

# Sports Data Analysis and Visualization

Code, data, visuals and the Tidyverse for journalists and other  
storytellers

By Matt Waite

July 29, 2019



# Contents



# Chapter 1

## Throwing cold water on hot takes

The 2018 season started out disastrously for the Nebraska Cornhuskers. The first game against a probably overmatched opponent? Called on account of an epic thunderstorm that plowed right over Memorial Stadium. The next game? Loss. The one following? Loss. The next four? All losses, after the fanbase was whipped into a hopeful frenzy by the hiring of Scott Frost, national title winning quarterback turned hot young coach come back home to save a mythical football program from the mediocrity it found itself mired in.

All that excitement lay in tatters.

On sports talk radio, on the sports pages and across social media and cafe conversations, one topic kept coming up again and again to explain why the team was struggling: Penalties. The team was just committing too many of them. In fact, six games and no wins into the season, they were dead last in the FBS penalty yards.

Worse yet for this line of reasoning? Nebraska won game 7, against Minnesota, committing only six penalties for 43 yards, just about half their average over the season. Then they won game 8 against FCS patsy Bethune Cookman, committing only five penalties for 35 yards. That's a whopping 75 yards less than when they were losing. See? Cut the penalties, win games screamed the radio show callers.

The problem? It's not true. Penalties might matter for a single drive. They may even throw a single game. But if you look at every top-level college football team since 2009, the number of penalty yards the team racks up means absolutely nothing to the total number of points they score. There's no relationship between them. Penalty yards have no discernible influence on points beyond just random noise.

Put this another way: If you were Scott Frost, and a major college football program was paying you \$5 million a year to make your team better, what should you focus on in practice? If you had growled at some press conference that you're going to work on penalties in practice until your team stops committing them, the results you'd get from all that wasted practice time would be impossible to separate from just random chance. You very well may reduce your penalty yards and still lose.

How do I know this? Simple statistics.

That's one of the three pillars of this book: Simple stats. The three pillars are:

1. Simple, easy to understand statistics ...
2. ... extracted using simple code ...
3. ... visualized simply to reveal new and interesting things in sports.

Do you need to be a math whiz to read this book? No. I'm not one either. What we're going to look at is pretty basic, but that's also why it's so powerful.

Do you need to be a computer science major to write code? Nope. I'm not one of those either. But anyone can think logically, and write simple code that is repeatable and replicable.

Do you need to be an artist to create compelling visuals? I think you see where this is going. No. I can barely draw stick figures, but I've been paid to make graphics in my career. With a little graphic design know how, you can create publication worthy graphics with code.

## 1.1 Requirements and Conventions

This book is all in the R statistical language. To follow along, you'll do the following:

1. Install the R language on your computer. Go to the R Project website, click download R and select a mirror closest to your location. Then download the version for your computer.
2. Install R Studio Desktop. The free version is great.

Going forward, you'll see passages like this:

```
install.packages("tidyverse")
```

Don't do it now, but that is code that you'll need to run in your R Studio. When you see that, you'll know what to do.

## 1.2 About this book

This book is the collection of class materials for the author's Sports Data Analysis and Visualization class at the University of Nebraska-Lincoln's College of

Journalism and Mass Communications. There's some things you should know about it:

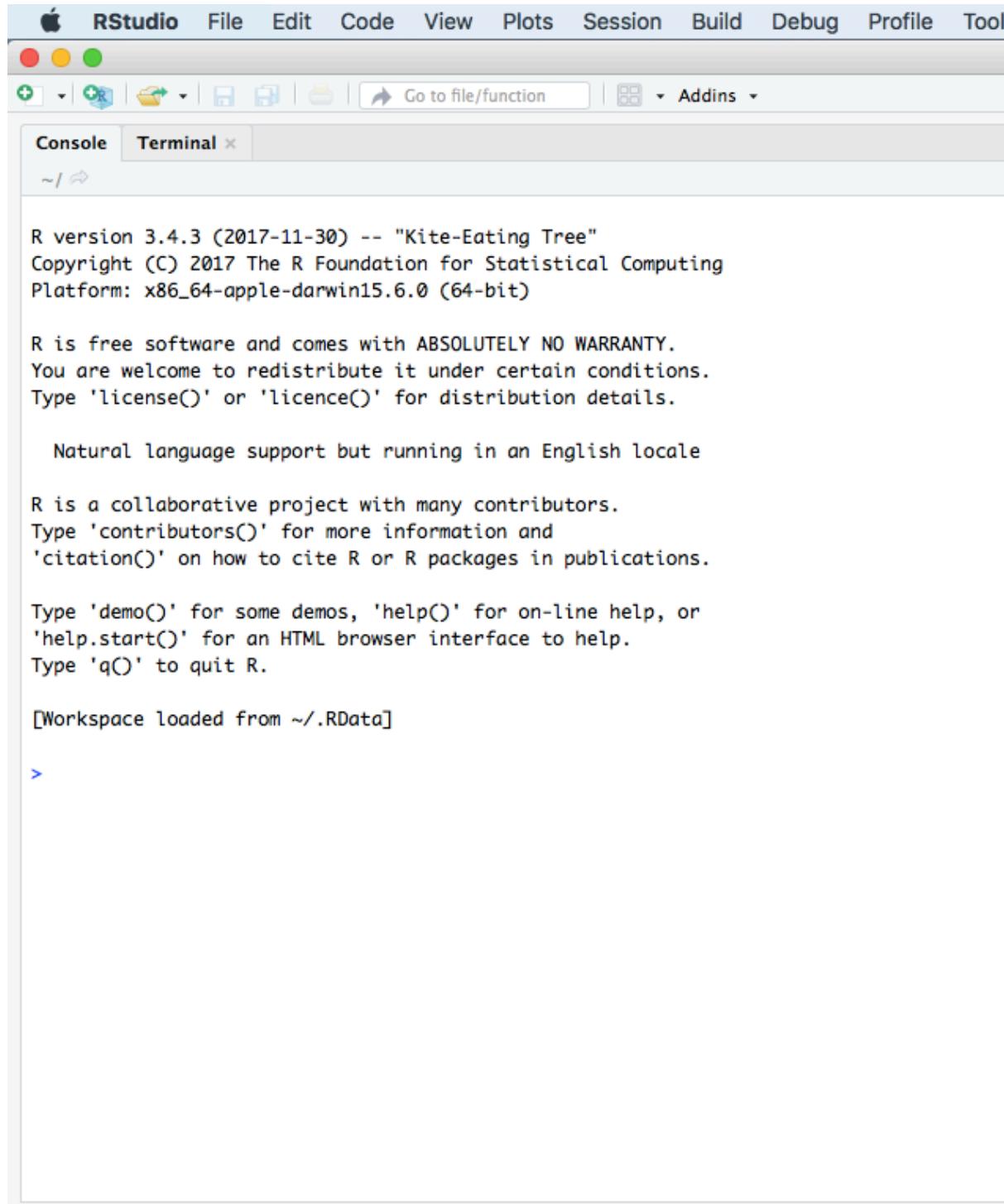
- It is free for students.
- The topics will remain the same but the text is going to be constantly tinkered with.
- What is the work of the author is copyright Matt Waite 2019.
- The text is Attribution-NonCommercial-ShareAlike 4.0 International Creative Commons licensed. That means you can share it and change it, but only if you share your changes with the same license and it cannot be used for commercial purposes. I'm not making money on this so you can't either.
- As such, the whole book – authored in Bookdown – is open sourced on Github. Pull requests welcomed!



# **Chapter 2**

## **The very basics**

R is a programming language, one specifically geared toward statistical analysis. Like all programming languages, it has certain built-in functions and you can interact with it in multiple ways. The first, and most basic, is the console.



The screenshot shows the RStudio interface with the R console tab selected. The console window displays the standard R startup message, including the version number (R version 3.4.3), copyright information, platform details, and various informational messages about the software's nature as free software and its collaborative development. The message concludes with the workspace loading status and a prompt for further interaction.

```
R version 3.4.3 (2017-11-30) -- "Kite-Eating Tree"
Copyright (C) 2017 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin15.6.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Workspace loaded from ~/.RData]

>
```

Think of the console like talking directly to R. It's direct, but it has some drawbacks and some quirks we'll get into later. For now, try typing this into the console and hit enter:

**2+2**

```
## [1] 4
```

Congrats, you've run some code. It's not very complex, and you knew the answer before hand, but you get the idea. We can compute things. We can also store things. **In programming languages, these are called variables.** We can assign things to variables using `<-`. And then we can do things with them. **The `<-` is a called an assignment operator.**

```
number <- 2
```

```
number * number
```

```
## [1] 4
```

Now assign a different number to the variable `number`. Try run `number * number` again. Get what you expected?

We can have as many variables as we can name. **We can even reuse them (but be careful you know you're doing that or you'll introduce errors).** Try this in your console.

```
firstnumber <- 1
secondnumber <- 2
```

```
(firstnumber + secondnumber) * secondnumber
```

```
## [1] 6
```

**We can store anything in a variable.** A whole table. An array of numbers. A single word. A whole book. All the books of the 18th century. They're really powerful. We'll explore them at length.

## 2.1 Adding libraries, part 1

The real strength of any given programming language is the external libraries that power it. The base language can do a lot, but it's the external libraries that solve many specific problems – even making the base language easier to use.

For this class, we're going to need several external libraries.

The first library we're going to use is called Swirl. So in the console, type `install.packages('swirl')` and hit enter. That installs swirl.

Now, to use the library, type `library(swirl)` and hit enter. That loads swirl.

Then type `swirl()` and hit enter. Now you're running swirl. Follow the directions on the screen. When you are asked, you want to install course 1 R Programming: The basics of programming in R. Then, when asked, you want to do option 1, R Programming, in that course.

When you are finished with the course – it will take just a few minutes – it will first ask you if you want credit on Coursera. You do not. Then type 0 to exit (it will not be very clear that's what you do when you are done).

## 2.2 Adding libraries, part 2

We'll mostly use two libraries for analysis – `dplyr` and `ggplot2`. To get them, and several other useful libraries, we can install a single collection of libraries called the tidyverse. Type this into your console:

```
install.packages('tidyverse')
```

**NOTE:** This is a pattern. You should always install libraries in the console.

Then, to help us with learning and replication, we're going to use R Notebooks. So we need to install that library. Type this into your console:

```
install.packages('rmarkdown')
```

## 2.3 Notebooks

For the rest of the class, we're going to be working in notebooks. In notebooks, you will both run your code and explain each step, much as I am doing here.

To start a notebook, you click on the green plus in the top left corner and go down to R Notebook. Do that now.

The screenshot shows the R Studio interface. A red arrow points to the 'New File' button in the top-left toolbar. Another red arrow points to the 'R Notebook' option in the dropdown menu that appears when the 'New File' button is clicked. The main workspace shows a portion of an R Markdown document with code and explanatory text.

```
35
36 ```{r}
37 install.packages('tidyverse')
38 ```
39 Then, to help us with learning and replication, we're going to
  need to install that library. Type this into your console:
40 ```{r}
41 install.packages('rmarkdown')
42 ```
43 You may have to quit and restart R Studio. I honestly can't re-
  member exactly why, but it's something to do with the knitr package
  and its dependencies.
44
45 #### Notebooks
```

You will see that the notebook adds a lot of text for you. It tells you how to work in notebooks – and you should read it. The most important parts are these:

To add text, simply type. To add code you can click on the *Insert* button on the toolbar or by pressing *Cmd+Option+I* on Mac or *Ctrl+Alt+I* on Windows.

Highlight all that text and delete it. You should have a blank document. This document is called a R Markdown file – it's a special form of text, one that you can style, and one you can include R in the middle of it. Markdown is a simple markup format that you can use to create documents. So first things first, let's give our notebook a big headline. Add this:

```
# My awesome notebook
```

Now, under that, without any markup, just type This is my awesome notebook.

Under that, you can make text bold by writing **It is \*\*really\*\* awesome.**

If you want it italics, just do this on the next line: *No, it's \_really\_ awesome. I swear.*

To see what it looks like without the markup, click the Preview or Knit button in the toolbar. That will turn your notebook into a webpage, with the formatting included.

Throughout this book, we're going to use this markdown to explain what we are doing and, more importantly, why we are doing it. Explaining your thinking is a vital part of understanding what you are doing.

That explanation, plus the code, is the real power of notebooks. To add a block of code, follow the instructions from above: click on the *Insert* button on the toolbar or by pressing *Cmd+Option+I* on Mac or *Ctrl+Alt+I* on Windows.

In that window, use some of the code from above and add two numbers together. To see it run, click the green triangle on the right. That runs the chunk. You should see the answer to your addition problem.

And that, just that, is the foundation you need to start this book.

# Chapter 3

## Data, structures and types

Data are everywhere (and data is plural of datum, thus the use of are in that statement). It surrounds you. Every time you use your phone, you are creating data. Lots of it. Your online life. Any time you buy something. It's everywhere. Sports, like life, is no different. Sports is drowning in data, and more comes along all the time.

In sports, and in this class, we'll be dealing largely with two kinds of data: event level data and summary data. It's not hard to envision event level data in sports. A pitch in baseball. A hit. A play in football. A pass in soccer. They are the events that make up the game. Combine them together – summarize them – and you'll have some notion of how the game went. What we usually see is summary data – who wants to scroll through 50 pitches to find out a player went 2-3 with a double and an RBI? Who wants to scroll through hundreds of pitches to figure out the Rays beat the Yankees?

To start with, we need to understand the shape of data.

EXERCISE: Try scoring a child's board game. For example, Chutes and Ladders. If you were placed in charge of analytics for the World Series of Chutes and Ladders, what is your event level data? What summary data do you keep? If you've got the game, try it.

### 3.1 Rows and columns

Data, oversimplifying it a bit, is information organized. Generally speaking, it's organized into rows and columns. Rows, generally, are individual elements. A team. A player. A game. Columns, generally, are components of the data, sometimes called variables. So if each row is a player, the first column might be their name. The second is their position. The third is their batting average. And so on.

G	Date	Opp	W/L	
1	2018-11-06	Mississippi Valley State	W	10
2	2018-11-11	Southeastern Louisiana	W	8
3	2018-11-14	Seton Hall	W	8
4	2018-11-19	N Missouri State	W	8
5	2018-11-20	N Texas Tech	L	5
6	2018-11-24	Western Illinois	W	7
7	2018-11-26	@ Clemson	W	6
8	2018-12-02	Rows		7
9	2018-12-05	@ Minnesota	L	7
10	2018-12-08	Creighton	W	9
11	2018-12-16	N Oklahoma State	W	7
12	2018-12-22	Cal State Fullerton	W	8
13	2018-12-29	Southwest Minnesota State	W	7
14	2019-01-02	@ Maryland	L	7
15	2019-01-06	@ Iowa	L	8
16	2019-01-10	Penn State	W	7
17	2019-01-14	@ Indiana	W	6
18	2019-01-17	Michigan State	L	6
19	2019-01-21	@ Rutgers	L	6

One of the critical components of data analysis, especially for beginners, is having a mental picture of your data. What does each row mean? What does each column in each row signify? How many rows do you have? How many columns?

## 3.2 Types

There are scores of data types in the world, and R has them. In this class, we're primarily going to be dealing with data frames, and each element of our data frames will have a data type.

Typically, they'll be one of four types of data:

- Numeric: a number, like the number of touchdown passes in a season or a batting average.
- Character: Text, like a name, a team, a conference.
- Date: Fully formed dates – 2019-01-01 – have a special date type. Elements of a date, like a year (ex. 2019) are not technically dates, so they'll appear as numeric data types.
- Logical: Rare, but every now and then we'll have a data type that's Yes or No, True or False, etc.

**Question:** Is a zip code a number? Is a jersey number a number? Trick question, because the answer is no. Numbers are things we do math on. If the thing you want is not something you're going to do math on – can you add two phone numbers together? – then make it a character type. If you don't, most every software system on the planet will drop leading zeros. For example, every zip code in Boston starts with 0. If you record that as a number, your zip code will become a four digit number, which isn't a zip code anymore.

## 3.3 A simple way to get data

One good thing about sports is that there's lots of interest in it. And that means there's outlets that put sports data on the internet. Now I'm going to show you a trick to getting it easily.

The site sports-reference.com takes NCAA (and other league) stats and puts them online. For instance, here's their page on Nebraska basketball's game logs, which you should open now.

Now, in a new tab, log into Google Docs/Drive and open a new spreadsheet. In the first cell of the first row, copy and paste this formula in:

```
=IMPORTHTML("https://www.sports-reference.com/cbb/schools/nebraska/2019-gamelogs.html", "table",
```

If it worked right, you've got the data from that page in a spreadsheet.

### 3.4 Cleaning the data

The first thing we need to do is recognize that we don't have data, really. We have the results of a formula. You can tell by putting your cursor on that field, where you'll see the formula again. This is where you'd look:

Screenshot of a Google Sheets spreadsheet titled "Untitled spreadsheet". The formula bar shows the formula =IMPORTHTML("https://www.sports-reference.com/cbb/schools/nebraska/2018.html", "table", 1). A red arrow points to the formula, and a red circle highlights the first row of the table.

	A	B	C	D	E
1					
2	G	Date		Opp	W/L
3	1	2018-11-06		Mississippi Valley	W
4	2	2018-11-11		Southeastern Louisiana	W
5	3	2018-11-14		Seton Hall	W
6	4	2018-11-19	N	Missouri State	W
7	5	2018-11-20	N	Texas Tech	L
8	6	2018-11-24		Western Illinois	W
9	7	2018-11-26	@	Clemson	W
10	8	2018-12-02		Illinois	W
11	9	2018-12-05	@	Minnesota	L
12	10	2018-12-08		Creighton	W
13	11	2018-12-16	N	Oklahoma State	W
14	12	2018-12-22		Cal State Fullerton	W
15	13	2018-12-29		Southwest Minnesota	W
16	14	2019-01-02	@	Maryland	L
17	15	2019-01-06	@	Iowa	L
18	16	2019-01-10		Penn State	W
19	17	2019-01-14	@	Indiana	W
20	18	2019-01-17		Michigan State	L

The solution is easy:

Edit > Select All or type command/control A Edit > Copy or type command/control c Edit > Paste Special > Values Only or type command/control shift v

You can verify that it worked by looking in that same row 1 column A, where you'll see the formula is gone.

Screenshot of the Google Sheets "Edit" menu open, showing the "Paste special" submenu.

The submenu includes:

- Paste values only
- Paste format only
- Paste all except borders
- Paste column width
- Paste formula only
- Paste data validation
- Paste conditional formats
- Paste transposed

The main menu bar shows: File, Edit (highlighted), View, Insert, Format, Data, Tools, Add-ons, Help, All.

The spreadsheet view shows rows 1 through 20. Row 1 contains the formula =IMPORT. Row 2 contains the letter G. Rows 3 through 20 are empty.

	A	B	C	D	E
1					
2	G				
3					
4					
5					
6					
7					
8					
9					
10					
11					
12					
13					
14		2019-01-02	@		Ma
15		2019-01-06	@		low
16		2019-01-10			Per
17		2019-01-14	@	Indiana	W
18		2019-01-17		Michigan State	L

Now you have data, but your headers are all wrong. You want your headers to be one line – not two, like they have. And the header names repeat – first for our team, then for theirs. So you have to change each header name to be UsORB or TeamORB and OpponentORB instead of just ORB.

After you've done that, note we have repeating headers. There's two ways to deal with that – you could just highlight it and go up to Edit > Delete Rows XX-XX depending on what rows you highlighted. That's the easy way with our data.

But what if you had hundreds of repeating headers like that? Deleting them would take a long time.

You can use sorting to get rid of anything that's not data. So click on Data > Sort Range. You'll want to check the "Data has header row" field. Then hit Sort.

Year	Yards	Pen./G	Yards/G
13	116	3.3	29
15	153	3	30.6
19	162	3.8	32.4
17	133	4.3	33.3
17	135		
13	136		
17	147		
24	189		
21	191		
19	192		
27	195		
21	198		
20	159		
20	160		
17	165		
18	165		
20	165		
32	207	6.4	41.4
22	208	4.4	41.6
22	210	4.4	42
18	172	4.5	43
28	215	5.6	43
23	172	5.8	43
26	221	5.2	44.2
19	178	4.8	44.5

Sort range from A1 to Z100

Data has header row

sort by

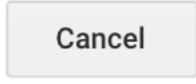
Year 

A → Z

Z → A

[+ Add another sort column](#)

 Sort

 Cancel

Now all you need to do is search through the data for where your junk data – extra headers, blanks, etc. – got sorted and delete it. After you've done that, you can export it for use in R. Go to File > Download as > Comma Separated Values. Remember to put it in the same directory as your R Notebook file so you can import the data easily.

# Chapter 4

## Aggregates

R is a statistical programming language that is purpose built for data analysis.

Base R does a lot, but there are a mountain of external libraries that do things to make R better/easier/more fully featured. We already installed the tidyverse – or you should have if you followed the instructions for the last assignment – which isn't exactly a library, but a collection of libraries. Together, they make up the tidyverse. Individually, they are extraordinarily useful for what they do. We can load them all at once using the tidyverse name, or we can load them individually. Let's start with individually.

The two libraries we are going to need for this assignment are `readr` and `dplyr`. The library `readr` reads different types of data in as a dataframe. For this assignment, we're going to read in csv data or Comma Separated Values data. That's data that has a comma between each column of data.

Then we're going to use `dplyr` to analyze it.

To use a library, you need to import it. Good practice – one I'm going to insist on – is that you put all your library steps at the top of your notebooks.

That code looks like this:

```
library(readr)
```

To load them both, you need to run that code twice:

```
library(readr)  
library(dplyr)
```

You can keep doing that for as many libraries as you need. I've seen notebooks with 10 or more library imports.

But the tidyverse has a neat little trick. We can load most of the libraries we'll need for the whole semester with one line:

```
library(tidyverse)
```

**From now on, if that's not the first line of your notebook, you're probably doing it wrong.**

## 4.1 Basic data analysis: Group By and Count

The first thing we need to do is get some data to work with. We do that by reading it in. In our case, we're going to read data from a csv file – a comma-separated values file.

The CSV file we're going to read from is a Basketball Reference page of advanced metrics for NBA players this season. The Sports Reference sites are a godsend of data, a trove of stuff, and we're going to use it a lot in this class.

So step 2, after setting up our libraries, is most often going to be importing data. In order to analyze data, we need data, so it stands to reason that this would be something we'd do very early.

The code looks *something* like this, but hold off copying it just yet:

```
nbaplayers <- read_csv("~/Box/SportsData/nbaadvancedplayers1920.csv")
```

Let's unpack that.

The first part – nbaplayers – is the name of your variable. A variable is just a name of a thing that stores stuff. In this case, our variable is a data frame, which is R's way of storing data (technically it's a tibble, which is the tidyverse way of storing data, but the differences aren't important and people use them interchangeably). **We can call this whatever we want.** I always want to name data frames after what is in it. In this case, we're going to import a dataset of NBA players. Variable names, by convention are one word all lower case. You can end a variable with a number, but you can't start one with a number.

The <- bit is the variable assignment operator. It's how we know we're assigning something to a word. Think of the arrow as saying “Take everything on the right of this arrow and stuff it into the thing on the left.” So we're creating an empty vessel called **nbaplayers** and stuffing all this data into it.

