-951,-952,-953,-954,-955,-956,-957,-958,-959,-960,-961,-962,-963,-964,-965,-966,-967,-968,-969,-970,-971,-972,-973,-974,-975,-976,-977,-978,-979,-980,-981,-982,-983,-984,-985,-986,-987,-988,-989,-990,-991,-992,-993,-994,-995,-996,-997,-998,-999,1000))
```

Use `scoreItems` function to score each respondent on the three scales of interest SJCa, SJCh, and DA – the default here in `scoreItems` is to takes the mean of the items (not additive though that is sometimes used) and also, it imputes missing values instead of dropping cases the `scoreItems` function calculates many things. At this stage, all we really want are the scores, so we include a line to only extract that info.

```
my_scales <- scoreItems(my_keys, seniorSurveyScales_df)
my_scores <- my_scales$scores
```

Now, if you view the first few rows of the `my.scores` vector using the `header – head()` command – it looks like we expect:

```
head(my_scores)
```

```
##      SJCa      SJCh  DA
## [1,] 2.75 3.000000 3.2
## [2,] 3.75 3.333333 4.2
## [3,] 3.00 3.000000 3.0
## [4,] 2.25 4.333333 3.6
## [5,] 3.00 3.333333 3.4
## [6,] 4.50 4.333333 3.4
```

Now, let's build a clean dataframe to prep for analysis - by clean in this case I mean that we have replaced item scores from the scales with their means and also that we have fixed the NAs that don't belong (for participation variables, in this dataset, the NAs should be 0s)

```
my_df <- data.frame(seniorSurvey_df[1:5],my_scores, seniorSurvey_df[19:25])
```

This is an old school method to replace NAs in specific columns

```
my_df$ParticipateServiceL[is.na(my_df$ParticipateServiceL)] <- 0
my_df$ParticipateCSservice[is.na(my_df$ParticipateCSservice)] <- 0
my_df$ParticipateStudyAbroadSemester[is.na(my_df$ParticipateStudyAbroadSemester)] <- 0
```

`my_dfParticipateInternCoop...[is.na(mydfParticipateInternCoop...)] <- 0` —
 — this variable read in clumsily named and I don't care about it right now so I'll skip

Here is An alternative method to replace NAs in specific columns:

```
my_df <- my_df %>%
  replace_na(list(ParticipateCSservice = 0, ParticipateStudyAbroadSemester = 0, ParticipateInternCoop = 0))
```

4.3.3 Preliminary analysis

At this point, we are ready for some analysis

Let's investigate correlations. What seems most obvious would just be to run `cor()` but, as we found out in class, this can cause us to run full speed ahead without considering assumptions

```
my_correlations <- my_df %>% select(SJCa,SJCh,DA) %>% cor()
print(my_correlations)
```

```
##           SJCa      SJCh      DA
## SJCa 1.0000000 0.2590211 0.3342276
## SJCh 0.2590211 1.0000000 0.2310703
## DA   0.3342276 0.2310703 1.0000000
```

Ok, so, it is important that we note that this ran correlations but R doesn't know that this was sample data and therefore that we are interested in statistical significance (or not) of these results AND that our data may need another method (e.g., non-parametric). `cor()` does have a way to run spearman instead.

```
my_spearman_correlations <- my_df %>% select(SJCa,SJCh,DA) %>% cor(method="spearman")
print(my_spearman_correlations)
```

```
##           SJCa      SJCh      DA
## SJCa 1.0000000 0.2727828 0.3148168
## SJCh 0.2727828 1.0000000 0.2340174
## DA   0.3148168 0.2340174 1.0000000
```

If we need p values though, we need to change to something else – `corr.test`

```
my_results <- corr.test(my_df$SJCa,my_df$DA)
```

Then we can pull out results from this list or print it. Let's do both.

4.3. ANOTHER WORKED EXAMPLE FOR CLEANING AND PRELIM ANALYSIS 65

```
print(my_results,short=FALSE)
```

```
## Call:corr.test(x = my_df$SJCa, y = my_df$DA)
## Correlation matrix
## [1] 0.33
## Sample Size
## [1] 1852
## Probability values  adjusted for multiple tests.
## [1] 0
##
## Confidence intervals based upon normal theory. To get bootstrapped values, try cor.ci
##      raw.lower raw.r raw.upper raw.p lower.adj upper.adj
## NA-NA      0.29 0.33      0.37    0      0.29      0.37
```

```
my_results$r # correlation coefficient
```

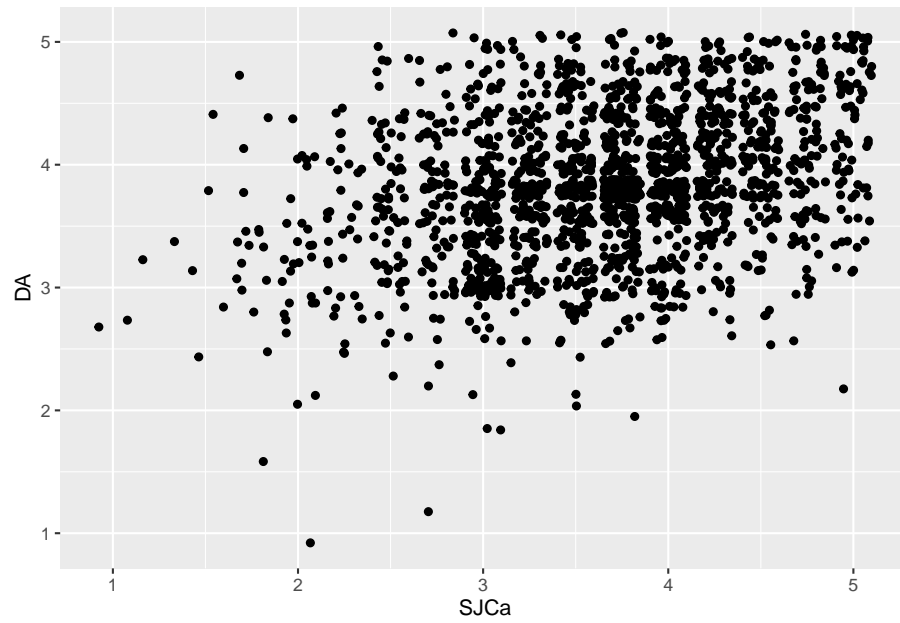
```
## [1] 0.3342276
```

```
my_results$p # p-value
```

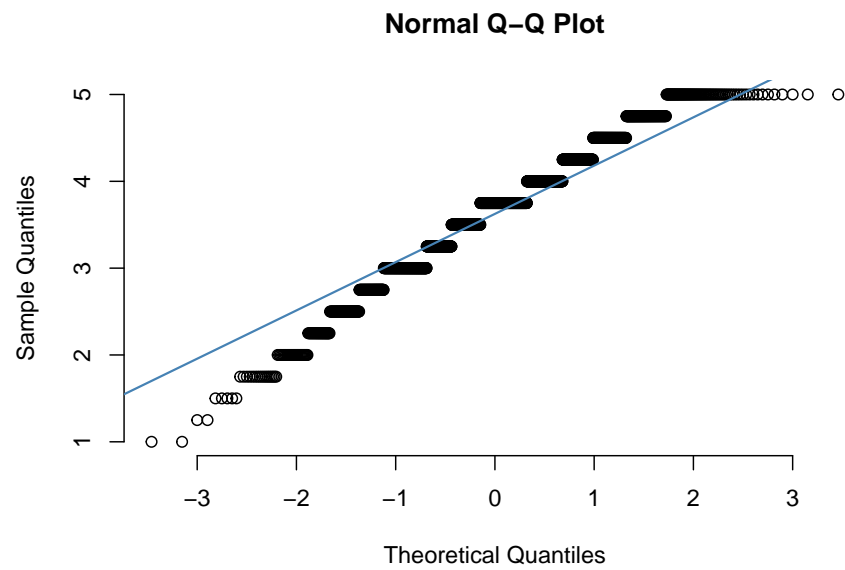
```
## [1] 1.433637e-49
```

Visually, we should be able to see this on a scatterplot. We are going to use `qplot` which stands for quickplot from within `ggplot`. It is useful and quicker for simple plotting than building up `ggplot` (though from the same package) we need to jitter my points (take `geom="jitter"` out if you want to see why)

```
qplot(SJCa,DA,data=my_df,geom="jitter")
```



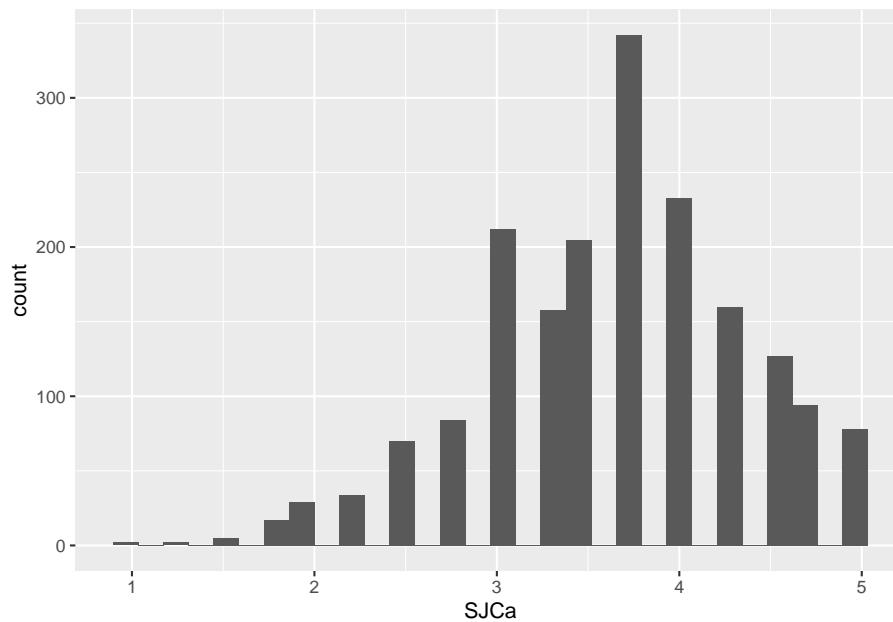
```
qqnorm(my_df$SJCa, frame = FALSE)
qqline(my_df$SJCa, col = "steelblue", lwd = 1.5)
```



4.3. ANOTHER WORKED EXAMPLE FOR CLEANING AND PRELIM ANALYSIS67

```
my_df %>% ggplot(aes(x = SJCa)) +  
  geom_histogram()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Other functions we used today in class were `describe()` and also the q-q plot creation to investigate normality assumption copying syntax from the Field, Miles, & Field book

Chapter 5

Week 5: Simple Regression

```
require(tidyverse)
require(psych)
require(kableExtra)

## Loading required package: kableExtra

##
## Attaching package: 'kableExtra'

## The following object is masked from 'package:dplyr':
##
##      group_rows

library(broom)
```

This week we will start learning about linear regression. In particular, we focus on simple regression. These kinds of models involve one predictor variable and one continuous outcome variable. Next week we will move to models with multiple regression, which involves - you guessed it - multiple predictor variables.

5.1 General Modeling Philosophy

The general approach is to model the outcome variable as a function of some predictor(s) plus an error term. Mathematically, this looks like:

$$outcome_i = model + error_i$$

where the i subscript refers to the i^{th} person in the sample.

5.1.0.1 Review of the normal distribution and standardizing variables

Just to review, let's think about normally distributed variables and the notion of centering and standardizing.

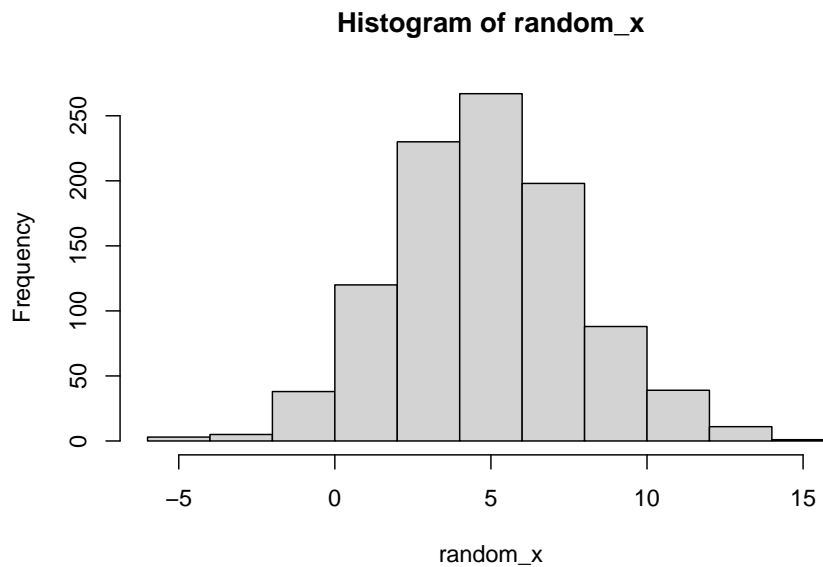
First, we will generate some data by drawing n random numbers from a normal distribution with a mean and standard deviation that we will specify.

```
mean <- 5
sd <- 3
n <- 1000

random_x <- rnorm(n = n, mean = mean, sd = sd)
```

We can then visualize those numbers

```
hist(random_x)
```

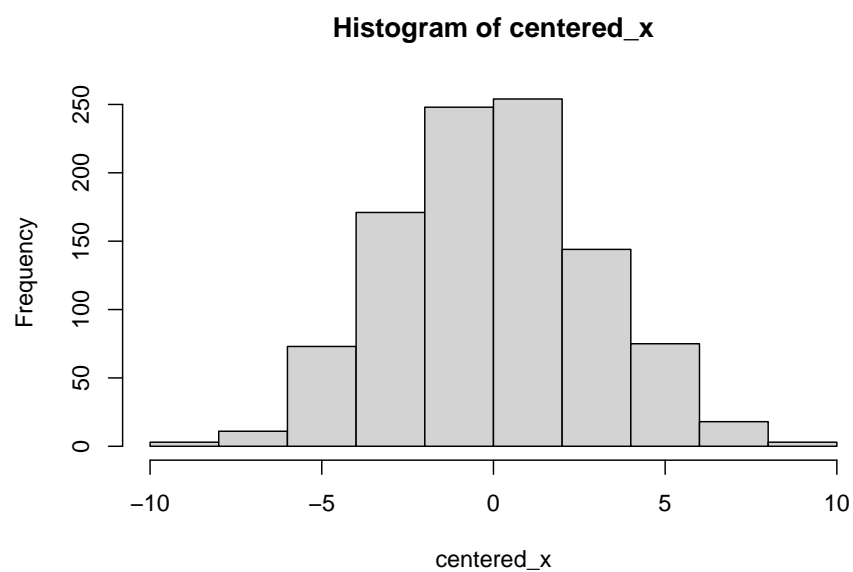


Now if we subtract the mean and plot the histogram, notice how the values have all basically shifted to the left along the x-axis.

```
sample_mean <- mean(random_x)

centered_x <- random_x - sample_mean
```

```
hist(centered_x)
```

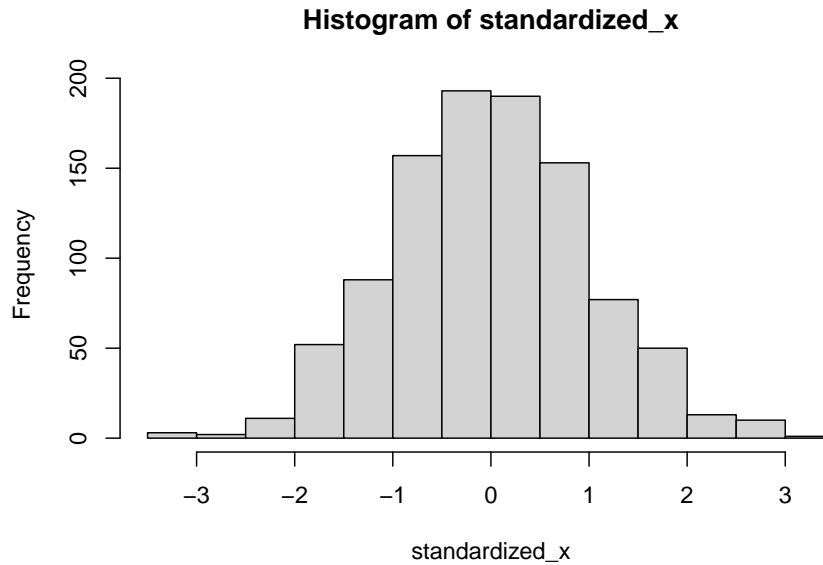


Finally, we can divide by the sample standard deviation, which should have the effect of either stretching or squishing the values along the x-axis (without changing their mean). Pay attention again to the values along the x-axis.

```
sample_sd <- sd(random_x)

standardized_x <- centered_x / sample_sd

hist(standardized_x)
```



This final plot should remind you have the standard normal plot (with mean 0 and standard deviation 1). This is noted as $x \sim \mathcal{N}(0, 1)$ and is read as “x is distributed according to a normal distribution with a mean of 0 and variance of 1”.

5.2 Data generation demo - one set sample size

The following is a demo from class, found in the week_5_demo.R file

```
#store the sample size that we want to use
samp_size <- 100
```

```
# uniformly sample X values (values for our predictor variable) from 0 to 20
x <- round(runif(n = samp_size, min = 0, max = 30), digits = 1) # this gives samp_size
```

Store the noise values for our different test models

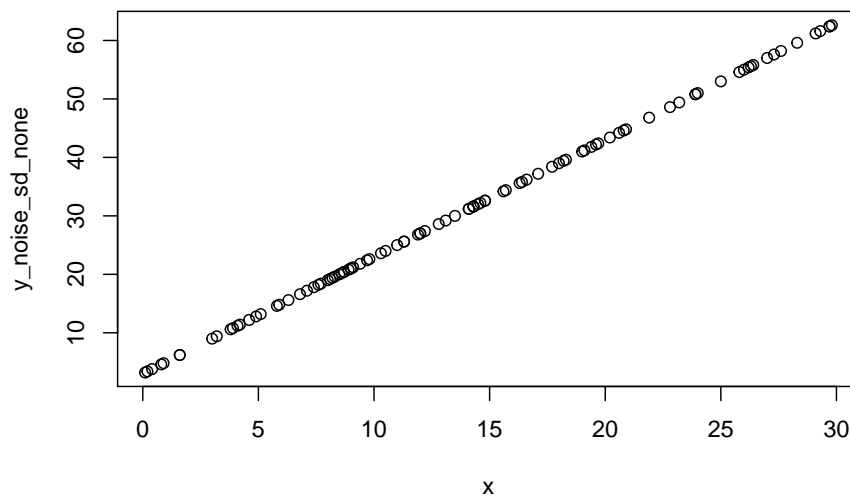
```
sd_min <- 2 # low noise
sd_med <- 6 # medium noise
sd_max <- 12 # high noise
```

Generate the outcome variable values under different amounts of noise (the rnorm() function is what is generating noise here)

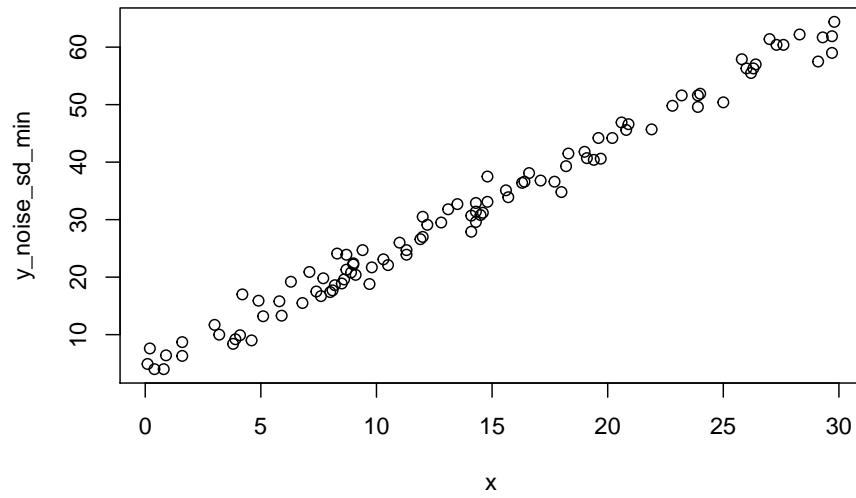

```
y_noise_sd_none <- 3 + 2*x # this is the true relationship without any noise  
y_noise_sd_min <- 3 + 2*x + round(x = rnorm(n = samp_size, mean = 0, sd = sd_min), digits = 1)  
y_noise_sd_med <- 3 + 2*x + round(x = rnorm(n = samp_size, mean = 0, sd = sd_med), digits = 1)  
y_noise_sd_max <- 3 + 2*x + round(x = rnorm(n = samp_size, mean = 0, sd = sd_max), digits = 1)
```

Typical step 1: visualize! Let's plot each of these x values vs y

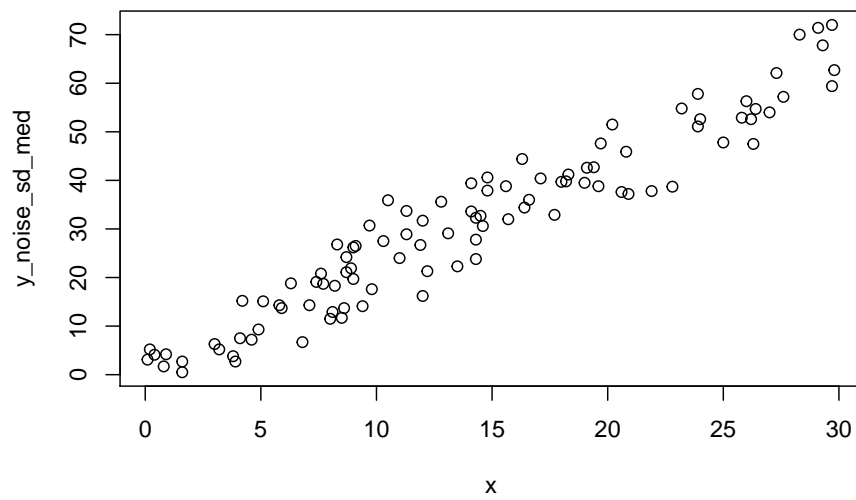
```
plot(x, y_noise_sd_none)
```



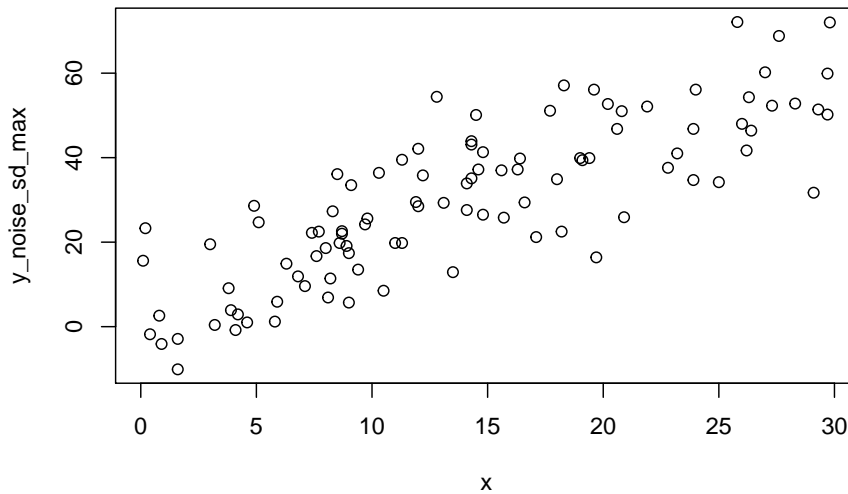
```
plot(x, y_noise_sd_min)
```



```
plot(x, y_noise_sd_med)
```



```
plot(x, y_noise_sd_max)
```



Let's put all of these vectors together into a data frame to make it easier to analyze later on. Note, this is not a vital step for conducting the simple regression

```
demo_df <- tibble("x" = x,
                  "y_noise_sd_none" = y_noise_sd_none,
                  "y_noise_sd_min" = y_noise_sd_min,
                  "y_noise_sd_med" = y_noise_sd_med,
                  "y_noise_sd_max" = y_noise_sd_max)
```

Check out what demo_df looks like

```
head(demo_df)
```

```
## # A tibble: 6 x 5
##       x y_noise_sd_none y_noise_sd_min y_noise_sd_med y_noise_sd_max
##   <dbl>         <dbl>         <dbl>         <dbl>         <dbl>
## 1  8.5           20           18.9           11.7           36.1
## 2 11.9           26.8           26.6           26.7           29.5
## 3 14.3           31.6           31.4           23.8           43.9
## 4  4.9           12.8           15.9            9.3           28.6
## 5 15.7           34.4           33.9            32           25.8
## 6  9.7           22.4           18.8           30.7           24.2
```

Order by increasing x value

```
demo_df <- demo_df %>%
  arrange(x)
```

Check out what the arrange() function did

```
head(demo_df)
```

```
## # A tibble: 6 x 5
##       x y_noise_sd_none y_noise_sd_min y_noise_sd_med y_noise_sd_max
##   <dbl>      <dbl>      <dbl>      <dbl>      <dbl>
## 1  0.1          3.2          4.9          3.1         15.6
## 2  0.2          3.4          7.6          5.2         23.3
## 3  0.4          3.8           4          4.1          -1.8
## 4  0.8          4.6          4.00         1.70          2.60
## 5  0.9          4.8          6.4          4.2          -4.1
## 6  1.6          6.2          8.7          2.7          -2.90
```

Let's make this a long df so that we can plot multiple standard deviation values together

```
demo_df_long <- demo_df %>%
  pivot_longer(cols = starts_with("y_noise"),
               names_to = "y_col",
               values_to = "y_val"
  )
```