The **read\_csv** bits are pretty obvious, except for one thing. What happens in the quote marks is the path to the data. In there, I have to tell R where it will find the data. The easiest thing to do, if you are confused about how to find your data, is to put your data in the same folder as your notebook (you'll have to save that notebook first). If you do that, then you just need to put the name of the file in there (nbaadvancedplayers1920.csv). In my case, I've got a folder called Box in my home directory (that's the ~ part), and in there is a folder called SportsData that has the file called nbaadvancedplayers1920.csv in it. Some people – insane people – leave the data in their downloads folder. The

data path then would be `~/Downloads/nameofthedatafilehere.csv` on PC or Mac.

**What you put in there will be different from mine.** So your first task is to import the data.

```
nbaplayers <- read_csv("data/nbaadvancedplayers1920.csv")
```

```
## 
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Player = col_character(),
##   Pos = col_character(),
##   Tm = col_character()
## )
## i Use `spec()` for the full column specifications.
```

Now we can inspect the data we imported. What does it look like? To do that, we use `head(nbaplayers)` to show the headers and **the first six rows of data**. If we wanted to see them all, we could just simply enter `mountainlions` and run it.

To get the number of records in our dataset, we run `nrow(nbaplayers)`

```
head(nbaplayers)

## # A tibble: 6 x 27
##       Rk Player Pos     Age Tm      G    MP    PER `TS%` `3PAr` `FTr` `ORB%` 
##   <dbl> <chr> <chr> <dbl> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1     1 Steve~ C        26 OKC     63 1680  20.5  0.604  0.006  0.421   14
## 2     2 Bam ~ PF       22 MIA     72 2417  20.3  0.598  0.018  0.484   8.5
## 3     3 LaMar~ C       34 SAS     53 1754  19.7  0.571  0.198  0.241   6.3
## 4     4 Kyle ~ PF      23 MIA     2   13   4.7   0.5    0     0     0     17.9
## 5     5 Nicke~ SG      21 NOP     47  591   8.9   0.473  0.5    0.139   1.6
## 6     6 Grays~ SG      24 MEM     38  718   12    0.609  0.562  0.179   1.2
## # ... with 15 more variables: `DRB%` <dbl>, `TRB%` <dbl>, `AST%` <dbl>,
## #   `STL%` <dbl>, `BLK%` <dbl>, `TOV%` <dbl>, `USG%` <dbl>, OWS <dbl>,
## #   DWS <dbl>, WS <dbl>, `WS/48` <dbl>, OBPM <dbl>, DBPM <dbl>, BPM <dbl>,
## #   VORP <dbl>

nrow(nbaplayers)

## [1] 651
```

Another way to look at `nrow` – we have 651 players from this season in our dataset.

What if we wanted to know how many players there were by position? To do that by hand, we'd have to take each of the 651 records and sort them into a

pile. We'd put them in groups and then count them.

`dplyr` has a **group by** function in it that does just this. A massive amount of data analysis involves grouping like things together at some point. So it's a good place to start.

So to do this, we'll take our dataset and we'll introduce a new operator: `%>%`. The best way to read that operator, in my opinion, is to interpret that as “and then do this.”

After we group them together, we need to count them. We do that first by saying we want to summarize our data (a count is a part of a summary). To get a summary, we have to tell it what we want. So in this case, we want a count. To get that, let's create a thing called `total` and set it equal to `n()`, which is `dplyr`'s way of counting something.

Here's the code:

```
nbaplayers %>%
  group_by(Pos) %>%
  summarise(
    total = n()
  )

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

## # A tibble: 9 x 2
##   Pos     total
##   <chr> <int>
## 1 C         111
## 2 C-PF      2
## 3 PF        135
## 4 PF-C      2
## 5 PG        111
## 6 SF        113
## 7 SF-PF     4
## 8 SF-SG     3
## 9 SG        170
```

So let's walk through that. We start with our dataset – `nbaplayers` – and then we tell it to group the data by a given field in the data which we get by looking at either the output of `head` or you can look in the environment where you'll see `nbaplayers`.

In this case, we wanted to group together positions, signified by the field name `Pos`. After we group the data, we need to count them up. In `dplyr`, we use `summarize` which can do more than just count things. Inside the parentheses in `summarize`, we set up the summaries we want. In this case, we just want a count of the positions: `total = n()`, says create a new field, called `total` and set it equal to `n()`, which might look weird, but it's common in stats. The

number of things in a dataset? Statisticians call in n. There are n number of players in this dataset. So n() is a function that counts the number of things there are.

And when we run that, we get a list of positions with a count next to them. But it's not in any order. So we'll add another And Then Do This %>% and use **arrange**. Arrange does what you think it does – it arranges data in order. By default, it's in ascending order – smallest to largest. But if we want to know the county with the most mountain lion sightings, we need to sort it in descending order. That looks like this:

```
nbaplayers %>%
  group_by(Pos) %>%
  summarise(
    total = n()
  ) %>% arrange(desc(total))

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

## # A tibble: 9 x 2
##   Pos     total
##   <chr>   <int>
## 1 SG        170
## 2 PF        135
## 3 SF        113
## 4 C         111
## 5 PG        111
## 6 SF-PF      4
## 7 SF-SG      3
## 8 C-PF       2
## 9 PF-C       2
```

So the most common position in the NBA? Shooting guard, followed by power forward.

We can, if we want, group by more than one thing. Which team has the most of a single position? To do that, we can group by the team – called Tm in the data – and position, or Pos in the data:

```
nbaplayers %>%
  group_by(Tm, Pos) %>%
  summarise(
    total = n()
  ) %>% arrange(desc(total))

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

## # A tibble: 159 x 3
##   Tm      Pos     total
##   <chr>   <chr>   <int>
## 1 BOS     SG        170
## 2 BOS     PF        135
## 3 BOS     SF        113
## 4 BOS     C         111
## 5 BOS     PG        111
## 6 BOS     SF-PF      4
## 7 BOS     SF-SG      3
## 8 BOS     C-PF       2
## 9 BOS     PF-C       2
```

```

##   Tm    Pos  total
##   <chr> <chr> <int>
## 1 TOT   PF     13
## 2 TOT   SG     13
## 3 SAC   PF      9
## 4 TOT   SF      9
## 5 BRK   SG      8
## 6 LAL   SG      8
## 7 TOT   PG      8
## 8 ATL   SG      7
## 9 BRK   SF      7
## 10 DAL  SG      7
## # ... with 149 more rows

```

So wait, what team is TOT?

Valuable lesson: whoever collects the data has opinions on how to solve problems. In this case, Basketball Reference, when a player get's traded, records stats for the player's first team, their second team, and a combined season total for a team called TOT, meaning Total. Is there a team abbreviated TOT? No. So ignore them here.

Sacramento has 9 power forward. Brooklyn has 8 shooting guards, as do the Lakers. You can learn a bit about how a team is assembled by looking at these simple counts.

## 4.2 Other aggregates: Mean and median

In the last example, we grouped some data together and counted it up, but there's so much more you can do. You can do multiple measures in a single step as well.

Sticking with our NBA player data, we can calculate any number of measures inside summarize. Here, we'll use R's built in mean and median functions to calculate ... well, you get the idea.

Let's look just at the number of minutes each position gets.

```

nbplayers %>%
  group_by(Pos) %>%
  summarise(
    count = n(),
    mean_minutes = mean(MP),
    median_minutes = median(MP)
  )

## `summarise()` ungrouping output (override with `.`groups` argument)
## # A tibble: 9 x 4

```

```

##   Pos   count mean_minutes median_minutes
##   <chr> <int>      <dbl>        <dbl>
## 1 C       111      891.        887
## 2 C-PF     2       316.        316.
## 3 PF      135      790.        567
## 4 PF-C     2      1548.       1548.
## 5 PG      111      944.        850
## 6 SF      113      877.        754
## 7 SF-PF    4       638.        286.
## 8 SF-SG    3      1211.       1688
## 9 SG      170      843.        654.

```

So there's 651 players in the data. Let's look at shooting guards. The average shooting guard plays 842 minutes and the median is 653.5 minutes.

Why?

Let's let sort help us.

```

nbaplayers %>% arrange(desc(MP))

## # A tibble: 651 x 27
##       Rk Player Pos   Age Tm      G   MP   PER `TS%` `3PAr`   FTr `ORB%` 
##   <dbl> <chr> <chr> <dbl> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1    323 CJ Mc~ SG     28 POR     70 2556  17 0.541  0.378 0.136  1.9
## 2     55 Devin~ SG     23 PHO     70 2512 20.6 0.618  0.31  0.397  1.3
## 3    198 James~ SG     30 HOU     68 2483 29.1 0.626  0.557 0.528  2.9
## 4     27 Harri~ PF     27 SAC     72 2482 13.3 0.574  0.338 0.337  3.4
## 5    297 Damia~ PG     29 POR     66 2474 26.9 0.627  0.5  0.384  1.4
## 6    204 Tobia~ PF     27 PHI     72 2469 17.2 0.556  0.304 0.184  3.1
## 7    479 P.J. ~ PF     34 HOU     72 2467  8.3 0.559  0.702 0.113  4.7
## 8    175 Shai ~ SG     21 OKC     70 2428 17.7 0.568  0.247 0.352  2.2
## 9     2 Bam A~ PF     22 MIA     72 2417 20.3 0.598  0.018 0.484  8.5
## 10   343 Donov~ SG     23 UTA     69 2364 18.8 0.558  0.352 0.24  2.6
## # ... with 641 more rows, and 15 more variables: `DRB%` <dbl>, `TRB%` <dbl>,
## # `AST%` <dbl>, `STL%` <dbl>, `BLK%` <dbl>, `TOV%` <dbl>, `USG%` <dbl>,
## # OWS <dbl>, DWS <dbl>, WS <dbl>, `WS/48` <dbl>, OBPM <dbl>, DBPM <dbl>,
## # BPM <dbl>, VORP <dbl>

```

The player with the most minutes on the floor is a shooting guard. Shooting guard is the most common position, so that means there's CJ McCollum rolling up 2,556 minutes in a season, and then there's Cleveland Cavalier's sensation J.P. Macura. Never heard of J.P. Macura? Might be because he logged one minute in one game this season.

That's a huge difference.

So when choosing a measure of the middle, you have to ask yourself – could I have extremes? Because a median won't be sensitive to extremes. It will be the

point at which half the numbers are above and half are below. The average or mean will be a measure of the middle, but if you have a bunch of pine riders and then one ironman superstar, the average will be wildly skewed.

### 4.3 Even more aggregates

There's a ton of things we can do in summarize – we'll work with more of them as the course progresses – but here's a few other questions you can ask.

Which position in the NBA plays the most minutes? And what is the highest and lowest minute total for that position? And how wide is the spread between minutes? We can find that with `sum` to add up the minutes to get the total minutes, `min` to find the minimum minutes, `max` to find the maximum minutes and `sd` to find the standard deviation in the numbers.

```
nbaplayers %>%
  group_by(Pos) %>%
  summarise(
    total = sum(MP),
    avgminutes = mean(MP),
    minminutes = min(MP),
    maxminutes = max(MP),
    stdev = sd(MP)) %>% arrange(desc(total))

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

## # A tibble: 9 x 6
##   Pos     total avgminutes minminutes maxminutes stdev
##   <chr>    <dbl>      <dbl>        <dbl>      <dbl> <dbl>
## 1 SG      143229      843.         1       2556  735.
## 2 PF      106654      790.         5       2482  719.
## 3 PG      104745      944.         8       2474  727.
## 4 SF      99109       877.        11      2316  709.
## 5 C       98914       891.         3       2336  619.
## 6 SF-SG    3633       1211        87      1858  977.
## 7 PF-C     3097       1548.        960     2137  832.
## 8 SF-PF    2553       638.        46      1936  873.
## 9 C-PF     633        316.        256     377   85.6
```

So again, no surprise, shooting guards spend the most minutes on the floor in the NBA. They average 842 minutes, but we noted why that's trouble. The minimum is the J.P. Macura Award, max is the Trailblazer's failing at load management, and the standard deviation is a measure of how spread out the data is. In this case, not the highest spread among positions, but pretty high. So you know you've got some huge minutes players and a bunch of bench players.

# Chapter 5

## Mutating data

One of the most common data analysis techniques is to look at change over time. The most common way of comparing change over time is through percent change. The math behind calculating percent change is very simple, and you should know it off the top of your head. The easy way to remember it is:

```
(new - old) / old
```

Or new minus old divided by old. Your new number minus the old number, the result of which is divided by the old number. To do that in R, we can use `dplyr` and `mutate` to calculate new metrics in a new field using existing fields of data.

So first we'll import the tidyverse so we can read in our data and begin to work with it.

```
library(tidyverse)
```

Now you'll need a common and simple dataset of total attendance at NCAA football games over the last few seasons.

You'll import it something like this.

```
attendance <- read_csv('data/attendance.csv')
```

```
##  
## -- Column specification -----  
## cols(  
##   Institution = col_character(),  
##   Conference = col_character(),  
##   `2013` = col_double(),  
##   `2014` = col_double(),  
##   `2015` = col_double(),  
##   `2016` = col_double(),  
##   `2017` = col_double(),
```

```
## `2018` = col_double()
## )
```

If you want to see the first six rows – handy to take a peek at your data – you can use the function `head`.

```
head(attendance)
```

```
## # A tibble: 6 x 8
##   Institution    Conference `2013` `2014` `2015` `2016` `2017` `2018`
##   <chr>          <chr>      <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Air Force      MWC        228562  168967  156158  177519  174924  166205
## 2 Akron          MAC        107101   55019  108588   62021  117416   92575
## 3 Alabama         SEC        710538  710736  707786  712747  712053  710931
## 4 Appalachian St. FBS Independent 149366     NA     NA     NA     NA     NA
## 5 Appalachian St. Sun Belt       NA  138995  128755  156916  154722  131716
## 6 Arizona         Pac-12      285713  354973  308355  338017  255791  318051
```

The code to calculate percent change is pretty simple. Remember, with `summarize`, we used `n()` to count things. With `mutate`, we use very similar syntax to calculate a new value using other values in our dataset. So in this case, we’re trying to do  $(\text{new-old})/\text{old}$ , but we’re doing it with fields. If we look at what we got when we did `head`, you’ll see there’s ‘2018’ as the new data, and we’ll use ‘2017’ as the old data. So we’re looking at one year. Then, to help us, we’ll use `arrange` again to sort it, so we get the fastest growing school over one year.

```
attendance %>% mutate(
  change = (`2018` - `2017`)/`2017`
)

## # A tibble: 150 x 9
##   Institution    Conference `2013` `2014` `2015` `2016` `2017` `2018`   change
##   <chr>          <chr>      <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Air Force      MWC        228562  168967  156158  177519  174924  166205 -0.0498
## 2 Akron          MAC        107101   55019  108588   62021  117416   92575 -0.212
## 3 Alabama         SEC        710538  710736  707786  712747  712053  710931 -0.00158
## 4 Appalachian ~ FBS Indepen~ 149366     NA     NA     NA     NA     NA NA
## 5 Appalachian ~ Sun Belt       NA  138995  128755  156916  154722  131716 -0.149
## 6 Arizona         Pac-12      285713  354973  308355  338017  255791  318051  0.243
## 7 Arizona St.    Pac-12      501509  343073  368985  286417  359660  291091 -0.191
## 8 Arkansas        SEC        431174  399124  471279  487067  442569  367748 -0.169
## 9 Arkansas St.   Sun Belt      149477  149163  138043  136200  119538  119001 -0.00449
## 10 Army West Po~ FBS Indepen~ 169781  171310  185946  163267  185543  190156  0.0249
## # ... with 140 more rows
```

What do we see right away? Do those numbers look like we expect them to? No. They’re a decimal expressed as a percentage. So let’s fix that by multiplying by 100.

```

attendance %>% mutate(
  change = ((`2018` - `2017`)/`2017`)*100
)

## # A tibble: 150 x 9
##   Institution   Conference `2013` `2014` `2015` `2016` `2017` `2018` change
##   <chr>         <chr>     <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Air Force     MWC        228562  168967  156158  177519  174924  166205 -4.98
## 2 Akron         MAC        107101   55019  108588   62021  117416   92575 -21.2 
## 3 Alabama        SEC        710538  710736  707786  712747  712053  710931 -0.158
## 4 Appalachian S~ FBS Indepen~ 149366      NA      NA      NA      NA      NA    NA
## 5 Appalachian S~ Sun Belt      NA  138995  128755  156916  154722  131716 -14.9 
## 6 Arizona        Pac-12      285713  354973  308355  338017  255791  318051  24.3 
## 7 Arizona St.    Pac-12      501509  343073  368985  286417  359660  291091 -19.1 
## 8 Arkansas       SEC        431174  399124  471279  487067  442569  367748 -16.9 
## 9 Arkansas St.   Sun Belt      149477  149163  138043  136200  119538  119001 -0.449
## 10 Army West Poi~ FBS Indepen~ 169781  171310  185946  163267  185543  190156  2.49
## # ... with 140 more rows

```

Now, does this ordering do anything for us? No. Let's fix that with arrange.

```

attendance %>% mutate(
  change = ((`2018` - `2017`)/`2017`)*100
) %>% arrange(desc(change))

## # A tibble: 150 x 9
##   Institution   Conference `2013` `2014` `2015` `2016` `2017` `2018` change
##   <chr>         <chr>     <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Ga. Southern  Sun Belt      NA  105510  124681  104095  61031  100814  65.2 
## 2 La.-Monroe    Sun Belt      85177   90540   58659   67057   49640   71048  43.1 
## 3 Louisiana     Sun Belt      129878  154652  129577  121346  78754  111303  41.3 
## 4 Hawaii        MWC        185931  192159  164031  170299  145463  205455  41.2 
## 5 Buffalo       MAC        136418  122418  110743  104957  80102  110280  37.7 
## 6 California    Pac-12      345303  286051  292797  279769  219290  300061  36.8 
## 7 UCF           AAC        252505  226869  180388  214814  257924  352148  36.5 
## 8 UTSA          C-USA      175282  165458  138048  138226  114104  148257  29.9 
## 9 Eastern Mich. MAC        20255   75127   29381  106064   73649   95632  29.8 
## 10 Louisville   ACC        NA  317829  294413  324391  276957  351755  27.0
## # ... with 140 more rows

```

So who had the most growth last year from the year before? Something going on at Georgia Southern.

## 5.1 A more complex example

There's metric in basketball that's easy to understand – shooting percentage. It's the number of shots made divided by the number of shots attempted. Simple,

right? Except it's a little too simple. Because what about three point shooters? They tend to be more valuable because the three point shot is worth more. What about players who get to the line? In shooting percentage, free throws are nowhere to be found.

Basketball nerds, because of these weaknesses, have created a new metric called True Shooting Percentage. True shooting percentage takes into account all aspects of a player's shooting to determine who the real shooters are.

Using `dplyr` and `mutate`, we can calculate true shooting percentage. So let's look at a new dataset, one of every college basketball player's season stats in 2018-19 season. It's a dataset of 5,386 players, and we've got 59 variables – one of them is True Shooting Percentage, but we're going to ignore that.

Import it like this:

```
players <- read_csv("data/players19.csv")

## Warning: Missing column names filled in: 'X1' [1]

##
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Team = col_character(),
##   Conference = col_character(),
##   Player = col_character(),
##   Class = col_character(),
##   Pos = col_character(),
##   Height = col_character(),
##   Hometown = col_character(),
##   `High School` = col_character(),
##   Summary = col_character()
## )
## i Use `spec()` for the full column specifications.
```

The basic true shooting percentage formula is  $(\text{Points} / (2 * (\text{FieldGoalAttempts} + (.44 * \text{FreeThrowAttempts})))) * 100$ . Let's talk that through. Points divided by a lot. It's really field goal attempts plus 44 percent of the free throw attempts. Why? Because that's about what a free throw is worth, compared to other ways to score. After adding those things together, you double it. And after you divide points by that number, you multiply the whole lot by 100.

In our data, we need to be able to find the fields so we can complete the formula. To do that, one way is to use the Environment tab in R Studio. In the Environment tab is a listing of all the data you've imported, and if you click the triangle next to it, it'll list all the field names, giving you a bit of information about each one.

The screenshot shows the RStudio interface with the Environment tab selected. The Global Environment pane displays the 'players' dataset, which has 5386 observations and 59 variables. The variables listed are X1, Team, Conference, Player, #, and Class. Below the environment pane, the Files tab is active, showing a file tree under the 'BookProjects' folder. Two red arrows point from the text in the sidebar to the 'File' and 'New Folder' icons in the Files tab.

ly field  
a free throw  
it. And

Environment History Connections

Import Dataset

Global Environment

players 5386 obs. of 59 variables

X1 : num 1 2 3 4 5 6 7 8 9 10 .

Team : chr "Youngstown State Pe"

Conference : chr "Horizon" "Hor"

Player : chr "Darius Quisenberr"

# : num 3 32 22 2 33 13 5 31 21

Class : chr "FR" "S0" "JR" "FR"

New Folder Delete Rename

Home > Box > BookProjects > Sport

Name

27-finishingtouches2.Rmd

28-assignments.Rmd

So what does True Shooting Percentage look like in code?

Let's think about this differently. Who had the best true shooting season last year?

```
players %>%
  mutate(trueshooting = (PTS/(2*(FGA + (.44*FTA)))*100) %>%
    arrange(desc(trueshooting))

## # A tibble: 5,386 x 60
##   X1 Team Conference Player `#` Class Pos Height Weight Hometown
##   <dbl> <chr> <chr>     <chr> <dbl> <chr> <chr> <chr> <dbl> <chr>
## 1 579 Texa~ Big 12 Drayt~     4 JR   G   6-0   156 Austin, ~
## 2 843 Ston~ AEC Nick ~     42 FR   F   6-7   240 Port Je~
## 3 1059 Sout~ Southland Patri~     22 SO   F   6-3   210 Folsom, ~
## 4 4269 Dayt~ A-10 Camro~     52 SO   G   5-7   160 Country~
## 5 4681 Cali~ Pac-12 David~     21 JR   G   6-4   185 Newbury~
## 6 326 Virg~ ACC Grant~     1 FR   G   <NA>  NA Charlott~
## 7 410 Vand~ SEC Mac H~     42 FR   G   6-6   182 Chattan~
## 8 1390 Sain~ A-10 Jack ~     31 JR   G   6-6   205 Mattoon~
## 9 2230 NJIT~ A-Sun Patri~     3 SO   G   5-9   160 West Or~
## 10 266 Wash~ Pac-12 Reaga~    34 FR   F   6-6   225 Santa A~
## # ... with 5,376 more rows, and 50 more variables: `High School` <chr>,
## #   Summary <chr>, Rk.x <dbl>, G <dbl>, GS <dbl>, MP <dbl>, FG <dbl>,
## #   FGA <dbl>, `FG%` <dbl>, `2P` <dbl>, `2PA` <dbl>, `2P%` <dbl>, `3P` <dbl>,
## #   `3PA` <dbl>, `3P%` <dbl>, FT <dbl>, FTA <dbl>, `FT%` <dbl>, ORB <dbl>,
## #   DRB <dbl>, TRB <dbl>, AST <dbl>, STL <dbl>, BLK <dbl>, TOV <dbl>, PF <dbl>,
## #   PTS <dbl>, Rk.y <dbl>, PER <dbl>, `TS%` <dbl>, `eFG%` <dbl>, `3PAr` <dbl>,
## #   FTr <dbl>, PProd <dbl>, `ORB%` <dbl>, `DRB%` <dbl>, `TRB%` <dbl>,
## #   `AST%` <dbl>, `STL%` <dbl>, `BLK%` <dbl>, `TOV%` <dbl>, `USG%` <dbl>,
## #   OWS <dbl>, DWS <dbl>, WS <dbl>, `WS/40` <dbl>, OBPM <dbl>, DBPM <dbl>,
## #   BPM <dbl>, trueshooting <dbl>
```

You'll be forgiven if you did not hear about Texas Longhorns shooting sensation Drayton Whiteside. He played in six games, took one shot and actually hit it. It happened to be a three pointer, which is one more three pointer than I've hit in college basketball. So props to him. Does that mean he had the best true shooting season in college basketball last year?

Not hardly.

We'll talk about how to narrow the pile and filter out data in the next chapter.

# Chapter 6

## Filters and selections

More often than not, we have more data than we want. Sometimes we need to be rid of that data. In `dplyr`, there's two ways to go about this: filtering and selecting.

**Filtering creates a subset of the data based on criteria.** All records where the count is greater than 10. All records that match “Nebraska”. Something like that.

**Selecting simply returns only the fields named.** So if you only want to see School and Attendance, you select those fields. When you look at your data again, you'll have two columns. If you try to use one of your columns that you had before you used `select`, you'll get an error.

Let's work with our football attendance data to show some examples.

First we'll need the tidyverse.

```
library(tidyverse)
```

Now import the data.

```
attendance <- read_csv('data/attendance.csv')

## 
## -- Column specification -----#
## cols(
##   Institution = col_character(),
##   Conference = col_character(),
##   `2013` = col_double(),
##   `2014` = col_double(),
##   `2015` = col_double(),
##   `2016` = col_double(),
##   `2017` = col_double(),
```

```
## `2018` = col_double()
## )
```

So, first things first, let's say we don't care about all this Air Force, Akron, Alabama crap and just want to see Dear Old Nebraska U. We do that with `filter` and then we pass it a condition.

Before we do that, a note about conditions. Most of the conditional operators you'll understand – greater than and less than are `>` and `<`. The tough one to remember is equal to. In conditional statements, equal to is `==` not `=`. If you haven't noticed, `=` is a variable assignment operator, not a conditional statement. So equal is `==` and NOT equal is `!=`.

So if you want to see Institutions equal to Nebraska, you do this:

```
attendance %>% filter(Institution == "Nebraska")

## # A tibble: 1 x 8
##   Institution Conference `2013` `2014` `2015` `2016` `2017` `2018`
##   <chr>        <chr>     <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Nebraska     Big Ten    727466  638744  629983  631402  628583  623240
```

Or if we want to see schools that had more than half a million people buy tickets to a football game in a season, we do the following. NOTE THE BACKTICKS.

```
attendance %>% filter(`2018` >= 500000)

## # A tibble: 17 x 8
##   Institution Conference `2013` `2014` `2015` `2016` `2017` `2018`
##   <chr>        <chr>     <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Alabama      SEC       710538  710736  707786  712747  712053  710931
## 2 Auburn       SEC       685252  612157  612157  695498  605120  591236
## 3 Clemson      ACC       574333  572262  588266  566787  565412  562799
## 4 Florida      SEC       524638  515001  630457  439229  520290  576299
## 5 Georgia      SEC       556476  649222  649222  556476  556476  649222
## 6 LSU          SEC       639927  712063  654084  708618  591034  705733
## 7 Michigan     Big Ten   781144  734364  771174  883741  669534  775156
## 8 Michigan St. Big Ten   506294  522765  522628  522666  507398  508088
## 9 Nebraska     Big Ten   727466  638744  629983  631402  628583  623240
## 10 Ohio St.    Big Ten   734528  744075  750705  750944  752464  713630
## 11 Oklahoma    Big 12    508334  510972  512139  521142  519119  607146
## 12 Penn St.    Big Ten   676112  711358  698590  701800  746946  738396
## 13 South Carolina SEC       576805  569664  472934  538441  550099  515396
## 14 Tennessee   SEC       669087  698276  704088  706776  670454  650887
## 15 Texas        Big 12    593857  564618  540210  587283  556667  586277
## 16 Texas A&M  SEC       697003  630735  725354  713418  691612  698908
## 17 Wisconsin   Big Ten   552378  556642  546099  476144  551766  540072
```

But what if we want to see all of the Power Five conferences? We *could* use conditional logic in our filter. The conditional logic operators are `|` for OR and

& for AND. NOTE: AND means all conditions have to be met. OR means any of the conditions work. So be careful about boolean logic.

```
attendance %>% filter(Conference == "Big 10" | Conference == "SEC" | Conference == "Pac-12" | Con
```

```
## # A tibble: 51 x 8
##   Institution  Conference `2013` `2014` `2015` `2016` `2017` `2018`
##   <chr>        <chr>     <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Alabama      SEC        710538  710736  707786  712747  712053  710931
## 2 Arizona      Pac-12    285713  354973  308355  338017  255791  318051
## 3 Arizona St.  Pac-12    501509  343073  368985  286417  359660  291091
## 4 Arkansas     SEC        431174  399124  471279  487067  442569  367748
## 5 Auburn       SEC        685252  612157  612157  695498  605120  591236
## 6 Baylor        Big 12    321639  280257  276960  275029  262978  248017
## 7 Boston College ACC       198035  239893  211433  192942  215546  263363
## 8 California   Pac-12    345303  286051  292797  279769  219290  300061
## 9 Clemson      ACC        574333  572262  588266  566787  565412  562799
## 10 Colorado    Pac-12   230778  226670  236331  279652  282335  274852
## # ... with 41 more rows
```

But that's a lot of repetitive code. And a lot of typing. And typing is the devil. So what if we could create a list and pass it into the filter? It's pretty simple.

We can create a new variable – remember variables can represent just about anything – and create a list. To do that we use the `c` operator, which stands for concatenate. That just means take all the stuff in the parenthesis after the `c` and bunch it into a list.

Note here: text is in quotes. If they were numbers, we wouldn't need the quotes.

```
powerfive <- c("SEC", "Big Ten", "Pac-12", "Big 12", "ACC")
```

Now with a list, we can use the `%in%` operator. It does what you think it does – it gives you data that matches things IN the list you give it.

```
attendance %>% filter(Conference %in% powerfive)
```

```
## # A tibble: 65 x 8
##   Institution  Conference `2013` `2014` `2015` `2016` `2017` `2018`
##   <chr>        <chr>     <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Alabama      SEC        710538  710736  707786  712747  712053  710931
## 2 Arizona      Pac-12    285713  354973  308355  338017  255791  318051
## 3 Arizona St.  Pac-12    501509  343073  368985  286417  359660  291091
## 4 Arkansas     SEC        431174  399124  471279  487067  442569  367748
## 5 Auburn       SEC        685252  612157  612157  695498  605120  591236
## 6 Baylor        Big 12    321639  280257  276960  275029  262978  248017
## 7 Boston College ACC       198035  239893  211433  192942  215546  263363
## 8 California   Pac-12    345303  286051  292797  279769  219290  300061
## 9 Clemson      ACC        574333  572262  588266  566787  565412  562799
```

```
## # 10 Colorado      Pac-12      230778 226670 236331 279652 282335 274852
## # ... with 55 more rows
```

## 6.1 Selecting data to make it easier to read

So now we have our Power Five list. What if we just wanted to see attendance from the most recent season and ignore all the rest? Select to the rescue.

```
attendance %>% filter(Conference %in% powerfive) %>% select(Institution, Conference, ^
```

```
## # A tibble: 65 x 3
##   Institution   Conference `2018`
##   <chr>         <chr>     <dbl>
## 1 Alabama        SEC       710931
## 2 Arizona        Pac-12    318051
## 3 Arizona St.   Pac-12    291091
## 4 Arkansas       SEC       367748
## 5 Auburn         SEC       591236
## 6 Baylor         Big 12    248017
## 7 Boston College ACC       263363
## 8 California     Pac-12    300061
## 9 Clemson        ACC       562799
## 10 Colorado      Pac-12   274852
## # ... with 55 more rows
```

If you have truly massive data, Select has tools to help you select fields that start \_with the same things or ends with a certain word. The documentation will guide you if you need those someday. For 90 plus percent of what we do, just naming the fields will be sufficient.

## 6.2 Using conditional filters to set limits

Let's return to the blistering season of Drayton Whiteside using our dataset of every college basketball player's season stats in 2018-19 season. How can we set limits in something like a question of who had the best season?

Let's get our Drayton Whiteside data from the previous chapter back up.

```
players <- read_csv("data/players19.csv")
```

```
## Warning: Missing column names filled in: 'X1' [1]
##
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Team = col_character(),
```

```

##   Conference = col_character(),
##   Player = col_character(),
##   Class = col_character(),
##   Pos = col_character(),
##   Height = col_character(),
##   Hometown = col_character(),
##   `High School` = col_character(),
##   Summary = col_character()
## )
## i Use `spec()` for the full column specifications.

players %>%
  mutate(trueshooting = (PTS/(2*(FGA + (.44*FTA))))*100) %>%
  arrange(desc(trueshooting))

## # A tibble: 5,386 x 60
##       X1 Team Conference Player  `#` Class Pos  Height Weight Hometown
##   <dbl> <chr> <chr>     <chr> <dbl> <chr> <chr> <chr> <dbl> <chr>
## 1    579 Texa~ Big 12     Drayt~     4 JR     G   6-0    156 Austin,~
## 2    843 Ston~ AEC        Nick ~    42 FR     F   6-7    240 Port Je~
## 3   1059 Sout~ Southland Patri~    22 SO     F   6-3    210 Folsom,~
## 4   4269 Dayt~ A-10      Camro~    52 SO     G   5-7    160 Country~
## 5   4681 Cali~ Pac-12    David~    21 JR     G   6-4    185 Newbury~
## 6    326 Virg~ ACC        Grant~     1 FR     G   <NA>   NA Charlott~
## 7    410 Vand~ SEC        Mac H~    42 FR     G   6-6    182 Chattan~
## 8   1390 Sain~ A-10      Jack ~    31 JR     G   6-6    205 Mattoon~
## 9   2230 NJIT~ A-Sun     Patri~     3 SO     G   5-9    160 West Or~
## 10   266 Wash~ Pac-12    Reaga~    34 FR     F   6-6    225 Santa A~
## # ... with 5,376 more rows, and 50 more variables: `High School` <chr>,
## #   Summary <chr>, Rk.x <dbl>, G <dbl>, GS <dbl>, MP <dbl>, FG <dbl>,
## #   FGA <dbl>, `FG%` <dbl>, `2P` <dbl>, `2PA` <dbl>, `2P%` <dbl>, `3P` <dbl>,
## #   `3PA` <dbl>, `3P%` <dbl>, FT <dbl>, FTA <dbl>, `FT%` <dbl>, ORB <dbl>,
## #   DRB <dbl>, TRB <dbl>, AST <dbl>, STL <dbl>, BLK <dbl>, TOV <dbl>, PF <dbl>,
## #   PTS <dbl>, Rk.y <dbl>, PER <dbl>, `TS%` <dbl>, `eFG%` <dbl>, `3PAr` <dbl>,
## #   FTr <dbl>, PProd <dbl>, `ORB%` <dbl>, `DRB%` <dbl>, `TRB%` <dbl>,
## #   `AST%` <dbl>, `STL%` <dbl>, `BLK%` <dbl>, `TOV%` <dbl>, `USG%` <dbl>,
## #   OWS <dbl>, DWS <dbl>, WS <dbl>, `WS/40` <dbl>, OBP <dbl>, DBP <dbl>,
## #   BPM <dbl>, trueshooting <dbl>

```

In most contests, like the batting title in Major League Baseball, there's a minimum number of X to qualify. In baseball, it's at bats. In basketball, it attempts. So let's set a floor and see how it changes. What if we said you had to have played 100 minutes in a season? The top players in college basketball play more than 1000 minutes in a season. So 100 is not that much. Let's try it and see.

```

players %>%
  mutate(trueshooting = (PTS/(2*(FGA + (.44*FTA)))*100) %>%
    arrange(desc(trueshooting)) %>%
    filter(MP > 100)

## # A tibble: 3,659 x 60
##       X1 Team Conference Player `#` Class Pos Height Weight Hometown
##   <dbl> <chr> <chr>     <chr> <dbl> <chr> <chr> <chr> <dbl> <chr>
## 1 4634 Cent~ Southland Jord~ 33 JR   G   6-1   185 Harris~ 
## 2 3623 Hart~ AEC      Max T~ 20 SR   G   6-5   200 Rye, NY  
## 3 2675 Mich~ Big Ten Thoma~ 15 FR   F   6-8   225 Clarkst~ 
## 4 5175 Litt~ Sun Belt Kris ~ 32 SO   F   6-8   194 Dewitt, ~
## 5 5205 Ariz~ Pac-12 De'Qu~ 32 SR   F   6-10  225 St. Tho~ 
## 6 4099 ETSU~ Southern Lucas~ 25 JR   C   7-0   220 De Lier~ 
## 7 3006 Loui~ Sun Belt Brand~ 0 SR   G   6-4   180 Hawthr~ 
## 8 570 Texa~ Big 12 Jaxso~ 10 FR   F   6-11  220 Lovelan~ 
## 9 1704 Pepp~ WCC      Victo~ 34 FR   C   6-9   200 Owerri, ~
## 10 4056 East~ MAC     Jalen~ 30 SO   F   6-9   215 Pasco, ~
## # ... with 3,649 more rows, and 50 more variables: `High School` <chr>,
## #   Summary <chr>, Rk.x <dbl>, G <dbl>, GS <dbl>, MP <dbl>, FG <dbl>,
## #   FGA <dbl>, `FG%` <dbl>, `2P` <dbl>, `2PA` <dbl>, `2P%` <dbl>, `3P` <dbl>,
## #   `3PA` <dbl>, `3P%` <dbl>, FT <dbl>, FTA <dbl>, `FT%` <dbl>, ORB <dbl>,
## #   DRB <dbl>, TRB <dbl>, AST <dbl>, STL <dbl>, BLK <dbl>, TOV <dbl>, PF <dbl>,
## #   PTS <dbl>, Rk.y <dbl>, PER <dbl>, `TS%` <dbl>, `eFG%` <dbl>, `3PAr` <dbl>,
## #   FTr <dbl>, PProd <dbl>, `ORB%` <dbl>, `DRB%` <dbl>, `TRB%` <dbl>,
## #   `AST%` <dbl>, `STL%` <dbl>, `BLK%` <dbl>, `TOV%` <dbl>, `USG%` <dbl>,
## #   OWS <dbl>, DWS <dbl>, WS <dbl>, `WS/40` <dbl>, OBPM <dbl>, DBPM <dbl>,
## #   BPM <dbl>, trueshooting <dbl>

```

Now you get Central Arkansas Bears Junior Jordan Grant, who played in 25 games and was on the floor for 152 minutes. So he played regularly. But in that time, he only attempted 16 shots, and made 68 percent of them. In other words, when he shot, he probably scored. He just rarely shot.

So is 100 minutes our level? Here's the truth – there's not really an answer here. We're picking a cutoff. If you can cite a reason for it and defend it, then it probably works.

### 6.3 Top list

One last little dplyr trick that's nice to have in the toolbox is a shortcut for selecting only the top values for your dataset. Want to make a Top 10 List? Or Top 25? Or Top Whatever You Want? It's easy.

So what are the top 10 Power Five schools by season attendance. All we're doing here is chaining commands together with what we've already got. We're

*filtering* by our list of Power Five conferences, we're *selecting* the three fields we need, now we're going to *arrange* it by total attendance and then we'll introduce the new function: `top_n`. The `top_n` function just takes a number. So we want a top 10 list? We do it like this:

```
attendance %>% filter(Conference %in% powerfive) %>% select(Institution, Conference, `2018`) %>%
```

```
## Selecting by 2018  
## # A tibble: 10 x 3  
##   Institution Conference `2018`  
##   <chr>       <chr>     <dbl>  
## 1 Michigan    Big Ten    775156  
## 2 Penn St.    Big Ten    738396  
## 3 Ohio St.   Big Ten    713630  
## 4 Alabama     SEC        710931  
## 5 LSU          SEC        705733  
## 6 Texas A&M  SEC        698908  
## 7 Tennessee   SEC        650887  
## 8 Georgia     SEC        649222  
## 9 Nebraska    Big Ten    623240  
## 10 Oklahoma   Big 12     607146
```

That's all there is to it. Just remember – for it to work correctly, you need to sort your data BEFORE you run `top_n`. Otherwise, you're just getting the first 10 values in the list. The function doesn't know what field you want the top values of. You have to do it.



# Chapter 7

## Transforming data

Sometimes long data needs to be wide, and sometimes wide data needs to be long. I'll explain.

You are soon going to discover that long before you can visualize data, **you need to have it in a form that the visualization library can deal with**. One of the ways that isn't immediately obvious is **how your data is cast**. Most of the data you will encounter will be **wide – each row will represent a single entity with multiple measures for that entity**. So think of states. Your row of your dataset could have the state name, population, average life expectancy and other demographic data.

But what if your visualization library needs one row for each measure? So state, data type and the data. Nebraska, Population, 1,929,000. That's one row. Then the next row is Nebraska, Average Life Expectancy, 76. That's the next row. That's where recasting your data comes in.

We can use a library called `tidyverse` to `pivot_longer` or `pivot_wider` the data, depending on what we need. We'll use a dataset of college football attendance to demonstrate.

First we need some libraries.

```
library(tidyverse)
```

Now we'll load the data.

```
attendance <- read_csv('data/attendance.csv')
```

```
##  
## -- Column specification -----  
## cols(  
##   Institution = col_character(),  
##   Conference = col_character(),
```

```

## `2013` = col_double(),
## `2014` = col_double(),
## `2015` = col_double(),
## `2016` = col_double(),
## `2017` = col_double(),
## `2018` = col_double()
## )
attendance

## # A tibble: 150 x 8
##   Institution Conference `2013` `2014` `2015` `2016` `2017` `2018`
##   <chr>        <chr>     <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Air Force    MWC      228562  168967  156158  177519  174924  166205
## 2 Akron        MAC      107101   55019  108588   62021  117416   92575
## 3 Alabama       SEC      710538  710736  707786  712747  712053  710931
## 4 Appalachian St. FBS Independent 149366     NA     NA     NA     NA     NA
## 5 Appalachian St. Sun Belt        NA  138995  128755  156916  154722  131716
## 6 Arizona       Pac-12     285713  354973  308355  338017  255791  318051
## 7 Arizona St.   Pac-12     501509  343073  368985  286417  359660  291091
## 8 Arkansas      SEC      431174  399124  471279  487067  442569  367748
## 9 Arkansas St.  Sun Belt     149477  149163  138043  136200  119538  119001
## 10 Army West Point FBS Independent 169781  171310  185946  163267  185543  190156
## # ... with 140 more rows

```

So as you can see, each row represents a school, and then each column represents a year. This is great for calculating the percent change – we can subtract a column from a column and divide by that column. But later, when we want to chart each school’s attendance over the years, we have to have each row be one team for one year. Nebraska in 2013, then Nebraska in 2014, and Nebraska in 2015 and so on.

To do that, we use `pivot_longer` because we’re making wide data long. Since all of the columns we want to make rows start with 20, we can use that in our `cols` directive. Then we give that column a name – Year – and the values for each year need a name too. Those are the attendance figure. We can see right away how this works.

```

attendance %>% pivot_longer(cols = starts_with("20"), names_to = "Year", values_to = "Attendance")

## # A tibble: 900 x 4
##   Institution Conference Year   Attendance
##   <chr>        <chr>   <chr>     <dbl>
## 1 Air Force    MWC     2013      228562
## 2 Air Force    MWC     2014      168967
## 3 Air Force    MWC     2015      156158
## 4 Air Force    MWC     2016      177519
## 5 Air Force    MWC     2017      174924

```

```

## 6 Air Force MWC 2018 166205
## 7 Akron MAC 2013 107101
## 8 Akron MAC 2014 55019
## 9 Akron MAC 2015 108588
## 10 Akron MAC 2016 62021
## # ... with 890 more rows

```

We've gone from 150 rows to 900, but that's expected when we have 6 years for each team.

## 7.1 Making long data wide

We can reverse this process using `pivot_wider`, which makes long data wide.

Why do any of this?

In some cases, you're going to be given long data and you need to calculate some metric using two of the years – a percent change for instance. So you'll need to make the data wide to do that. You might then have to re-lengthen the data now with the percent change. Some project require you to do all kinds of flexing like this. It just depends on the data.

So let's take what we made above and turn it back into wide data.

```

longdata <- attendance %>% pivot_longer(cols = starts_with("20"), names_to = "Year", values_to =
longdata

## # A tibble: 900 x 4
##   Institution Conference Year Attendance
##   <chr>        <chr>    <chr>     <dbl>
## 1 Air Force    MWC      2013     228562
## 2 Air Force    MWC      2014     168967
## 3 Air Force    MWC      2015     156158
## 4 Air Force    MWC      2016     177519
## 5 Air Force    MWC      2017     174924
## 6 Air Force    MWC      2018     166205
## 7 Akron        MAC      2013     107101
## 8 Akron        MAC      2014     55019
## 9 Akron        MAC      2015     108588
## 10 Akron       MAC      2016     62021
## # ... with 890 more rows

```

To `pivot_wider`, we just need to say where our column names are coming from – the Year – and where the data under it should come from – Attendance.

```

longdata %>% pivot_wider(names_from = Year, values_from = Attendance)

## # A tibble: 150 x 8
##   Institution Conference 2013 2014 2015 2016 2017 2018
##   <chr>        <chr>    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Air Force    MWC      228562 168967 156158 177519 174924 166205
## 2 Akron        MAC      107101 55019 108588 62021  NA     NA
## 3 Akron        MAC      107101 55019 108588 62021  NA     NA
## 4 Akron        MAC      107101 55019 108588 62021  NA     NA
## 5 Akron        MAC      107101 55019 108588 62021  NA     NA
## 6 Akron        MAC      107101 55019 108588 62021  NA     NA
## 7 Akron        MAC      107101 55019 108588 62021  NA     NA
## 8 Akron        MAC      107101 55019 108588 62021  NA     NA
## 9 Akron        MAC      107101 55019 108588 62021  NA     NA
## 10 Akron       MAC      107101 55019 108588 62021  NA     NA
## # ... with 140 more rows

```

```

##   Institution    Conference `2013` `2014` `2015` `2016` `2017` `2018`
##   <chr>          <chr>      <dbl>  <dbl>  <dbl>  <dbl>  <dbl>  <dbl>
## 1 Air Force      MWC        228562 168967 156158 177519 174924 166205
## 2 Akron          MAC        107101 55019 108588 62021 117416 92575
## 3 Alabama         SEC        710538 710736 707786 712747 712053 710931
## 4 Appalachian St. FBS Independent 149366 NA NA NA NA NA
## 5 Appalachian St. Sun Belt       NA 138995 128755 156916 154722 131716
## 6 Arizona          Pac-12     285713 354973 308355 338017 255791 318051
## 7 Arizona St.     Pac-12     501509 343073 368985 286417 359660 291091
## 8 Arkansas         SEC        431174 399124 471279 487067 442569 367748
## 9 Arkansas St.    Sun Belt     149477 149163 138043 136200 119538 119001
## 10 Army West Point FBS Independent 169781 171310 185946 163267 185543 190156
## # ... with 140 more rows

```

And just like that, we're back.