Again, check on what this did

```
head(demo_df_long)
```

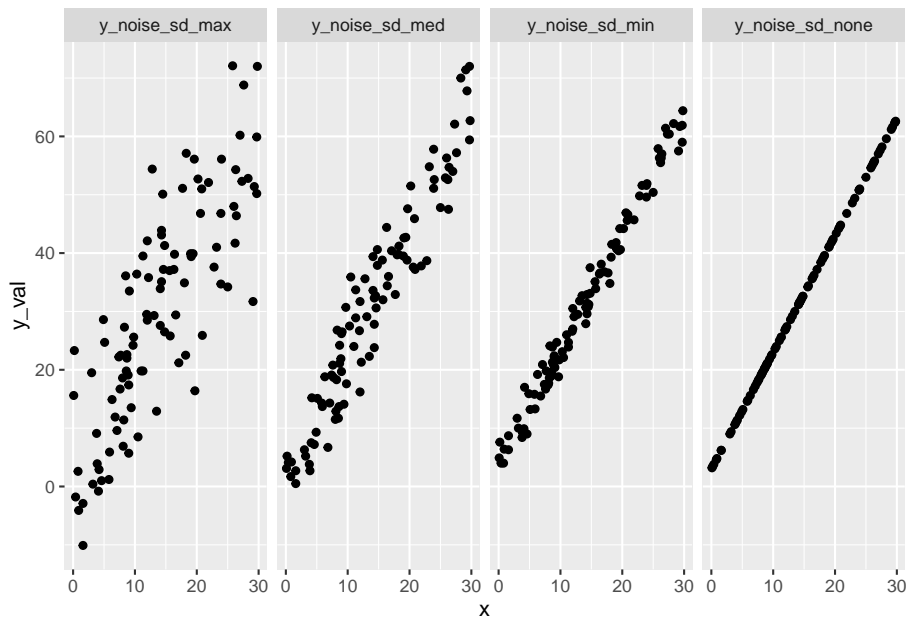
```
## # A tibble: 6 x 3
##       x y_col          y_val
##   <dbl> <chr>      <dbl>
## 1  0.1 y_noise_sd_none  3.2
## 2  0.1 y_noise_sd_min  4.9
## 3  0.1 y_noise_sd_med  3.1
## 4  0.1 y_noise_sd_max 15.6
## 5  0.2 y_noise_sd_none  3.4
## 6  0.2 y_noise_sd_min  7.6
```

Let's add in a column to note whether the value is from the min, med, max, or zero sd (noise) model

```
demo_df_long <- demo_df_long %>%
  mutate(sd_val = case_when(str_detect(y_col, "sd_none") ~ 0,
                             str_detect(y_col, "sd_min") ~ sd_min,
                             str_detect(y_col, "sd_med") ~ sd_med,
                             str_detect(y_col, "sd_max") ~ sd_max))
```

Use `facet_grid` to separate the plots out by

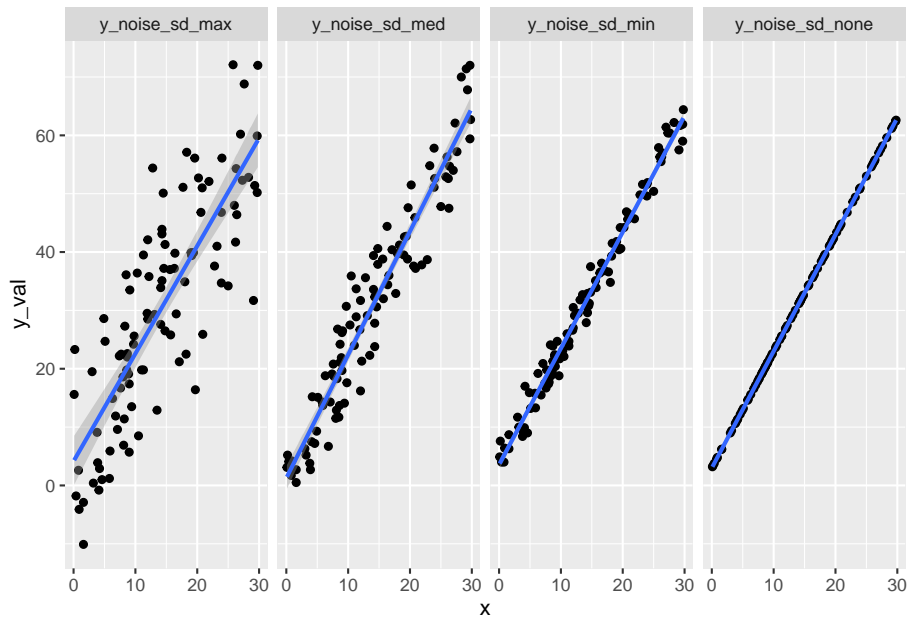
```
demo_df_long %>%
  ggplot(aes(x = x, y = y_val)) +
  geom_point() +
  facet_grid(.~y_col)
```



You can also automatically add in a line with the `geom_smooth()` function

```
demo_df_long %>%
  ggplot(aes(x = x, y = y_val)) +
  geom_point() +
  geom_smooth(method = "lm") +
  facet_grid(.~y_col)
```

```
## `geom_smooth()` using formula 'y ~ x'
```



Now we can create a linear model for the data with minimum noise with the following command:

```
fit_demo_min <- lm(y_noise_sd_min ~ x)
```

...and we can look at the summary of the model with:

```
summary(fit_demo_min)
```

```
##
## Call:
## lm(formula = y_noise_sd_min ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.6315 -1.4674  0.0561  1.4041  5.0728
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.5564     0.4064   8.751 6.19e-14 ***
## x             1.9931     0.0250  79.726 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## Residual standard error: 2.053 on 98 degrees of freedom
## Multiple R-squared:  0.9848, Adjusted R-squared:  0.9847
## F-statistic: 6356 on 1 and 98 DF,  p-value: < 2.2e-16
```

We can also look at model results with the `glance()` function from the `broom` package

```
broom::glance(fit_demo_min)
```

```
## # A tibble: 1 x 12
##   r.squared adj.r.squared sigma statistic p.value    df logLik   AIC   BIC
##   <dbl>      <dbl> <dbl>    <dbl>    <dbl> <dbl> <dbl> <dbl> <dbl>
## 1    0.985        0.985  2.05    6356. 6.25e-91     1 -213.  432.  439.
## # ... with 3 more variables: deviance <dbl>, df.residual <int>, nobs <int>
```

We can create models for the med and max sd values as well and take a look at those with the `summary()` function once again

```
fit_demo_med <- lm(y_noise_sd_med ~ x)
summary(fit_demo_med)
```

```
##
## Call:
## lm(formula = y_noise_sd_med ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.7657  -3.8200   0.1044   3.8614  12.4352
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.26899    1.03129    1.23   0.221
## x            2.11389    0.06344   33.32 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.209 on 98 degrees of freedom
## Multiple R-squared:  0.9189, Adjusted R-squared:  0.9181
## F-statistic: 1110 on 1 and 98 DF,  p-value: < 2.2e-16
```

```
fit_demo_max <- lm(y_noise_sd_max ~ x)
summary(fit_demo_max)
```

```
##
## Call:
## lm(formula = y_noise_sd_max ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -26.1504  -7.7075  -0.2646   7.7050  26.6308
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   4.1472     2.1243   1.952  0.0538 .
## x             1.8455     0.1307  14.123 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 10.73 on 98 degrees of freedom
## Multiple R-squared:  0.6706, Adjusted R-squared:  0.6672
## F-statistic: 199.5 on 1 and 98 DF,  p-value: < 2.2e-16
```

Notice the increase in the standard error of the coefficient estimates as the noise in y values went up

From a programming perspective, this was not very efficient because I just copied, pasted, and corrected these values. There is a better way to do this using lists (see below)

Let's do some fancy stuff to make multiple models at once rather than having to write new lines for each model *Some of these ideas are taken from the R4DS book chapter 25

```
test_nest <- demo_df_long %>% nest(data = -sd_val)
```

```
linear_model <- function(df) {
  lm(y_val ~ x, data = df)
}
```

```
models <- map(test_nest$data, linear_model)
```

```
summary(models[[2]])
```

```
##
## Call:
## lm(formula = y_val ~ x, data = df)
##
```



```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.6315 -1.4674  0.0561  1.4041  5.0728
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.5564     0.4064   8.751 6.19e-14 ***
## x             1.9931     0.0250  79.726 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.053 on 98 degrees of freedom
## Multiple R-squared:  0.9848, Adjusted R-squared:  0.9847
## F-statistic: 6356 on 1 and 98 DF,  p-value: < 2.2e-16
```

```
summary(models[[3]])
```

```
##
## Call:
## lm(formula = y_val ~ x, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.7657 -3.8200  0.1044  3.8614  12.4352
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   1.26899     1.03129   1.23   0.221
## x             2.11389     0.06344  33.32 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.209 on 98 degrees of freedom
## Multiple R-squared:  0.9189, Adjusted R-squared:  0.9181
## F-statistic: 1110 on 1 and 98 DF,  p-value: < 2.2e-16
```

```
summary(models[[4]])
```

```
##
## Call:
## lm(formula = y_val ~ x, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -26.1504 -7.7075 -0.2646 7.7050 26.6308
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.1472     2.1243   1.952  0.0538 .
## x            1.8455     0.1307  14.123 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 10.73 on 98 degrees of freedom
## Multiple R-squared:  0.6706, Adjusted R-squared:  0.6672
## F-statistic: 199.5 on 1 and 98 DF, p-value: < 2.2e-16
```

We can also store the models as new columns in the nested dataframe

```
test_nest <- test_nest %>%
  mutate(model = map(data, linear_model))
```

Finally, we can unnest the models to make it easier to compare them with each other in a data frame

```
test_nest <- test_nest %>%
  mutate(glance = map(model, broom::glance)) %>%
  unnest(glance)
```

```
## Warning in summary.lm(x): essentially perfect fit: summary may be unreliable
```

```
## Warning in summary.lm(x): essentially perfect fit: summary may be unreliable
```

5.3 Data generation demo - one set sample size;

The change from the past demo is that we are now sampling from integer values rather than continuous for the predictor

Store the sample size that we want to use

```
samp_size <- 200
```

Instead of sampling uniformly from 0 to 20, this is to sample integers from 40 to 100 uniformly. We take “samp_size” number of samples. Replace = TRUE means we can get the same x value multiple times

```
x <- sample(x = c(60:100), size = samp_size, replace = TRUE)
```

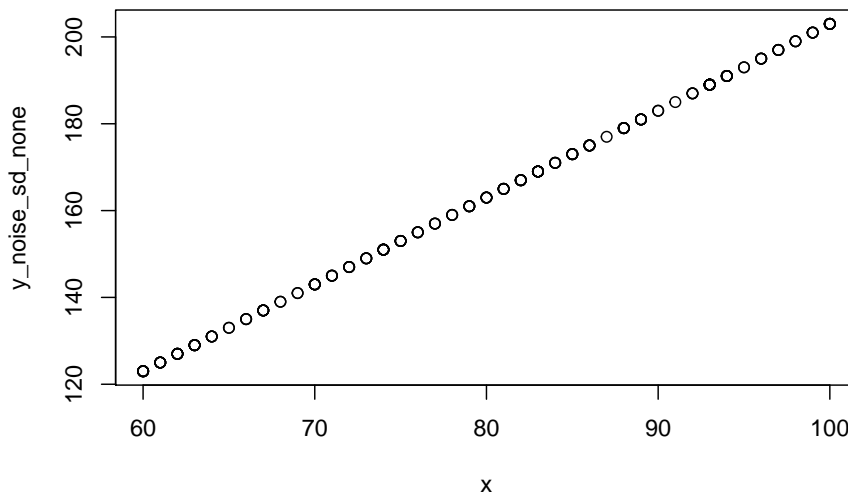
As before, store the noise values for our different test models

```
sd_min <- 2
sd_med <- 6
sd_max <- 12

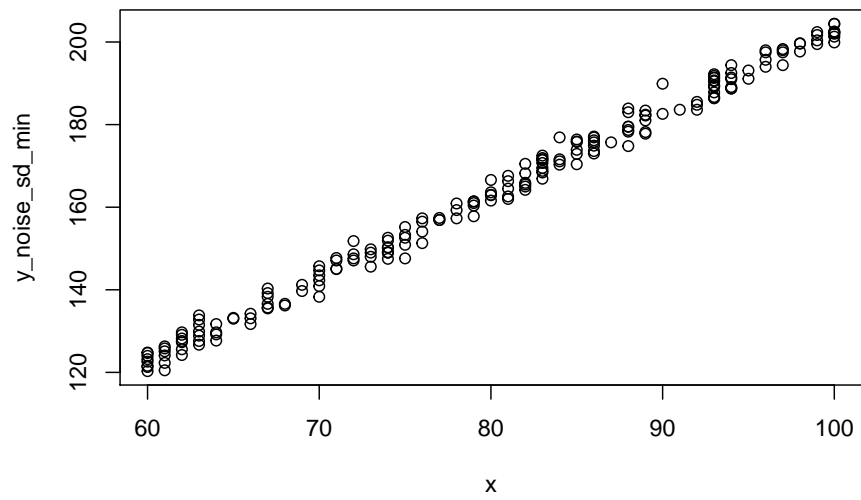
y_noise_sd_none <- 3 + 2*x
y_noise_sd_min <- 3 + 2*x + round(x = rnorm(n = samp_size, mean = 0, sd = sd_min), digits = 1)
y_noise_sd_med <- 3 + 2*x + round(x = rnorm(n = samp_size, mean = 0, sd = sd_med), digits = 1)
y_noise_sd_max <- 3 + 2*x + round(x = rnorm(n = samp_size, mean = 0, sd = sd_max), digits = 1)
```

Typical step 1: visualize! Let's plot each of these x values vs y

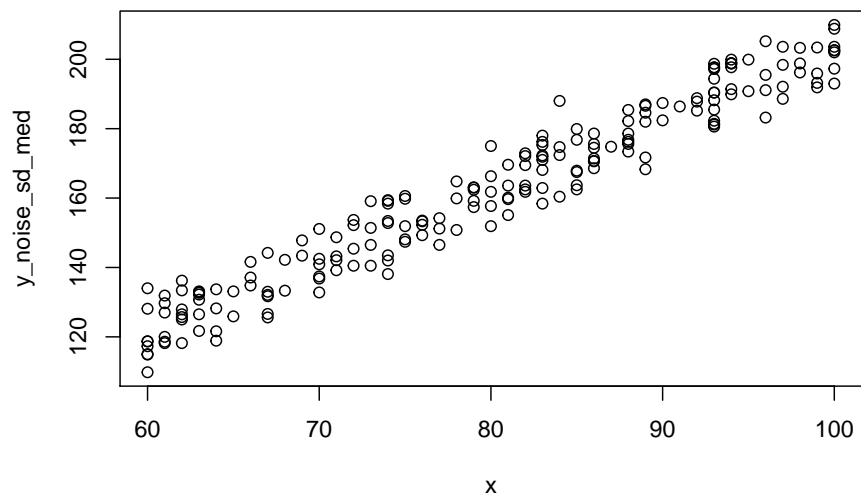
```
plot(x, y_noise_sd_none)
```



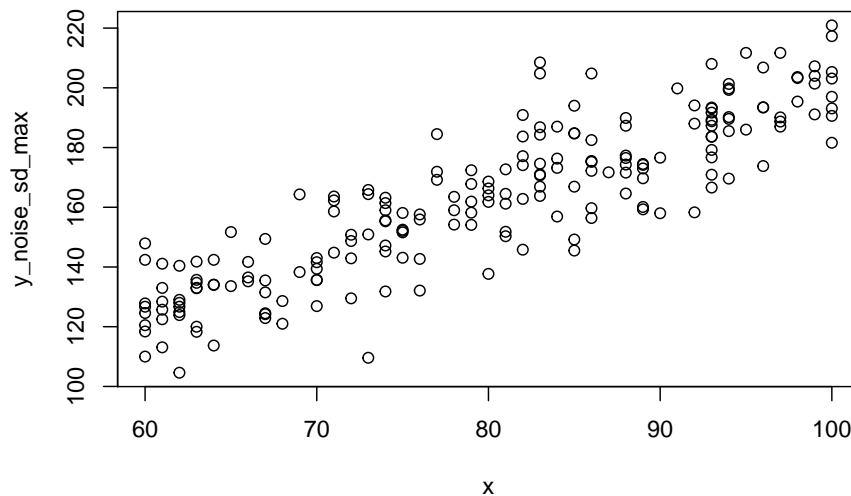
```
plot(x, y_noise_sd_min)
```



```
plot(x, y_noise_sd_med)
```



```
plot(x, y_noise_sd_max)
```



Let's put all of these vectors together into a data frame to make it easier to analyze later on. Note, this is not a vital step for conducting the simple regression

```
demo_df <- tibble("x" = x,
                  "y_noise_sd_none" = y_noise_sd_none,
                  "y_noise_sd_min" = y_noise_sd_min,
                  "y_noise_sd_med" = y_noise_sd_med,
                  "y_noise_sd_max" = y_noise_sd_max)
```

Order by increasing x value

```
demo_df <- demo_df %>%
  arrange(x)
```

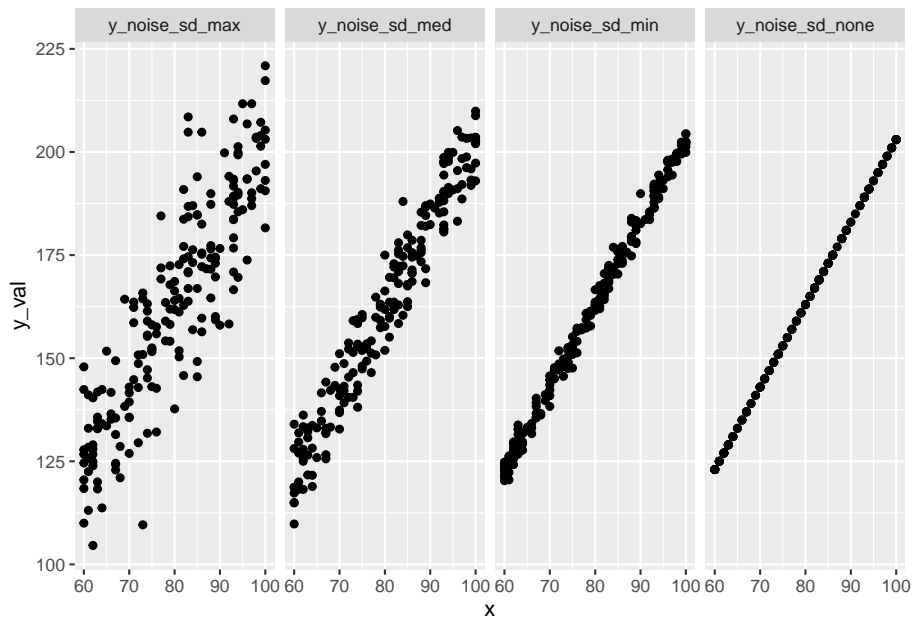
Let's make this a long df so that we can plot multiple standard deviation values together

```
demo_df_long <- demo_df %>%
  pivot_longer(cols = starts_with("y_noise"),
               names_to = "y_col",
               values_to = "y_val"
  )
```

```
demo_df_long <- demo_df_long %>%
  mutate(sd_val = case_when(str_detect(y_col, "sd_none") ~ 0,
                             str_detect(y_col, "sd_min") ~ sd_min,
                             str_detect(y_col, "sd_med") ~ sd_med,
                             str_detect(y_col, "sd_max") ~ sd_max))
```

And visualize the data, faceting by different noise

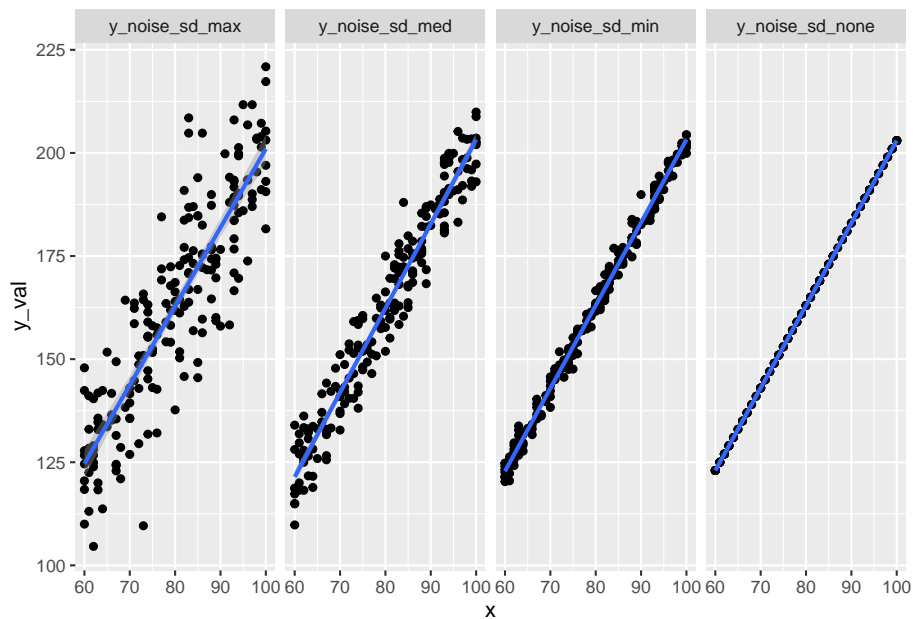
```
demo_df_long %>%
  ggplot(aes(x = x, y = y_val)) +
  geom_point() +
  facet_grid(.~y_col)
```



And add in a line with `geom_smooth(method = 'lm')`

```
demo_df_long %>%
  ggplot(aes(x = x, y = y_val)) +
  geom_point() +
  geom_smooth(method = "lm") +
  facet_grid(.~y_col)
```

```
## `geom_smooth()` using formula 'y ~ x'
```



Create a linear model and look at the summary.

```
fit_demo_min <- lm(y_noise_sd_min ~ x)
summary(fit_demo_min)
```

```
##
## Call:
## lm(formula = y_noise_sd_min ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5.3835 -1.5211  0.0379  1.3254  6.7636
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.21905    0.99404   2.232  0.0267 *
## x            2.01019    0.01229 163.530 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.104 on 198 degrees of freedom
## Multiple R-squared:  0.9927, Adjusted R-squared:  0.9926
## F-statistic: 2.674e+04 on 1 and 198 DF,  p-value: < 2.2e-16
```

We can also look at model results with the `glance()` function from the `broom`

package

```
broom::glance(fit_demo_min)
```

```
## # A tibble: 1 x 12
##   r.squared adj.r.squared sigma statistic    p.value    df logLik   AIC   BIC
##   <dbl>      <dbl> <dbl>      <dbl>      <dbl> <dbl> <dbl> <dbl> <dbl>
## 1     0.993        0.993  2.10      26742. 3.27e-213     1 -432.  869.  879.
## # ... with 3 more variables: deviance <dbl>, df.residual <int>, nobs <int>
```

We can create models for the med and max sd values as well and take a look at those with the `summary()` function once again.

```
fit_demo_med <- lm(y_noise_sd_med ~ x)
summary(fit_demo_med)
```

```
##
## Call:
## lm(formula = y_noise_sd_med ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.4632  -4.6878  -0.4005   4.6437  17.4622
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.24876     2.87289  -0.435   0.664
## x             2.04508     0.03553  57.564 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.08 on 198 degrees of freedom
## Multiple R-squared:  0.9436, Adjusted R-squared:  0.9433
## F-statistic: 3314 on 1 and 198 DF, p-value: < 2.2e-16
```

```
fit_demo_max <- lm(y_noise_sd_max ~ x)
summary(fit_demo_max)
```

```
##
## Call:
## lm(formula = y_noise_sd_max ~ x)
##
## Residuals:
```



```
##      Min      1Q  Median      3Q      Max
## -39.867 -6.497   0.018   7.027  39.951
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 10.16895    5.87444   1.731   0.085 .
## x           1.90820    0.07264  26.268 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.43 on 198 degrees of freedom
## Multiple R-squared:  0.777, Adjusted R-squared:  0.7759
## F-statistic: 690 on 1 and 198 DF, p-value: < 2.2e-16
```

Notice the increase in the standard error of the coefficient estimates as the noise in y values went up

From a programming perspective, this was not very efficient because I just copied, pasted, and corrected these values. There is a better way to do this using lists (see below)

Let's do some fancy stuff to make multiple models at once rather than having to write new lines for each model *Some of these ideas are taken from the R4DS book chapter 25

```
test_nest <- demo_df_long %>% nest(data = -sd_val)
```

```
linear_model <- function(df) {
  lm(y_val ~ x, data = df)
}
```

```
models <- map(test_nest$data, linear_model)
```

```
summary(models[[2]])
```

```
##
## Call:
## lm(formula = y_val ~ x, data = df)
##
## Residuals:
##      Min      1Q  Median      3Q      Max
## -5.3835 -1.5211  0.0379  1.3254  6.7636
##
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.21905    0.99404   2.232   0.0267 *
## x           2.01019    0.01229 163.530   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.104 on 198 degrees of freedom
## Multiple R-squared:  0.9927, Adjusted R-squared:  0.9926
## F-statistic: 2.674e+04 on 1 and 198 DF,  p-value: < 2.2e-16
```

```
summary(models[[3]])
```

```
##
## Call:
## lm(formula = y_val ~ x, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.4632  -4.6878  -0.4005   4.6437  17.4622
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.24876    2.87289  -0.435   0.664
## x           2.04508    0.03553  57.564   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.08 on 198 degrees of freedom
## Multiple R-squared:  0.9436, Adjusted R-squared:  0.9433
## F-statistic: 3314 on 1 and 198 DF,  p-value: < 2.2e-16
```

```
summary(models[[4]])
```

```
##
## Call:
## lm(formula = y_val ~ x, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -39.867  -6.497   0.018   7.027  39.951
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 10.16895    5.87444   1.731   0.085 .
```

5.4. DATA GENERATION WITH THREE DIFFERENT SAMPLE SIZES⁹¹

```
## x          1.90820    0.07264  26.268    <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.43 on 198 degrees of freedom
## Multiple R-squared:  0.777, Adjusted R-squared:  0.7759
## F-statistic:   690 on 1 and 198 DF,  p-value: < 2.2e-16
```

We can also store the models as new columns in the nested dataframe

```
test_nest <- test_nest %>%
  mutate(model = map(data, linear_model))
```

Finally, we can unnest the models to make it easier to compare them with each other in a data frame

```
test_nest <- test_nest %>%
  mutate(glance = map(model, broom::glance)) %>%
  unnest(glance)
```

5.4 Data generation with three different sample sizes

Let's run the same demo but now have three different sample sizes - 10, 50, and 500

First, store the sample sizes we want to use

```
samp_sizes <- c(10, 50, 500)
```

Next, create a bookkeeping column for ourselves to keep track of which sample size the future values will come from

```
samp_size_col <- rep(x = c(10, 50, 500), times = samp_sizes)
```

Calculate the total number of values we will need from the three samples combined

```
tot_samp_size <- sum(samp_sizes)
```

Sample uniformly from 0 to 20

```
x <- round(x = runif(n = tot_samp_size, min = 0, max = 20), digits = 1)
```