## 7.2 Why this matters

This matters because certain visualization types need wide or long data. A significant hurdle you will face for the rest of the semester is getting the data in the right format for what you want to do.

So let me walk you through an example using this data.

Let's look at Nebraska's attendance over the time period. In order to do that, I need long data because that's what the charting library, `ggplot2`, needs. You're going to learn a lot more about `ggplot` later.

```
nebraska <- longdata %>% filter(Institution == "Nebraska")
```

Now that we have long data for just Nebraska, we can chart it.

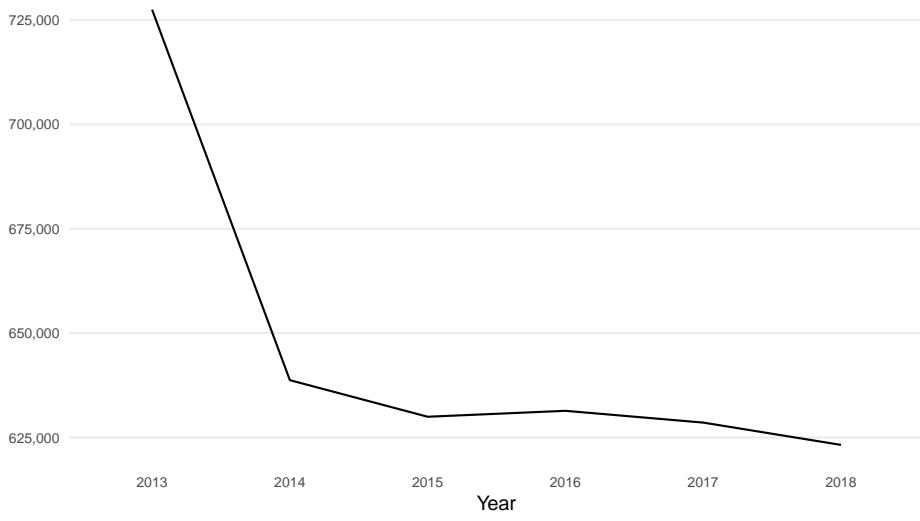
```

ggplot(nebraska, aes(x=Year, y=Attendance, group=1)) +
  geom_line() +
  scale_y_continuous(labels = scales::comma) +
  labs(x="Year", y="Attendance", title="We'll all stick together?", subtitle="It's not",
       theme_minimal() +
       theme(
         plot.title = element_text(size = 16, face = "bold"),
         axis.title = element_text(size = 10),
         axis.title.y = element_blank(),
         axis.text = element_text(size = 7),
         axis.ticks = element_blank(),
         panel.grid.minor = element_blank(),
         panel.grid.major.x = element_blank(),
         legend.position="bottom"
       )

```

## We'll all stick together?

It's not as bad as you think -- they widened the seats, cutting the number.



Source: NCAA | By Matt Waite



# Chapter 8

## Significance tests

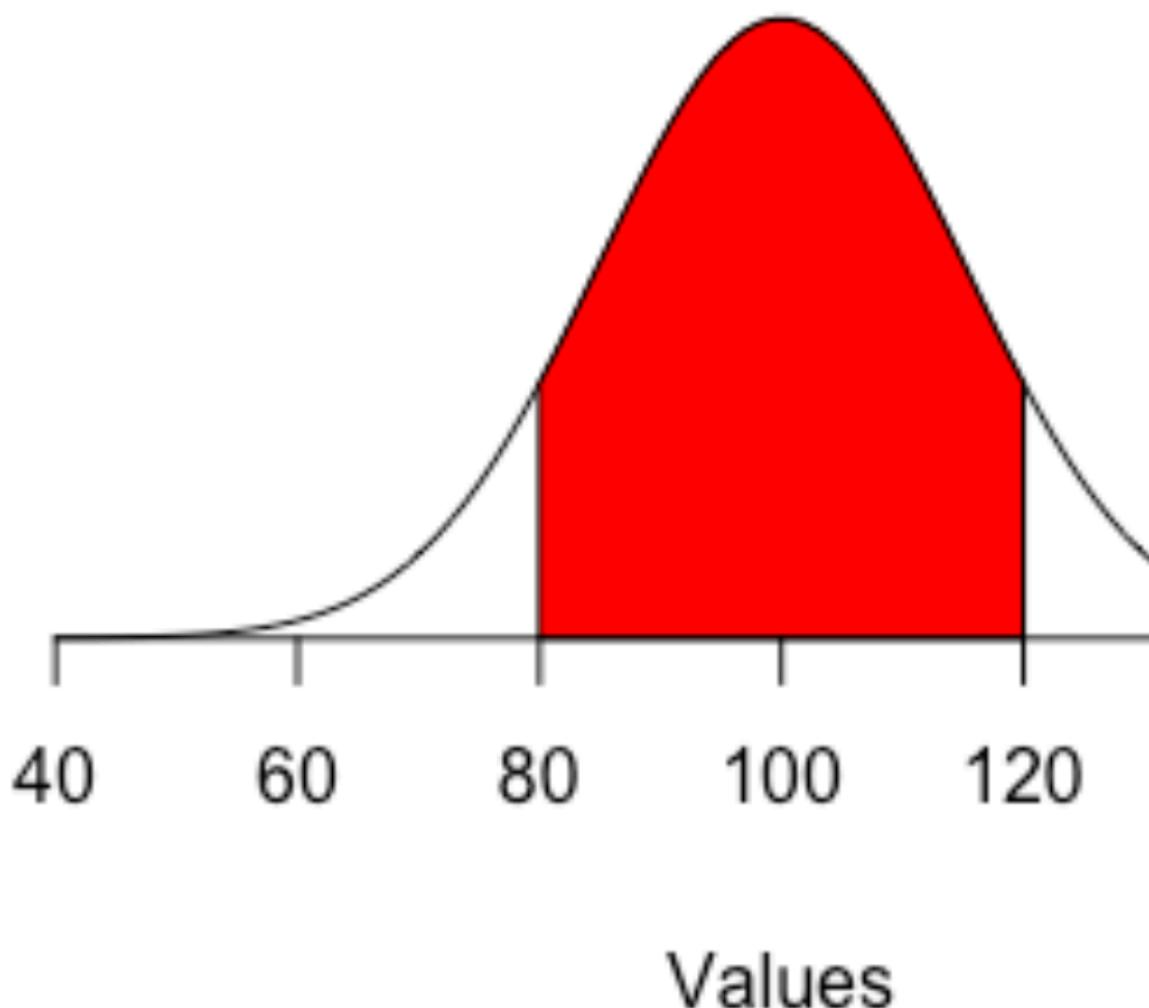
Now that we've worked with data a little, it's time to start asking more probing questions of our data. One of the most probing questions we can ask – one that so few sports journalists ask – is if the difference between this thing and the normal thing is real.

We have a perfect natural experiment going on in sports right now to show how significance tests work. The NBA, to salvage a season and get to the playoffs, put their players in a bubble – more accurately a hotel complex at Disney World in Orlando – and had them play games without fans.

So are the games different from other regular season games that had fans?

To answer this, we need to understand that a significance test is a way to determine if two numbers are *significantly* different from each other. Generally speaking, we're asking if a subset of data – a sample – is different from the total data pool – the population. Typically, this relies on data being in a normal distribution.

## Normal Distribution



If it is, then we know certain things about it. Like the mean – the average – will be a line right at the peak of cases. And that 66 percent of cases will be in that red area – the first standard deviation.

A significance test will determine if a sample taken from that group is different from the total.

Significance testing involves stating a hypothesis. In our case, our hypothesis is that there is a difference between bubble games without people and regular games with people.

In statistics, the **null hypothesis** is the opposite of your hypothesis. In this case, that there is no difference between fans and no fans.

What we're driving toward is a metric called a p-value, which is the probability that you'd get your sample mean *if the null hypothesis is true*. So in our case, it's the probability we'd see the numbers we get if there was no difference between fans and no fans. If that probability is below .05, then we consider the difference significant and we reject the null hypothesis.

So let's see. We'll need a log of every game last NBA season. In this data, there's a field called COVID, which labels the game as a regular game or a bubble game.

Load the tidyverse.

```
library(tidyverse)
```

And import the data.

```
logs <- read_csv("data/nbabubble.csv")
```

```
## 
## -- Column specification -----#
## cols(
##   .default = col_double(),
##   Season = col_character(),
##   Conference = col_character(),
##   Team = col_character(),
##   Date = col_date(format = ""),
##   HomeAway = col_character(),
##   Opponent = col_character(),
##   W_L = col_character(),
##   COVID = col_character()
## )
## i Use `spec()` for the full column specifications.
```

First, let's just look at scoring. Here's a theory: fans make players nervous. The screaming makes players tense up, and tension makes for bad shooting. An alternative to this: screaming fans make you defend harder. So my hypothesis is that not only is the scoring different, it's lower.

First things first, let's create a new field, called **totalpoints** and add the two scores together. We'll need this, so we're going to make this a new dataframe called **points**.

```
points <- logs %>% mutate(totalpoints = TeamScore + OpponentScore )
```

Typically speaking, with significance tests, the process involves creating two different means and then running a bunch of formulas on them. R makes this easy by giving you a `t.test` function, which does all the work for you. What we have to tell it is what is the value we are testing, over which groups, and from what data. It looks like this:

```
t.test(totalpoints ~ COVID, data=points)

##
## Welch Two Sample t-test
##
## data: totalpoints by COVID
## t = -5.232, df = 206.88, p-value = 4.099e-07
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -11.64698 -5.27178
## sample estimates:
## mean in group With Fans mean in group Without Fans
## 222.8929 231.3523
```

Now let's talk about the output. I prefer to read these bottom up. So at the bottom, it says that the mean number of points score in an NBA game With Fans is 222.89. The mean scored in games Without Fans is 231.35. That means teams are scoring almost 8.5 points MORE without fans on average.

But, some games are defenseless track meets, some games are defensive slugfests. We learned that averages can be skewed by extremes. So the next thing we need to look at is the p-value. Remember, this is the probability that we'd get this sample mean – the without fans mean – if there was no difference between fans and no fans.

The probability?  $4.099e-07$  or  $4.099 \times 10^{-7}$ . Don't remember your scientific notation? That's  $.00000004099$ . The decimal, seven zeros and the number.

Remember, if the probability is below .05, then we determine that this number is statistically significant. We'll talk more about statistical significance soon, but in this case, statistical significance means that our hypothesis is correct: points are different without fans than with. And since our hypothesis is correct, we *reject the null hypothesis* and we can confidently say that bubble teams are scoring more than they were when fans packed arenas.

## 8.1 Accepting the null hypothesis

So what does it look like when your hypothesis is wrong?

Let's test another thing that may have been impacted by bubble games: home court advantage. If you're the home team, but you're not at home, does it affect you? It has to, right? Your fans aren't there. Home and away are just positions on the scoreboard. It can't matter, can it?

My hypothesis is that home court is no longer an advantage, and the home team will score less relative to the away team.

First things first: We need to make a dataframe where Team is the home team. And then we'll create a differential between the home team and away team. If home court is an advantage, the differential should average out to be positive – the home team scores more than the away team.

```
homecourt <- logs %>% filter(is.na(HomeAway) == TRUE) %>% mutate(differential = TeamScore - Oppon
```

Now let's test it.

```
t.test(differential ~ COVID, data=homecourt)

##
## Welch Two Sample t-test
##
## data: differential by COVID
## t = 0.36892, df = 107.84, p-value = 0.7129
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -2.301628 3.354268
## sample estimates:
## mean in group With Fans mean in group Without Fans
## 2.174047 1.647727
```

So again, start at the bottom. With Fans, the home team averages 2.17 more points than the away team. Without fans, they average 1.64 more.

If you are a bad sportswriter or a hack sports talk radio host, you look at this and scream "the bubble killed home court!"

But two things: first, the home team is STILL, on average, scoring more than the away team on the whole.

And two: Look at the p-value. It's .7129. Is that less than .05? No, no it is not. So that means we have to **accept the null hypothesis** that there is no difference between fans and no fans when it comes to the difference between the home team and the away team's score.

Now, does this mean that the bubble hasn't impacted the magic of home court? Not necessarily. What it's saying is that the variance between one and the other is too large to be able to say that they're different. It could just be random noise that's causing the difference, and so it's not real. More to the point, it's saying that this metric isn't capable of telling you that there's no home court in the bubble.

We're going to be analyzing these bubble games for *years* trying to find the true impact of fans.

# Chapter 9

## Correlations and regression

Throughout sports, you will find no shortage of opinions. From people yelling at their TV screens to an entire industry of people paid to have opinions, there are no shortage of reasons why this team sucks and that player is great. They may have their reasons, but a better question is, does that reason really matter?

Can we put some numbers behind that? Can we prove it or not?

This is what we're going to start to answer. And we'll do it with correlations and regressions.

First, we need data from the 2020 college football season.

Then load the tidyverse.

```
library(tidyverse)
```

Now import the data.

```
correlations <- read_csv("data/footballlogs20.csv")
```

```
## 
## -- Column specification -----  
## cols(  
##   .default = col_double(),  
##   Date = col_date(format = ""),  
##   HomeAway = col_character(),  
##   Opponent = col_character(),  
##   Result = col_character(),  
##   TeamFull = col_character(),  
##   TeamURL = col_character(),  
##   Outcome = col_character(),  
##   Team = col_character(),  
##   Conference = col_character()
```

```
## )
## i Use `spec()` for the full column specifications.
```

To do this, we need all FBS college football teams and their season stats from last year. How much, over the course of a season, does a thing matter? That's the question you're going to answer.

In our case, we want to know how much does a team's accumulated penalties influence the number of points they score in a season? How much difference can we explain in points with penalties?

We're going to use two different methods here and they're closely related. Correlations – specifically the Pearson Correlation Coefficient – is a measure of how related two numbers are in a linear fashion. In other words – if our X value goes up one, what happens to Y? If it also goes up 1, that's a perfect correlation. X goes up 1, Y goes up 1. Every time. Correlation coefficients are a number between 0 and 1, with zero being no correlation and 1 being perfect correlation **if our data is linear**. We'll soon go over scatterplots to visually determine if our data is linear, but for now, we have a hypothesis: More penalties are bad. Penalties hurt. So if a team gets lots of them, they should have worse outcomes than teams that get few of them. That is an argument for a linear relationship between them.

But is there one?

We're going to create a new dataframe called `newcorrelations` that takes our data that we imported and adds a column called `differential` because we don't have separate offense and defense penalties, and then we'll use correlations to see how related those two things are.

```
newcorrelations <- correlations %>%
  mutate(
    differential = TeamScore - OpponentScore,
    TotalPenalties = Penalties+DefPenalties,
    TotalPenaltyYards = PenaltyYds+DefPenaltyYds
  )
```

In R, there is a `cor` function, and it works much the same as `mean` or `median`. So we want to see if `differential` is correlated with `TotalPenaltyYards`, which is the yards of penalties a team gets in a game. We do that by referencing `differential` and `TotalPenaltyYards` and specifying we want a `pearson` correlation. The number we get back is the correlation coefficient.

```
newcorrelations %>% summarise(correlation = cor(differential, TotalPenaltyYards, method = "pearson"))

## # A tibble: 1 x 1
##   correlation
##       <dbl>
## 1     -0.00655
```

So on a scale of -1 to 1, where 0 means there's no relationship at all and 1 or -1 means a perfect relationship, penalty yards and whether or not the team scores more points than it give up are at -0.0065. You could say they're .7 percent related toward the negative – more penalties, the lower your differential. Another way to say it? They're 99.3 percent not related.

What about the number of penalties instead of the yards?

```
newcorrelations %>%
  summarise(correlation = cor(differential, TotalPenalties, method="pearson"))

## # A tibble: 1 x 1
##   correlation
##       <dbl>
## 1     0.000676
```

So wait, what does this all mean?

It means that when you look at every game in college football, the number of penalties and penalty yards does have a negative impact on the score difference between your team and the other team. But the relationship between penalties, penalty yards and the difference between scores is barely anything at all. Like 99 percent plus not related.

Normally, at this point, you'd quit while you were ahead. A correlation coefficient that shows there's no relationship between two things means stop. It's pointless to go on. But let's beat a dead horse a bit for the sake of talk radio callers who want to complain about undisciplined football teams.

Enter regression. Regression is how we try to fit our data into a line that explains the relationship the best. Regressions will help us predict things as well – if we have a team that has so many penalties, what kind of point differential could we expect? So regressions are about prediction, correlations are about description. Correlations describe a relationship. Regressions help us predict what that relationship means and what it might look like in the real world. Specifically, it tells us how much of the change in a dependent variable can be explained by the independent variable.

Another thing regressions do is give us some other tools to evaluate if the relationship is real or not.

Here's an example of using linear modeling to look at penalty yards. Think of the ~ character as saying "is predicted by". The output looks like a lot, but what we need is a small part of it.

```
fit <- lm(differential ~ TotalPenaltyYards, data = newcorrelations)
summary(fit)

## 
## Call:
```

```

## lm(formula = differential ~ TotalPenaltyYards, data = newcorrelations)
##
## Residuals:
##   Min     1Q Median     3Q    Max
## -67.02 -14.68    0.24  14.04  64.98
##
## Coefficients:
##                   Estimate Std. Error t value Pr(>|t|)
## (Intercept)      1.226995  1.817342  0.675   0.500
## TotalPenaltyYards -0.003383  0.015816 -0.214   0.831
##
## Residual standard error: 21.46 on 1068 degrees of freedom
## Multiple R-squared:  4.284e-05, Adjusted R-squared:  -0.0008935
## F-statistic: 0.04575 on 1 and 1068 DF, p-value: 0.8307

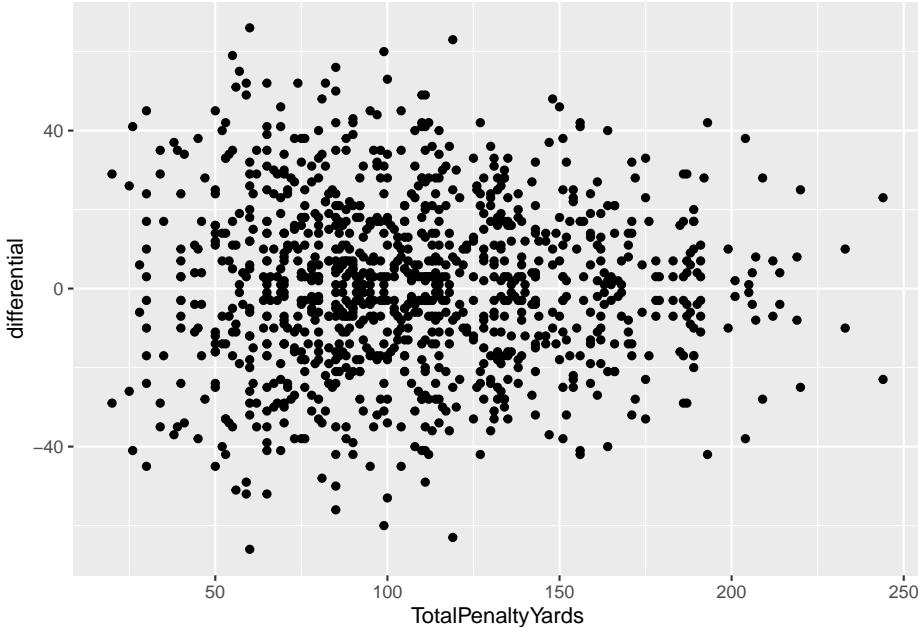
```

There's three things we need here:

1. First we want to look at the p-value. It's at the bottom right corner of the output. In the case of Total Penalty Yards, the p-value is .8307. The threshold we're looking for here is .05. If it's less than .05, then the relationship is considered to be *statistically significant*. Significance here does not mean it's a big deal. It means it's not random. That's it. Just that. Not random. So in our case, the relationship between total penalty yards and a team's aggregate point differential are **not statistically significant**. The differences in score difference and penalty yards could be completely random. This is another sign we should just stop with this.
2. Second, we look at the Adjusted R-squared value. It's right above the p-value. Adjusted R-squared is a measure of how much of the difference between teams aggregate point values can be explained by penalty yards. Our correlation coefficient said they're .7 percent related to each other, but penalty yard's ability to explain the difference between teams? About .08 percent. That's ... not much. It's really nothing. Again, we should quit.
3. The third thing we can look at, and we only bother if the first two are meaningful, is the coefficients. In the middle, you can see the (Intercept) is 1.226995 and the TotalPenaltyYards coefficient is -0.003383. Remember high school algebra? Remember learning the equation of a line? Remember swearing that learning  $y=mx+b$  is stupid because you'll never need it again? Surprise. It's useful again. In this case, we could try to predict a team's score differential in a game – will they score more than they give up – by using  $y=mx+b$ . In this case, y is the aggregate score, m is -0.003383 and b is 1.226995. So we would multiply a team's total penalty yards by -0.003383 and then add 1.226995 to it. The result would tell you what the total aggregate score in the game would be, according to our model. Chance that your even close with this? About .08 percent. In other words, you've got a 99.92 percent chance of being completely wrong. Did I say we should quit? Yeah.

So penalty yards are totally meaningless to the outcome of a game.

You can see the problem in a graph. On the X axis is penalty yards, on the y is aggregate score. If these elements had a strong relationship, we'd see a clear pattern moving from right to left, sloping down. On the left would be the teams with few penalties and a positive point differential. On right would be teams with high penalty yards and negative point differentials. Do you see that below?



## 9.1 A more predictive example

So we've **firmlly** established that penalties aren't predictive. But what is?

So instead of looking at penalty yards, let's make a new metric: Net Yards. Can we predict the score differential by looking at the yards a team gained minus the yards they gave up.

```
regressions <- newcorrelations %>% mutate(NetYards = OffensiveYards - DefYards)
```

First, let's look at the correlation coefficient.

```
regressions %>%
  summarise(correlation = cor(differential, NetYards, method="pearson"))

## # A tibble: 1 x 1
##   correlation
##       <dbl>
## 1     0.806
```

Answer: 81 percent. Not a perfect relationship, but very good. But how meaningful is that relationship and how predictive is it?

```
net <- lm(differential ~ NetYards, data = regressions)
summary(net)

##
## Call:
## lm(formula = differential ~ NetYards, data = regressions)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -49.379  -8.520   0.059   8.487  48.748
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 0.315632  0.388994  0.811   0.417
## NetYards    0.102473  0.002307 44.427  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.72 on 1068 degrees of freedom
## Multiple R-squared:  0.6489, Adjusted R-squared:  0.6486
## F-statistic: 1974 on 1 and 1068 DF,  p-value: < 2.2e-16
```

First we check p-value. See that e-16? That means scientific notation. That means our number is 2.2 times 10 to the -16 power. So -.0000000000000022. That's sixteen zeros between the decimal and 22. Is that less than .05? Uh, yeah. So this is really, really, really not random. But anyone who has watched a game of football knows this is true. It makes intuitive sense.

Second, Adjusted R-squared: 0.6486. So we can predict a whopping 65 percent of the difference in the score differential by simply looking at the net yards the team has.

Third, the coefficients: In this case, our  $y=mx+b$  formula looks like  $y = 0.102473x + 0.315632$ . So if we were applying this, let's look at Nebraska's 26-20 loss to Iowa in 2020. Nebraska's net yards that game? 16. That's right – we outgained them.