Store the standard deviations for the min, med, and max models

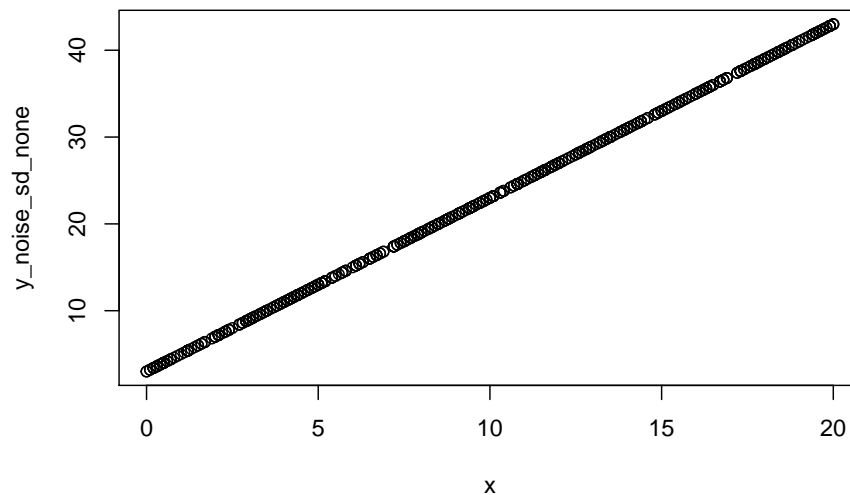
```
sd_min <- 2  
sd_med <- 6  
sd_max <- 12
```

Calculate the y values for the different scenarios where there is no noise up to max noise

```
y_noise_sd_none <- 3 + 2*x  
y_noise_sd_min <- 3 + 2*x + round(x = rnorm(n = tot_samp_size, mean = 0, sd = sd_min),  
y_noise_sd_med <- 3 + 2*x + round(x = rnorm(n = tot_samp_size, mean = 0, sd = sd_med),  
y_noise_sd_max <- 3 + 2*x + round(x = rnorm(n = tot_samp_size, mean = 0, sd = sd_max),
```

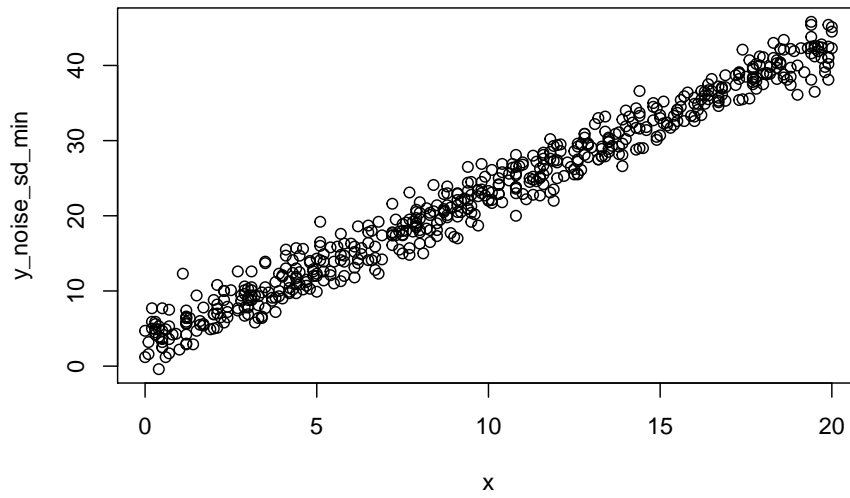
Typical step 1: visualize! Let's plot each of these x values vs y

```
plot(x, y_noise_sd_none)
```

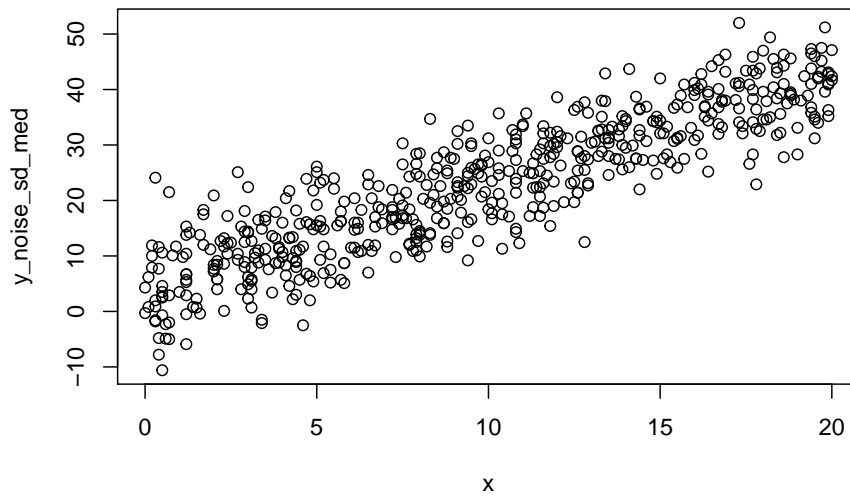


```
plot(x, y_noise_sd_min)
```

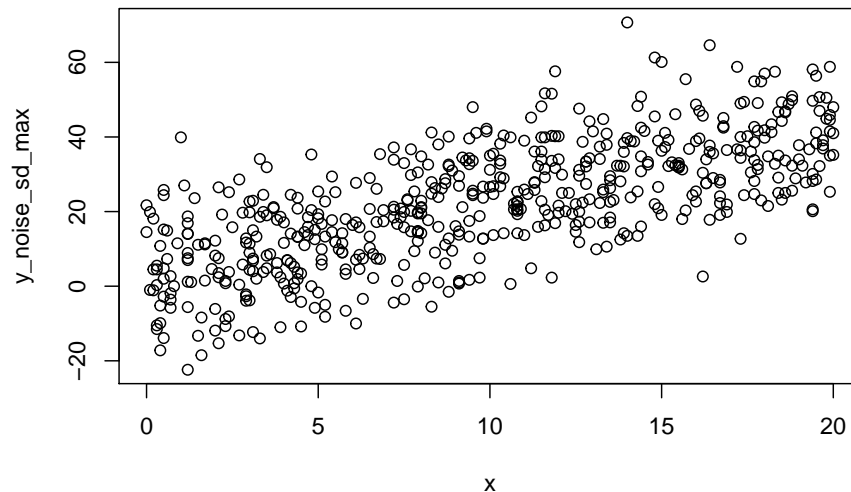
5.4. DATA GENERATION WITH THREE DIFFERENT SAMPLE SIZES93



```
plot(x, y_noise_sd_med)
```



```
plot(x, y_noise_sd_max)
```



Can we calculate the correlations between x and these different y values? (pro tip: yes)

Let's put all of these vectors together into a data frame to make it easier to analyze later on Note, this is not a vital step for conducting the simple regression

```
demo_df <- tibble("n" = samp_size_col,
                  "x" = x,
                  "y_noise_sd_none" = y_noise_sd_none,
                  "y_noise_sd_min" = y_noise_sd_min,
                  "y_noise_sd_med" = y_noise_sd_med,
                  "y_noise_sd_max" = y_noise_sd_max)
```

Order by increasing x value

```
demo_df <- demo_df %>%
  arrange(n, x)
```

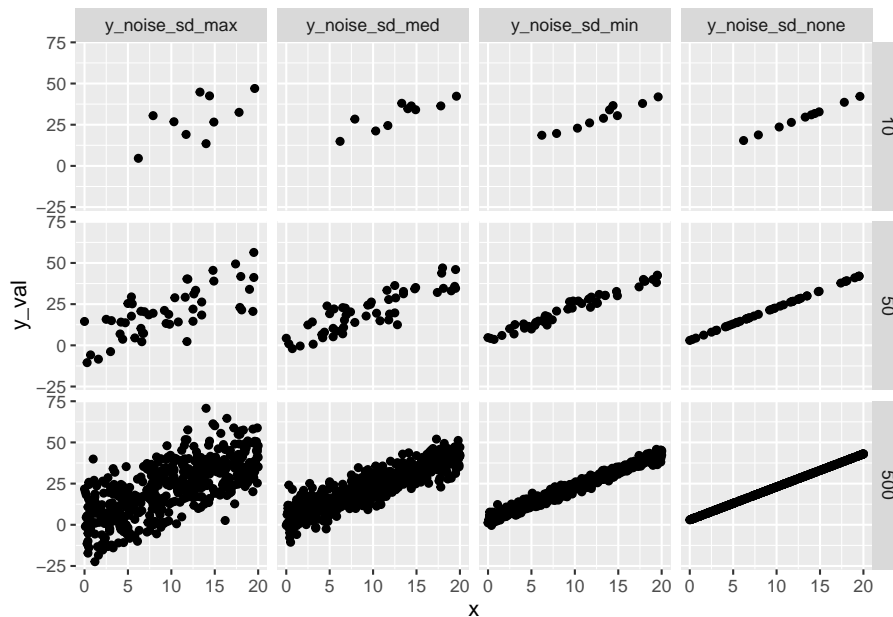
Let's make this a long df so that we can plot multiple standard deviation values together

5.4. DATA GENERATION WITH THREE DIFFERENT SAMPLE SIZES 95

```
demo_df_long <- demo_df %>%
  pivot_longer(cols = starts_with("y_noise"),
               names_to = "y_col",
               values_to = "y_val"
  )

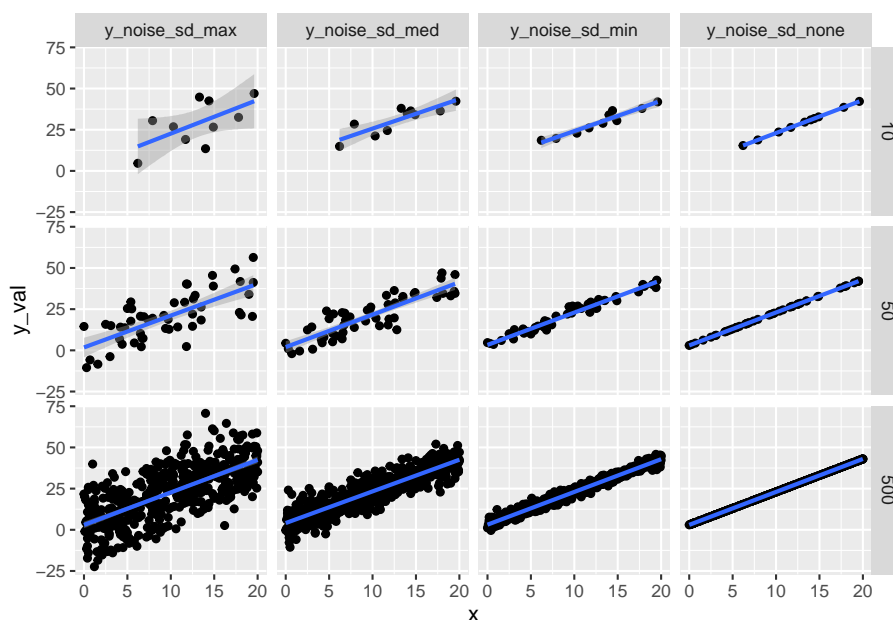
demo_df_long <- demo_df_long %>%
  mutate(sd_val = case_when(str_detect(y_col, "sd_none") ~ 0,
                             str_detect(y_col, "sd_min") ~ sd_min,
                             str_detect(y_col, "sd_med") ~ sd_med,
                             str_detect(y_col, "sd_max") ~ sd_max))
```

```
demo_df_long %>%
  ggplot(aes(x = x, y = y_val)) +
  geom_point() +
  facet_grid(n~y_col)
```



```
demo_df_long %>%
  ggplot(aes(x = x, y = y_val)) +
  geom_point() +
  geom_smooth(method = "lm") +
  facet_grid(n~y_col)
```

`geom_smooth()` using formula 'y ~ x'



```
fit_demo_min <- lm(y_noise_sd_min ~ x)
summary(fit_demo_min)
```

```
##
## Call:
## lm(formula = y_noise_sd_min ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5.3503 -1.2768 -0.0035  1.2450  6.9759
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.14043    0.15889   19.77  <2e-16 ***
## x            1.98512    0.01397  142.12  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.91 on 558 degrees of freedom
## Multiple R-squared:  0.9731, Adjusted R-squared:  0.9731
## F-statistic: 2.02e+04 on 1 and 558 DF, p-value: < 2.2e-16
```

We can also look at model results with the `glance()` function from the `broom` package

5.4. DATA GENERATION WITH THREE DIFFERENT SAMPLE SIZES97

```
broom::glance(fit_demo_min)
```

```
## # A tibble: 1 x 12
##   r.squared adj.r.squared sigma statistic p.value    df logLik   AIC   BIC
##   <dbl>      <dbl> <dbl>      <dbl>    <dbl> <dbl> <dbl> <dbl> <dbl>
## 1    0.973      0.973  1.91      20198.      0      1 -1156. 2318. 2331.
## # ... with 3 more variables: deviance <dbl>, df.residual <int>, nobs <int>
```

We can create models for the med and max sd values as well and take a look at those with the `summary()` function once again

```
fit_demo_med <- lm(y_noise_sd_med ~ x)
summary(fit_demo_med)
```

```
##
## Call:
## lm(formula = y_noise_sd_med ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16.1057  -4.0478   0.1628   3.9796  19.6085
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.91276    0.48368   8.09 3.75e-15 ***
## x             1.92913    0.04252  45.37 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.814 on 558 degrees of freedom
## Multiple R-squared:  0.7867, Adjusted R-squared:  0.7863
## F-statistic: 2058 on 1 and 558 DF, p-value: < 2.2e-16
```

```
fit_demo_max <- lm(y_noise_sd_max ~ x)
summary(fit_demo_max)
```

```
##
## Call:
## lm(formula = y_noise_sd_max ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -32.189  -7.761   0.154   7.975  40.250
```

```
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.83362    0.96027   2.951  0.0033 **
## x           1.97258    0.08442  23.366  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11.54 on 558 degrees of freedom
## Multiple R-squared:  0.4946, Adjusted R-squared:  0.4937
## F-statistic: 546 on 1 and 558 DF, p-value: < 2.2e-16
```

Notice the increase in the standard error of the coefficient estimates as the noise in y values went up

From a programming perspective, this was not very efficient because I just copied, pasted, and corrected these values. There is a better way to do this using lists (see below)

Let's do some fancy stuff to make multiple models at once rather than having to write new lines for each model *Some of these ideas are taken from the R4DS book chapter 25

```
test_nest <- demo_df_long %>% nest(data = -c(sd_val, n))

linear_model <- function(df) {
  lm(y_val ~ x, data = df)
}

models <- map(test_nest$data, linear_model)

summary(models[[2]])
```

```
##
## Call:
## lm(formula = y_val ~ x, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.6922 -1.3073 -0.6152  1.0899  4.3290
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept)   5.7397     2.4885   2.307   0.05 *
```

5.4. DATA GENERATION WITH THREE DIFFERENT SAMPLE SIZES99

```
## x          1.8425      0.1831  10.063  8.1e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.276 on 8 degrees of freedom
## Multiple R-squared:  0.9268, Adjusted R-squared:  0.9176
## F-statistic: 101.3 on 1 and 8 DF,  p-value: 8.102e-06
```

```
summary(models[[3]])
```

```
##
## Call:
## lm(formula = y_val ~ x, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5.0241 -3.7944 -0.5031  2.5688  6.4852
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   7.7300     4.9571   1.559  0.15753
## x             1.7955     0.3647   4.923  0.00116 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.534 on 8 degrees of freedom
## Multiple R-squared:  0.7518, Adjusted R-squared:  0.7208
## F-statistic: 24.23 on 1 and 8 DF,  p-value: 0.00116
```

```
summary(models[[4]])
```

```
##
## Call:
## lm(formula = y_val ~ x, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -17.317  -6.780  -1.251   9.328  15.416
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    2.155     12.626   0.171  0.8687
## x              2.047     0.929   2.204  0.0587 .
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11.55 on 8 degrees of freedom
## Multiple R-squared:  0.3777, Adjusted R-squared:    0.3
## F-statistic: 4.856 on 1 and 8 DF,  p-value: 0.05865
```

We can also store the models as new columns in the nested dataframe

```
test_nest <- test_nest %>%
  mutate(model = map(data, linear_model))
```

Finally, we can unnest the models to make it easier to compare them with each other in a data frame

```
test_nest <- test_nest %>%
  mutate(glance = map(model, broom::glance)) %>%
  unnest(glance)
```

```
## Warning in summary.lm(x): essentially perfect fit: summary may be unreliable
## Warning in summary.lm(x): essentially perfect fit: summary may be unreliable
## Warning in summary.lm(x): essentially perfect fit: summary may be unreliable
## Warning in summary.lm(x): essentially perfect fit: summary may be unreliable
```

And look at the different models by just calling the data frame

```
test_nest
```

```
## # A tibble: 12 x 16
##       n sd_val data  model r.squared adj.r.squared  sigma statistic  p.value
##   <dbl> <dbl> <lis> <lis>    <dbl>        <dbl>    <dbl>    <dbl>    <dbl>
## 1     10      0 <tib~ <lm>      1          1 2.80e-15  7.90e31 2.87e-125
## 2     10      2 <tib~ <lm>  0.927      0.918 2.28e+ 0  1.01e 2 8.10e- 6
## 3     10      6 <tib~ <lm>  0.752      0.721 4.53e+ 0  2.42e 1 1.16e- 3
## 4     10     12 <tib~ <lm>  0.378      0.300 1.15e+ 1  4.86e 0 5.87e- 2
## 5     50      0 <tib~ <lm>      1          1  4.48e-15  3.00e32 0.
## 6     50      2 <tib~ <lm>  0.965      0.965 2.08e+ 0  1.34e 3 9.56e- 37
## 7     50      6 <tib~ <lm>  0.754      0.749 6.31e+ 0  1.47e 2 3.07e- 16
## 8     50     12 <tib~ <lm>  0.521      0.511 1.04e+ 1  5.22e 1 3.34e- 9
## 9    500      0 <tib~ <lm>      1          1  2.95e-14  7.79e31 0.
## 10   500      2 <tib~ <lm>  0.974      0.974 1.89e+ 0  1.87e 4 0.
```

5.4. DATA GENERATION WITH THREE DIFFERENT SAMPLE SIZES101

```
## 11  500      6 <tib~ <lm>      0.790      0.790 5.78e+ 0   1.87e 3 7.39e-171
## 12  500     12 <tib~ <lm>      0.492      0.491 1.17e+ 1   4.83e 2 2.51e- 75
## # ... with 7 more variables: df <dbl>, logLik <dbl>, AIC <dbl>, BIC <dbl>,
## #   deviance <dbl>, df.residual <int>, nobs <int>
```


Chapter 6

Week 6: Regression II

This week we will be discussing multiple regression. We will use the child aggression data set as a running example for this.

6.1 Explore the child aggression data set

```
ca_df <- read.table("ChildAggression.dat", header = TRUE)
```

Method 1 for quickly getting summary statistics for the data you have, using the built-in `summary()` function

```
summary(ca_df)
```

##	Aggression	Television	Computer_Games
##	Min. : -1.295608	Min. : -1.46012	Min. : -1.1538345
##	1st Qu.: -0.174279	1st Qu.: -0.18206	1st Qu.: -0.1687007
##	Median : -0.005548	Median : -0.01247	Median : -0.0001997
##	Mean : -0.005011	Mean : -0.02758	Mean : 0.0103812
##	3rd Qu.: 0.149611	3rd Qu.: 0.14983	3rd Qu.: 0.1881810
##	Max. : 1.178823	Max. : 0.98162	Max. : 1.6175039
##	Sibling_Aggression	Diet	Parenting_Style
##	Min. : -1.433127	Min. : -1.28490	Min. : -4.46041
##	1st Qu.: -0.156414	1st Qu.: -0.16136	1st Qu.: -0.58008
##	Median : 0.008459	Median : 0.00934	Median : 0.02736
##	Mean : 0.008275	Mean : 0.01162	Mean : 0.00000
##	3rd Qu.: 0.185136	3rd Qu.: 0.18708	3rd Qu.: 0.51784
##	Max. : 1.103671	Max. : 1.22383	Max. : 3.99326

Method 2, using the `describe()` function from the `psych` package

```
psych::describe(ca_df)
```

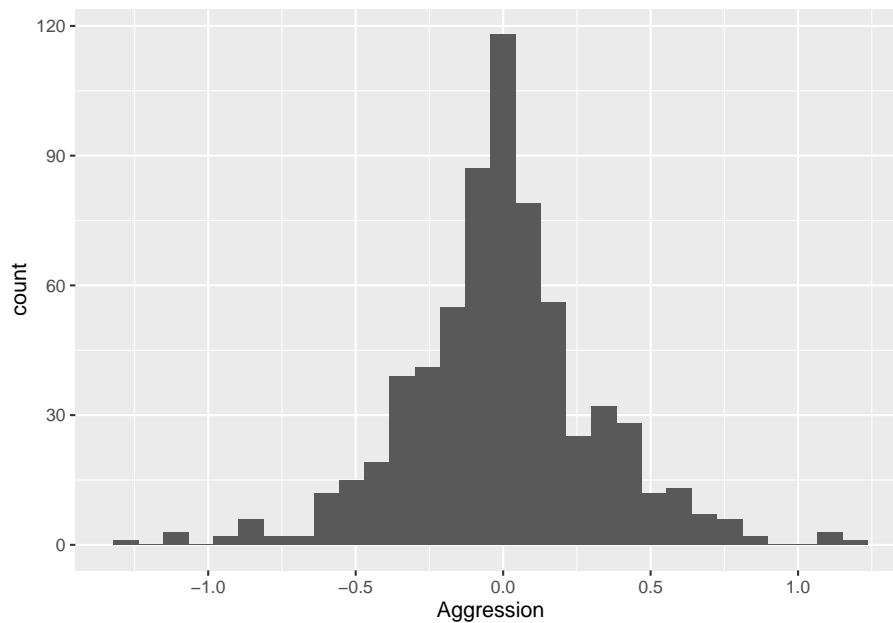
```
##              vars    n mean  sd median trimmed  mad   min  max range
## Aggression          1 666 -0.01 0.32  -0.01  -0.01 0.24 -1.30  1.18  2.47
## Television          2 666 -0.03 0.31  -0.01  -0.02 0.25 -1.46  0.98  2.44
## Computer_Games      3 666  0.01 0.34   0.00   0.00 0.27 -1.15  1.62  2.77
## Sibling_Aggression   4 666  0.01 0.33   0.01   0.01 0.26 -1.43  1.10  2.54
## Diet                5 666  0.01 0.34   0.01   0.01 0.26 -1.28  1.22  2.51
## Parenting_Style      6 666  0.00 1.00   0.03   0.01 0.82 -4.46  3.99  8.45
##              skew kurtosis   se
## Aggression      -0.02     1.61 0.01
## Television       -0.36     1.34 0.01
## Computer_Games    0.25     1.60 0.01
## Sibling_Aggression -0.17     1.40 0.01
## Diet             -0.12     1.51 0.01
## Parenting_Style  -0.22     1.67 0.04
```

Spot check: where do the “se” values come from?

Now that we have seen some of the numbers, let’s try to visualize some of these data. Let’s try to visualize the data while we’re at it

```
ca_df %>%
  ggplot(aes(x = Aggression)) +
  geom_histogram()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

We could do this for each of our variables or we could try something a little fancier by using the `pivot_longer()` function to make a longer dataframe and plot everything at once

First, create the long dataframe. After running this command, it's a good idea to view the new data frame to make sure this did what you intended

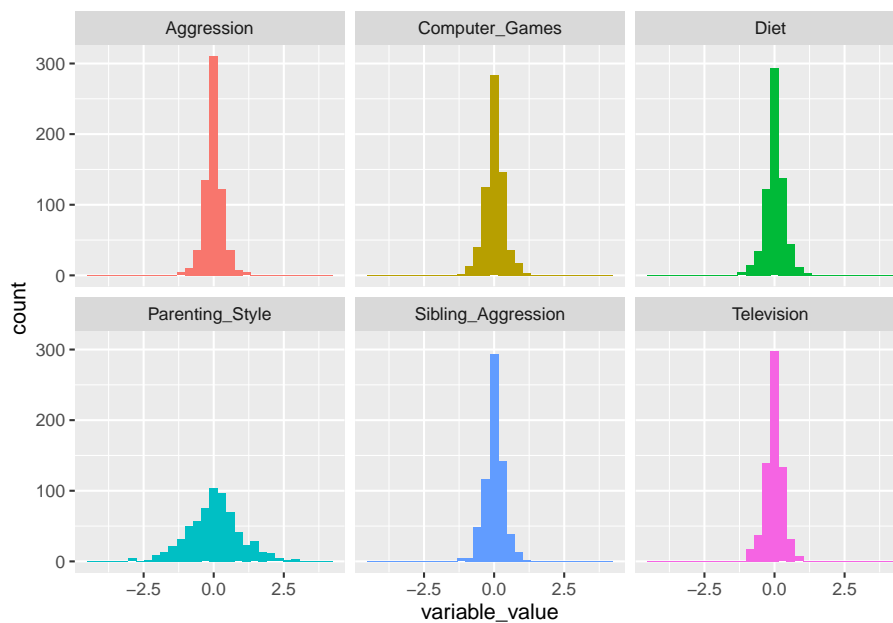
```
ca_df_long <- ca_df %>%
  pivot_longer(cols = Aggression:Parenting_Style, names_to = "variable_name", values_to = "variable_value")
```

After making sure everything looks in order, run the same plot command with the addition of a facet.

Note that the `fill = variable_name` tells R to color the plot with different colors by `variable_name` and the `theme(...)` tells R to get rid of the legend that automatically shows up

```
ca_df_long %>%
  ggplot(aes(x = variable_value, fill = variable_name)) +
  geom_histogram() +
  facet_wrap(variable_name ~.) +
  theme(legend.position = "none")
```

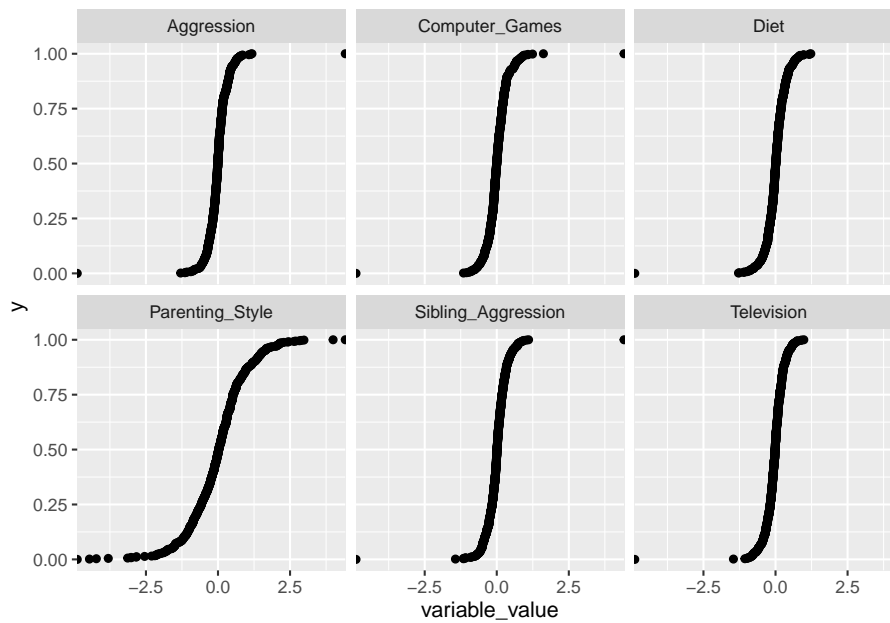
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



From this we should be able to see that everything is normally distributed and appears to be standardized

Let's try a different kind of visualization: the empirical cumulative distribution function (eCDF) (sounds fancy, but it's not bad)