`(0.102473*16)+0.315632`

`## [1] 1.9552`

So by our model, Nebraska should have won by 1.96 points. Some games are closer than others. But when you can explain 65 percent of the difference, this is the kind of result you get. What would improve the model? Using more data to start. And using more inputs.

# Chapter 10

## Multiple regression

Last chapter, we looked at correlations and linear regression to predict how one element of a game would predict the score. But we know that a single variable, in all but the rarest instances, is not going to be that predictive. We need more than one. Enter multiple regression. Multiple regression lets us add – wait for it – multiple predictors to our equation to help us get a better fit to reality.

That presents it's own problems. So let's get set up. The dataset we'll use is all college football games since the 2011 season.

We need the tidyverse.

```
library(tidyverse)
```

And the data.

```
logs <- read_csv("data/footballlogs1120.csv")
```

```
##  
## -- Column specification -----  
## cols(  
##   .default = col_double(),  
##   Date = col_date(format = ""),  
##   HomeAway = col_character(),  
##   Opponent = col_character(),  
##   Result = col_character(),  
##   TeamFull = col_character(),  
##   TeamURL = col_character(),  
##   Outcome = col_character(),  
##   Team = col_character(),  
##   Conference = col_character()  
## )  
## i Use `spec()` for the full column specifications.
```

One way to show how successful a football team was for a game is to show the differential between the team's score and the opponent's score. Score a lot more than the opponent = good, score a lot less than the opponent = bad. And, relatively speaking, the more the better. So let's create that differential. Let's also get our net yardage stat back. And because we'll need it later, let's add the turnover margin.

```
logs <- logs %>% mutate(
  Differential = TeamScore - OpponentScore,
  NetYards = OffensiveYards - DefYards,
  TurnoverMargin = DefTotalTurnovers - TotalTurnovers)
```

The linear model code we used before is pretty straight forward. Its `field` is predicted by `field`. Here's a simple linear model that looks at predicting a team's point differential by looking at their net yards.

```
yards <- lm(Differential ~ NetYards, data=logs)
summary(yards)

##
## Call:
## lm(formula = Differential ~ NetYards, data = logs)
##
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -52.052  -8.736  -0.004   8.753  64.486 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 0.4918375  0.1052811  4.672 3.01e-06 ***
## NetYards    0.1046592  0.0005882 177.920  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 13.12 on 15605 degrees of freedom
## Multiple R-squared:  0.6698, Adjusted R-squared:  0.6698 
## F-statistic: 3.166e+04 on 1 and 15605 DF,  p-value: < 2.2e-16
```

Remember: There's a lot here, but only some of it we care about. What is the Adjusted R-squared value? What's the p-value and is it less than .05? In this case, we can predict 67 percent of the difference in differential with the net yardage in the game.

To add more predictors to this mix, we merely add them. But it's not that simple, as you'll see in a moment. So first, let's look at adding how well the other team shot to our prediction model:

```

model1 <- lm(Differential ~ NetYards + TurnoverMargin, data=logs)
summary(model1)

##
## Call:
## lm(formula = Differential ~ NetYards + TurnoverMargin, data = logs)
##
## Residuals:
##    Min      1Q  Median      3Q     Max 
## -39.031 -6.983 -0.025  6.932 40.789 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 0.4608403  0.0846844  5.442 5.35e-08 ***
## NetYards     0.0966321  0.0004811 200.864 < 2e-16 ***
## TurnoverMargin 4.1968337  0.0454801  92.278 < 2e-16 ***
## ---      
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
##
## Residual standard error: 10.55 on 15604 degrees of freedom
## Multiple R-squared:  0.7864, Adjusted R-squared:  0.7864 
## F-statistic: 2.872e+04 on 2 and 15604 DF,  p-value: < 2.2e-16

```

First things first: What is the adjusted R-squared?

Second: what is the p-value and is it less than .05?

Third: Compare the residual standard error. We went from 13.12 to 10.55. The meaning of this is both really opaque and also simple – by adding data, we reduced the amount of error in our model. Residual standard error is the total distance between what our model would predict and what we actually have in the data. So lots of residual error means the distance between reality and our model is wider. So the width of our predictive range in this example shrank while we improved the amount of the difference we could predict. That's good, and not always going to be the case.

One of the more difficult things to understand about multiple regression is the issue of multicollinearity. What that means is that there is significant correlation overlap between two variables – the two are related to each other as well as to the target output – and all you are doing by adding both of them is adding error with no real value to the R-squared. In pure statistics, we don't want any multicollinearity at all. Violating that assumption limits the applicability of what you are doing. So if we have some multicollinearity, it limits our scope of application to college football. We can't say this will work for every football league and level everywhere. What we need to do is see how correlated each value is to each other and throw out ones that are highly co-correlated.

So to find those, we have to create a correlation matrix that shows us how

each value is correlated to our outcome variable, but also with each other. We can do that in the `Hmisc` library. We install that in the console with

```
install.packages("Hmisc")
```

```
library(Hmisc)
```

We can pass in every numeric value to the `Hmisc` library and get a correlation matrix out of it, but since we have a large number of values – and many of them character values – we should strip that down and reorder them. So that's what I'm doing here. I'm saying give me differential first, and then columns 9-24, and then 26-41. Why the skip? There's a blank column in the middle of the data – a remnant of the scraper I used.

```
simplelogs <- logs %>% select_if(is.numeric) %>% select(-Game) %>% select(Differential
```

Before we proceed, what we're looking to do is follow the Differential column down, looking for correlation values near 1 or -1. Correlations go from -1, meaning perfect negative correlation, to 0, meaning no correlation, to 1, meaning perfect positive correlation. So we're looking for numbers near 1 or -1 for their predictive value. BUT: We then need to see if that value is also highly correlated with something else. If it is, we have a decision to make.

We get our correlation matrix like this:

```
cormatrix <- rcorr(as.matrix(simplelogs))  
cormatrix$r
```

	Differential	NetYards	TurnoverMargin	PassingCmp
## Differential	1.0000000	8.184193e-01	0.4837849198	0.03118170
## NetYards	0.81841934	1.000000e+00	0.1808173764	0.19290045
## TurnoverMargin	0.48378492	1.808174e-01	1.0000000000	-0.09233553
## PassingCmp	0.03118170	1.929004e-01	-0.0923355289	1.00000000
## PassingAtt	-0.19686950	-7.097271e-03	-0.1939446793	0.88589623
## PassingPct	0.42410345	4.230815e-01	0.1567432526	0.50717498
## PassingYds	0.22322068	3.660003e-01	-0.0255596419	0.80654722
## PassingTD	0.41571492	3.776967e-01	0.1458497043	0.42231692
## RushingAtt	0.36621257	4.169451e-01	0.1942339975	-0.36911655
## RushingYds	0.52945141	5.517793e-01	0.1959559925	-0.32669337
## RushingAvg	0.49317675	4.892094e-01	0.1512425735	-0.19464168
## RushingTD	0.57578332	4.734649e-01	0.2627916039	-0.13305246
## OffensivePlays	0.14922605	3.856779e-01	-0.0098582568	0.53304628
## OffensiveYards	0.58963784	7.222939e-01	0.1308007817	0.39936463
## OffenseAvg	0.61303088	6.273595e-01	0.1640078840	0.15135101
## FirstDownPass	0.17742548	3.296124e-01	-0.0519043178	0.85823184
## FirstDownRush	0.46104465	5.138962e-01	0.1438888221	-0.24594648
## FirstDownPen	-0.00567194	-1.419963e-02	-0.0212171554	0.15771675
## FirstDownTotal	0.46569860	6.160784e-01	0.0606678892	0.48947643

## Penalties	-0.01446172	6.431274e-02	0.0245826177	0.13394910
## PenaltyYds	0.01760113	9.250198e-02	0.0375348365	0.12542719
## Fumbles	-0.14420105	-9.215913e-05	-0.4544647395	0.01753247
## Interceptions	-0.34740859	-1.828408e-01	-0.5658983150	0.09952907
## TotalTurnovers	-0.35226029	-1.357271e-01	-0.7176732980	0.08533780
## TeamScore	0.78106738	6.372500e-01	0.3770822640	0.18334268
## OpponentScore	-0.76964262	-6.319258e-01	-0.3731455129	0.13856056
## DefPassingCmp	-0.06473817	-2.227403e-01	0.0828218918	0.08346001
## DefPassingAtt	0.15048839	-3.476179e-02	0.1819010221	0.07599800
## DefPassingPct	-0.42046253	-4.221770e-01	-0.1547820852	0.03694525
## DefPassingYds	-0.25733994	-3.976691e-01	0.0170156333	0.12715093
## DefPassingTD	-0.41930782	-3.805155e-01	-0.1479955551	0.11207117
## DefRushingAtt	-0.36480427	-4.156720e-01	-0.1927488957	-0.03564161
## DefRushingYds	-0.52342461	-5.478198e-01	-0.1942127993	0.01864717
## DefRushingAvg	-0.48630318	-4.871370e-01	-0.1469244778	0.04905028
## DefRushingTD	-0.56152960	-4.645077e-01	-0.2584394301	0.06194592
## DefPlays	-0.19091439	-4.210217e-01	0.0001525553	0.04182614
## DefYards	-0.59891323	-7.299795e-01	-0.1317978120	0.11616839
## DefAvg	-0.61875849	-6.395677e-01	-0.1615742932	0.10878889
## DefFirstDownPass	-0.21709433	-3.658307e-01	0.0425861829	0.09123305
## DefFirstDownRush	-0.45863262	-5.137995e-01	-0.1442285028	-0.01962705
## DefFirstDownPen	-0.02107814	-9.546899e-03	0.0146372877	0.05162783
## DefFirstDownTotal	-0.48687095	-6.329252e-01	-0.0664607512	0.06425565
## DefPenalties	0.01317935	-6.452130e-02	-0.0302781991	0.09772103
## DefPenaltyYds	-0.02890210	-1.011338e-01	-0.0461923015	0.10806661
## DefFumbles	0.15069928	1.689405e-03	0.4611567538	-0.04092384
## DefInterceptions	0.33233883	1.665064e-01	0.5702003054	-0.02870947
## DefTotalTurnovers	0.34523008	1.250021e-01	0.7239648943	-0.04795582
## Season	-0.00363621	-1.569119e-03	-0.0028876028	-0.03091173
	PassingAtt	PassingPct	PassingYds	PassingTD
## Differential	-0.196869502	0.42410345	0.223220680	0.415714923
## NetYards	-0.007097271	0.42308150	0.366000270	0.377696725
## TurnoverMargin	-0.193944679	0.15674325	-0.025559642	0.145849704
## PassingCmp	0.885896233	0.50717498	0.806547223	0.422316922
## PassingAtt	1.000000000	0.09315536	0.687271795	0.272813521
## PassingPct	0.093155364	1.000000000	0.474576058	0.394989670
## PassingYds	0.687271795	0.47457606	1.000000000	0.628178519
## PassingTD	0.272813521	0.39498967	0.628178519	1.000000000
## RushingAtt	-0.467518367	0.03613210	-0.263107487	-0.059376574
## RushingYds	-0.449975567	0.10449207	-0.200609277	0.066350138
## RushingAvg	-0.303577976	0.13444998	-0.082888644	0.150499189
## RushingTD	-0.292374841	0.24175736	0.011342402	-0.029015586
## OffensivePlays	0.553796697	0.12658649	0.435433701	0.215293539
## OffensiveYards	0.207675438	0.46417599	0.653171610	0.558994355
## OffenseAvg	-0.083538631	0.47899455	0.504666780	0.524444297
## FirstDownPass	0.746157690	0.47052335	0.883657791	0.521551760

```

## FirstDownRush      -0.367019503  0.12282841  -0.161359790  0.067439565
## FirstDownPen       0.197351786 -0.01284499  0.129006081  0.102871859
## FirstDownTotal     0.327399378  0.43330605  0.563311167  0.456860192
## Penalties          0.141768049  0.03463586  0.151550273  0.082259734
## PenaltyYds         0.119567124  0.05599523  0.151131721  0.100160624
## Fumbles            0.027478866 -0.01783733  0.008811146  -0.048018859
## Interceptions      0.254940410 -0.23627286  0.017618719  -0.114458545
## TotalTurnovers     0.207169253 -0.18699980  0.018846333  -0.116391065
## TeamScore          -0.035933316  0.44835810  0.431051906  0.618539629
## OpponentScore      0.272001348 -0.20669709  0.090668157  -0.019512592
## DefPassingCmp      0.083570904  0.02206021  0.126357936  0.091773507
## DefPassingAtt      0.057944813  0.05692377  0.145063141  0.144169653
## DefPassingPct      0.077374243 -0.07337905  -0.002052903 -0.081413409
## DefPassingYds      0.156497317 -0.02196601  0.155256065  0.075376864
## DefPassingTD        0.171067592 -0.08316878  0.097982769  0.035053345
## DefRushingAtt      0.083993629 -0.21380037  -0.037088960 -0.051053515
## DefRushingYds      0.119747148 -0.17947974  -0.010300250 -0.071383520
## DefRushingAvg       0.114462808 -0.11505322  0.008649510 -0.065672712
## DefRushingTD        0.155716440 -0.15511709  0.028151375 -0.033926688
## DefPlays            0.135316859 -0.14277058  0.108624632  0.094751037
## DefYards            0.215484942 -0.15211246  0.117075612  0.007116432
## DefAvg              0.175755717 -0.10704691  0.066058087 -0.058840155
## DefFirstDownPass    0.109208886 -0.01119101  0.123142671  0.065395375
## DefFirstDownRush   -0.062313958 -0.15434159  -0.034360788 -0.068080474
## DefFirstDownPen     0.049336590  0.02691594  0.078576992  0.072627007
## DefFirstDownTotal   0.134090419 -0.11025611  0.083106670  0.016190934
## DefPenalties        0.144271935 -0.04301213  0.092139354  0.087539127
## DefPenaltyYds       0.162654164 -0.05390508  0.091809187  0.087458617
## DefFumbles          -0.035268306 -0.01779875  -0.037573349 0.043415306
## DefInterceptions    -0.067393814  0.06886121  0.008612120  0.088384660
## DefTotalTurnovers   -0.073067543  0.03967676  -0.018006148  0.093981546
## Season              -0.042815333  0.01153492  -0.006117554  0.007749168
##                               RushingAtt  RushingYds  RushingAvg  RushingTD
## Differential        0.366212575  0.529451415  0.493176753  0.575783321
## NetYards             0.416945120  0.551779347  0.489209446  0.473464850
## TurnoverMargin       0.194233997  0.195955992  0.151242574  0.262791604
## PassingCmp          -0.369116550 -0.326693369 -0.194641684 -0.133052460
## PassingAtt          -0.467518367 -0.449975567 -0.303577976 -0.292374841
## PassingPct           0.036132096  0.104492066  0.134449978  0.241757360
## PassingYds          -0.263107487 -0.200609277 -0.082888644  0.011342402
## PassingTD            -0.059376574  0.066350138  0.150499189 -0.029015586
## RushingAtt           1.000000000  0.736583943  0.364284984  0.490658949
## RushingYds           0.736583943  1.000000000  0.871565174  0.695903189
## RushingAvg           0.364284984  0.871565174  1.000000000  0.607637826
## RushingTD             0.490658949  0.695903189  0.607637826  1.000000000
## OffensivePlays       0.477140312  0.246460859  0.041324341  0.171497768

```

```

## OffensiveYards      0.356671646  0.610784556  0.606659400  0.547046394
## OffenseAvg         0.138425296  0.576907524  0.708081627  0.545667695
## FirstDownPass      -0.235040292 -0.195601164 -0.092469432 -0.002819467
## FirstDownRush       0.788796974  0.868884675  0.659404370  0.595246572
## FirstDownPen        -0.003204555 -0.066910301 -0.080759132 -0.001136082
## FirstDownTotal      0.404815319  0.474808460  0.392505703  0.432561158
## Penalties          -0.024921198 -0.003651045  0.017975757 -0.025539221
## PenaltyYds         -0.002199279  0.033168959  0.054284078  0.009108827
## Fumbles             0.025578601  -0.028118399 -0.058261942 -0.058709124
## Interceptions      -0.193975079  -0.216449107 -0.172564128 -0.233234180
## TotalTurnovers     -0.127168356  -0.179027989 -0.166216652 -0.211526460
## TeamScore           0.343571991  0.571429522  0.560340983  0.698724429
## OpponentScore       -0.223024981  -0.246049036 -0.200498698 -0.188554772
## DefPassingCmp       -0.044763075  0.002351343  0.033307036  0.043038468
## DefPassingAtt       0.073608528  0.098937001  0.093487708  0.130716692
## DefPassingPct       -0.216287969  -0.180720848 -0.115717537 -0.152791861
## DefPassingYds       -0.050233964  -0.030415762 -0.010078082  0.005118324
## DefPassingTD         -0.059010431  -0.078502621 -0.070786506 -0.044702528
## DefRushingAtt       -0.410216642  -0.276569387 -0.117694938 -0.196242310
## DefRushingYds       -0.279183263  -0.225077345 -0.131739440 -0.196349300
## DefRushingAvg        -0.121621052  -0.134136582 -0.108134174 -0.151165356
## DefRushingTD         -0.194448775  -0.192176011 -0.145311374 -0.156635569
## DefPlays             -0.309019722  -0.159689121 -0.017255978 -0.053600005
## DefYards              -0.249498265  -0.193040997 -0.106779574 -0.142960915
## DefAvg                -0.119690358  -0.146487656 -0.127188735 -0.152049693
## DefFirstDownPass     -0.090092408  -0.047490894 -0.008586427 -0.007229665
## DefFirstDownRush     -0.317355086  -0.220166727 -0.100294586 -0.173767081
## DefFirstDownPen      -0.047381988  -0.006469058  0.026743015  0.011628275
## DefFirstDownTotal    -0.301942200  -0.194065174 -0.073946954 -0.126211617
## DefPenalties          -0.044234601  -0.066155832 -0.057171981  0.001478081
## DefPenaltyYds        -0.060368138  -0.088133097 -0.076423381 -0.023191416
## DefFumbles            0.071037391  0.023590761 -0.006904102  0.081934337
## DefInterceptions      0.143215339  0.118930804  0.076394767  0.153598788
## DefTotalTurnovers    0.152731898  0.103830416  0.052362318  0.167539382
## Season                -0.006470762  0.015090607  0.029310251  0.010761615
## OffensivePlays      OffensiveYards  OffenseAvg   FirstDownPass
## Differential         0.1492260492  0.589637838  0.61303088  0.177425485
## NetYards              0.3856779392  0.722293918  0.62735954  0.329612380
## TurnoverMargin        -0.0098582568  0.130800782  0.16400788 -0.051904318
## PassingCmp            0.5330462787  0.399364632  0.15135101  0.858231839
## PassingAtt            0.5537966967  0.207675438 -0.08353863  0.746157690
## PassingPct            0.1265864872  0.464175994  0.47899455  0.470523352
## PassingYds             0.4354337011  0.653171610  0.50466678  0.883657791
## PassingTD              0.2152935388  0.558994355  0.52444430  0.521551760
## RushingAtt             0.4771403121  0.356671646  0.13842530 -0.235040292
## RushingYds             0.2464608593  0.610784556  0.57690752 -0.195601164

```

```

## RushingAvg      0.0413243412   0.606659400   0.70808163   -0.092469432
## RushingTD       0.1714977683   0.547046394   0.54566769   -0.002819467
## OffensivePlays 1.0000000000   0.542424397   0.04733582   0.520413497
## OffensiveYards  0.5424243966   1.000000000   0.85379007   0.563011577
## OffenseAvg      0.0473358159   0.853790068   1.000000000  0.350092632
## FirstDownPass   0.5204134972   0.563011577   0.35009263   1.000000000
## FirstDownRush   0.3781141523   0.541165081   0.41380414   -0.176035283
## FirstDownPen    0.1931821714   0.052549729   -0.04818553   0.127923002
## FirstDownTotal  0.7067980430   0.822273393   0.54765028   0.632578575
## Penalties       0.1174670579   0.119665050   0.06686068   0.095769976
## PenaltyYds     0.1167981247   0.147785771   0.10159079   0.096191798
## Fumbles         0.0514118771   -0.014611934   -0.04910112   -0.002669702
## Interceptions  0.0707420861   -0.153058437   -0.22086304   0.047918056
## TotalTurnovers  0.0861769432   -0.123142627   -0.19605167   0.033805622
## TeamScore       0.2878976199   0.790058468   0.76105515   0.348506805
## OpponentScore  0.0603404859   -0.116896356   -0.18323754   0.078087301
## DefPassingCmp   0.0409196811   0.103943303   0.09421622   0.091467148
## DefPassingAtt   0.1269410280   0.193714572   0.15018634   0.099017964
## DefPassingPct   -0.1268359214   -0.141363954   -0.09631045   0.008201755
## DefPassingYds   0.1082675145   0.101973122   0.04818640   0.125262682
## DefPassingTD    0.1144859698   0.018515887   -0.05340694   0.090586531
## DefRushingAtt   -0.3028923478   -0.243742972   -0.10829113   -0.078858613
## DefRushingYds   -0.1439228873   -0.182292260   -0.13449228   -0.026930969
## DefRushingAvg   -0.0007670721   -0.096687847   -0.11851074   0.011869282
## DefRushingTD    -0.0283493049   -0.125784441   -0.13842291   0.018320394
## DefPlays        -0.1565479001   -0.035633936   0.04743620   0.024332221
## DefYards        -0.0207822877   -0.054582037   -0.06201003   0.080510479
## DefAvg          0.0619862494   -0.059835995   -0.11515218   0.073494860
## DefFirstDownPass 0.0237109843   0.062820458   0.05355860   0.094812163
## DefFirstDownRush -0.2369763215   -0.197943110   -0.09623351   -0.048360863
## DefFirstDownPen  0.0044182172   0.058507968   0.06481826   0.059488255
## DefFirstDownTotal -0.1511006437   -0.082828263   -0.01566227   0.043655644
## DefPenalties    0.1017644325   0.023336190   -0.02925259   0.074673629
## DefPenaltyYds   0.1048427505   0.006082620   -0.05169993   0.076942358
## DefFumbles       0.0318497808   -0.012133968   -0.02853415   -0.048511328
## DefInterceptions 0.0678983530   0.098884830   0.08018058   -0.012722149
## DefTotalTurnovers 0.0712216676   0.065699798   0.04112451   -0.040989258
## Season          -0.0486606092   0.006719496   0.03847022   -0.021594636
## FirstDownRush   FirstDownPen  FirstDownTotal  Penalties
## Differential   0.4610446515   -0.0056719400   0.465698595   -0.014461721
## NetYards        0.5138962499   -0.0141996290   0.616078422   0.064312738
## TurnoverMargin  0.1438888221   -0.0212171554   0.060667889   0.024582618
## PassingCmp      -0.2459464797   0.1577167545   0.489476431   0.133949104
## PassingAtt      -0.3670195030   0.1973517860   0.327399378   0.141768049
## PassingPct      0.1228284090   -0.0128449897   0.433306049   0.034635855
## PassingYds      -0.1613597905   0.1290060812   0.563311167   0.151550273