```
ca_df_long %>%
  ggplot(aes(x = variable_value)) +
  stat_ecdf(geom = "point") +
  facet_wrap(variable_name ~.)
```

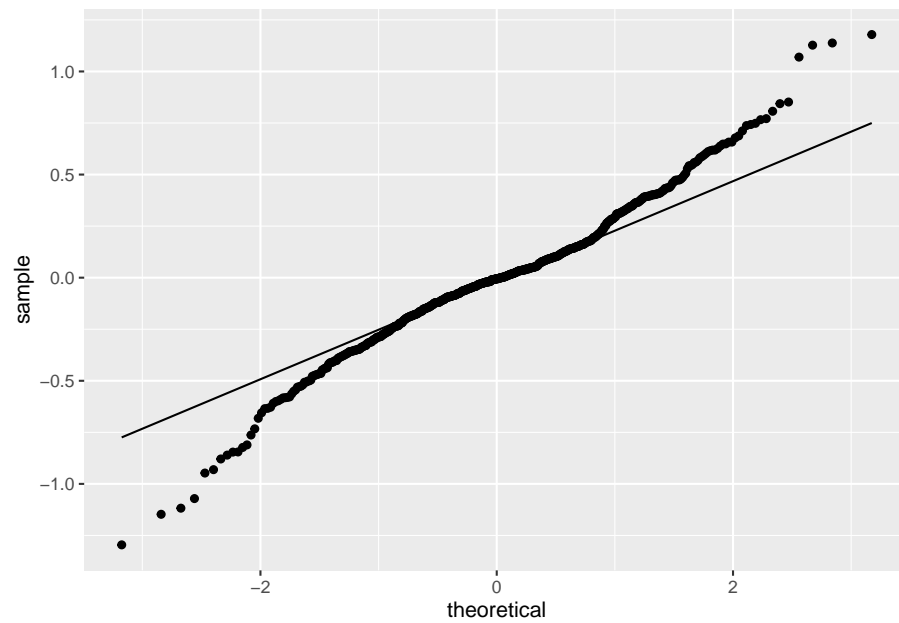


It looks like a bunch of distorted s shapes. In order to understand what is going on here, we should take a brief detour into the world of cumulative distribution functions (CDFs), eCDFs, and eventually QQ plots. This detour is at the end of the markdown file.

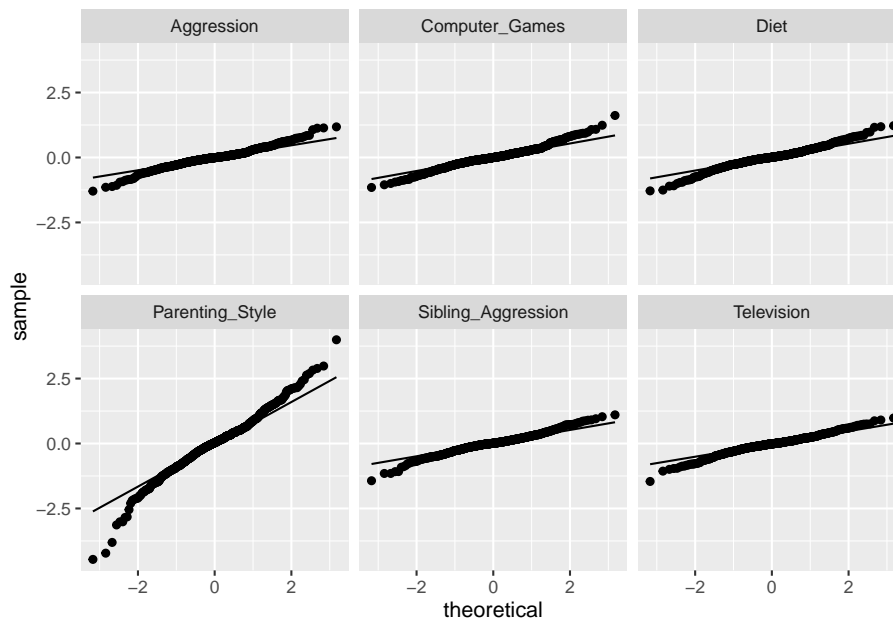
Now we can try a third way of visualizing our data - QQ plots!

As before, we can plot each variable with separate calls or plot them all together with the long data frame

```
ca_df %>%
  ggplot(aes(sample=Aggression)) +
  stat_qq() +
  stat_qq_line()
```



```
ca_df_long %>%  
  ggplot(aes(sample=variable_value)) +  
    stat_qq() +  
    stat_qq_line() +  
    facet_wrap(variable_name ~ .)
```



6.2 Multiple regression

Let's create a model of child aggression as a function of parenting style and sibling aggression

```
model_fam <- lm(Aggression ~ Parenting_Style + Sibling_Aggression, data = ca_df)
```

Run the summary on the model

```
summary(model_fam)
```

```
##
## Call:
## lm(formula = Aggression ~ Parenting_Style + Sibling_Aggression,
##     data = ca_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.09755 -0.17180  0.00092  0.15405  1.23037
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)      -0.005784    0.012065   -0.479    0.632
## Parenting_Style    0.061984    0.012257    5.057 5.51e-07 ***
## Sibling_Aggression 0.093409    0.037505    2.491    0.013 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3113 on 663 degrees of freedom
## Multiple R-squared:  0.05325,    Adjusted R-squared:  0.05039
## F-statistic: 18.64 on 2 and 663 DF,  p-value: 1.325e-08
```

This will be helpful later when we want to plot the residuals of the model.

Let's create another model of child aggression as a function of computer games and television.

```
model_screens <- lm(Aggression ~ Computer_Games + Television, data = ca_df)
summary(model_screens)
```

```
##
## Call:
## lm(formula = Aggression ~ Computer_Games + Television, data = ca_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.05234 -0.15191 -0.00512  0.15156  1.24062
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -0.002879   0.012132  -0.237  0.812526
## Computer_Games  0.153874   0.035845   4.293 2.03e-05 ***
## Television     0.135263   0.039546   3.420 0.000664 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3115 on 663 degrees of freedom
## Multiple R-squared:  0.05149,    Adjusted R-squared:  0.04862
## F-statistic: 17.99 on 2 and 663 DF,  p-value: 2.455e-08
```

Now let's see what the full model looks like with all five predictors

```
model_all <- lm(Aggression ~ Parenting_Style + Sibling_Aggression + Diet + Computer_Games +
               data = ca_df)
summary(model_all)
```

```
##
## Call:
## lm(formula = Aggression ~ Parenting_Style + Sibling_Aggression +
##     Diet + Computer_Games + Television, data = ca_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.12629 -0.15253 -0.00421  0.15222  1.17669
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -0.004988   0.011983  -0.416 0.677350
## Parenting_Style    0.056648   0.014557   3.891 0.000110 ***
## Sibling_Aggression 0.081684   0.038780   2.106 0.035550 *
## Diet           -0.109054   0.038076  -2.864 0.004315 **
## Computer_Games    0.142161   0.036920   3.851 0.000129 ***
## Television        0.032916   0.046057   0.715 0.475059
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3071 on 660 degrees of freedom
## Multiple R-squared:  0.08258,    Adjusted R-squared:  0.07563
## F-statistic: 11.88 on 5 and 660 DF,  p-value: 5.025e-11
```

We can get the output from the `summary()` function but as a tibble instead:

```
broom::tidy(model_all, conf.int = TRUE)
```

```
## # A tibble: 6 x 7
##   term                estimate std.error statistic  p.value conf.low conf.high
##   <chr>              <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)      -0.00499   0.0120   -0.416  0.677    -0.0285   0.0185
## 2 Parenting_Style    0.0566    0.0146    3.89  0.000110  0.0281   0.0852
## 3 Sibling_Aggression 0.0817    0.0388    2.11  0.0355    0.00554   0.158
## 4 Diet             -0.109    0.0381   -2.86  0.00432  -0.184   -0.0343
## 5 Computer_Games     0.142    0.0369    3.85  0.000129  0.0697   0.215
## 6 Television         0.0329    0.0461    0.715  0.475    -0.0575   0.123
```

We can also get summary statistics like the R^2 , F-statistics, and AIC from the `glance()` function in the Broom package

```
broom::glance(model_all)
```

```
## # A tibble: 1 x 12
```

```
##   r.squared adj.r.squared sigma statistic p.value    df logLik   AIC   BIC
##   <dbl>         <dbl> <dbl>         <dbl>    <dbl> <dbl> <dbl> <dbl>
## 1    0.0826         0.0756 0.307          11.9 5.02e-11     5 -156.  325.  357.
## # ... with 3 more variables: deviance <dbl>, df.residual <int>, nobs <int>
```

Let's stop here to interpret what the model output is telling us (specifically, the coefficient estimates, the standard errors, the p values, and the R^2 values).

Notice the Television p value and standard errors (especially how the estimate is so close to zero, and its 95% confidence interval crosses 0)

We can also have R handle factors automatically in a regression model for us.

First, we will generate a new factor for the child aggression data - the handedness of the child (left vs right).

This code generates a random handedness for each of the 666 students and assigns it to a new column in `ca_df`

```
handed_levels <- c("left", "right")
handed_vector <- sample(handed_levels, size = nrow(ca_df), replace = TRUE, prob = c(0.4, 0.6))
ca_df$Handedness <- handed_vector
```

Check to make sure we have all the predictor names correct to add into the linear model.

```
names(ca_df)
```

```
## [1] "Aggression"      "Television"      "Computer_Games"
## [4] "Sibling_Aggression" "Diet"           "Parenting_Style"
## [7] "Handedness"
```

```
model_all_3 <- lm(Aggression ~ Parenting_Style + Sibling_Aggression + Television + Computer_Games + Diet + Handedness, data = ca_df)
```

Check the model summary and notice how the Handedness predictor variable was handled.

```
summary(model_all_3)
```

```
##
## Call:
## lm(formula = Aggression ~ Parenting_Style + Sibling_Aggression +
##     Television + Computer_Games + Diet + Handedness, data = ca_df)
##
## Residuals:
```



```
##      Min      1Q   Median      3Q      Max
## -1.12485 -0.15230 -0.00332  0.15254  1.17596
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -0.006228   0.019320  -0.322  0.747302
## Parenting_Style  0.056528   0.014642   3.861  0.000124 ***
## Sibling_Aggression 0.081905   0.038903   2.105  0.035636 *
## Television     0.033078   0.046134   0.717  0.473636
## Computer_Games  0.142117   0.036951   3.846  0.000132 ***
## Diet          -0.109070   0.038106  -2.862  0.004339 **
## Handednessright  0.002019   0.024671   0.082  0.934809
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3073 on 659 degrees of freedom
## Multiple R-squared:  0.08259,    Adjusted R-squared:  0.07424
## F-statistic: 9.888 on 6 and 659 DF,  p-value: 1.847e-10
```

```
lm(Aggression ~ ., data=ca_df)
```

```
##
## Call:
## lm(formula = Aggression ~ ., data = ca_df)
##
## Coefficients:
##      (Intercept)      Television      Computer_Games  Sibling_Aggression
##      -0.006228         0.033078         0.142117         0.081905
##           Diet      Parenting_Style      Handednessright
##      -0.109070         0.056528         0.002019
```

6.2.1 Assumption testing

Note: A lot of the things we are doing here can be accomplished with the `augment()` function in the `broom` package

Check for correlation between adjacent residual terms using Durbin-Watson test

```
dwt(model_all)
```

```
## lag Autocorrelation D-W Statistic p-value
## 1      0.04218005      1.912808  0.282
## Alternative hypothesis: rho != 0
```

Check for multicollinearity.

```
car::vif(model_all)
```

```
##      Parenting_Style Sibling_Aggression      Diet      Computer_Games
##      1.494296      1.132618      1.160466      1.122719
##      Television
##      1.435525
```

Now let's go back and check for influential cases following the procedure Field does in the textbook.

```
ca_df$residuals <- resid(model_all) # notice how this is the same as looking at View(m
#model_all$residuals
ca_df$standardized.residuals <- rstandard(model_all)
ca_df$studentized.residuals <- rstudent(model_all)
ca_df$cooks.distance <- cooks.distance(model_all)
ca_df$dfbeta <- dfbeta(model_all)
ca_df$dfbet <- dfbet(model_all)
ca_df$leverage <- hatvalues(model_all)
ca_df$covariance.ratios <- covratio(model_all)
```

Now that we have all of the extra model fit statistics calculated, we can add a new column using a different method than what he does in the textbook (but this accomplishes the same task).

We are using a combination of `mutate()` and `case_when()` instead of the traditional logical test subsetting that we see in the book on pp. 289-290

```
ca_df <- ca_df %>%
  mutate(large.residual = case_when(standardized.residuals > 2 | standardized.residuals
    abs(standardized.residuals) <= 2 ~ FALSE))
```

Now let's just look at the observations where there was a large residual (absolute value > 2)

```
ca_df %>% filter(large.residual == TRUE) %>% head()
```

```
##      Aggression Television Computer_Games Sibling_Aggression      Diet
## 1  0.7711534 -0.03287184    0.70991822    0.576836667 -0.0229903
## 2 -0.9309839 -0.14695730    0.58503832    0.143914123  0.1341567
## 3  0.8437696 -0.53337187    0.62674291    0.006193792 -0.0542531
## 4 -0.8604237 -0.32734996   -0.10228524    0.088454410  0.3403472
## 5  0.7374707  0.05104538    0.06804261    0.742478946  0.3607529
```

```
## 6 0.6158280 -0.17945254 -0.14381635 0.337002793 -0.5992525
## Parenting_Style Handedness residuals standardized.residuals
## 1 -1.2481665 left 0.6973813 2.289855
## 2 -1.0470961 left -0.9421375 -3.081966
## 3 0.2386578 right 0.7572738 2.484122
## 4 -0.6988965 left -0.7606374 -2.483667
## 5 0.9243197 left 0.6574374 2.152417
## 6 -1.0334894 right 0.6128345 2.006389
## studentized.residuals cooks.distance dfbeta.(Intercept)
## 1 2.297263 0.014610960 0.0010048327
## 2 -3.102033 0.014474124 -0.0013336489
## 3 2.493925 0.015185553 0.0008375093
## 4 -2.493465 0.005595170 -0.0010247662
## 5 2.158374 0.008347004 0.0008799517
## 6 2.011011 0.007252602 0.0009296294
## dfbeta.Parenting_Style dfbeta.Sibling_Aggression dfbeta.Diet
## 1 -0.0022102310 0.0053377408 -0.0015602060
## 2 0.0022877988 -0.0008058492 -0.0011872151
## 3 0.0016068696 0.0002508274 -0.0028943101
## 4 0.0006636867 -0.0011984511 -0.0042173119
## 5 0.0008756410 0.0069052922 0.0014075292
## 6 -0.0005401040 0.0047888699 -0.0054530647
## dfbeta.Computer_Games dfbeta.Television dffit leverage
## 1 0.0067966200 0.0013345865 0.2970419 0.016444209
## 2 -0.0079629954 -0.0003344591 -0.2966131 0.009060145
## 3 0.0072143563 -0.0098994502 0.3030411 0.014550253
## 4 0.0015183740 0.0029152017 -0.1839467 0.005412783
## 5 -0.0015476206 -0.0021915865 0.2244095 0.010694476
## 6 -0.0004611714 -0.0010521222 0.2090845 0.010694131
## covariance.ratios large.residual
## 1 0.9780649 TRUE
## 2 0.9335336 TRUE
## 3 0.9679182 TRUE
## 4 0.9590456 TRUE
## 5 0.9778335 TRUE
## 6 0.9832825 TRUE
```

There's a problem here because we don't preserve observation numbers. In general, it might be a good idea to give each observation an ID with something like

```
dim(ca_df)[1]
```

```
## [1] 666
```

```
ca_df <- ca_df %>%  
  mutate(ID = seq(1:dim(ca_df)[1]))
```

In the tidyverse, there is an even simpler way to do this with the `rowid_to_column()` function!

```
ca_df <- ca_df %>%  
  rowid_to_column(var = "participant_id")
```

Now we can run the filter again and see which observations have large residuals

```
ca_df %>% filter(large.residual == TRUE) %>% view()
```

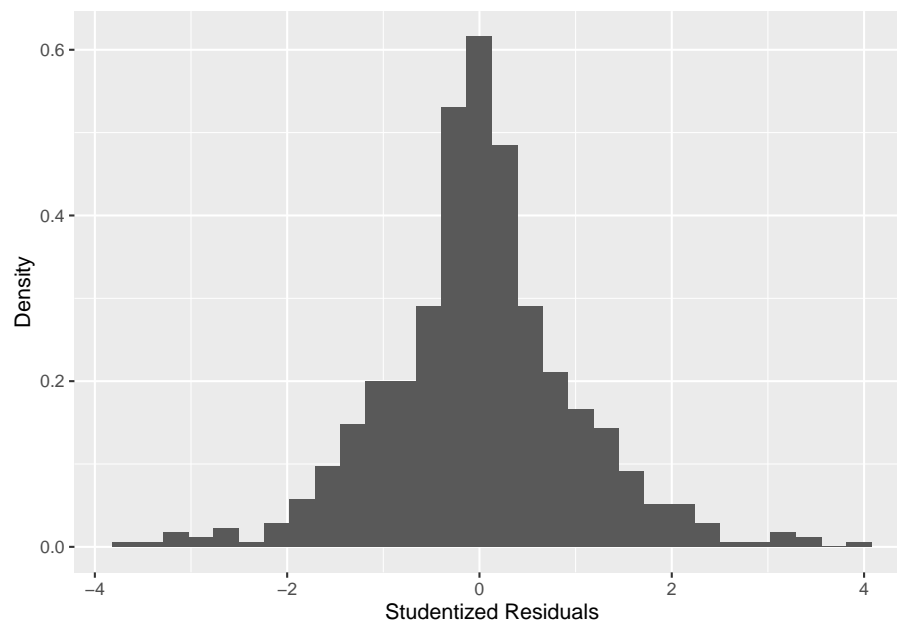
Plot the residuals against the fitted values

```
ca_df$fitted <- model_all$fitted.values
```

First, plot the distribution of the residuals by themselves

```
histogram <- ca_df %>%  
  ggplot(aes(x = studentized.residuals)) +  
  geom_histogram(aes(y = ..density..)) +  
  labs(x = "Studentized Residuals",  
       y = "Density")  
histogram
```

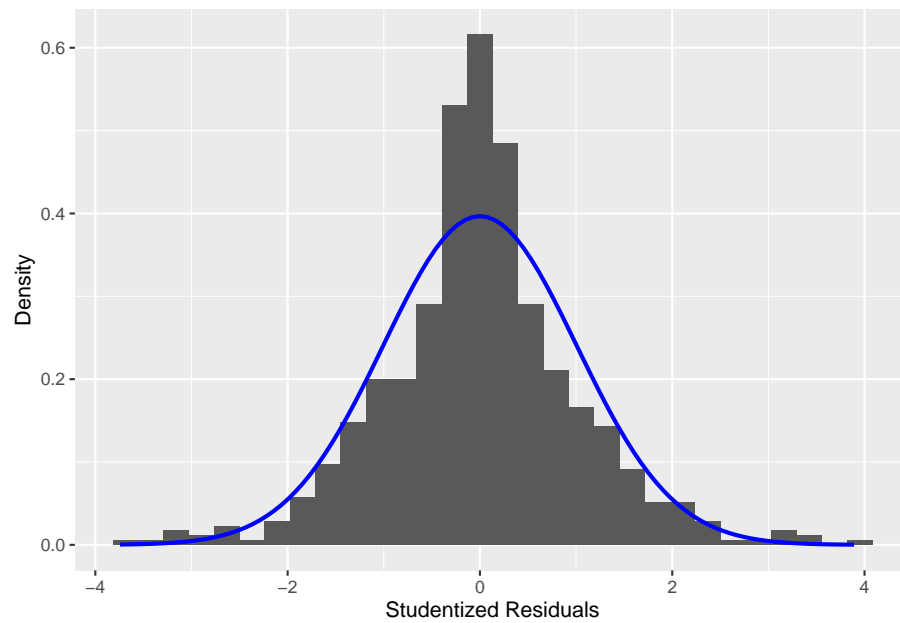
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Now let's add a normal density on top of that for comparison

```
histogram + stat_function(fun = dnorm,  
                           args = list(mean = mean(ca_df$studentized.residuals, na.rm = TRUE),  
                                       sd = sd(ca_df$studentized.residuals, na.rm = TRUE)),  
                           color = "blue", size = 1)
```

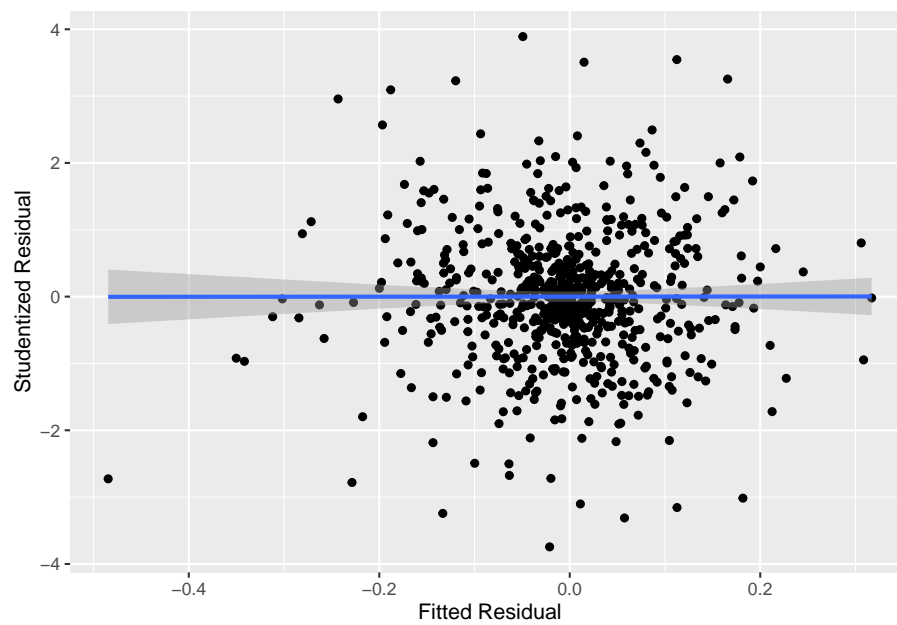
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Second, plot the residuals against the fitted values to make sure there are no systematic patterns

```
ca_df %>%  
  ggplot(aes(x = fitted, y = studentized.residuals)) +  
  geom_point() +  
  geom_smooth(method = "lm") +  
  labs(x = "Fitted Residual",  
       y = "Studentized Residual")
```

```
## `geom_smooth()` using formula 'y ~ x'
```

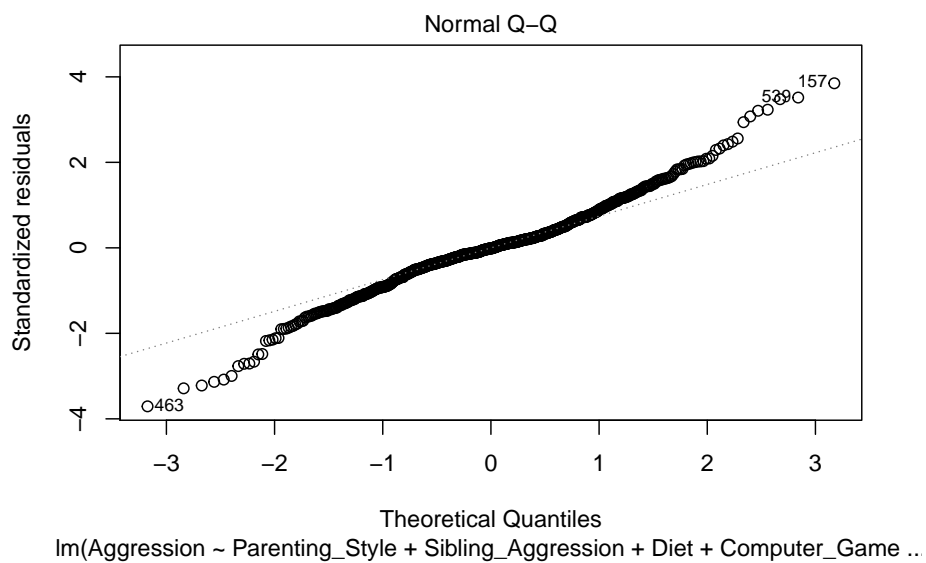
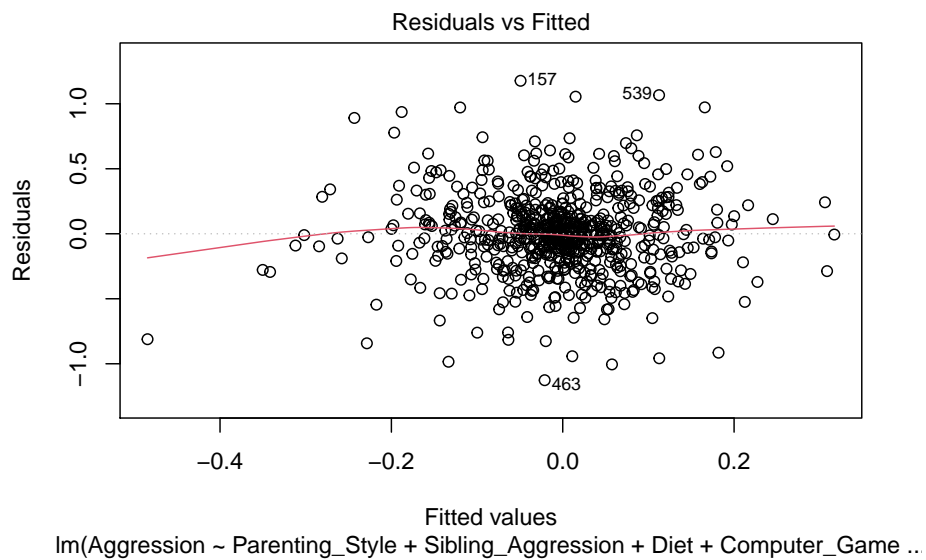


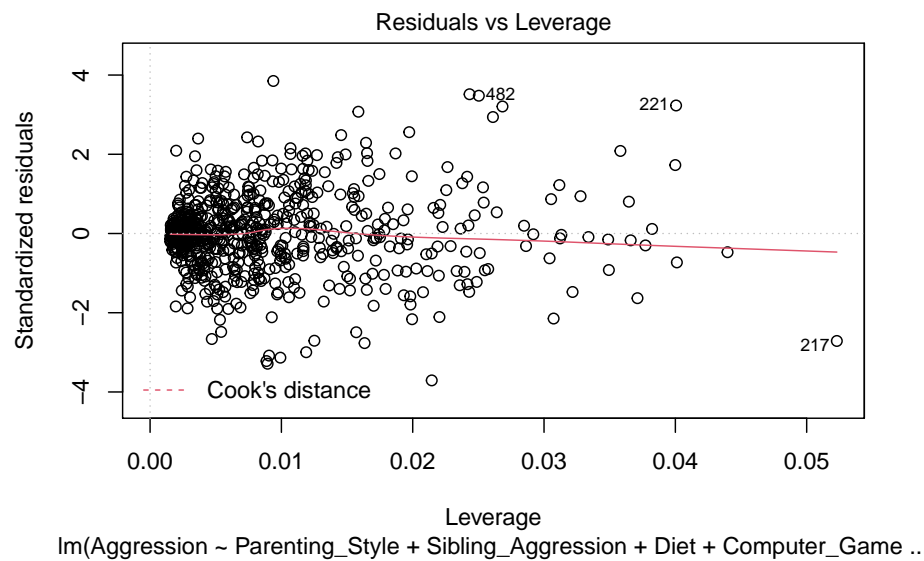
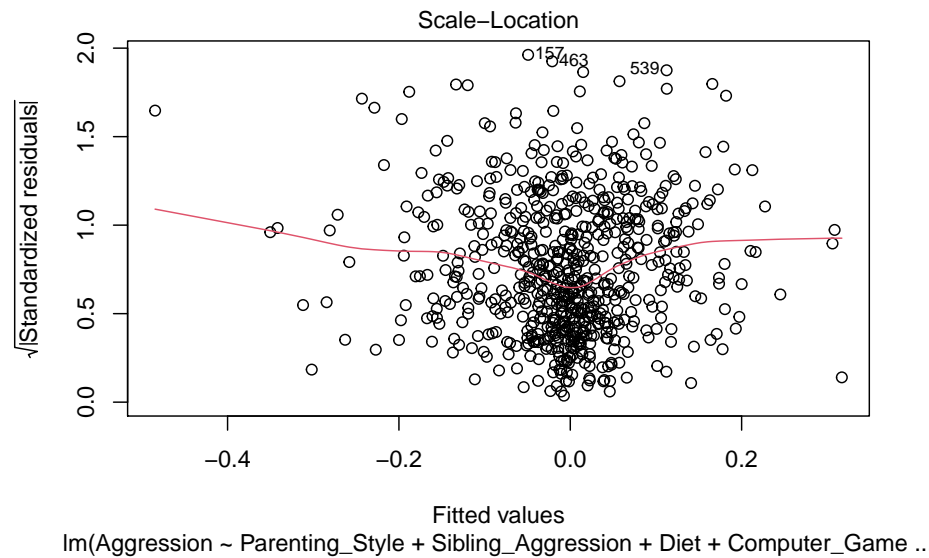
We can also skip all the fancy code from the book and use the `plot()` function on the model.

R knows that this means we want to plot four fit plots

```
par(mfrow=c(2,2)) # this changes the arrangement for the plots to show two rows and two columns
```

```
plot(model_all)
```





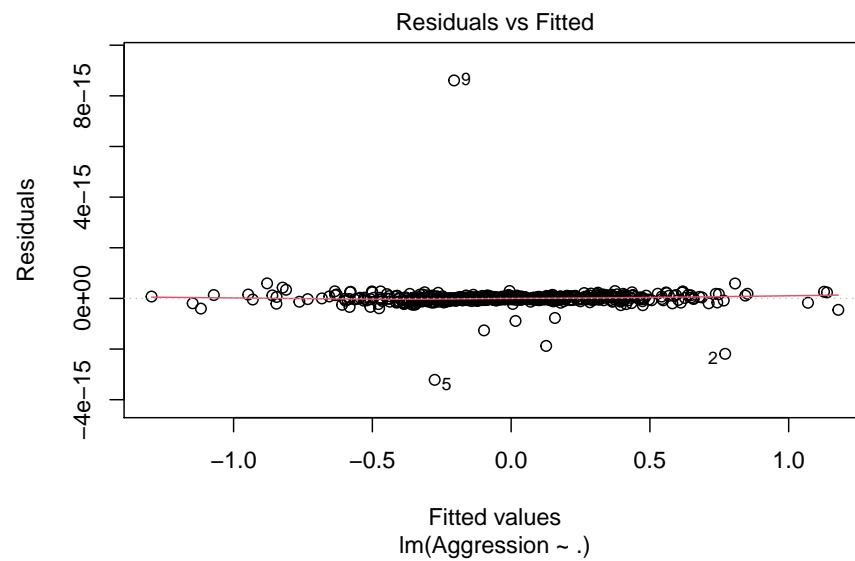
```
par(mfrow=c(1,1)) # this changes the plotting arrangement back to the default of one plot at a time
```

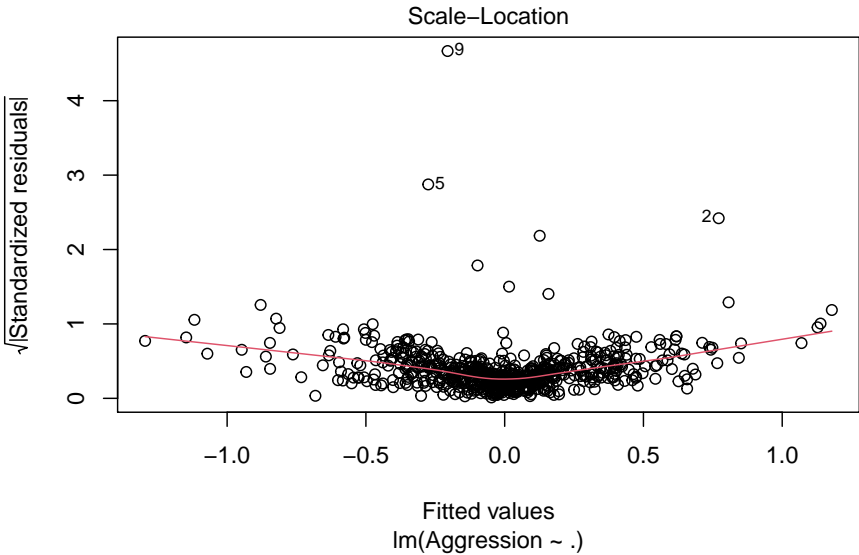
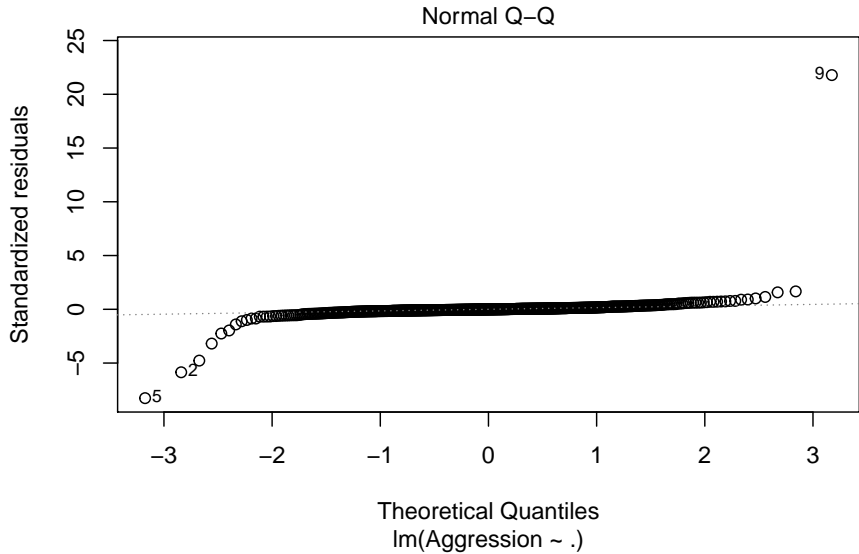
Note that we could also create a model that uses all our covariates as predictors with the following call

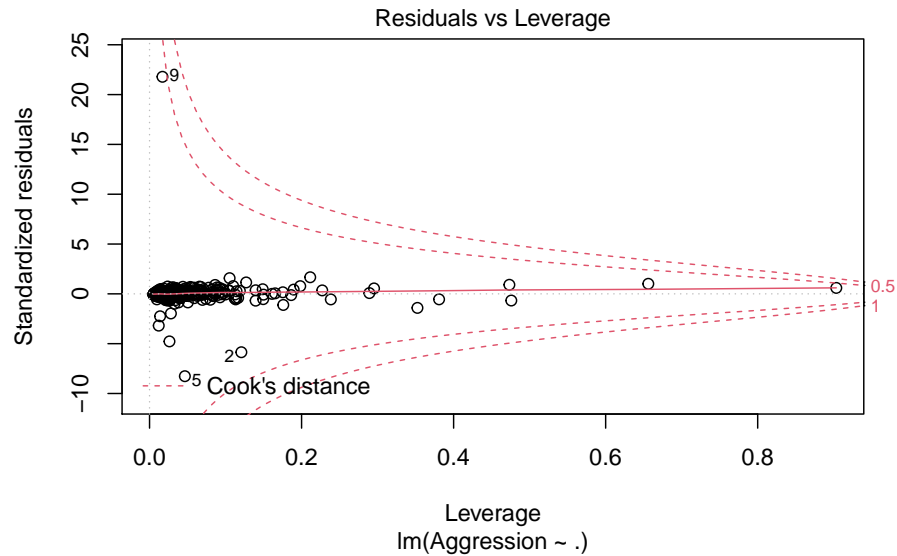
```
model_all_2 <- lm(Aggression ~ ., data=ca_df)
```

Inspect the results to make sure it is giving the same results as before

```
plot(model_all_2)
```







```
summary(model_all_2)
```

```
##
## Call:
## lm(formula = Aggression ~ ., data = ca_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.216e-15 -4.010e-17 -2.600e-18  4.070e-17  8.607e-15
##
## Coefficients: (2 not defined because of singularities)
##              Estimate Std. Error  t value Pr(>|t|)
## (Intercept)   -4.988e-03  2.150e-15 -2.320e+12 < 2e-16 ***
## participant_id  3.761e-19  8.163e-20  4.607e+00  4.93e-06 ***
## Television     3.292e-02  6.086e-17  5.409e+14 < 2e-16 ***
## Computer_Games 1.422e-01  4.883e-17  2.911e+15 < 2e-16 ***
## Sibling_Aggression 8.168e-02  5.123e-17  1.594e+15 < 2e-16 ***
## Diet          -1.091e-01  5.053e-17 -2.158e+15 < 2e-16 ***
## Parenting_Style 5.665e-02  1.934e-17  2.929e+15 < 2e-16 ***
## Handednessright 2.481e-17  3.228e-17  7.690e-01    0.442
## residuals      1.000e+00  4.111e-12  2.432e+11 < 2e-16 ***
## standardized.residuals 8.540e-13  2.457e-12  3.480e-01    0.728
## studentized.residuals -1.487e-16  5.938e-15 -2.500e-02    0.980
## cooks.distance  5.031e-15  5.584e-15  9.010e-01    0.368
```

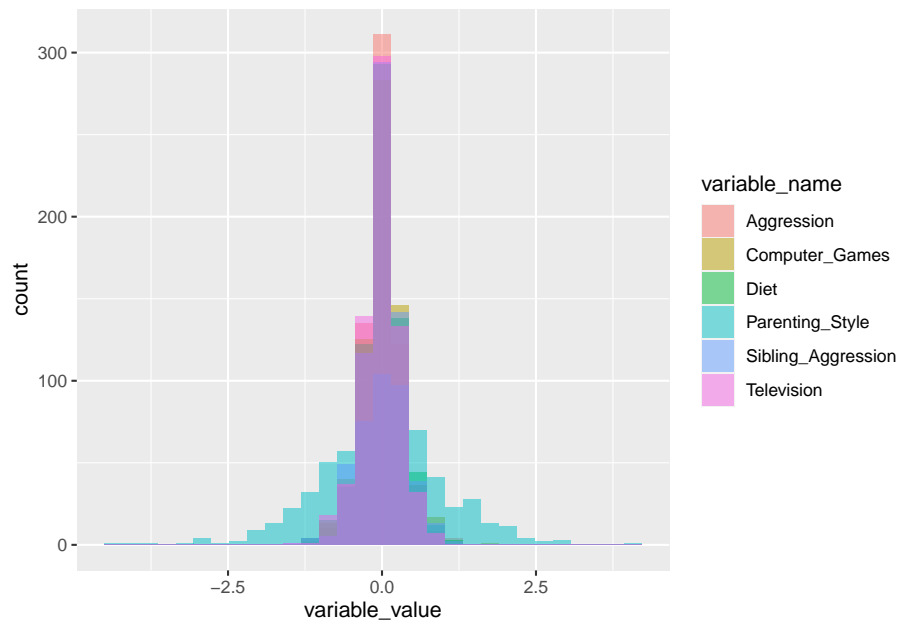
```
## dfbeta(Intercept)      -8.887e-10  2.593e-09 -3.430e-01  0.732
## dfbetaParenting_Style   3.083e-14  2.986e-14  1.033e+00  0.302
## dfbetaSibling_Aggression -7.369e-12  2.146e-11 -3.430e-01  0.731
## dfbetaDiet              -1.032e-11  3.013e-11 -3.430e-01  0.732
## dfbetaComputer_Games    -9.233e-12  2.692e-11 -3.430e-01  0.732
## dfbetaTelevision        2.452e-11  7.150e-11  3.430e-01  0.732
## dffit                  -2.383e-15  5.561e-15 -4.280e-01  0.668
## leverage               -3.056e-15  3.512e-15 -8.700e-01  0.384
## covariance.ratios       1.874e-16  2.140e-15  8.800e-02  0.930
## large.residualTRUE      -1.521e-16  1.215e-16 -1.251e+00  0.211
## ID                      NA          NA          NA      NA
## fitted                  NA          NA          NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.985e-16 on 644 degrees of freedom
## Multiple R-squared:    1, Adjusted R-squared:    1
## F-statistic: 2.034e+31 on 21 and 644 DF, p-value: < 2.2e-16
```

Check how `augment()` works and which parts of the above it replicates

```
augmented_model_all <- augment(model_all)
```

```
ca_df_long %>%
  ggplot(aes(x = variable_value, fill = variable_name)) +
  geom_histogram(position = "identity", alpha = 0.5)
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



6.2.2 Detour: Cumulative Distribution Function (CDF), empirical CDF, and QQ Plot Discussion

Our goal here is to visualize the cumulative distribution function for a normal distribution as a stepping stone toward understanding QQ plots. We will use plots of the empirical cumulative distribution function (eCDF)

As a way to visualize our data. Although we have not discussed eCDFs before, they are another tool you can use to visualize the distribution of quantitative data.

We will also look at the effect of sample size on the kinds of plots that we may make when visualizing the distribution of data.

First, let's generate some data to help us see what these concepts are all about.

Note that We are generating data to make things easier on ourselves - we know the underlying data-generating process. So that lets us form an expectation of what *should* happen. These kinds of toy examples can be helpful for building an intuition about these concepts that we can then use in our own work.

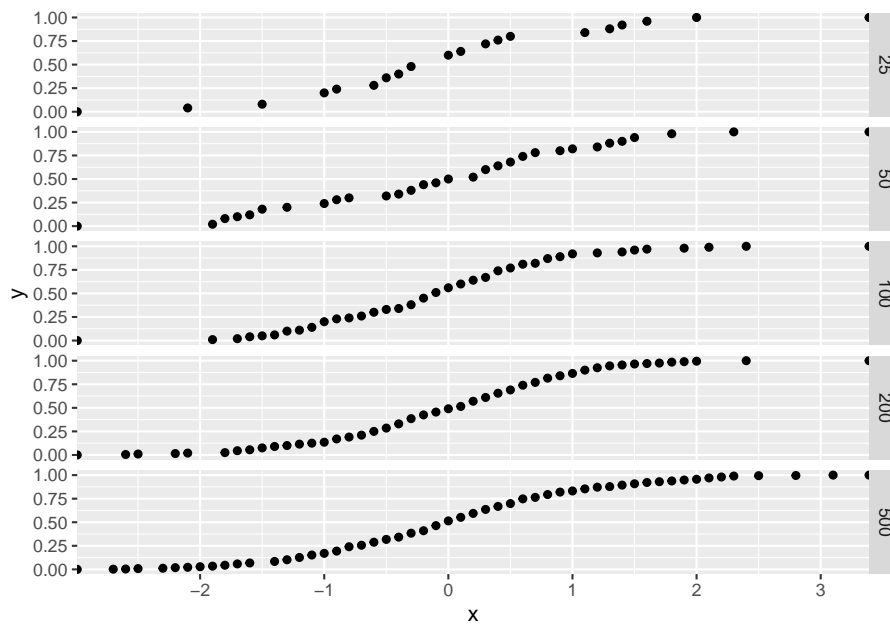
```
df_1 <- tibble(x = round(rnorm(n = 25, mean = 0, sd = 1), digits = 1), samp_size = rep(
df_2 <- tibble(x = round(rnorm(n = 50, mean = 0, sd = 1), digits = 1), samp_size = rep(
df_3 <- tibble(x = round(rnorm(n = 100, mean = 0, sd = 1), digits = 1), samp_size = rep(
df_4 <- tibble(x = round(rnorm(n = 200, mean = 0, sd = 1), digits = 1), samp_size = rep(
df_5 <- tibble(x = round(rnorm(n = 500, mean = 0, sd = 1), digits = 1), samp_size = rep(
```

Combine the generated data together into one data frame to make it easy to plot everything in one ggplot call.

```
sample_df <- bind_rows(df_1, df_2, df_3, df_4, df_5)
```

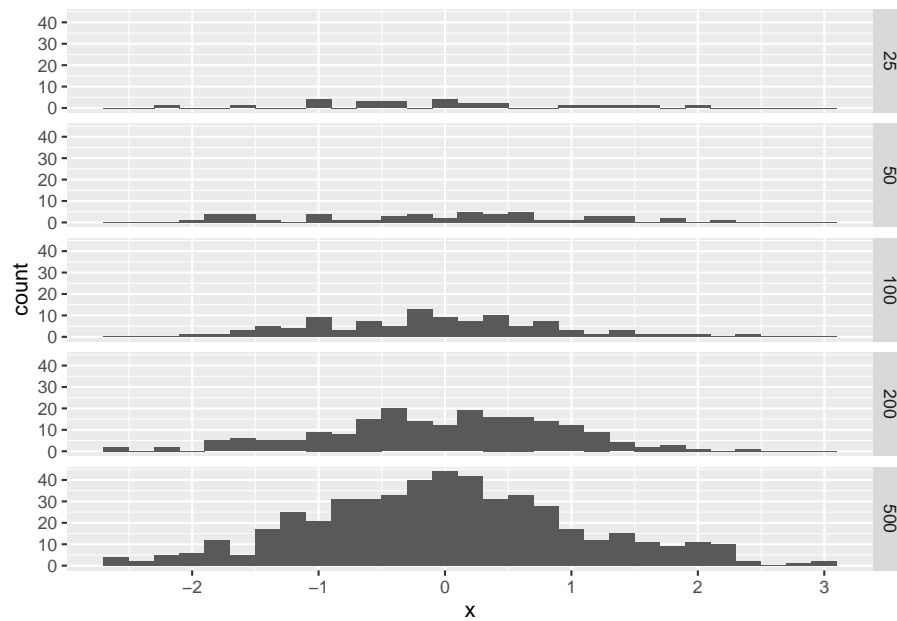
Look at the empirical CDFs

```
sample_df %>%
  ggplot(aes(x)) +
  stat_ecdf(geom = "point") +
  facet_grid(samp_size~.)
```



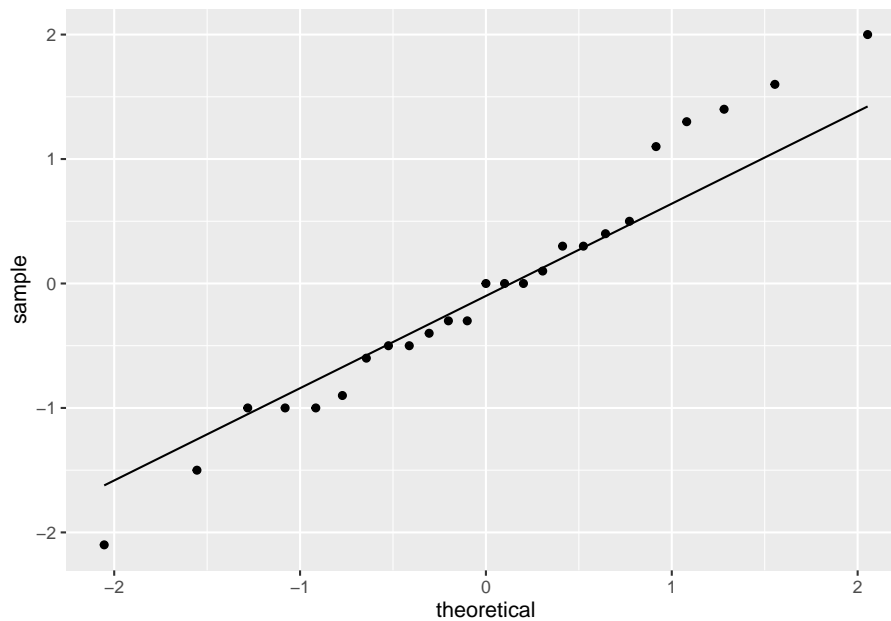
```
sample_df %>%
  ggplot(aes(x)) +
  geom_histogram() +
  facet_grid(samp_size~.)
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



In class, this should lead to a discussion about how to find quantiles now let's try to look at some QQ plots.

```
df_1 %>%  
  ggplot(aes(sample=x)) +  
  stat_qq() +  
  stat_qq_line()
```

Let's try to wrap our heads around what is going on by looking at sorted versions of `df_1` values

```
head(df_1 %>% arrange(x))
```

```
## # A tibble: 6 x 2
##       x samp_size
##   <dbl>   <dbl>
## 1 -2.1     25
## 2 -1.5     25
## 3 -1       25
## 4 -1       25
## 5 -1       25
## 6 -0.9     25
```

Let's also look at the other end of the dataframe. We can do this either using `arrange(desc(x))` or `tail` instead of `head`. Confirm that they're basically two ways of accomplishing the same objective.

```
head(df_1 %>% arrange(desc(x)))
```

```
## # A tibble: 6 x 2
##       x samp_size
##   <dbl>   <dbl>
```

```
## 1  2      25
## 2  1.6    25
## 3  1.4    25
## 4  1.3    25
## 5  1.1    25
## 6  0.5    25
```

```
tail(df_1 %>% arrange(x))
```

```
## # A tibble: 6 x 2
##       x samp_size
##   <dbl>   <dbl>
## 1  0.5     25
## 2  1.1     25
## 3  1.3     25
## 4  1.4     25
## 5  1.6     25
## 6  2       25
```

End detour

6.2.3 Generating data for student happiness exercise.

We are going to student happiness scores on a scale from 0 to 100 as a function of their time spent outdoors, time spent on Zoom, class standing, and department.

```
pop_size <- 5000
sample_size <- 200

pop_df <- tibble(standing = rep(c("undergrad", "masters", "phd"),
                                times = c(pop_size/2, pop_size/4, pop_size/4)))

pop_df$discipline <- rep(c("mechanical", "civil", "electrical"),
                          length.out = pop_size)

pop_df$min_outdoors <- round(runif(n = pop_size, min = 0, max = 300), 0)

pop_df$min_zoom <- round(runif(n = pop_size, min = 0, max = 400), 0)

pop_df <- pop_df %>%
  mutate(undergrad = case_when(standing == "undergrad" ~ 1,
                                TRUE ~ 0),
         masters = case_when(standing == "masters" ~ 1,
```

```

      TRUE ~ 0),
  phd = case_when(standing == "phd" ~ 1,
      TRUE ~ 0),
  civil = case_when(discipline == "civil" ~ 1,
      TRUE ~ 0),
  mechanical = case_when(discipline == "mechanical" ~ 1,
      TRUE ~ 0),
  electrical = case_when(discipline == "electrical" ~ 1,
      TRUE ~ 0))

```

With the predictors generated, we can now generate the happiness scores.

```

#pop_df <- pop_df %>%
# mutate(happiness = 50 + -.2*(min_zoom+rnorm(n = 1, min = 10, max = 30)) + .1 * min_outdoors +

b_0 <- 50
b_out <- 0.03
b_zoom <- -0.05
b_civ <- 6
b_mech <- 2
b_ele <- -3
b_under <- 5
b_masters <- -2
b_phd <- 10

pop_df$happiness <- b_0 + b_out*(pop_df$min_outdoors + rnorm(pop_size, 0, 40)) +
  b_zoom*(pop_df$min_zoom + rnorm(pop_size, 0, 40)) +
  b_civ*(pop_df$civil + rnorm(pop_size, 0, 2)) +
  b_mech*pop_df$mechanical +
  b_ele*pop_df$electrical + b_under * pop_df$undergrad +
  b_masters*pop_df$masters +
  b_phd*pop_df$phd + round(rnorm(n = pop_size, mean = 0, sd = 10), 0)

```

We will save the full population dataset. That is what we will start with in class.

```
pop_df %>% write_csv("student_happiness.csv")
```

Now we create a sample of students by sampling from the population dataframe.

```
samp_df <- sample_n(pop_df, size = 200)
```

On the modeling side, we can first start with a very simple model with only one predictor.

```
happiness_zoom <- lm(happiness ~ min_zoom, data = samp_df)
summary(happiness_zoom)
```

```
##
## Call:
## lm(formula = happiness ~ min_zoom, data = samp_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -42.871  -9.296  -0.308   10.089   41.709
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 63.584166   2.248145   28.28 < 2e-16 ***
## min_zoom    -0.065520   0.009677   -6.77 1.42e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 15.08 on 198 degrees of freedom
## Multiple R-squared:  0.188, Adjusted R-squared:  0.1839
## F-statistic: 45.84 on 1 and 198 DF, p-value: 1.425e-10
```

Next, we can look at what happens with a more complex model with multiple predictors. Pay attention to how the R^2 value changes.