```

## PassingTD	0.0674395650	0.1028718590	0.456860192	0.082259734
## RushingAtt	0.7887969735	-0.0032045552	0.404815319	-0.024921198
## RushingYds	0.8688846753	-0.0669103014	0.474808460	-0.003651045
## RushingAvg	0.6594043700	-0.0807591323	0.392505703	0.017975757
## RushingTD	0.5952465724	-0.0011360821	0.432561158	-0.025539221
## OffensivePlays	0.3781141523	0.1931821714	0.706798043	0.117467058
## OffensiveYards	0.5411650814	0.0525497292	0.822273393	0.119665050
## OffenseAvg	0.4138041377	-0.0481855308	0.547650277	0.066860679
## FirstDownPass	-0.1760352833	0.1279230020	0.632578575	0.095769976
## FirstDownRush	1.0000000000	-0.0640698685	0.578677075	-0.036810209
## FirstDownPen	-0.0640698685	1.0000000000	0.271006893	0.131796886
## FirstDownTotal	0.5786770748	0.2710068926	1.0000000000	0.073197638
## Penalties	-0.0368102087	0.1317968858	0.073197638	1.0000000000
## PenaltyYds	-0.0038110643	0.1331976516	0.098139564	0.903246349
## Fumbles	-0.0127589089	-0.0001733341	-0.010522510	-0.005346654
## Interceptions	-0.1823861332	0.0214575248	-0.091526136	0.027625989
## TotalTurnovers	-0.1436894309	0.0158077753	-0.074806518	0.016995007
## TeamScore	0.4922799315	0.0696961188	0.632897568	0.045390014
## OpponentScore	-0.2196973570	0.0801568972	-0.083219329	0.069089389
## DefPassingCmp	-0.0351952073	0.0550827638	0.054621111	0.102064433
## DefPassingAtt	0.0425376517	0.0504467878	0.115086837	0.148129510
## DefPassingPct	-0.1563389121	0.0361246166	-0.098512229	-0.036848531
## DefPassingYds	-0.0530266913	0.0854845045	0.073302813	0.096286575
## DefPassingTD	-0.0724142119	0.0819003858	0.033416927	0.089942486
## DefRushingAtt	-0.3153770802	-0.0437639450	-0.299959508	-0.045386998
## DefRushingYds	-0.2196474882	0.0049716253	-0.181281160	-0.069403770
## DefRushingAvg	-0.1024397880	0.0402209526	-0.059266250	-0.059746435
## DefRushingTD	-0.1691695555	0.0296305740	-0.103644138	0.001729374
## DefPlays	-0.2514420854	0.0090586645	-0.165515581	0.103930020
## DefYards	-0.2071478834	0.0724343129	-0.076869161	0.025406067
## DefAvg	-0.1077475877	0.0813642808	-0.006762158	-0.026690341
## DefFirstDownPass	-0.0684673287	0.0650428178	0.030148892	0.078466818
## DefFirstDownRush	-0.2246581052	-0.0164089821	-0.213167195	-0.080014867
## DefFirstDownPen	-0.0271049617	0.0709464960	0.034896710	0.544815636
## DefFirstDownTotal	-0.2243526184	0.0464600265	-0.123036459	0.124940137
## DefPenalties	-0.0775132732	0.5421750750	0.124255326	0.190344623
## DefPenaltyYds	-0.0976956398	0.6585481016	0.137953942	0.193217609
## DefFumbles	0.0007492654	-0.0219681886	-0.039437448	0.017507901
## DefInterceptions	0.0855201980	-0.0006859365	0.051837528	0.054737340
## DefTotalTurnovers	0.0641260141	-0.0147851433	0.012948409	0.052108799
## Season	0.0087181647	0.0662421947	0.003573008	0.029358706
##				
##				
## Differential	0.017601130	-1.442010e-01	-0.347408586	-0.352260285
## NetYards	0.092501983	-9.215913e-05	-0.182840770	-0.135727097
## TurnoverMargin	0.037534837	-4.544647e-01	-0.565898315	-0.717673298
## PassingCmp	0.125427191	1.753247e-02	0.099529068	0.085337798

## PassingAtt	0.119567124	2.747887e-02	0.254940410	0.207169253
## PassingPct	0.055995229	-1.783733e-02	-0.236272865	-0.186999804
## PassingYds	0.151131721	8.811146e-03	0.017618719	0.018846333
## PassingTD	0.100160624	-4.801886e-02	-0.114458545	-0.116391065
## RushingAtt	-0.002199279	2.557860e-02	-0.193975079	-0.127168356
## RushingYds	0.033168959	-2.811840e-02	-0.216449107	-0.179027989
## RushingAvg	0.054284078	-5.826194e-02	-0.172564128	-0.166216652
## RushingTD	0.009108827	-5.870912e-02	-0.233234180	-0.211526460
## OffensivePlays	0.116798125	5.141188e-02	0.070742086	0.086176943
## OffensiveYards	0.147785771	-1.461193e-02	-0.153058437	-0.123142627
## OffenseAvg	0.101590794	-4.910112e-02	-0.220863037	-0.196051667
## FirstDownPass	0.096191798	-2.669702e-03	0.047918056	0.033805622
## FirstDownRush	-0.003811064	-1.275891e-02	-0.182386133	-0.143689431
## FirstDownPen	0.133197652	-1.733341e-04	0.021457525	0.015807775
## FirstDownTotal	0.098139564	-1.052251e-02	-0.091526136	-0.074806518
## Penalties	0.903246349	-5.346654e-03	0.027625989	0.016995007
## PenaltyYds	1.000000000	-1.403761e-02	0.015917475	0.002612761
## Fumbles	-0.014037613	1.000000e+00	0.020686856	0.670582481
## Interceptions	0.015917475	2.068686e-02	1.000000000	0.755548472
## TotalTurnovers	0.002612761	6.705825e-01	0.755548472	1.000000000
## TeamScore	0.081291125	-9.221149e-02	-0.276837519	-0.265831719
## OpponentScore	0.055514238	1.318612e-01	0.261769161	0.280630915
## DefPassingCmp	0.111411452	-3.973445e-02	-0.017382571	-0.038933078
## DefPassingAtt	0.162686876	-3.109202e-02	-0.058249476	-0.063593263
## DefPassingPct	-0.043181735	-2.298935e-02	0.079458063	0.043886543
## DefPassingYds	0.096147185	-3.794686e-02	0.023141664	-0.007693062
## DefPassingTD	0.093042717	3.917292e-02	0.099335433	0.099373662
## DefRushingAtt	-0.059428915	7.073269e-02	0.145548041	0.154342125
## DefRushingYds	-0.086487888	2.161863e-02	0.125655719	0.107400997
## DefRushingAvg	-0.073680859	-9.993266e-03	0.083883084	0.055685410
## DefRushingTD	-0.016742606	7.668263e-02	0.161302037	0.169930092
## DefPlays	0.105229463	3.512879e-02	0.077938067	0.080847102
## DefYards	0.012497218	-1.430737e-02	0.112722764	0.074264925
## DefAvg	-0.043324322	-3.375146e-02	0.094214745	0.047783145
## DefFirstDownPass	0.080096704	-4.961707e-02	0.001262556	-0.031573932
## DefFirstDownRush	-0.095022978	1.225369e-03	0.092104959	0.069144207
## DefFirstDownPen	0.656398393	-1.908304e-02	0.005624566	-0.008330444
## DefFirstDownTotal	0.140738712	-3.865811e-02	0.066245867	0.023823953
## DefPenalties	0.193097182	2.205344e-02	0.053183584	0.053912026
## DefPenaltyYds	0.197512296	2.200282e-02	0.064072314	0.061958232
## DefFumbles	0.016184547	2.463791e-02	-0.012823924	0.006628298
## DefInterceptions	0.061722386	-5.156149e-03	-0.074209524	-0.058441483
## DefTotalTurnovers	0.056446760	1.217256e-02	-0.063555262	-0.039181745
## Season	0.043715876	-7.685922e-02	-0.041631219	-0.081250800
## Differential		TeamScore OpponentScore DefPassingCmp DefPassingAtt		
	0.781067377	-0.769642620	-0.064738171	0.150488394

## NetYards	0.637250008	-0.631925770	-0.222740348	-0.034761786
## TurnoverMargin	0.377082264	-0.373145513	0.082821892	0.181901022
## PassingCmp	0.183342682	0.138560556	0.083460009	0.075998001
## PassingAtt	-0.035933316	0.272001348	0.083570904	0.057944813
## PassingPct	0.448358095	-0.206697090	0.022060210	0.056923769
## PassingYds	0.431051906	0.090668157	0.126357936	0.145063141
## PassingTD	0.618539629	-0.019512592	0.091773507	0.144169653
## RushingAtt	0.343571991	-0.223024981	-0.044763075	0.073608528
## RushingYds	0.571429522	-0.246049036	0.002351343	0.098937001
## RushingAvg	0.560340983	-0.200498698	0.033307036	0.093487708
## RushingTD	0.698724429	-0.188554772	0.043038468	0.130716692
## OffensivePlays	0.287897620	0.060340486	0.040919681	0.126941028
## OffensiveYards	0.790058468	-0.116896356	0.103943303	0.193714572
## OffenseAvg	0.761055147	-0.183237536	0.094216225	0.150186343
## FirstDownPass	0.348506805	0.078087301	0.091467148	0.099017964
## FirstDownRush	0.492279931	-0.219697357	-0.035195207	0.042537652
## FirstDownPen	0.069696119	0.080156897	0.055082764	0.050446788
## FirstDownTotal	0.632897568	-0.083219329	0.054621111	0.115086837
## Penalties	0.045390014	0.069089389	0.102064433	0.148129510
## PenaltyYds	0.081291125	0.055514238	0.1114111452	0.162686876
## Fumbles	-0.092211487	0.131861239	-0.039734445	-0.031092018
## Interceptions	-0.276837519	0.261769161	-0.017382571	-0.058249476
## TotalTurnovers	-0.265831719	0.280630915	-0.038933078	-0.063593263
## TeamScore	1.000000000	-0.202449142	0.110635407	0.235566171
## OpponentScore	-0.202449142	1.000000000	0.214646894	0.004853732
## DefPassingCmp	0.110635407	0.214646894	1.000000000	0.888534051
## DefPassingAtt	0.235566171	0.004853732	0.888534051	1.000000000
## DefPassingPct	-0.205209988	0.449588246	0.527936650	0.124019463
## DefPassingYds	0.056396074	0.461238083	0.811002597	0.696116964
## DefPassingTD	-0.026902967	0.630074932	0.431985918	0.288984139
## DefRushingAtt	-0.220249045	0.346909797	-0.357667664	-0.457042860
## DefRushingYds	-0.246450675	0.568876938	-0.303784505	-0.424394812
## DefRushingAvg	-0.202890701	0.555199375	-0.170640401	-0.274715155
## DefRushingTD	-0.184304923	0.692177201	-0.107643356	-0.260693452
## DefPlays	0.027517830	0.327538771	0.543938457	0.561465097
## DefYards	-0.139273423	0.796847454	0.424314893	0.241629847
## DefAvg	-0.199723885	0.766161290	0.189980837	-0.036500573
## DefFirstDownPass	0.040803814	0.382180089	0.858417270	0.750367056
## DefFirstDownRush	-0.221952586	0.492314907	-0.220058651	-0.338567902
## DefFirstDownPen	0.060041899	0.094446672	0.166792838	0.202250647
## DefFirstDownTotal	-0.112838954	0.648164723	0.507510150	0.352847040
## DefPenalties	0.066777296	0.047608838	0.129228123	0.138563162
## DefPenaltyYds	0.045061806	0.091400033	0.124944592	0.119897735
## DefFumbles	0.141183882	-0.091979574	0.015863025	0.025781372
## DefInterceptions	0.249951833	-0.265625593	0.094022164	0.243583489
## DefTotalTurnovers	0.277739467	-0.257430498	0.080273452	0.198013119

```

## Season          0.006049668  0.011888084 -0.028328659 -0.041701025
## DefPassingPct -0.420462531 -0.257339941 -0.41930782 -0.3648042744
## Differential   -0.422176950 -0.397669138 -0.38051554 -0.4156719928
## NetYards        -0.154782085  0.017015633 -0.14799556 -0.1927488957
## TurnoverMargin  0.036945250  0.127150934  0.11207117 -0.0356416142
## PassingCmp      0.077374243  0.156497317  0.17106759  0.0839936289
## PassingAtt       0.073379047 -0.021966011 -0.08316878 -0.2138003701
## PassingPct       -0.002052903  0.155256065  0.09798277 -0.0370889603
## PassingYds       -0.081413409  0.075376864  0.03505334 -0.0510535147
## PassingTD        -0.216287969 -0.050233964 -0.05901043 -0.4102166420
## RushingAtt       -0.180720848 -0.030415762 -0.07850262 -0.2765693874
## RushingYds       -0.115717537 -0.010078082 -0.07078651 -0.1176949379
## RushingAvg        -0.152791861  0.005118324 -0.04470253 -0.1962423099
## RushingTD         -0.126835921  0.108267515  0.11448597 -0.3028923478
## OffensivePlays   -0.141363954  0.101973122  0.01851589 -0.2437429725
## OffensiveYards   -0.096310447  0.048186399 -0.05340694 -0.1082911321
## OffenseAvg        -0.008201755  0.125262682  0.09058653 -0.0788586132
## FirstDownPass     -0.156338912 -0.053026691 -0.07241421 -0.3153770802
## FirstDownRush     0.036124617  0.085484505  0.08190039 -0.0437639450
## FirstDownPen      -0.098512229  0.073302813  0.03341693 -0.2999595083
## FirstDownTotal    -0.036848531  0.096286575  0.08994249 -0.0453869976
## Penalties         -0.043181735  0.096147185  0.09304272 -0.0594289149
## PenaltyYds        -0.022989345 -0.037946859  0.03917292  0.0707326851
## Fumbles           0.079458063  0.023141664  0.09933543  0.1455480406
## Interceptions    0.043886543 -0.007693062  0.09937366  0.1543421251
## TotalTurnovers   -0.205209988  0.056396074 -0.02690297 -0.2202490448
## TeamScore         0.449588246  0.461238083  0.63007493  0.3469097966
## OpponentScore    0.527936650  0.811002597  0.43198592 -0.3576676643
## DefPassingCmp    0.124019463  0.696116964  0.28898414 -0.4570428600
## DefPassingAtt     0.000000000  0.489075819  0.39650239  0.0334096923
## DefPassingPct     0.489075819  1.000000000  0.63385166 -0.2418870888
## DefPassingYds     0.396502392  0.633851663  1.00000000 -0.0473965023
## DefPassingTD       0.033409692 -0.241887089 -0.04739650  1.0000000000
## DefRushingAtt     0.102635384 -0.172013159  0.07771343  0.7391873224
## DefRushingYds     0.133036964 -0.057509274  0.15729694  0.3720422494
## DefRushingAvg      0.238808298  0.038341710 -0.01416026  0.4883122699
## DefRushingTD       0.153529419  0.461801976  0.24103860  0.4794022253
## DefPlays          0.469912225  0.674929292  0.56768532  0.3592633482
## DefYards          0.483874653  0.534876301  0.53421323  0.1519570904
## DefAvg            0.485967332  0.888603335  0.53272207 -0.2149967952
## DefFirstDownPass   0.123910001 -0.129646601  0.08022841  0.7888528350
## DefFirstDownRush   0.004410449  0.141000183  0.11566582 -0.0002455711
## DefFirstDownPen    0.442541414  0.587380037  0.47044073  0.4033934413
## DefFirstDownTotal  0.034579594  0.147205592  0.08036485 -0.0238558553
## DefPenalties       0.058957963  0.152557145  0.10315158  0.0025346468

```

## DefFumbles	-0.018162723	0.007005965	-0.05021171	0.0263769198
## DefInterceptions	-0.224175953	0.016451788	-0.10924136	-0.1892407555
## DefTotalTurnovers	-0.178623380	0.016794868	-0.11391813	-0.1236824864
## Season	0.017373437	-0.003226989	0.01101258	-0.0065777142
##	DefRushingYds	DefRushingAvg	DefRushingTD	DefPlays
## Differential	-0.523424607	-0.4863031759	-0.561529600	-0.1909143856
## NetYards	-0.547819825	-0.4871370004	-0.464507699	-0.4210217236
## TurnoverMargin	-0.194212799	-0.1469244778	-0.258439430	0.0001525553
## PassingCmp	0.018647173	0.0490502813	0.061945922	0.0418261399
## PassingAtt	0.119747148	0.1144628084	0.155716440	0.1353168586
## PassingPct	-0.179479743	-0.1150532210	-0.155117094	-0.1427705801
## PassingYds	-0.010300250	0.0086495098	0.028151375	0.1086246317
## PassingTD	-0.071383520	-0.0656727117	-0.033926688	0.0947510369
## RushingAtt	-0.279183263	-0.1216210521	-0.194448775	-0.3090197221
## RushingYds	-0.225077345	-0.1341365822	-0.192176011	-0.1596891208
## RushingAvg	-0.131739440	-0.1081341745	-0.145311374	-0.0172559779
## RushingTD	-0.196349300	-0.1511653556	-0.156635569	-0.0536000054
## OffensivePlays	-0.143922887	-0.0007670721	-0.028349305	-0.1565479001
## OffensiveYards	-0.182292260	-0.0966878468	-0.125784441	-0.0356339359
## OffenseAvg	-0.134492281	-0.1185107448	-0.138422906	0.0474361958
## FirstDownPass	-0.026930969	0.0118692816	0.018320394	0.0243322208
## FirstDownRush	-0.219647488	-0.1024397880	-0.169169555	-0.2514420854
## FirstDownPen	0.004971625	0.0402209526	0.029630574	0.0090586645
## FirstDownTotal	-0.181281160	-0.0592662502	-0.103644138	-0.1655155810
## Penalties	-0.069403770	-0.0597464352	0.001729374	0.1039300196
## PenaltyYds	-0.086487888	-0.0736808591	-0.016742606	0.1052294630
## Fumbles	0.021618628	-0.0099932661	0.076682629	0.0351287926
## Interceptions	0.125655719	0.0838830841	0.161302037	0.0779380669
## TotalTurnovers	0.107400997	0.0556854104	0.169930092	0.0808471024
## TeamScore	-0.246450675	-0.2028907014	-0.184304923	0.0275178301
## OpponentScore	0.568876938	0.5551993748	0.692177201	0.3275387714
## DefPassingCmp	-0.303784505	-0.1706404012	-0.107643356	0.5439384572
## DefPassingAtt	-0.424394812	-0.2747151550	-0.260693452	0.5614650973
## DefPassingPct	0.102635384	0.1330369637	0.238808298	0.1535294189
## DefPassingYds	-0.172013159	-0.0575092741	0.038341710	0.4618019756
## DefPassingTD	0.077713426	0.1572969436	-0.014160258	0.2410385975
## DefRushingAtt	0.739187322	0.3720422494	0.488312270	0.4794022253
## DefRushingYds	1.000000000	0.8730947621	0.688245075	0.2689666214
## DefRushingAvg	0.873094762	1.0000000000	0.600041582	0.0750924781
## DefRushingTD	0.688245075	0.6000415822	1.000000000	0.1970839405
## DefPlays	0.268966621	0.0750924781	0.197083940	1.0000000000
## DefYards	0.610787296	0.6079029094	0.546346178	0.5726537077
## DefAvg	0.579906813	0.7047277389	0.543402451	0.1054439997
## DefFirstDownPass	-0.165502747	-0.0636257511	0.025657171	0.5403468891
## DefFirstDownRush	0.868314472	0.6624669099	0.587340983	0.3998567473
## DefFirstDownPen	-0.055707347	-0.0663814432	0.015665968	0.1993277518

```

## DefFirstDownTotal  0.480393124  0.4048150894  0.437620119  0.7234459497
## DefPenalties     -0.002806677  0.0189733608  -0.021067247  0.1145227495
## DefPenaltyYds    0.039258125  0.0602355258  0.018572174  0.1206585558
## DefFumbles        -0.026783262  -0.0547759457  -0.054120831  0.0499778396
## DefInterceptions -0.208131868  -0.1613333457  -0.224866874  0.0642770948
## DefTotalTurnovers -0.172283106  -0.1556512003  -0.202500002  0.0803063695
## Season            0.015352689  0.0291418446  0.016225977  -0.0472650922
##                               DefYards          DefAvg DefFirstDownPass DefFirstDownRush
## Differential      -0.598913227  -0.618758489  -0.217094325  -0.458632617
## NetYards           -0.729979536  -0.639567669  -0.365830735  -0.513799549
## TurnoverMargin    -0.131797812  -0.161574293  0.042586183  -0.144228503
## PassingCmp         0.116168390  0.108788890  0.091233051  -0.019627052
## PassingAtt         0.215484942  0.175755717  0.109208886  0.062313958
## PassingPct         -0.152112461  -0.107046906  -0.011191012  -0.154341585
## PassingYds         0.117075612  0.066058087  0.123142671  -0.034360788
## PassingTD          0.007116432  -0.058840155  0.065395375  -0.068080474
## RushingAtt         -0.249498265  -0.119690358  -0.090092408  -0.317355086
## RushingYds         -0.193040997  -0.146487656  -0.047490894  -0.220166727
## RushingAvg         -0.106779574  -0.127188735  -0.008586427  -0.100294586
## RushingTD          -0.142960915  -0.152049693  -0.007229665  -0.173767081
## OffensivePlays    -0.020782288  0.061986249  0.023710984  -0.236976322
## OffensiveYards    -0.054582037  -0.059835995  0.062820458  -0.197943110
## OffenseAvg         -0.062010033  -0.115152176  0.053558602  -0.096233511
## FirstDownPass      0.080510479  0.073494860  0.094812163  -0.048360863
## FirstDownRush      -0.207147883  -0.107747588  -0.068467329  -0.224658105
## FirstDownPen       0.072434313  0.081364281  0.065042818  -0.016408982
## FirstDownTotal     -0.076869161  -0.006762158  0.030148892  -0.213167195
## Penalties          0.025406067  -0.026690341  0.078466818  -0.080014867
## PenaltyYds         0.012497218  -0.043324322  0.080096704  -0.095022978
## Fumbles             -0.014307373  -0.033751463  -0.049617071  0.001225369
## Interceptions      0.112722764  0.094214745  0.001262556  0.092104959
## TotalTurnovers     0.074264925  0.047783145  -0.031573932  0.069144207
## TeamScore          -0.139273423  -0.199723885  0.040803814  -0.221952586
## OpponentScore      0.796847454  0.766161290  0.382180089  0.492314907
## DefPassingCmp      0.424314893  0.189980837  0.858417270  -0.220058651
## DefPassingAtt      0.241629847  -0.036500573  0.750367056  -0.338567902
## DefPassingPct      0.469912225  0.483874653  0.485967332  0.123910001
## DefPassingYds      0.674929292  0.534876301  0.888603335  -0.129646601
## DefPassingTD        0.567685318  0.534213234  0.532722072  0.080228407
## DefRushingAtt      0.359263348  0.151957090  -0.214996795  0.788852835
## DefRushingYds      0.610787296  0.579906813  -0.165502747  0.868314472
## DefRushingAvg      0.607902909  0.704727739  -0.063625751  0.662466910
## DefRushingTD        0.546346178  0.543402451  0.025657171  0.587340983
## DefPlays            0.572653708  0.105444000  0.540346889  0.399856747
## DefYards            1.000000000  0.864462231  0.590268014  0.546201809
## DefAvg              0.864462231  1.000000000  0.390009234  0.423536375