```
happiness_all <- lm(happiness ~ min_outdoors + min_zoom + civil + mechanical + electrical +
summary(happiness_all)
```

```
##
## Call:
## lm(formula = happiness ~ min_outdoors + min_zoom + civil + mechanical +
##      electrical + undergrad + masters + phd, data = samp_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -45.992 -10.389  -0.302    8.511   37.301
##
## Coefficients: (2 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 62.94592    3.76431  16.722 < 2e-16 ***
## min_outdoors  0.02556    0.01241   2.060  0.04076 *
## min_zoom     -0.06954    0.00964  -7.214 1.21e-11 ***
## civil         4.66774    2.57563   1.812  0.07150 .
## mechanical    0.96157    2.66793   0.360  0.71893
```

```
## electrical      NA      NA      NA      NA
## undergrad    -4.60262    2.61826   -1.758   0.08035 .
## masters      -7.79982    2.96814   -2.628   0.00928 **
## phd           NA      NA      NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14.76 on 193 degrees of freedom
## Multiple R-squared:  0.2418, Adjusted R-squared:  0.2183
## F-statistic: 10.26 on 6 and 193 DF,  p-value: 7.479e-10
```

An alternative way of writing the full model

```
happiness_all <- lm(happiness ~ min_outdoors + min_zoom + discipline + standing, data = samp_df)
summary(happiness_all)
```

```
##
## Call:
## lm(formula = happiness ~ min_outdoors + min_zoom + discipline +
##     standing, data = samp_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -45.992 -10.389  -0.302   8.511  37.301
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    59.81384     3.61923   16.527 < 2e-16 ***
## min_outdoors     0.02556     0.01241    2.060  0.04076 *
## min_zoom        -0.06954     0.00964   -7.214 1.21e-11 ***
## disciplineelectrical -4.66774     2.57563   -1.812  0.07150 .
## disciplinemechanical -3.70617     2.51991   -1.471  0.14299
## standingphd       7.79982     2.96814    2.628  0.00928 **
## standingundergrad    3.19720     2.53864    1.259  0.20940
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14.76 on 193 degrees of freedom
## Multiple R-squared:  0.2418, Adjusted R-squared:  0.2183
## F-statistic: 10.26 on 6 and 193 DF,  p-value: 7.479e-10
```


Chapter 7

Week 7: Logistic Regression

To study how to run a logistic regression model in R, we are going to create simulate some data.

NOTE: Simulating data is a good way to get a sense of how the model works since you already know the data-generate process that you are trying to characterize with the model. This means you already know the “ground truth”. In practice, we do not have this luxury when we run our studies - it is this “ground truth” that we’re looking for. Nonetheless, it can be helpful to remove one piece of uncertainty when learning the technical aspects by using data that are already characterized because you made them.

We will proceed through several rounds of data generation to see how the model may vary.

7.1 Round 1 - No systematic variation in outcomes

```
disciplines <- c("civil", "mechanical", "electrical", "systems")
disciplines_prob <- c(0.25, 0.3, 0.35, 0.1)

know_engineer <- c("immediate_fam", "distant_fam", "friend", "none")
know_engineer_prob <- c(0.2, 0.4, 0.3, 0.1)

persistence <- c("yes", "no")
persistence_prob <- c(0.8, 0.2)
```

Now we will generate the actual sample of 500 students

```
samp_size <- 500

disc_samp <- sample(x = disciplines, size = samp_size, prob = disciplines_prob, replace = TRUE)
know_samp <- sample(x = know_engineer, size = samp_size, prob = know_engineer_prob, replace = TRUE)
pers_samp <- sample(x = persistence, size = samp_size, prob = persistence_prob, replace = TRUE)

samp_df <- tibble(discipline = disc_samp,
                  know_eng = know_samp,
                  persist = pers_samp,
                  gpa = round(rnorm(n = samp_size, mean = 3, sd = 0.3), 2))
```

Up to now, we have simulated the data collection process. This is the point where we would typically be cleaning the data and starting our analysis

First, add in a new binary column for the persistence variable (coding “yes” as 1 and “no” as 0)

```
samp_df <- samp_df %>%
  mutate(persist_bin = case_when(persist == "yes" ~ 1,
                                  persist == "no" ~ 0))
```

Second, get a sense of the distributions of some of the variables in our model

```
str(samp_df)
```

```
## tibble [500 x 5] (S3: tbl_df/tbl/data.frame)
## $ discipline : chr [1:500] "mechanical" "electrical" "civil" "electrical" ...
## $ know_eng   : chr [1:500] "friend" "friend" "distant_fam" "none" ...
## $ persist    : chr [1:500] "yes" "yes" "yes" "no" ...
## $ gpa        : num [1:500] 3.1 3.37 3.32 2.97 3.26 3.46 3.16 3.18 3.06 3.41 ...
## $ persist_bin: num [1:500] 1 1 1 0 1 0 1 1 1 1 ...
```

```
table(samp_df$persist)
```

```
##
## no yes
## 85 415
```

```
describe(samp_df)
```



```
## Warning in describe(samp_df): NAs introduced by coercion
## Warning in describe(samp_df): NAs introduced by coercion
## Warning in describe(samp_df): NAs introduced by coercion
## Warning in FUN(newX[, i], ...): no non-missing arguments to min; returning Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to min; returning Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to min; returning Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to max; returning -Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to max; returning -Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to max; returning -Inf

##          vars   n mean   sd median trimmed  mad  min  max range  skew
## discipline*   1 500  NaN   NA     NA      NaN   NA  Inf -Inf  -Inf   NA
## know_eng*     2 500  NaN   NA     NA      NaN   NA  Inf -Inf  -Inf   NA
## persist*      3 500  NaN   NA     NA      NaN   NA  Inf -Inf  -Inf   NA
## gpa           4 500 3.02 0.31   3.02   3.02 0.31 2.02 4.09  2.07 -0.05
## persist_bin   5 500 0.83 0.38   1.00   0.91 0.00 0.00 1.00  1.00 -1.75
##          kurtosis   se
## discipline*      NA   NA
## know_eng*        NA   NA
## persist*         NA   NA
## gpa              0.21 0.01
## persist_bin      1.07 0.02

table(samp_df$discipline)

##
##      civil electrical mechanical      systems
##      128          176          147          49

xtabs(~ discipline + persist_bin, data = samp_df)

##          persist_bin
## discipline    0    1
##   civil      23 105
##   electrical 29 147
##   mechanical 23 124
##   systems    10  39
```

Third, model the outcome (persistence) as a function of three predictor variables (discipline, knowing an engineering, and gpa)

```
model <- glm(persist_bin ~ discipline + know_eng + gpa, data = samp_df, family = binomial)
summary(model)
```

```
##
## Call:
## glm(formula = persist_bin ~ discipline + know_eng + gpa, family = binomial(),
##      data = samp_df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.1661    0.4728    0.6094    0.6462    0.8378
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      0.802358   1.174157   0.683   0.4944
## disciplineelectrical 0.110177   0.309131   0.356   0.7215
## disciplinemechanical 0.158249   0.325697   0.486   0.6271
## disciplinesystems  -0.101610   0.426971  -0.238   0.8119
## know_engfriend      0.571152   0.329111   1.735   0.0827
## know_engimmediate_fam 0.004592   0.313112   0.015   0.9883
## know_engnone        -0.345766   0.366554  -0.943   0.3455
## gpa                 0.206253   0.383967   0.537   0.5912
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 455.89  on 499  degrees of freedom
## Residual deviance: 449.32  on 492  degrees of freedom
## AIC: 465.32
##
## Number of Fisher Scoring iterations: 4
```

```
tidy(model)
```

```
## # A tibble: 8 x 5
##   term                estimate std.error statistic p.value
##   <chr>              <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)        0.802      1.17     0.683    0.494
## 2 disciplineelectrical 0.110      0.309    0.356    0.722
## 3 disciplinemechanical 0.158      0.326    0.486    0.627
```

```
## 4 disciplinesystems      -0.102      0.427    -0.238    0.812
## 5 know_engfriend        0.571      0.329     1.74     0.0827
## 6 know_engimmediate_fam 0.00459    0.313    0.0147    0.988
## 7 know_engnone         -0.346     0.367   -0.943     0.346
## 8 gpa                   0.206     0.384    0.537     0.591
```

Since we generated the data without any relationships between the predictors and the binary outcome, we do not expect any of the predictor variables to be statistically significant. Sometimes there will be a significant predictor, but that emphasizes the idea p values are not necessarily the most reliable indicator of importance.

7.2 Round 2 - Systematic variation in outcomes as a function of discipline

Now, to introduce some systematic variation, we will change the probabilities of persistence (the outcome variable we are modeling) as a function of major but not as a function of the other two predictors

```
disciplines <- c("civil", "mechanical", "electrical", "systems")
disciplines_prob <- c(0.25, 0.3, 0.35, 0.1)

know_engineer <- c("immediate_fam", "distant_fam", "friend", "none")
know_engineer_prob <- c(0.2, 0.4, 0.3, 0.1)
```

We generate the actual data here.

```
samp_size <- 5000

disc_samp_2 <- sample(x = disciplines, size = samp_size, prob = disciplines_prob, replace = TRUE)
know_samp_2 <- sample(x = know_engineer, size = samp_size, prob = know_engineer_prob, replace = TRUE)
#pers_samp_2 <- sample(x = persistence, size = samp_size, prob = persistence_prob, replace = TRUE)
```

Now we will combine our data into one dataframe.

```
samp_df_2 <- tibble(discipline = disc_samp_2,
                    know_eng = know_samp_2,
                    gpa = round(rnorm(n = samp_size, mean = 3, sd = 0.3), 2))
```

Now we want to have some different outcomes whose probabilities vary by discipline. We'll create a new column called `persist_prob` that describes the probability of persisting from year one to year two (e.g., 0.6 means there is a 0.6 prob of a student persisting)

```
samp_df_2 <- samp_df_2 %>%
  mutate(persist_prob = case_when(discipline == "civil" ~ 0.6,
                                   discipline == "mechanical" ~ 0.7,
                                   discipline == "electrical" ~ 0.8,
                                   discipline == "systems" ~ 0.9))
```

We will create a vector that samples depending on the value of the persistence probability at that index. That value varies depending on the discipline for that student at that index value in the vector.

```
persist_outcome <- modify(.x = samp_df_2$persist_prob, .f = ~rbinom(n = 1, size = 1, p =
```

Now we add that persistence outcome column to our dataframe

```
samp_df_2$persist_bin <- persist_outcome
```

Up to now, we have simulated the data collection process. This is the point where we would typically be cleaning the data and starting our analysis,

```
str(samp_df_2)
```

```
## tibble [5,000 x 5] (S3: tbl_df/tbl/data.frame)
## $ discipline : chr [1:5000] "systems" "civil" "electrical" "civil" ...
## $ know_eng   : chr [1:5000] "friend" "distant_fam" "immediate_fam" "distant_fam"
## $ gpa        : num [1:5000] 2.51 3.41 3.11 3.43 2.78 2.95 3.28 2.94 2.9 2.67 ...
## $ persist_prob: num [1:5000] 0.9 0.6 0.8 0.6 0.8 0.7 0.7 0.6 0.7 0.7 ...
## $ persist_bin : num [1:5000] 1 0 1 0 0 1 1 1 0 1 ...
```

Let's check on the distribution of persistence by major. `xtabs()` is a function that creates a contingency table (more on that in 2 weeks)

```
xtabs(~ persist_bin + discipline, data=samp_df_2)
```

```
##           discipline
## persist_bin civil electrical mechanical systems
##           0   486           353           497           40
##           1   745          1399          1001          479
```

Or we can use `describe()`

7.2. ROUND 2 - SYSTEMATIC VARIATION IN OUTCOMES AS A FUNCTION OF DISCIPLINE 141

```
describe(samp_df_2)
```

```
## Warning in describe(samp_df_2): NAs introduced by coercion
## Warning in describe(samp_df_2): NAs introduced by coercion

## Warning in FUN(newX[, i], ...): no non-missing arguments to min; returning Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to min; returning Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to max; returning -Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to max; returning -Inf

##          vars      n mean  sd median trimmed  mad  min  max range  skew
## discipline*    1 5000  NaN   NA    NA      NaN   NA  Inf -Inf  -Inf   NA
## know_eng*      2 5000  NaN   NA    NA      NaN   NA  Inf -Inf  -Inf   NA
## gpa            3 5000  3.00 0.30   3.0    3.00 0.30 1.65 4.11  2.46 -0.01
## persist_prob   4 5000  0.73 0.10   0.7    0.73 0.15 0.60 0.90  0.30  0.06
## persist_bin    5 5000  0.72 0.45   1.0    0.78 0.00 0.00 1.00  1.00 -1.01
##          kurtosis    se
## discipline*      NA   NA
## know_eng*        NA   NA
## gpa              0.23 0.00
## persist_prob    -1.02 0.00
## persist_bin     -0.99 0.01
```

Or even use `table()`

```
table(samp_df_2$discipline)
```

```
##
##      civil electrical mechanical      systems
##      1231      1752      1498      519
```

Now, model the outcome (persistence) as a function of three predictor variables (discipline, knowing an engineering, and gpa)

```
model_2 <- glm(persist_bin ~ discipline + know_eng + gpa, data = samp_df_2, family = binomial())
```

And examine the model output with either `summary()`...

```
summary(model_2)
```

```
##
## Call:
## glm(formula = persist_bin ~ discipline + know_eng + gpa, family = binomial(),
##      data = samp_df_2)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.3096  -1.3462   0.6673   0.8955   1.0749
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.8906195  0.3302760   2.697 0.007005 **
## disciplineelectrical 0.9502611  0.0834057  11.393 < 2e-16 ***
## disciplinemechanical 0.2710469  0.0801181   3.383 0.000717 ***
## disciplinesystems    2.0517523  0.1746544  11.748 < 2e-16 ***
## know_engfriend      0.0005848  0.0785944   0.007 0.994064
## know_engimmediate_fam -0.0946710  0.0890747  -1.063 0.287860
## know_engnone        -0.0657999  0.1150084  -0.572 0.567233
## gpa                -0.1454997  0.1070743  -1.359 0.174189
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 5883.6  on 4999  degrees of freedom
## Residual deviance: 5594.5  on 4992  degrees of freedom
## AIC: 5610.5
##
## Number of Fisher Scoring iterations: 5
```

```
...or tidy()
```

```
tidy(model_2)
```

```
## # A tibble: 8 x 5
##   term                estimate std.error statistic  p.value
##   <chr>              <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)        0.891       0.330      2.70  7.01e- 3
## 2 disciplineelectrical 0.950       0.0834    11.4  4.52e-30
## 3 disciplinemechanical 0.271       0.0801     3.38  7.17e- 4
## 4 disciplinesystems    2.05        0.175     11.7  7.27e-32
## 5 know_engfriend      0.000585    0.0786    0.00744 9.94e- 1
```


Now, since we want to look at the potential effect of gpa on persistence, we create a bookkeeping column for gpa_mean for students who do and do not persist to year two.

```
samp_df_3 <- samp_df_3 %>%
  mutate(gpa_mean = case_when(persist == "yes" ~ 3.4,
                              persist == "no" ~ 3.0))
```

Simulate the gpa data

```
gpa_vec <- round(modify(.x = samp_df_3$gpa_mean, .f = ~rnorm(n = 1, mean = .x, sd = .2)
```

Add the simulated data back to our data frame

```
samp_df_3$gpa <- gpa_vec
```

Up to now, we have simulated the data collection process. This is the point where we would typically be cleaning the data and starting our analysis

Check the structure of the dataframe to make sure it looks as expected

```
str(samp_df_3)
```

```
## tibble [500 x 6] (S3: tbl_df/tbl/data.frame)
## $ discipline : chr [1:500] "mechanical" "civil" "mechanical" "mechanical" ...
## $ know_eng   : chr [1:500] "distant_fam" "immediate_fam" "distant_fam" "friend" ..
## $ persist    : chr [1:500] "yes" "no" "yes" "yes" ...
## $ persist_bin: num [1:500] 1 0 1 1 1 1 1 1 1 0 ...
## $ gpa_mean   : num [1:500] 3.4 3 3.4 3.4 3.4 3.4 3.4 3.4 3.4 3 ...
## $ gpa        : num [1:500] 3.39 2.98 3.52 3.47 3.15 3.02 3.59 3.61 3.34 2.9 ...
```

Let's check on the distribution of persistence by major. xtabs() is a function that creates a contingency table (more on that in 2 weeks)

```
xtabs(~ persist_bin + discipline, data=samp_df_3)
```

```
##           discipline
## persist_bin civil electrical mechanical systems
##           0    31         24         34         5
##           1    90        120        145        51
```

...or with describe()...

7.3. ROUND 3 - SYSTEMATIC VARIATION IN OUTCOMES AS A FUNCTION OF DISCIPLINE AND GPA145

```
describe(samp_df_3)
```

```
## Warning in describe(samp_df_3): NAs introduced by coercion
## Warning in describe(samp_df_3): NAs introduced by coercion
## Warning in describe(samp_df_3): NAs introduced by coercion

## Warning in FUN(newX[, i], ...): no non-missing arguments to min; returning Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to min; returning Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to min; returning Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to max; returning -Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to max; returning -Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to max; returning -Inf

##           vars    n mean   sd median trimmed  mad  min  max range  skew
## discipline*    1 500  NaN   NA     NA      NaN   NA  Inf -Inf  -Inf    NA
## know_eng*      2 500  NaN   NA     NA      NaN   NA  Inf -Inf  -Inf    NA
## persist*       3 500  NaN   NA     NA      NaN   NA  Inf -Inf  -Inf    NA
## persist_bin    4 500 0.81 0.39   1.00   0.89 0.00 0.00 1.00 1.00 -1.59
## gpa_mean       5 500 3.32 0.16   3.40   3.36 0.00 3.00 3.40 0.40 -1.59
## gpa            6 500 3.32 0.25   3.33   3.33 0.24 2.61 4.06 1.45 -0.26
##           kurtosis    se
## discipline*        NA   NA
## know_eng*          NA   NA
## persist*           NA   NA
## persist_bin       0.54 0.02
## gpa_mean          0.54 0.01
## gpa               0.01 0.01
```

...or with table().

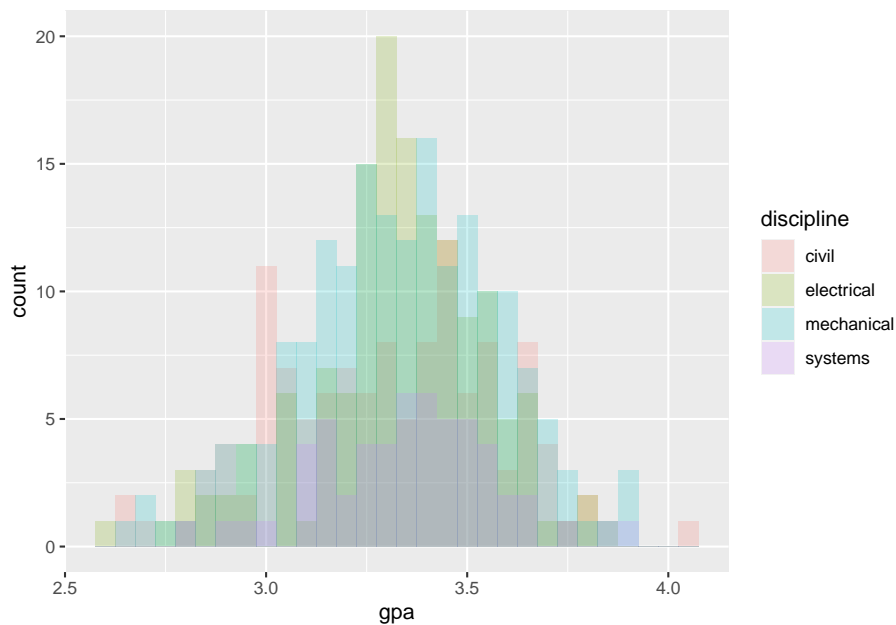
```
table(samp_df_3$discipline)
```

```
##
##      civil electrical mechanical  systems
##      121          144          179       56
```

Check the distribution of the gpa values by discipline, just to make sure

```
samp_df_3 %>%
  ggplot(aes(x = gpa, fill = discipline)) +
  geom_histogram(alpha = 0.2, position = "identity")
```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.

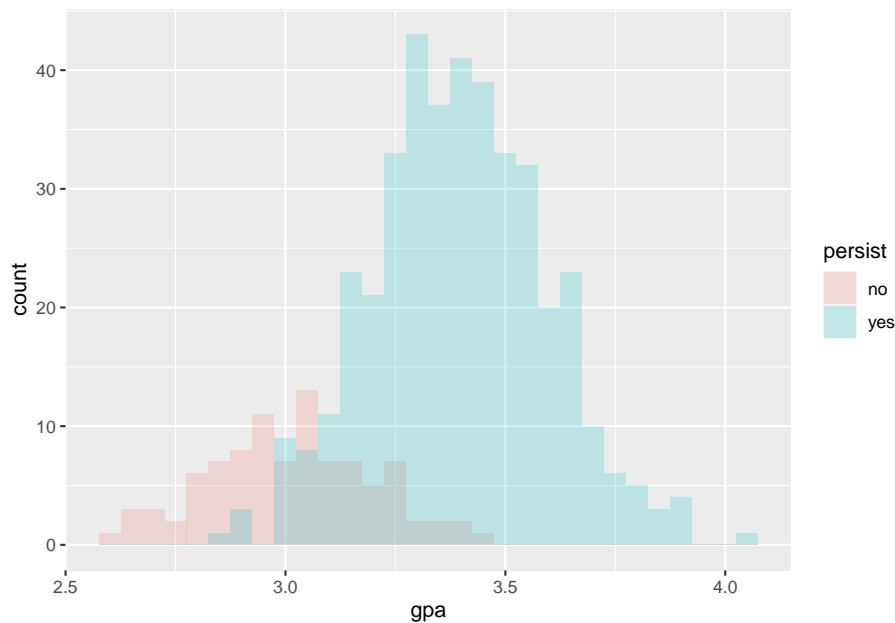


be sure to put alpha before position

```
samp_df_3 %>%
  ggplot(aes(x = gpa, fill = persist)) +
  geom_histogram(alpha = 0.2, position = "identity")
```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.

7.3. ROUND 3 - SYSTEMATIC VARIATION IN OUTCOMES AS A FUNCTION OF DISCIPLINE AND GPA147



Now, model the outcome (persistence) as a function of three predictor variables (discipline, knowing an engineering, and gpa)

```
model_3 <- glm(persist_bin ~ discipline + know_eng + gpa, data = samp_df_3, family = binomial())
```

And examine the results with `summary()` or `tidy()`

```
summary(model_3)
```

```
##
## Call:
## glm(formula = persist_bin ~ discipline + know_eng + gpa, family = binomial(),
##      data = samp_df_3)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.71674   0.05747   0.18177   0.41143   2.22106
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -32.297462   3.419792  -9.444  <2e-16 ***
## disciplineelectrical    0.007722   0.439815   0.018   0.986
## disciplinemechanical    0.110320   0.403782   0.273   0.785
## disciplinesystems     1.075912   0.656979   1.638   0.101
```

```
## know_engfriend      0.320038    0.399080    0.802    0.423
## know_engimmediate_fam -0.021159    0.428994   -0.049    0.961
## know_engnone        0.332457    0.583245    0.570    0.569
## gpa                 10.458703    1.083815    9.650   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 483.31  on 499  degrees of freedom
## Residual deviance: 258.80  on 492  degrees of freedom
## AIC: 274.8
##
## Number of Fisher Scoring iterations: 6
```

```
tidy(model_3)
```

```
## # A tibble: 8 x 5
##   term                estimate std.error statistic  p.value
##   <chr>              <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)      -32.3      3.42    -9.44   3.58e-21
## 2 disciplineelectrical  0.00772    0.440    0.0176 9.86e- 1
## 3 disciplinemechanical  0.110     0.404    0.273  7.85e- 1
## 4 disciplinesystems     1.08     0.657    1.64   1.01e- 1
## 5 know_engfriend       0.320     0.399    0.802  4.23e- 1
## 6 know_engimmediate_fam -0.0212    0.429   -0.0493 9.61e- 1
## 7 know_engnone         0.332     0.583    0.570  5.69e- 1
## 8 gpa                 10.5      1.08     9.65   4.92e-22
```

7.4 Round 4 - GPA and persistence vary by discipline

```
disciplines <- c("civil", "mechanical", "electrical", "systems")
disciplines_prob <- c(0.25, 0.3, 0.35, 0.1)

know_engineer <- c("immediate_fam", "distant_fam", "friend", "none")
know_engineer_prob <- c(0.2, 0.4, 0.3, 0.1)
```

Simulate the data for disciplines, knowing an engineer, and the persistence outcome

```
samp_size <- 500

student_id <- seq(samp_size)
disc_samp_4 <- sample(x = disciplines, size = samp_size, prob = disciplines_prob, replace = TRUE)
know_samp_4 <- sample(x = know_engineer, size = samp_size, prob = know_engineer_prob, replace = TRUE)
pers_samp_4 <- sample(x = persistence, size = samp_size, prob = persistence_prob, replace = TRUE)
```

Combine these all together in `tibble()`.

```
samp_df_4 <- tibble(sid = student_id,
                    discipline = disc_samp_4,
                    know_eng = know_samp_4)
```

Up to now, we have simulated the data collection process. This is the point where we would typically be cleaning the data and starting out analysis

Start the data analysis for logistic regression here

Now, since we want to look at the potential effect of gpa on persistence, we create a bookkeeping column for `gpa_mean` for students who do and do not persist to year two.

```
samp_df_4 <- samp_df_4 %>%
  mutate(gpa_mean = case_when(discipline == "civil" ~ 3.0,
                              discipline == "electrical" ~ 3.15,
                              discipline == "mechanical" ~ 3.3,
                              discipline == "systems" ~ 3.45))
```

Simulate the gpa data.

```
gpa_vec <- round(modify(.x = samp_df_4$gpa_mean, .f = ~rnorm(n = 1, mean = .x, sd = .1)), 2)
```

Add the simulated data back to our data frame

```
samp_df_4$gpa <- gpa_vec
```

We will create a vector that samples depending on the value of the persistence probability at that index.

That value varies depending on the discipline for that student at that index value in the vector.

Now we want to have some different outcomes whose probabilities vary by discipline. We'll create a new column called.

`persist_prob` that describes the probability of persisting from year one to year two (e.g., 0.6 means there is a 0.6 prob of a student persisting)

```
samp_df_4 <- samp_df_4 %>%
  mutate(persist_prob = case_when(discipline == "civil" ~ 0.6,
                                   discipline == "mechanical" ~ 0.7,
                                   discipline == "electrical" ~ 0.8,
                                   discipline == "systems" ~ 0.9))
```

```
persist_outcome <- modify(.x = samp_df_4$persist_prob, .f = ~rbinom(n = 1, size = 1, p = persist_prob))
```

Now we add that persistence outcome column to our dataframe

```
samp_df_4$persist_bin <- persist_outcome
```

Up to now, we have simulated the data collection process. This is the point where we would typically be cleaning the data and starting our analysis

Check the structure of the dataframe to make sure it looks as expected.

```
str(samp_df_4)
```

```
## tibble [500 x 7] (S3: tbl_df/tbl/data.frame)
## $ sid      : int [1:500] 1 2 3 4 5 6 7 8 9 10 ...
## $ discipline : chr [1:500] "civil" "mechanical" "electrical" "systems" ...
## $ know_eng  : chr [1:500] "friend" "distant_fam" "distant_fam" "friend" ...
## $ gpa_mean  : num [1:500] 3 3.3 3.15 3.45 3 3 3.15 3.45 3.15 3.3 ...
## $ gpa       : num [1:500] 3.04 3.17 3 3.51 3.17 3.04 3.38 3.6 3.1 3.14 ...
## $ persist_prob: num [1:500] 0.6 0.7 0.8 0.9 0.6 0.6 0.8 0.9 0.8 0.7 ...
## $ persist_bin : num [1:500] 0 1 1 1 1 0 1 1 1 0 ...
```

Let's check on the distribution of persistence by major. `xtabs()` is a function that creates a contingency table (more on that in 2 weeks).

```
xtabs(~ persist_bin + discipline, data=samp_df_4)
```

```
##           discipline
## persist_bin civil electrical mechanical systems
##           0    46          33          47         3
##           1    64         145         118        44
```

```
describe(samp_df_4)
```

```
## Warning in describe(samp_df_4): NAs introduced by coercion
```

```
## Warning in describe(samp_df_4): NAs introduced by coercion
```

```
## Warning in FUN(newX[, i], ...): no non-missing arguments to min; returning Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to min; returning Inf

## Warning in FUN(newX[, i], ...): no non-missing arguments to max; returning -Inf
## Warning in FUN(newX[, i], ...): no non-missing arguments to max; returning -Inf

##           vars   n  mean    sd median trimmed   mad  min    max  range
## sid           1 500 250.50 144.48 250.50  250.50 185.32 1.00 500.00 499.00
## discipline*   2 500   NaN    NA    NA    NaN    NA   Inf  -Inf  -Inf
## know_eng*     3 500   NaN    NA    NA    NaN    NA   Inf  -Inf  -Inf
## gpa_mean      4 500   3.19   0.14   3.15   3.19   0.22 3.00   3.45   0.45
## gpa           5 500   3.20   0.18   3.20   3.20   0.19 2.74   3.61   0.87
## persist_prob  6 500   0.73   0.09   0.70   0.73   0.15 0.60   0.90   0.30
## persist_bin   7 500   0.74   0.44   1.00   0.80   0.00 0.00   1.00   1.00
##           skew kurtosis   se
## sid           0.00   -1.21 6.46
## discipline*   NA      NA   NA
## know_eng*     NA      NA   NA
## gpa_mean      0.11   -0.87 0.01
## gpa           0.01   -0.46 0.01
## persist_prob  0.04   -0.92 0.00
## persist_bin -1.10   -0.79 0.02
```

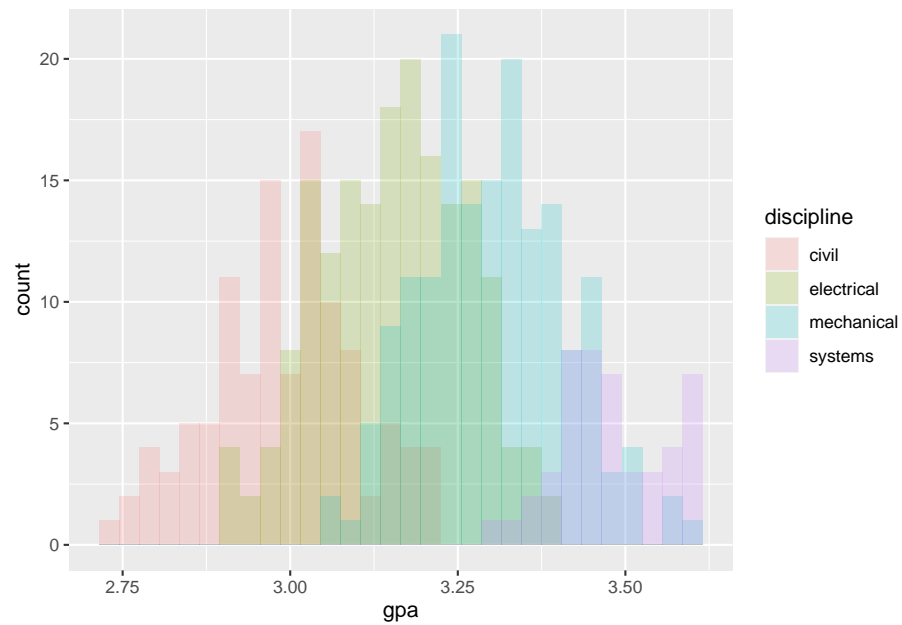
```
table(samp_df_4$discipline)
```

```
##
##      civil electrical mechanical  systems
##      110         178         165         47
```

Check the distribution of the gpa values by discipline, just to make sure.

```
samp_df_4 %>%
  ggplot(aes(x = gpa, fill = discipline)) +
  geom_histogram(alpha = 0.2, position = "identity")
```

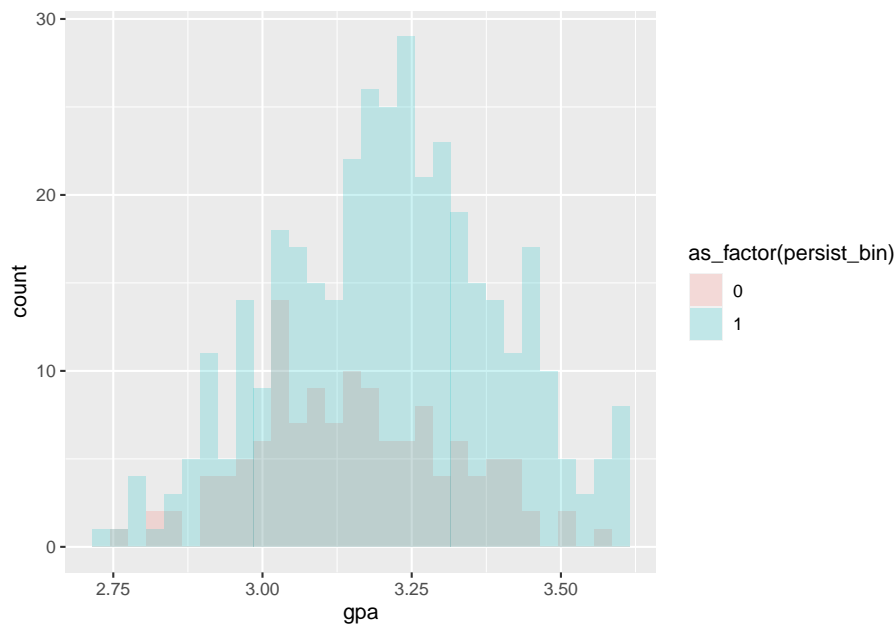
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



- be sure to put alpha before position

```
samp_df_4 %>%
  ggplot(aes(x = gpa, fill = as_factor(persist_bin))) +
  geom_histogram(alpha = 0.2, position = "identity")
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

Now, model the outcome (persistence) as a function of three predictor variables (discipline, knowing an engineering, and gpa)

```
model_4 <- glm(persist_bin ~ discipline + know_eng + gpa, data = samp_df_4, family = binomial())
summary(model_4)
```

```
##
## Call:
## glm(formula = persist_bin ~ discipline + know_eng + gpa, family = binomial(),
##      data = samp_df_4)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.4476  -1.2197   0.6448   0.8259   1.1700
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -0.5614     2.9799  -0.188 0.850568
## disciplineelectrical  1.0919     0.3176   3.438 0.000586 ***
## disciplinemechanical  0.5202     0.4033   1.290 0.197078
## disciplinesystems    2.2443     0.7931   2.830 0.004660 **
## know_engfriend    -0.2586     0.2568  -1.007 0.313871
## know_engimmediate_fam -0.4132     0.2885  -1.432 0.152023
## know_engnone      -0.1745     0.3881  -0.450 0.653007
## gpa              0.3543     0.9955   0.356 0.721925
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 570.95  on 499  degrees of freedom
## Residual deviance: 537.38  on 492  degrees of freedom
## AIC: 553.38
##
## Number of Fisher Scoring iterations: 5
```

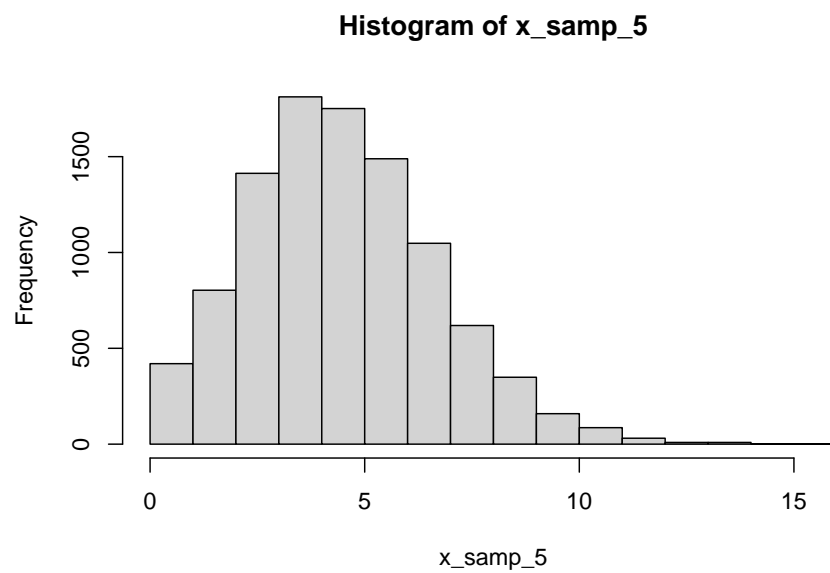
```
tidy(model_4)
```

```
## # A tibble: 8 x 5
##   term                estimate std.error statistic  p.value
##   <chr>                <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)        -0.561       2.98     -0.188  0.851
## 2 disciplineelectrical  1.09       0.318      3.44  0.000586
## 3 disciplinemechanical  0.520       0.403      1.29  0.197
## 4 disciplinesystems    2.24       0.793      2.83  0.00466
## 5 know_engfriend      -0.259       0.257     -1.01  0.314
## 6 know_engimmediate_fam -0.413       0.288     -1.43  0.152
## 7 know_engnone        -0.175       0.388     -0.450 0.653
## 8 gpa                 0.354       0.996      0.356 0.722
```

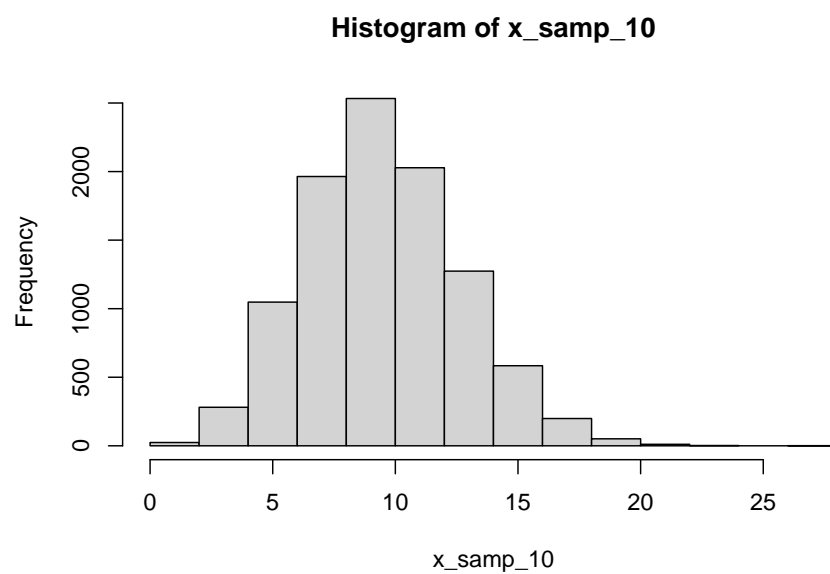
7.4.1 Interlude looking at poisson distribution

```
n <- 10000
set.seed(123)

x_samp_5 <- rpois(n, lambda = 5)
hist(x_samp_5)
```



```
x_samp_10 <- rpois(n, lambda = 10)
hist(x_samp_10)
```

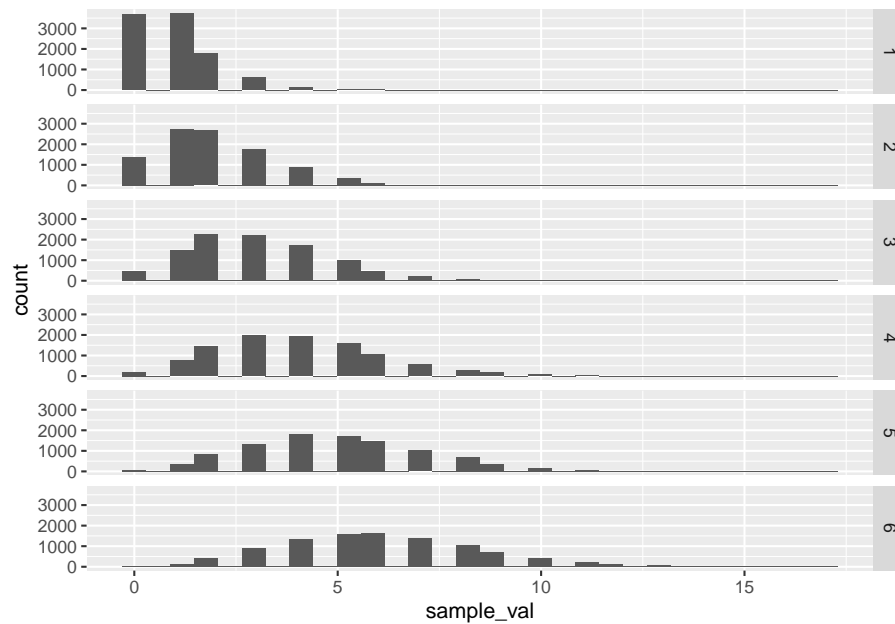


```
param_vect <- rep(c(1, 2, 3, 4, 5, 6), each = n)
samp_vect <- modify(.x = param_vect, .f = ~ rpois(n = 1, lambda = .x))
```

```
samp_df <- tibble(param_val = param_vect,
                  sample_val = samp_vect)
```

```
samp_df %>%
  ggplot(aes(x = sample_val)) +
  geom_histogram() +
  facet_grid(param_val ~ .)
```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.



Chapter 8

Week 8: Comparing two means (t-tests)

This week we focus on comparing two means. This can be the means of an outcome variable collected in two discrete groups (e.g., first-year and second-year students) or the same group across two time points (e.g., pre-test and post-test).

The following contains a number of demos to illustrate how and when we would use t-tests.

8.1 Demo 1 - Comparing salary data for chemical engineering and environmental engineering

First, let's generate some data so that we know the underlying data-generating process that we use statistical tests to characterize.

We'll have a sample of 100 participants.

```
N <- 100

student_id <- seq(N)

group_size <- N/2

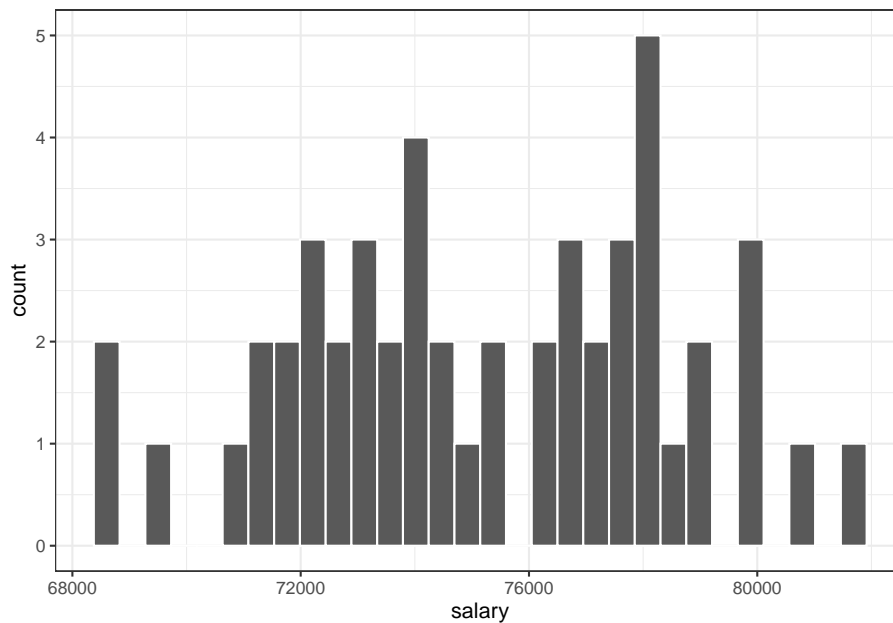
chem_eng_df <- tibble(id = seq(group_size),
                      salary = round(rnorm(n = group_size, mean = 75000, sd = 3000), 2),
```

```

                                discipline = "chemical")
chem_eng_df %>%
  ggplot(aes(x= salary)) +
  geom_histogram(color = "white") +
  theme_bw()

```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.



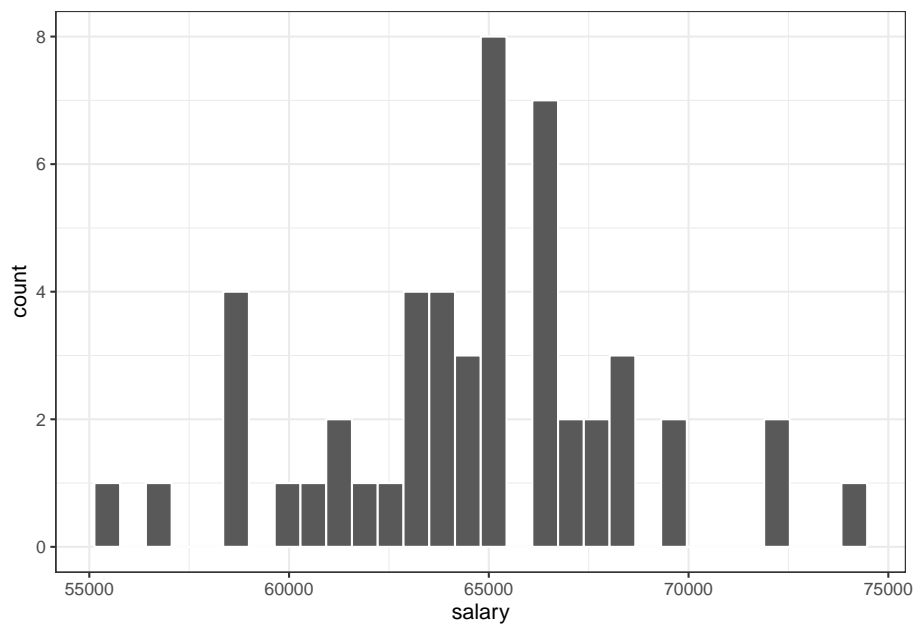
```

env_eng_df <- tibble(id = seq(group_size),
                      salary = round(rnorm(n = group_size, mean = 65000, sd = 4000), 2),
                      discipline = "environmental")
env_eng_df %>%
  ggplot(aes(x= salary)) +
  geom_histogram(color = "white") +
  theme_bw()

```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.

8.1. DEMO 1 - COMPARING SALARY DATA FOR CHEMICAL ENGINEERING AND ENVIRONMENTAL ENGINEERING

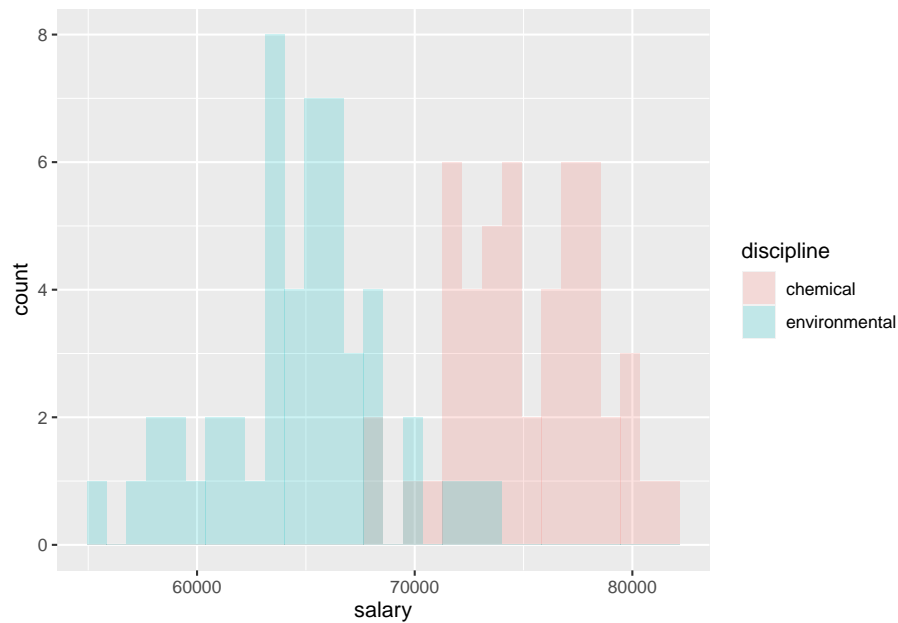


Combine the two disciplinary dataframes together.

```
salary_comb_long <- bind_rows(chem_eng_df, env_eng_df)

salary_comb_long %>%
  ggplot(aes(x = salary, fill = discipline)) +
  geom_histogram(alpha = 0.2, position = "identity")
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Independent samples t-test using the long format approach (NB: the long format is a good practice - it keeps the grouping variable in one column and the continuous outcome variable in a separate column).

```
salary_test <- t.test(salary ~ discipline, data = salary_comb_long, paired = FALSE)
salary_test
```

```
##
##  Welch Two Sample t-test
##
## data: salary by discipline
## t = 14.982, df = 95.279, p-value < 2.2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  9159.492 11957.530
## sample estimates:
##      mean in group chemical mean in group environmental
##                75218.01                64659.50
```

Use the `tidy()` function from the broom package to generate a tibble with the output of the `t.test()`.

```
tidy(salary_test)
```


8.1. DEMO 1 - COMPARING SALARY DATA FOR CHEMICAL ENGINEERING AND ENVIRONMENTAL ENGINEERING

```
## # A tibble: 1 x 10
##   estimate estimate1 estimate2 statistic p.value parameter conf.low conf.high
##   <dbl>      <dbl>      <dbl>      <dbl>    <dbl>      <dbl>      <dbl>      <dbl>
## 1  10559.    75218.    64659.    15.0 8.65e-27    95.3    9159.    11958.
## # ... with 2 more variables: method <chr>, alternative <chr>
```

Remember that a t-test is similar to a linear regression model with a binary predictor variable and continuous outcome variable.

In this case, the model is something like $salary_i = \beta_0 + \beta_1 * discipline_i$

Notice how the value of the t-statistic (t value in the linear model summary table) is close to the t value calculated with the t-test!

```
salary_lm <- lm(salary ~ discipline, data = salary_comb_long)
summary(salary_lm)
```

```
##
## Call:
## lm(formula = salary ~ discipline, data = salary_comb_long)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9496.0 -2100.7   239.4  2288.4  9202.1
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      75218.0      498.3   150.94  <2e-16 ***
## disciplineenvironmental -10558.5      704.7   -14.98  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3524 on 98 degrees of freedom
## Multiple R-squared:  0.6961, Adjusted R-squared:  0.693
## F-statistic: 224.5 on 1 and 98 DF, p-value: < 2.2e-16
```

Calculate the effect size of discipline on salary.

```
(t <- salary_test$statistic[[1]]) # wrapping this in parentheses prints out the value of t at the
```

```
## [1] 14.98229
```

```
(df <- salary_test$parameter[[1]])
```

```
## [1] 95.27933
```

```
(r <- sqrt(t^2/(t^2 + df)))
```

```
## [1] 0.8378649
```

```
round(r,3)
```

```
## [1] 0.838
```

```
r^2
```

```
## [1] 0.7020177
```

We now have the results of an independent samples t-test and the effect size.

Think about what would happen if we changed the sample size or if we changed the underlying distributions that we used to generate the salary data.

We can also run the t-test using a wide format approach

```
salary_comb_wide <- salary_comb_long %>%
  pivot_wider(names_from = discipline, values_from = salary)
```

```
salary_test_2 <- t.test(salary_comb_wide$chemical, salary_comb_wide$environmental, pair
salary_test_2
```

```
##
## Welch Two Sample t-test
##
## data: salary_comb_wide$chemical and salary_comb_wide$environmental
## t = 14.982, df = 95.279, p-value < 2.2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## 9159.492 11957.530
## sample estimates:
## mean of x mean of y
## 75218.01 64659.50
```

8.2 Demo 2 - Student SAT scores

Let's try a second example. This time we'll simulate data we might have where a dependent samples t-test would be appropriate.

We are going to simulate the scenario where we have student SAT scores on both the verbal section and math section.