```

## DefFirstDownPass	0.590268014	0.390009234	1.000000000	-0.141850255
## DefFirstDownRush	0.546201809	0.423536375	-0.141850255	1.000000000
## DefFirstDownPen	0.071605026	-0.025614875	0.137431352	-0.052058450
## DefFirstDownTotal	0.831958681	0.576278354	0.653641009	0.584664115
## DefPenalties	0.116217881	0.068228282	0.091887049	-0.035947615
## DefPenaltyYds	0.152027850	0.109505111	0.099033975	0.001989852
## DefFumbles	-0.014430707	-0.047710501	-0.002813518	-0.011352217
## DefInterceptions	-0.142677056	-0.206269402	0.042542473	-0.176151476
## DefTotalTurnovers	-0.115548882	-0.184499527	0.029829434	-0.138458223
## Season	0.008906115	0.041294557	-0.019898843	0.009579210
##		DefFirstDownPen	DefFirstDownTotal	DefPenalties
## Differential	-0.0210781403	-0.486870948	0.013179353	-0.028902103
## NetYards	-0.0095468986	-0.632925216	-0.064521305	-0.101133844
## TurnoverMargin	0.0146372877	-0.066460751	-0.030278199	-0.046192301
## PassingCmp	0.0516278342	0.064255648	0.097721027	0.108066606
## PassingAtt	0.0493365901	0.134090419	0.144271935	0.162654164
## PassingPct	0.0269159393	-0.110256110	-0.043012134	-0.053905077
## PassingYds	0.0785769921	0.083106670	0.092139354	0.091809187
## PassingTD	0.0726270069	0.016190934	0.087539127	0.087458617
## RushingAtt	-0.0473819882	-0.301942200	-0.044234601	-0.060368138
## RushingYds	-0.0064690578	-0.194065174	-0.066155832	-0.088133097
## RushingAvg	0.0267430155	-0.073946954	-0.057171981	-0.076423381
## RushingTD	0.0116282754	-0.126211617	0.001478081	-0.023191416
## OffensivePlays	0.0044182172	-0.151100644	0.101764432	0.104842750
## OffensiveYards	0.0585079679	-0.082828263	0.023336190	0.006082620
## OffenseAvg	0.0648182595	-0.015662267	-0.029252587	-0.051699931
## FirstDownPass	0.0594882546	0.043655644	0.074673629	0.076942358
## FirstDownRush	-0.0271049617	-0.224352618	-0.077513273	-0.097695640
## FirstDownPen	0.0709464960	0.046460027	0.542175075	0.658548102
## FirstDownTotal	0.0348967103	-0.123036459	0.124255326	0.137953942
## Penalties	0.5448156356	0.124940137	0.190344623	0.193217609
## PenaltyYds	0.6563983927	0.140738712	0.193097182	0.197512296
## Fumbles	-0.0190830400	-0.038658110	0.022053436	0.022002823
## Interceptions	0.0056245655	0.066245867	0.053183584	0.064072314
## TotalTurnovers	-0.0083304443	0.023823953	0.053912026	0.061958232
## TeamScore	0.0600418989	-0.112838954	0.066777296	0.045061806
## OpponentScore	0.0944466717	0.648164723	0.047608838	0.091400033
## DefPassingCmp	0.1667928379	0.507510150	0.129228123	0.124944592
## DefPassingAtt	0.2022506473	0.352847040	0.138563162	0.119897735
## DefPassingPct	0.0044104488	0.442541414	0.034579594	0.058957963
## DefPassingYds	0.1410001827	0.587380037	0.147205592	0.152557145
## DefPassingTD	0.1156658207	0.470440732	0.080364847	0.103151581
## DefRushingAtt	-0.0002455711	0.403393441	-0.023855855	0.002534647
## DefRushingYds	-0.0557073473	0.480393124	-0.002806677	0.039258125
## DefRushingAvg	-0.0663814432	0.404815089	0.018973361	0.060235526
## DefRushingTD	0.0156659684	0.437620119	-0.021067247	0.018572174

## DefPlays	0.1993277518	0.723445950	0.114522749	0.120658556
## DefYards	0.0716050265	0.831958681	0.116217881	0.152027850
## DefAvg	-0.0256148747	0.576278354	0.068228282	0.109505111
## DefFirstDownPass	0.1374313523	0.653641009	0.091887049	0.099033975
## DefFirstDownRush	-0.0520584497	0.584664115	-0.035947615	0.001989852
## DefFirstDownPen	1.00000000000	0.282939522	0.131267111	0.134211435
## DefFirstDownTotal	0.2829395221	1.000000000	0.070161291	0.103058264
## DefPenalties	0.1312671113	0.070161291	1.000000000	0.900744066
## DefPenaltyYds	0.1342114354	0.103058264	0.900744066	1.000000000
## DefFumbles	-0.0026003776	-0.010026151	-0.008397144	-0.017131676
## DefInterceptions	0.0194055745	-0.087683835	0.020716236	0.008367985
## DefTotalTurnovers	0.0127507987	-0.071764094	0.009959436	-0.004905025
## Season	0.0678208894	0.005517816	0.027121670	0.041720288
##		DefFumbles	DefInterceptions	DefTotalTurnovers
## Differential	0.1506992798	0.3323388335	0.345230083	-0.003636210
## NetYards	0.0016894054	0.1665064148	0.125002132	-0.001569119
## TurnoverMargin	0.4611567538	0.5702003054	0.723964894	-0.002887603
## PassingCmp	-0.0409238359	-0.0287094744	-0.047955818	-0.030911733
## PassingAtt	-0.0352683058	-0.0673938145	-0.073067543	-0.042815333
## PassingPct	-0.0177987463	0.0688612134	0.039676760	0.011534916
## PassingYds	-0.0375733487	0.0086121204	-0.018006148	-0.006117554
## PassingTD	0.0434153059	0.0883846599	0.093981546	0.007749168
## RushingAtt	0.0710373911	0.1432153386	0.152731898	-0.006470762
## RushingYds	0.0235907609	0.1189308042	0.103830416	0.015090607
## RushingAvg	-0.0069041024	0.0763947673	0.052362318	0.029310251
## RushingTD	0.0819343368	0.1535987879	0.167539382	0.010761615
## OffensivePlays	0.0318497808	0.0678983530	0.071221668	-0.048660609
## OffensiveYards	-0.0121339680	0.0988848301	0.065699798	0.006719496
## OffenseAvg	-0.0285341525	0.0801805815	0.041124510	0.038470215
## FirstDownPass	-0.0485113283	-0.0127221493	-0.040989258	-0.021594636
## FirstDownRush	0.0007492654	0.0855201980	0.064126014	0.008718165
## FirstDownPen	-0.0219681886	-0.0006859365	-0.014785143	0.066242195
## FirstDownTotal	-0.0394374483	0.0518375280	0.012948409	0.003573008
## Penalties	0.0175079007	0.0547373404	0.052108799	0.029358706
## PenaltyYds	0.0161845471	0.0617223857	0.056446760	0.043715876
## Fumbles	0.0246379075	-0.0051561488	0.012172560	-0.076859221
## Interceptions	-0.0128239235	-0.0742095243	-0.063555262	-0.041631219
## TotalTurnovers	0.0066282979	-0.0584414828	-0.039181745	-0.081250800
## TeamScore	0.1411838824	0.2499518329	0.277739467	0.006049668
## OpponentScore	-0.0919795740	-0.2656255930	-0.257430498	0.011888084
## DefPassingCmp	0.0158630249	0.0940221643	0.080273452	-0.028328659
## DefPassingAtt	0.0257813722	0.2435834885	0.198013119	-0.041701025
## DefPassingPct	-0.0181627225	-0.2241759530	-0.178623380	0.017373437
## DefPassingYds	0.0070059647	0.0164517885	0.016794868	-0.003226989
## DefPassingTD	-0.0502117148	-0.1092413599	-0.113918126	0.011012582
## DefRushingAtt	0.0263769198	-0.1892407555	-0.123682486	-0.006577714

```

## DefRushingYds      -0.0267832623   -0.2081318682   -0.172283106   0.015352689
## DefRushingAvg     -0.0547759457   -0.1613333457   -0.155651200   0.029141845
## DefRushingTD       -0.0541208313   -0.2248668737   -0.202500002   0.016225977
## DefPlays           0.0499778396    0.0642770948    0.080306370   -0.047265092
## DefYards            -0.0144307070   -0.1426770564   -0.115548882   0.008906115
## DefAvg              -0.0477105007   -0.2062694023   -0.184499527   0.041294557
## DefFirstDownPass   -0.0028135177   0.0425424733    0.029829434   -0.019898843
## DefFirstDownRush   -0.0113522169   -0.1761514759   -0.138458223   0.009579210
## DefFirstDownPen    -0.0026003776   0.0194055745    0.012750799   0.067820889
## DefFirstDownTotal  -0.0100261507   -0.0876838353   -0.071764094   0.005517816
## DefPenalties        -0.0083971436   0.0207162360    0.009959436   0.027121670
## DefPenaltyYds      -0.0171316758   0.0083679848   -0.004905025   0.041720288
## DefFumbles          1.0000000000    0.0248441007   0.668277532   -0.077380805
## DefInterceptions   0.0248441007   1.0000000000    0.760285188   -0.046159898
## DefTotalTurnovers  0.6682775322   0.7602851879    1.0000000000  -0.084630780
## Season              -0.0773808053   -0.0461598980   -0.084630780   1.0000000000

```

Notice right away – NetYards is highly correlated. But NetYards's also highly correlated with RushingYards, OffensiveYards and DefYards. And that makes sense: those things all feed into NetYards. Including all of these measures would be pointless – they would add error without adding much in the way of predictive power.

**Your turn:** What else do you see? What other values have predictive power and aren't co-correlated?

We can add more just by simply adding them. Let's add the average yard per play for both offense and defense. They're correlated to NetYards, but not as much as you might expect.

```

model2 <- lm(Differential ~ NetYards + TurnoverMargin + DefAvg + OffenseAvg, data=logs)
summary(model2)

##
## Call:
## lm(formula = Differential ~ NetYards + TurnoverMargin + DefAvg +
##     OffenseAvg, data = logs)
##
## Residuals:
##      Min      1Q      Median      3Q      Max 
## -38.254  -6.255  -0.002   6.229  37.507 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 0.4960303  0.4292118  1.156   0.248    
## NetYards     0.0547465  0.0008007  68.376  <2e-16 ***  
## TurnoverMargin 3.8806793  0.0410484  94.539  <2e-16 ***  
## DefAvg      -3.9374905  0.0738754 -53.299  <2e-16 *** 
## 
```

```

## OffenseAvg      3.9152803  0.0729586  53.664    <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 9.45 on 15601 degrees of freedom
##   (1 observation deleted due to missingness)
## Multiple R-squared:  0.8287, Adjusted R-squared:  0.8287
## F-statistic: 1.887e+04 on 4 and 15601 DF,  p-value: < 2.2e-16

```

Go down the list:

What is the Adjusted R-squared now? What is the p-value and is it less than .05? What is the Residual standard error?

The final thing we can do with this is predict things. Look at our coefficients table. See the Estimates? We can build a formula from that, same as we did with linear regressions.

How does this apply in the real world? Let's pretend for a minute that you are Scott Frost, and you have a mess on your hands. Your job is to win conference titles. To do that, we need to know what attributes of a team should we emphasize. We can do that by looking at what previous Big Ten conference champions looked like.

So if our goal is to predict a conference champion team, we need to know what those teams did. Here's the regular season conference champions in this dataset.

```

logs %>%
  filter(Team == "Ohio State" & Season == 2020 | Team == "Ohio State" & Season == 2019
  summarise(
    meanNetYards = mean(NetYards),
    meanTurnoverMargin = mean(TurnoverMargin),
    meanDefAvg = mean(DefAvg),
    meanOffenseAvg = mean(OffenseAvg)
  )

## # A tibble: 1 x 4
##   meanNetYards meanTurnoverMargin meanDefAvg meanOffenseAvg
##       <dbl>           <dbl>        <dbl>         <dbl>
## 1      196.          0.676        5.04        6.91

```

Now it's just plug and chug.

```
(0.0547465*195.8824) + (3.8806793*0.6764706) + (-3.9374905*5.044118) + (3.9152803*6.90
```

```
## [1] 21.03389
```

So a team with those numbers is going to average scoring 21 more points per game than their opponent. Sound like Ohio State in the last three years?

How does that compare to Nebraska this season?

```

logs %>%
  filter(
    Team == "Nebraska" & Season == 2020
  ) %>%
  summarise(
    meanNetYards = mean(NetYards),
    meanTurnoverMargin = mean(TurnoverMargin),
    meanDefAvg = mean(DefAvg),
    meanOffenseAvg = mean(OffenseAvg)
  )

## # A tibble: 1 x 4
##   meanNetYards meanTurnoverMargin meanDefAvg meanOffenseAvg
##       <dbl>           <dbl>        <dbl>          <dbl>
## 1         5            -1.38        5.44         5.54
## (0.0547465*5) + (3.8806793*-1.375) + (-3.9374905*5.4375) + (3.9152803*5.5375) + 0.4960303
## [1] -4.295411

```

By this model, it predicted we would average being outscored by our opponents by 4.3 points over the season. Reality? We were outscored by 6.25 on average.



# Chapter 11

## Residuals

When looking at a linear model of your data, there's a measure you need to be aware of called residuals. The residual is the distance between what the model predicted and what the real outcome is. Take our model at the end of the correlation and regression chapter. Our model predicted Nebraska, given a 5 net yardage margin would beat Iowa by 1.96 points. They lost by 6. So our residual is -7.96.

Residuals can tell you several things, but most important is if a linear model is the right model for your data. If the residuals appear to be random, then a linear model is appropriate. If they have a pattern, it means something else is going on in your data and a linear model isn't appropriate.

Residuals can also tell you who is underperforming and overperforming the model. And the more robust the model – the better your r-squared value is – the more meaningful that label of under or overperforming is.

Let's go back to our net yards model.

Then load the tidyverse.

```
library(tidyverse)

logs <- read_csv("data/footballlogs20.csv")

## 
## -- Column specification --
## cols(
##   .default = col_double(),
##   Date = col_date(format = ""),
##   HomeAway = col_character(),
##   Opponent = col_character(),
##   Result = col_character(),
```

```

##   TeamFull = col_character(),
##   TeamURL = col_character(),
##   Outcome = col_character(),
##   Team = col_character(),
##   Conference = col_character()
## )
## i Use `spec()` for the full column specifications.

```

First, let's make the columns we'll need.

```
residualmodel <- logs %>% mutate(differential = TeamScore - OpponentScore, NetYards =
```

Now let's create our model.

```

fit <- lm(differential ~ NetYards, data = residualmodel)
summary(fit)

##
## Call:
## lm(formula = differential ~ NetYards, data = residualmodel)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -49.379  -8.520   0.059   8.487  48.748
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 0.315632  0.388994  0.811   0.417
## NetYards    0.102473  0.002307 44.427  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.72 on 1068 degrees of freedom
## Multiple R-squared:  0.6489, Adjusted R-squared:  0.6486
## F-statistic: 1974 on 1 and 1068 DF,  p-value: < 2.2e-16

```

We've seen this output before, but let's review because if you are using scatter-plots to make a point, you should do this. First, note the Min and Max residual at the top. A team has underperformed the model by 39 points (!), and a team has overperformed it by 39 points (!!). The median residual, where half are above and half are below, is just slightly above the fit line. Close here is good.

Next: Look at the Adjusted R-squared value. What that says is that 72 percent of a team's scoring output can be predicted by their net yards.

Last: Look at the p-value. We are looking for a p-value smaller than .05. At .05, we can say that our correlation didn't happen at random. And, in this case, it REALLY didn't happen at random. But if you know a little bit about football, it doesn't surprise you that the more you outgain your opponent, the more you

win by. It's an intuitive result.

What we want to do now is look at those residuals. We want to add them to our individual game records. We can do that by creating two new fields – predicted and residuals – to our dataframe like this:

```
residualmodel$predicted <- predict(fit)
residualmodel$residuals <- residuals(fit)
```

Now we can sort our data by those residuals. Sorting in descending order gives us the games where teams overperformed the model. To make it easier to read, I'm going to use select to give us just the columns we need to see.

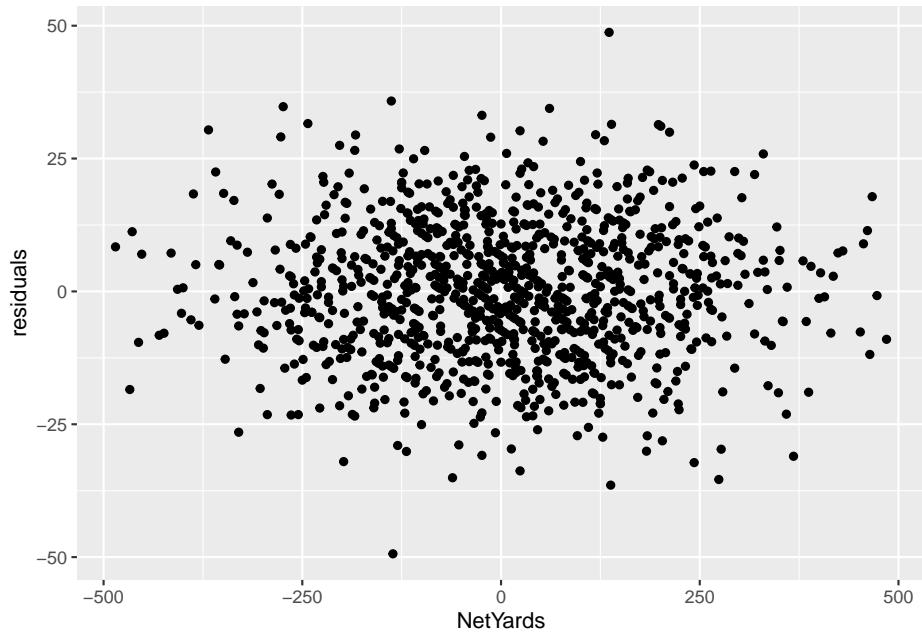
```
residualmodel %>% arrange(desc(residuals)) %>% select(Team, Opponent, Result, NetYards, residuals)
```

```
## # A tibble: 1,070 x 5
##   Team        Opponent     Result  NetYards residuals
##   <chr>       <chr>      <chr>    <dbl>     <dbl>
## 1 Arizona State Arizona W (70-7) 136 48.7
## 2 Kentucky      Mississippi State W (24-2) -138 35.8
## 3 Mississippi State Vanderbilt W (24-17) -274 34.8
## 4 Kansas State  Kansas W (55-14)  61 34.4
## 5 Boise State   Colorado State W (52-21) -24 33.1
## 6 Texas         Oklahoma State W (41-34) -243 31.6
## 7 Army          Mercer W (49-3)  139 31.4
## 8 Notre Dame    South Florida W (52-0)  198 31.4
## 9 BYU           North Alabama W (66-14) 201 31.1
## 10 Rutgers      Nebraska L (21-28) -368 30.4
## # ... with 1,060 more rows
```

So looking at this table, what you see here are the teams who scored more than their net yards would indicate. One of them should jump off the page at you.

Remember Nebraska vs Rutgers? We won and everyone was happy and rel\*ieved the season was over? We outgained Rutgers by **368 yards** in that game and won by 7. Our model predicted Nebraska should have won that game by **37** points. We should have blown Rutgers out of their own barn. But Rutgers isn't as hard done as Arizona, which should have lost by 48, but ended up losing by 63.

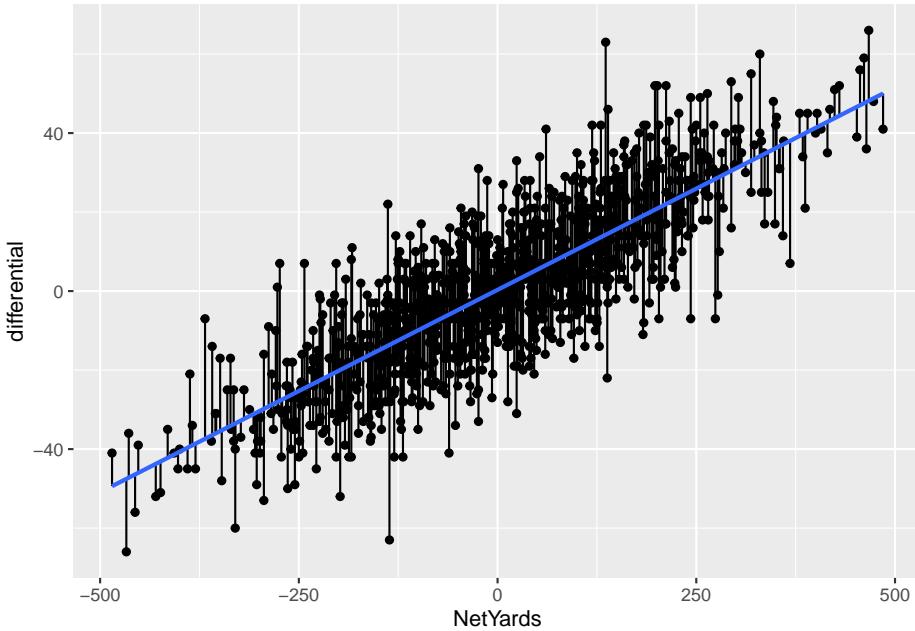
But, before we can bestow any validity on this model, we need to see if this linear model is appropriate. We've done that some looking at our p-values and R-squared values. But one more check is to look at the residuals themselves. We do that by plotting the residuals with the predictor. We'll get into plotting soon, but for now just seeing it is enough.



The lack of a shape here – the seemingly random nature – is a good sign that a linear model works for our data. If there was a pattern, that would indicate something else was going on in our data and we needed a different model.

Another way to view your residuals is by connecting the predicted value with the actual value.

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



The blue line here separates underperformers from overperformers.

## 11.1 Penalties

Now let's look at it where it doesn't work: Penalties.

```
penalties <- logs %>%
  mutate(
    differential = TeamScore - OpponentScore,
    TotalPenalties = Penalties+DefPenalties,
    TotalPenaltyYards = PenaltyYds+DefPenaltyYds
  )

pfit <- lm(differential ~ TotalPenaltyYards, data = penalties)
summary(pfit)

##
## Call:
## lm(formula = differential ~ TotalPenaltyYards, data = penalties)
## 
## Residuals:
##     Min      1Q  Median      3Q     Max 
## -67.02 -14.68   0.24  14.04  64.98 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept)  14.04     1.04   13.46  <2e-16 ***
## TotalPenaltyYards  0.24     0.01    24.00  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```

## (Intercept)      1.226995  1.817342  0.675   0.500
## TotalPenaltyYards -0.003383  0.015816 -0.214   0.831
##
## Residual standard error: 21.46 on 1068 degrees of freedom
## Multiple R-squared:  4.284e-05, Adjusted R-squared:  -0.0008935
## F-statistic: 0.04575 on 1 and 1068 DF, p-value: 0.8307

```

So from top to bottom:

- Our min and max go from -67 to positive 65
- Our adjusted R-squared is ... -0.0008935. Not much at all.
- Our p-value is ... 0.8307, which is more than than .05.