```
samp_size <- 1000

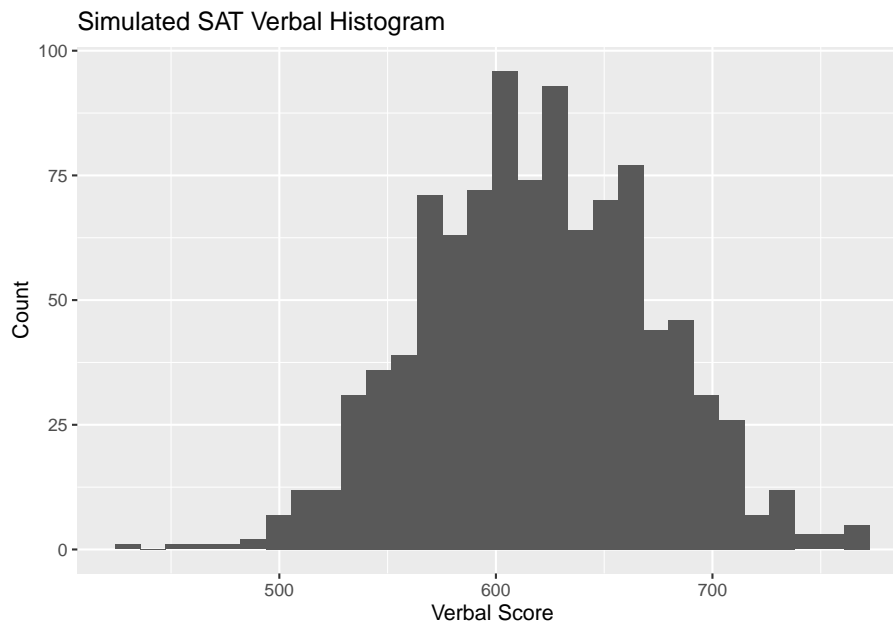
student_id <- seq(samp_size)
sat_verbal <- round(rnorm(n = samp_size, mean = 620, sd = 55), 0)
sat_math <- round(rnorm(n = samp_size, mean = 700, sd = 40), 0)

sat_verbal_capped <- modify(.x = sat_verbal, .f = ~ min(800, .x))
sat_math_capped <- modify(.x = sat_math, .f = ~min(800, .x))

sat_df <- tibble(id = student_id,
                 verbal = sat_verbal_capped,
                 math = sat_math_capped)

sat_df %>% ggplot(aes(x = verbal)) +
  geom_histogram() +
  labs(title = "Simulated SAT Verbal Histogram",
       x = "Verbal Score",
       y = "Count")
```

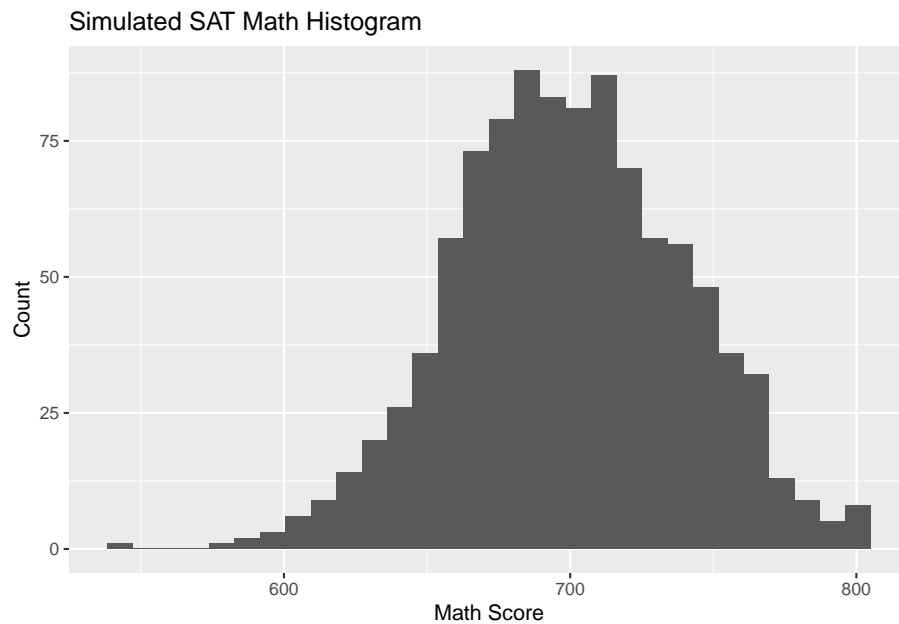
`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.



```
sat_df %>% ggplot(aes(x = math)) +
  geom_histogram() +
```

```
labs(title = "Simulated SAT Math Histogram",
     x = "Math Score",
     y = "Count")
```

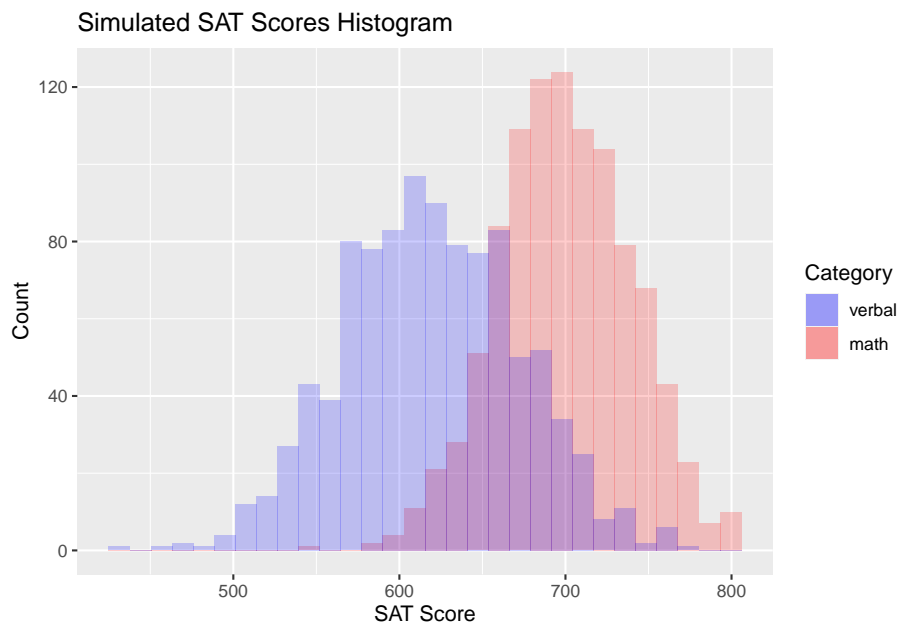
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Alternatively, we could plot these together by either making a long data set or just overlaying histograms and modifying transparency (the alpha parameter).

```
sat_df %>% ggplot() +
  geom_histogram(aes(x = math, fill = "red"), alpha = 0.2) +
  labs(title = "Simulated SAT Scores Histogram",
       x = "SAT Score",
       y = "Count") +
  geom_histogram(aes(x = verbal, fill = "blue"), alpha = 0.2) +
  scale_fill_manual(name = "Category", values = c("blue", "red"), labels = c("verbal",
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Looking at the histogram, we can see that there appears to be a systematic difference in the scores. Running a t-test will help us see if that initial observation is statistically correct.

Since we will be comparing verbal and math scores for each student, this will be a dependent samples (or paired samples) t-test.

```
sat_t_test <- t.test(sat_df$verbal, sat_df$math, paired = TRUE) # don't forget to add: paired = TRUE
# now look at the results
sat_t_test
```

```
##
## Paired t-test
##
## data: sat_df$verbal and sat_df$math
## t = -38.699, df = 999, p-value < 2.2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -84.23739 -76.10661
## sample estimates:
## mean of the differences
## -80.172
```

We see there was a statistically significant difference, but now let's try changing the distributions and see what happens.

First, we'll change make the means very close together and see how sensitive the test is to that

```
samp_size <- 1000

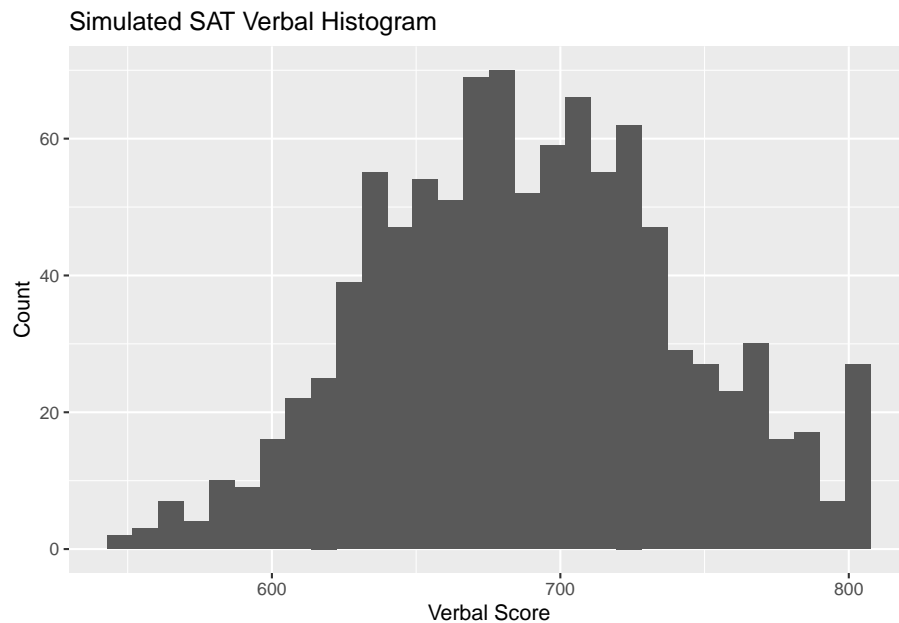
student_id_2 <- seq(samp_size)
sat_verbal_2 <- round(rnorm(n = samp_size, mean = 690, sd = 55), 0)
sat_math_2 <- round(rnorm(n = samp_size, mean = 700, sd = 40), 0)

# let's cap the simulated scores at 800
sat_verbal_2_capped <- modify(.x = sat_verbal_2, .f = ~ min(800, .x))
sat_math_2_capped <- modify(.x = sat_math_2, .f = ~ min(800, .x))

sat_df_2 <- tibble(id = student_id,
                   verbal = sat_verbal_2_capped,
                   math = sat_math_2_capped)

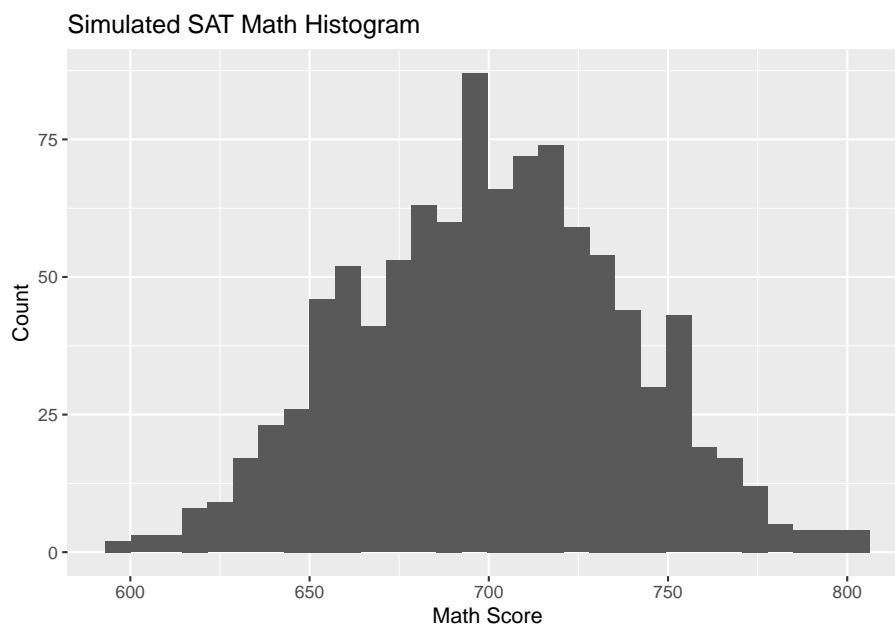
sat_df_2 %>% ggplot(aes(x = verbal)) +
  geom_histogram() +
  labs(title = "Simulated SAT Verbal Histogram",
       x = "Verbal Score",
       y = "Count")
```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.



```
sat_df_2 %>% ggplot(aes(x = math)) +
  geom_histogram() +
  labs(title = "Simulated SAT Math Histogram",
       x = "Math Score",
       y = "Count")
```

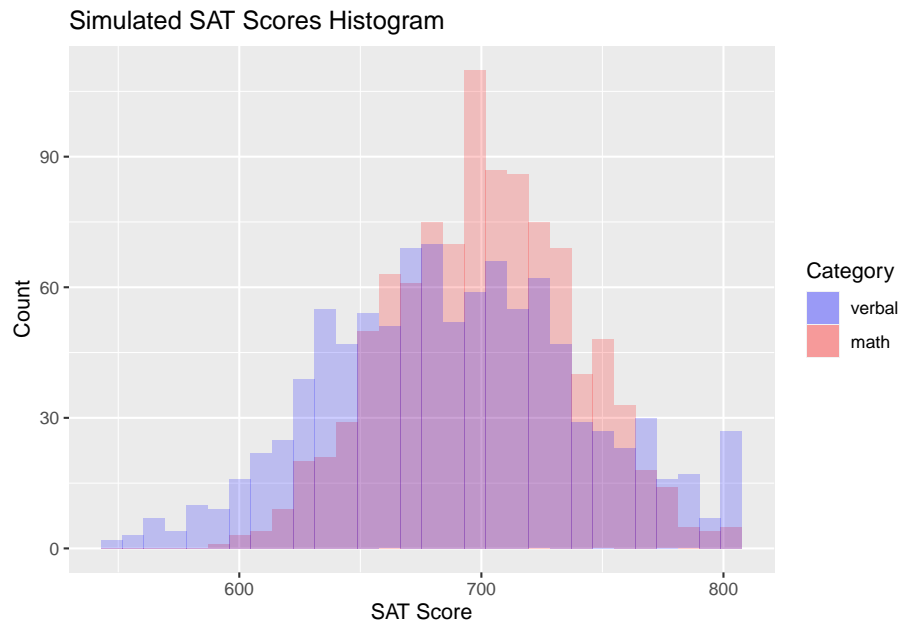
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Alternatively, we could plot these together by either making a long data set or just overlaying histograms and modifying transparency (the alpha parameter).

```
sat_df_2 %>% ggplot() +
  geom_histogram(aes(x = math, fill = "red"), alpha = 0.2) +
  labs(title = "Simulated SAT Scores Histogram",
       x = "SAT Score",
       y = "Count") +
  geom_histogram(aes(x = verbal, fill = "blue"), alpha = 0.2) +
  scale_fill_manual(name = "Category", values = c("blue", "red"), labels = c("verbal", "math"))
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Looking at the histogram, we can see that there appears to be a systematic difference in the scores. Running a t-test will help us see if that initial observation is statistically correct.

Since we will be comparing verbal and math scores for each student, this will be a dependent samples (or paired samples) t-test

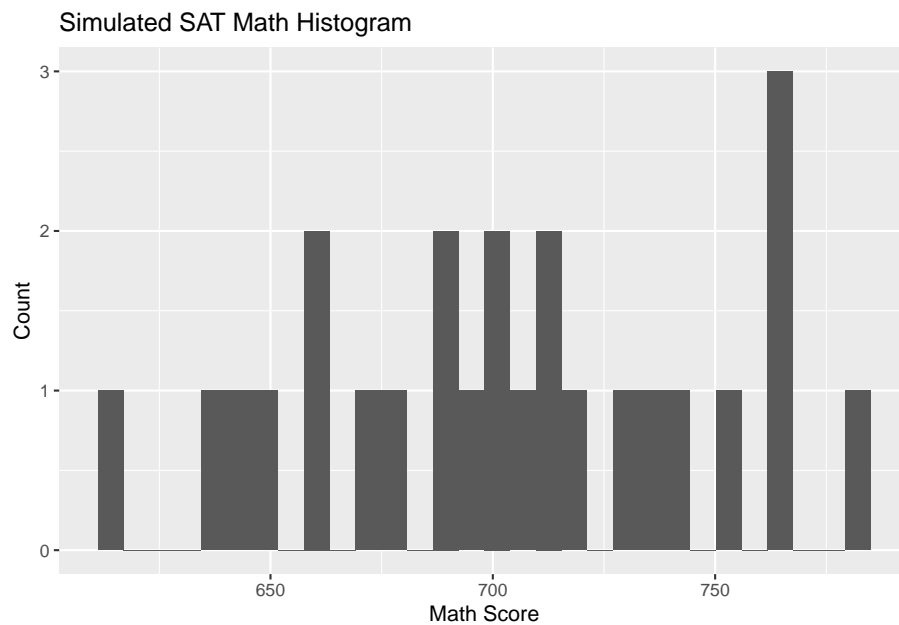
```
sat_t_test_2 <- t.test(sat_df_2$verbal, sat_df_2$math, paired = TRUE) # don't forget t
# now look at the results
sat_t_test_2
```

```
##
## Paired t-test
##
## data: sat_df_2$verbal and sat_df_2$math
## t = -5.7463, df = 999, p-value = 1.21e-08
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -15.724989 -7.719011
## sample estimates:
## mean of the differences
## -11.722
```

Let's try a third test - this time we'll change the sample size to be 25 students rather than 1000


```
sat_df_3 %>% ggplot(aes(x = math)) +
  geom_histogram() +
  labs(title = "Simulated SAT Math Histogram",
       x = "Math Score",
       y = "Count")
```

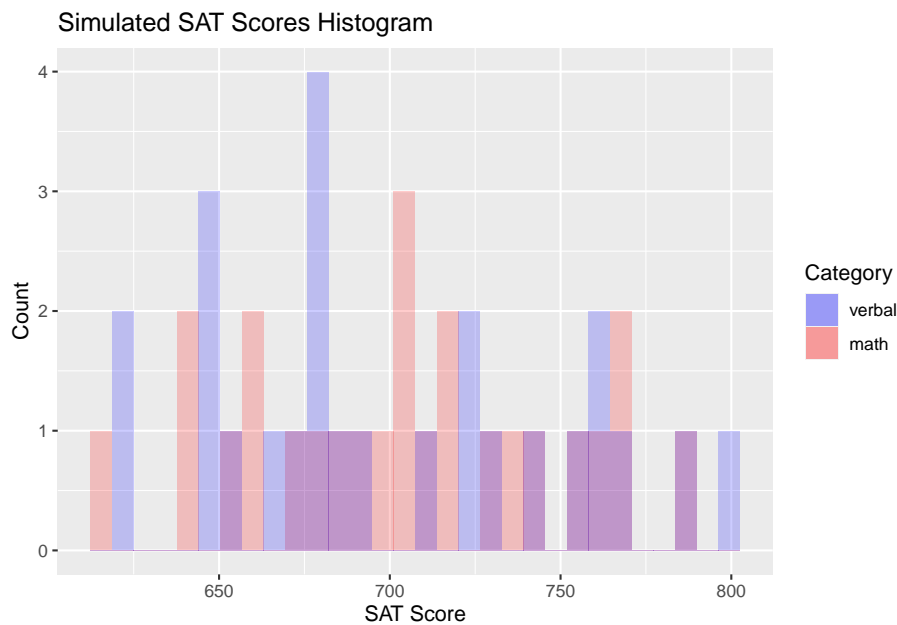
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Alternatively, we could plot these together by either making a long data set or just overlaying histograms and modifying transparency (the alpha parameter)

```
sat_df_3 %>% ggplot() +
  geom_histogram(aes(x = math, fill = "red"), alpha = 0.2) +
  labs(title = "Simulated SAT Scores Histogram",
       x = "SAT Score",
       y = "Count") +
  geom_histogram(aes(x = verbal, fill = "blue"), alpha = 0.2) +
  scale_fill_manual(name = "Category", values = c("blue", "red"), labels = c("verbal",
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Looking at the histogram, we can see that there appears to be a systematic difference in the scores. Running a t-test will help us see if that initial observation is statistically correct. Since we will be comparing verbal and math scores for each student, this will be a dependent samples (or paired samples) t-test

```
sat_t_test_3 <- t.test(sat_df_3$verbal, sat_df_3$math, paired = TRUE) # don't forget to add: paired
# now look at the results
sat_t_test_3
```

```
##
## Paired t-test
##
## data: sat_df_3$verbal and sat_df_3$math
## t = -0.14479, df = 24, p-value = 0.8861
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -30.50846 26.50846
## sample estimates:
## mean of the differences
## -2
```

Now we see that the statistically significant difference is gone. How about if we go back to having a large difference in the means

```

samp_size <- 25

student_id_4 <- seq(samp_size)
sat_verbal_4 <- round(rnorm(n = samp_size, mean = 620, sd = 55), 0)
sat_math_4 <- round(rnorm(n = samp_size, mean = 700, sd = 40), 0)

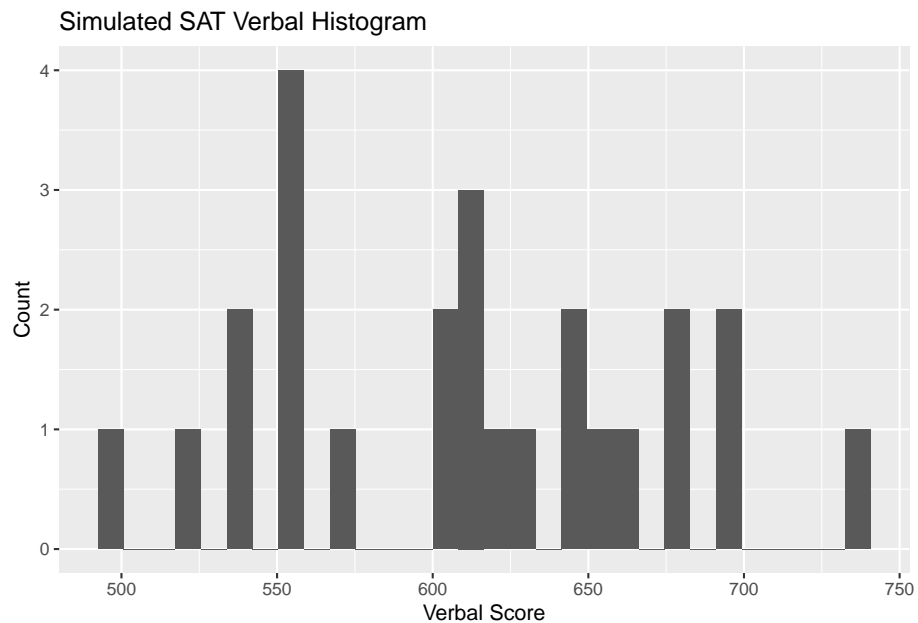
#cap the scores at 800
sat_verbal_4_capped <- modify(.x = sat_verbal_4, .f = ~ min(800, .x))
sat_math_4_capped <- modify(.x = sat_math_4, .f = ~ min(800, .x))

sat_df_4 <- tibble(id = student_id_4,
                   verbal = sat_verbal_4_capped,
                   math = sat_math_4_capped)

# visualize the distributions of scores
sat_df_4 %>% ggplot(aes(x = verbal)) +
  geom_histogram() +
  labs(title = "Simulated SAT Verbal Histogram",
       x = "Verbal Score",
       y = "Count")

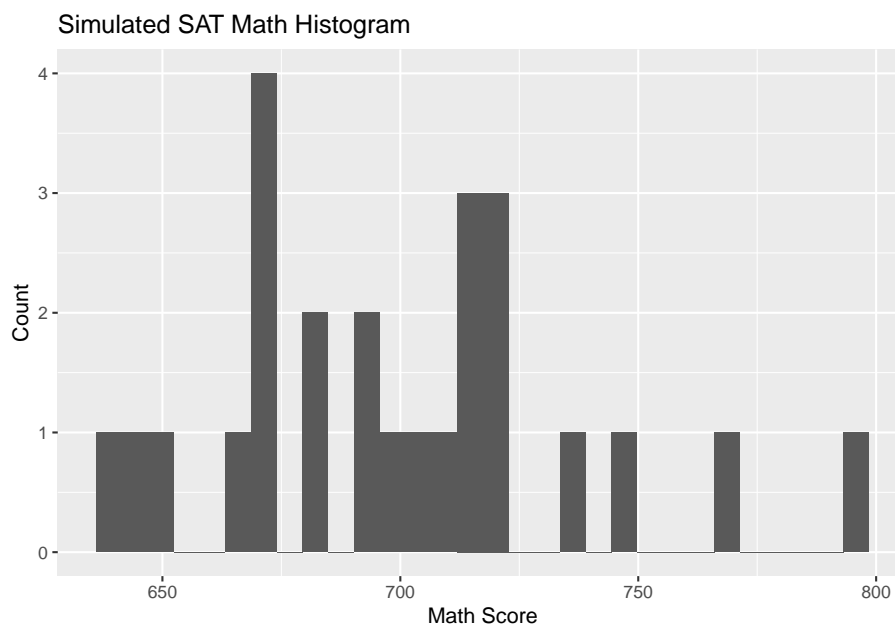
```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.



```
sat_df_4 %>% ggplot(aes(x = math)) +
  geom_histogram() +
  labs(title = "Simulated SAT Math Histogram",
       x = "Math Score",
       y = "Count")
```

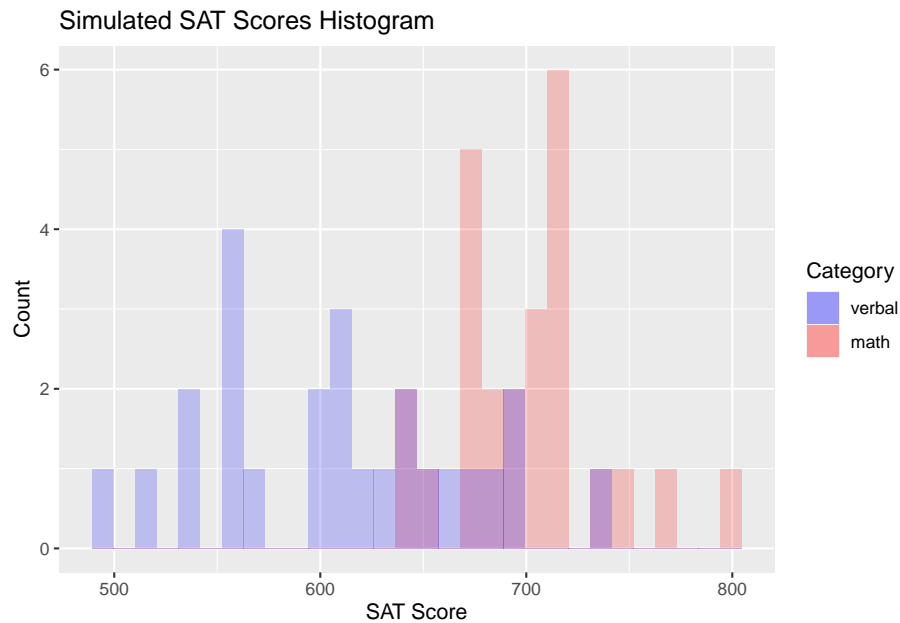
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Alternatively, we could plot these together by either making a long data set or just overlaying histograms and modifying transparency (the alpha parameter)

```
sat_df_4 %>% ggplot() +
  geom_histogram(aes(x = math, fill = "red"), alpha = 0.2) +
  labs(title = "Simulated SAT Scores Histogram",
       x = "SAT Score",
       y = "Count") +
  geom_histogram(aes(x = verbal, fill = "blue"), alpha = 0.2) +
  scale_fill_manual(name = "Category", values = c("blue", "red"), labels = c("verbal", "math"))
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Looking at the histogram, we can see that there appears to be a systematic difference in the scores. Running a t-test will help us see if that initial observation is statistically correct. Since we will be comparing verbal and math scores for each student, this will be a dependent samples (or paired samples) t-test

```
sat_t_test_4 <- t.test(sat_df_4$math, sat_df_4$verbal, paired = TRUE) # don't forget t
# now look at the results
sat_t_test_4
```

```
##
## Paired t-test
##
## data: sat_df_4$math and sat_df_4$verbal
## t = 7.1392, df = 24, p-value = 2.226e-07
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## 63.8111 115.7089
## sample estimates:
## mean of the differences
## 89.76
```

We should see that the t-test will pick up differences between means if there is a sufficiently large difference.