So what we can say about this model is that it's statistically insignificant and utterly meaningless. Normally, we'd stop right here – why bother going forward with a predictive model that isn't predictive? But let's do it anyway.

```

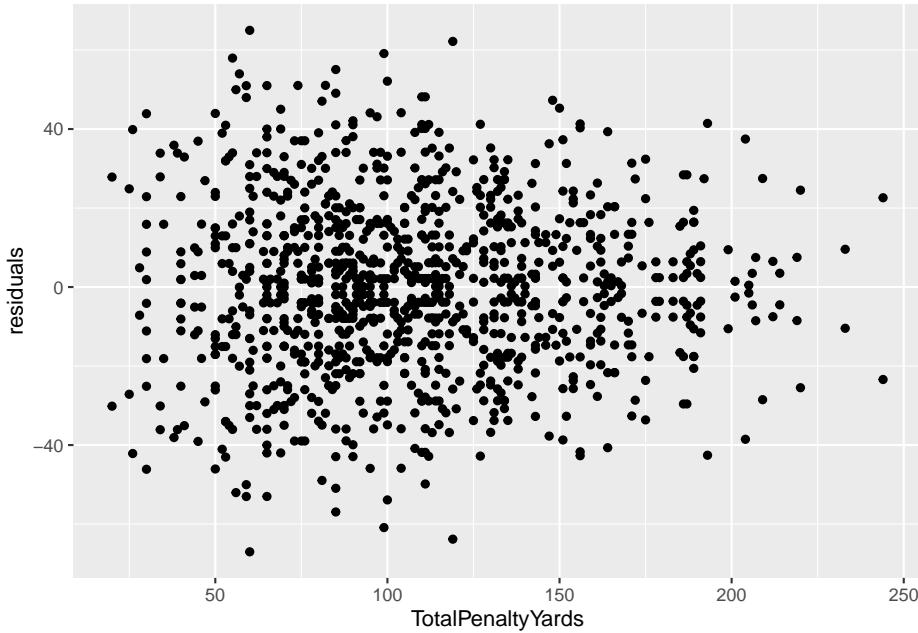
penalties$predicted <- predict(pfit)
penalties$residuals <- residuals(pfit)

penalties %>% arrange(desc(residuals)) %>% select(Team, Opponent, Result, TotalPenaltyYards, residuals)

## # A tibble: 1,070 x 5
##   Team        Opponent     Result TotalPenaltyYards residuals
##   <chr>       <chr>      <chr>          <dbl>      <dbl>
## 1 Clemson    Georgia Tech W (73-7)        60        65.0
## 2 Arizona State Arizona    W (70-7)        119       62.2
## 3 Alabama     Kentucky    W (63-3)        99        59.1
## 4 Marshall    Eastern Kentucky W (59-0)        55        58.0
## 5 Texas       Texas-El Paso W (59-3)        85        55.1
## 6 Pitt         Austin Peay  W (55-0)        57        54.0
## 7 Oklahoma    Kansas     W (62-9)        100       52.1
## 8 Wake Forest Campbell   W (66-14)       82        51.1
## 9 BYU          North Alabama W (66-14)       74        51.0
## 10 Notre Dame South Florida W (52-0)        65        51.0
## # ... with 1,060 more rows

```

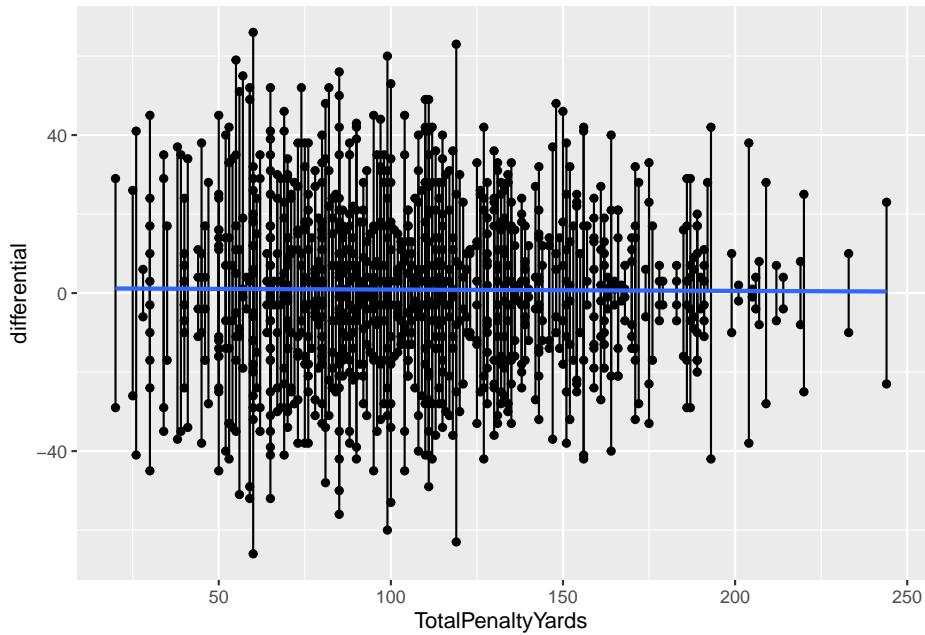
First, note all of the biggest misses here are all blowout games. The worst games of the season, the worst being Clemson vs Georgia Tech. The model missed that differential by ... 65 points. The margin of victory? 66 points. In other words, this model is terrible. But let's look at it anyway.



Well ... it actually says that a linear model is appropriate. Which an important lesson – just because your residual plot says a linear model works here, that doesn't say your linear model is good. There are other measures for that, and you need to use them.

Here's the segment plot of residuals – you'll see some really long lines. That's a bad sign. Another bad sign? A flat fit line. It means there's no relationship between these two things. Which we already know.

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



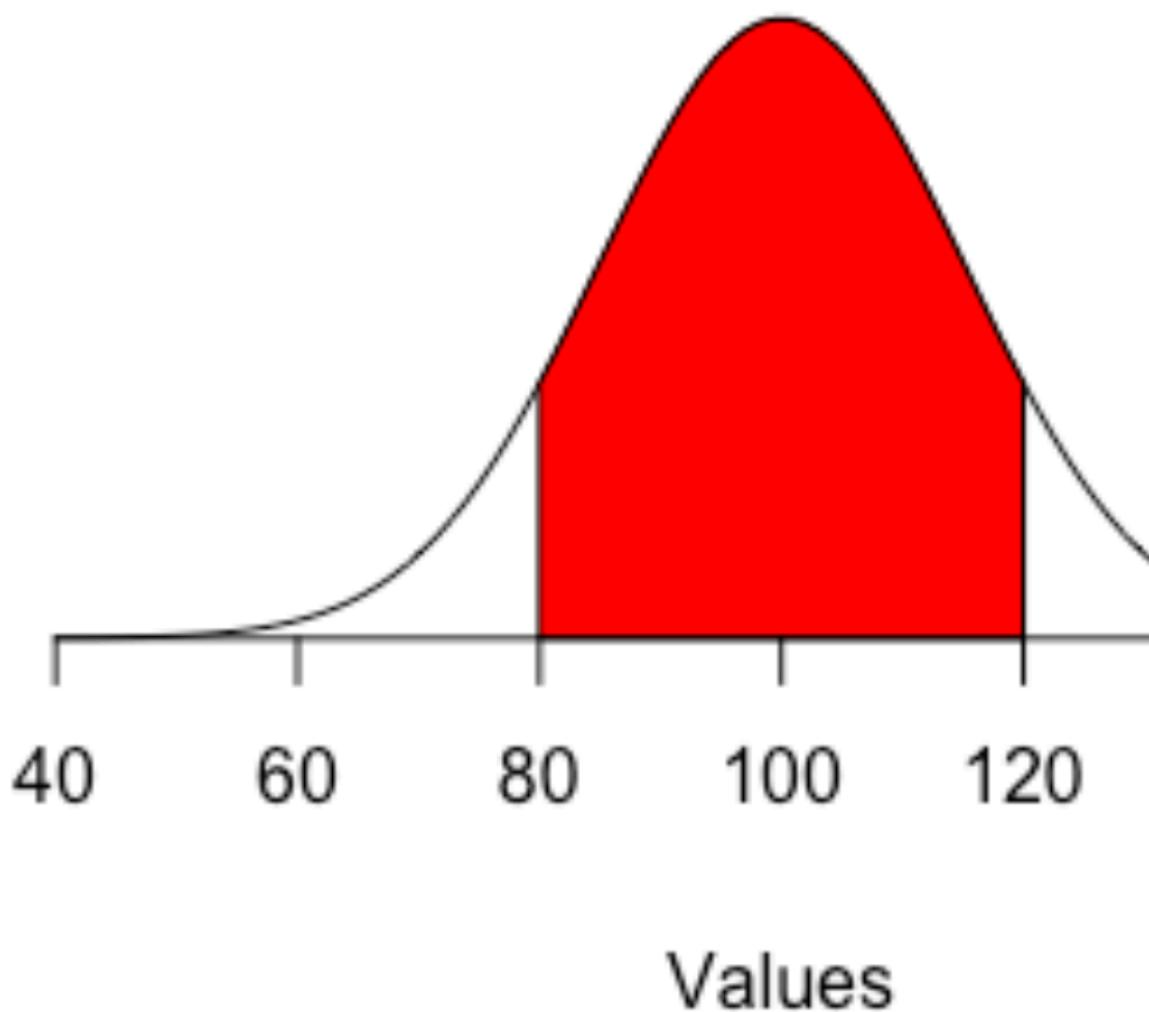
# **Chapter 12**

## **Z-scores**

Z-scores are a handy way to standardize numbers so you can compare things across groupings or time. In this class, we may want to compare teams by year, or era. We can use z-scores to answer questions like who was the greatest X of all time, because a z-score can put them in context to their era.

A z-score is a measure of how a particular stat is from the mean. It's measured in standard deviations from that mean. A standard deviation is a measure of how much variation – how spread out – numbers are in a data set. What it means here, with regards to z-scores, is that zero is perfectly average. If it's 1, it's one standard deviation above the mean, and 34 percent of all cases are between 0 and 1.

## Normal Distribution



If you think of the normal distribution, it means that 84.3 percent of all cases are below that 1. If it were -1, it would mean the number is one standard deviation below the mean, and 84.3 percent of cases would be above that -1. So if you have numbers with z-scores of 3 or even 4, that means that number is waaaaay above the mean.

So let's use last year's Nebraska basketball team, which if haven't been paying attention to current events, was not good at basketball.

## 12.1 Calculating a Z score in R

For this we'll need the logs of all college basketball games last season.

Load the tidyverse.

```
library(tidyverse)
```

And load the data.

```
gamelogs <- read_csv("data/logs20.csv")
```

```
## 
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Date = col_date(format = ""),
##   HomeAway = col_character(),
##   Opponent = col_character(),
##   W_L = col_character(),
##   Blank = col_logical(),
##   Team = col_character(),
##   Conference = col_character(),
##   season = col_character()
## )
## i Use `spec()` for the full column specifications.
```

The first thing we need to do is select some fields we think represent team quality and a few things to help us keep things straight. So I'm going to pick shooting percentage, rebounding and the opponent version of the same two:

```
teamquality <- gamelogs %>%
  select(Conference, Team, TeamFGPCT, TeamTotalRebounds, OpponentFGPCT, OpponentTotalRebounds)
```

And since we have individual game data, we need to collapse this into one record for each team. We do that with ... group by.

```
teamtotals <- teamquality %>%
  group_by(Conference, Team) %>%
  summarise(
    FGAvg = mean(TeamFGPCT),
    ReboundAvg = mean(TeamTotalRebounds),
    OppFGAvg = mean(OpponentFGPCT),
    OffRebAvg = mean(OpponentTotalRebounds)
  )
```

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

To calculate a z-score in R, the easiest way is to use the `scale` function in base R. To use it, you use `scale(Fieldname, center=TRUE, scale=TRUE)`. The center and scale indicate if you want to subtract from the mean and if you want to divide by the standard deviation, respectively. We do.

When we have multiple z-scores, it's pretty standard practice to add them together into a composite score. That's what we're doing at the end here with `TotalZscore`. Note: We have to invert OppZscore and OppRebZScore by multiplying it by a negative 1 because the lower someone's opponent number is, the better.

```
teamzscore <- teamtotals %>%
  mutate(
    FGzscore = as.numeric(scale(FGAvg, center = TRUE, scale = TRUE)),
    RebZscore = as.numeric(scale(ReboundAvg, center = TRUE, scale = TRUE)),
    OppZscore = as.numeric(scale(OppFGAvg, center = TRUE, scale = TRUE)) * -1,
    OppRebZScore = as.numeric(scale(OffRebAvg, center = TRUE, scale = TRUE)) * -1,
    TotalZscore = FGzscore + RebZscore + OppZscore + OppRebZScore
  )
```

So now we have a dataframe called `teamzscore` that has 353 basketball teams with Z scores. What does it look like?

```
head(teamzscore)

## # A tibble: 6 x 11
## # Groups:   Conference [1]
##   Conference Team FGAvg ReboundAvg OppFGAvg OffRebAvg FGzscore RebZscore
##   <chr>      <chr>   <dbl>     <dbl>     <dbl>     <dbl>     <dbl>
## 1 A-10       Davi~  0.454     31.1     0.437    30.4     0.505    -0.619
## 2 A-10       Dayt~  0.525     32.5     0.413    29.0     2.59     0.0352
## 3 A-10       Duqu~  0.444     32.4     0.427    32.4     0.216    -0.0168
## 4 A-10       Ford~  0.384     30.0     0.402    33.9    -1.53     -1.13
## 5 A-10       Geor~  0.424     33.8     0.440    30.5    -0.358     0.620
## 6 A-10       Geor~  0.422     30.5     0.452    32.7    -0.410    -0.904
## # ... with 3 more variables: OppZscore <dbl>, OppRebZScore <dbl>,
## #   TotalZscore <dbl>
```

A way to read this – a team at zero is precisely average. The larger the positive number, the more exceptional they are. The larger the negative number, the more truly terrible they are.

So who are the best teams in the country?

```
teamzscore %>% arrange(desc(TotalZscore))
```

```
## # A tibble: 353 x 11
## # Groups:   Conference [32]
```

```

##   Conference Team   FGAvg ReboundAvg OppFGAvg OffRebAvg FGzscore RebZscore
##   <chr>     <chr> <dbl>      <dbl>    <dbl>      <dbl>    <dbl>      <dbl>
## 1 Big West  UC-I~  0.473      36.6    0.390      27.1    1.60      2.23
## 2 Big 12   Kans~  0.482      35.9    0.378      29.0    2.36      1.13
## 3 WCC      Gonz~  0.517      37.4    0.424      28.2    1.73      1.90
## 4 Southland Step~  0.490      34.2    0.427      26.6    1.76      1.05
## 5 Big Ten   Mich~  0.460      37.7    0.382      29.6    1.38      1.55
## 6 OVC      Murr~  0.477      35.3    0.401      29.2    1.31      1.36
## 7 Summit    Sout~  0.492      35.5    0.423      31.3    1.58      1.52
## 8 A-10     Dayt~  0.525      32.5    0.413      29.0    2.59      0.0352
## 9 A-10     Sain~  0.457      37.4    0.403      30.5    0.598      2.21
## 10 ACC     Loui~  0.457      36.6    0.392      29.8    1.11      1.37
## # ... with 343 more rows, and 3 more variables: OppZscore <dbl>,
## #   OppRebZScore <dbl>, TotalZscore <dbl>

```

Don't sleep on the Anteaters! Would have been a tough out at the tournament that never happened.

But closer to home, how is Nebraska doing.

```

teamzscore %>%
  filter(Conference == "Big Ten") %>%
  arrange(desc(TotalZscore))

## # A tibble: 14 x 11
## # Groups:   Conference [1]
##   Conference Team   FGAvg ReboundAvg OppFGAvg OffRebAvg FGzscore RebZscore
##   <chr>     <chr> <dbl>      <dbl>    <dbl>      <dbl>    <dbl>
## 1 Big Ten   Mich~  0.460      37.7    0.382      29.6    1.38      1.55
## 2 Big Ten   Rutg~  0.449      37     0.385      31.1    0.727      1.22
## 3 Big Ten   Ohio~  0.447      33.6    0.400      28.4    0.592     -0.393
## 4 Big Ten   Illi~  0.444      36.1    0.418      29.1    0.439      0.779
## 5 Big Ten   Indi~  0.445      35.1    0.419      29.4    0.480      0.306
## 6 Big Ten   Mary~  0.419      36.1    0.401      31.9    -0.952      0.794
## 7 Big Ten   Mich~  0.463      33.0    0.428      31.9    1.56     -0.682
## 8 Big Ten   Penn~  0.432      35.6    0.411      34.2    -0.237      0.550
## 9 Big Ten   Minn~  0.426      35.5    0.411      33     -0.560      0.520
## 10 Big Ten  Iowa~  0.452      34.2    0.430      32.4    0.918     -0.104
## 11 Big Ten  Purd~  0.418      33.8    0.410      29.3    -1.02     -0.271
## 12 Big Ten  Wisc~  0.426      31.3    0.410      32.0    -0.587     -1.49
## 13 Big Ten  Nort~  0.417      30.5    0.422      34.8    -1.12     -1.84
## 14 Big Ten  Nebr~  0.408      32.4    0.453      42.2    -1.62     -0.947
## # ... with 3 more variables: OppZscore <dbl>, OppRebZScore <dbl>,
## #   TotalZscore <dbl>

```

So, as we can see, with our composite Z Score, Nebraska is ... not good. Not good at all.

## 12.2 Writing about z-scores

The great thing about z-scores is that they make it very easy for you, the sports analyst, to create your own measures of who is better than who. The downside: Only a small handful of sports fans know what the hell a z-score is.

As such, you should try as hard as you can to avoid writing about them.

If the word z-score appears in your story or in a chart, you need to explain what it is. “The ranking uses a statistical measure of the distance from the mean called a z-score” is a good way to go about it. You don’t need a full stats textbook definition, just a quick explanation. And keep it simple.

**Never use z-score in a headline.** Write around it. Away from it. Z-score in a headline is attention repellent. You won’t get anyone to look at it. So “Tottenham tops in z-score” bad, “Tottenham tops in the Premiere League” good.

# Chapter 13

## Clustering

One common effort in sports is to classify teams and players – who are this players peers? What teams are like this one? Who should we compare a player to? Truth is, most sports commentators use nothing more sophisticated than looking at a couple of stats or use the “eye test” to say a player is like this or that.

There's better ways.

In this chapter, we're going to use a method that sounds advanced but it's really quite simple called k-means clustering. It's based on the concept of the k-nearest neighbor algorithm. You're probably already scared. Don't be.

Imagine two dots on a scatterplot. If you took a ruler out and measured the distance between those dots, you'd know how far apart they are. In math, that's called the Euclidean distance. It's just the space between them in numbers. Where k-nearest neighbor comes in, you have lots of dots and you want to measure the distance between all of them. What does k-means clustering do? It lumps them into groups based on the average distance between them. Players who are good on offense but bad on defense are over here, good offense good defense are over there. And using the Euclidean distance between them, we can decide who is in and who is out of those groups.

For this exercise, I want to look at Cam Mack, Nebraska's point guard and probably the most interesting player on Fred Hoiberg's first team. This was Mack's first – only? – year in major college basketball. I believe Mack could have been one of the best players Nebraska ever had, but it didn't work out. So who does Cam Mack compare to?

To answer this, we'll use k-means clustering.

First thing we do is load some libraries and set a seed, so if we run this repeatedly, our random numbers are generated from the same base. If you don't have the

```
cluster library, just add it on the console with install.packages("cluster")
library(tidyverse)
library(cluster)

set.seed(1234)
```

I've gone and scraped stats for every player in that season.

Now load that data.

```
players <- read_csv("data/players20.csv")
```

```
##
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Team = col_character(),
##   Player = col_character(),
##   Class = col_character(),
##   Pos = col_character(),
##   Height = col_character(),
##   Hometown = col_character(),
##   `High School` = col_character(),
##   Summary = col_character()
## )
## i Use `spec()` for the full column specifications.
```

To cluster this data properly, we have some work to do.

First, it won't do to have players who haven't played, so we can use filter to find anyone with greater than 0 minutes played. Next, Cam Mack is a guard, so let's just look at guards. Third, we want to limit the data to things that make sense to look at for Cam Mack – things like shooting, three point shooting, assists, turnovers and points.

```
playersselected <- players %>%
  filter(MP>0) %>% filter(Pos == "G") %>%
  select(Player, Team, Pos, MP, `FG%`, `3P%`, AST, TOV, PTS) %>%
  na.omit()
```

Now, k-means clustering doesn't work as well with data that can be on different scales. So comparing a percentage to a count metric – shooting percentage to points – would create chaos because shooting percentages are a fraction of 1 and points, depending on when they are in the season, could be quite large. So we have to scale each metric – put them on a similar basis using the distance from the max value as our guide. Also, k-means clustering won't work with text data, so we need to create a dataframe that's just the numbers, but scaled. We can do that with another select, and using mutate\_all with the scale function. The

`na.omit()` means get rid of any blanks, because they too will cause errors.

```
playersscaled <- playersselected %>%
  select(MP, `FG%`, `3P%`, AST, TOV, PTS) %>%
  mutate_all(scale) %>%
  na.omit()
```

With k-means clustering, we decide how many clusters we want. Most often, researchers will try a handful of different cluster numbers and see what works. But there are methods for finding the optimal number. One method is called the Elbow method. One implementation of this, borrowed from the University of Cincinnati's Business Analytics program, does this quite nicely with a graph that will help you decide for yourself.

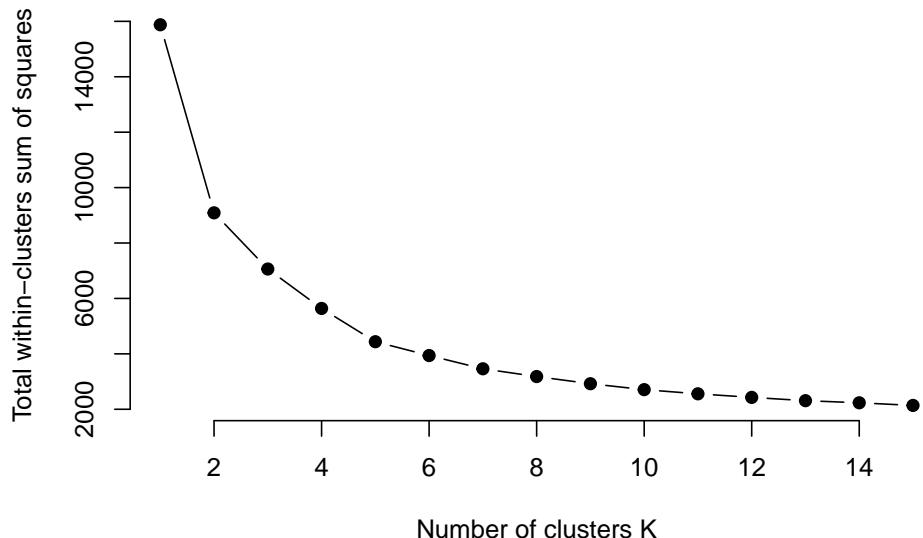
All you need to do in this code is change out the data frame – `playersscaled` in this case – and run it.

```
# function to compute total within-cluster sum of square
wss <- function(k) {
  kmeans(playersscaled, k, nstart = 10 )$tot.withinss
}

# Compute and plot wss for k = 1 to k = 15
k.values <- 1:15

# extract wss for 2-15 clusters
wss_values <- map_dbl(k.values, wss)

## Warning: did not converge in 10 iterations
plot(k.values, wss_values,
      type="b", pch = 19, frame = FALSE,
      xlab="Number of clusters K",
      ylab="Total within-clusters sum of squares")
```



The Elbow method – so named because you’re looking for the “elbow” where the line flattens out. In this case, it looks like a K of 5 is ideal. So let’s try that. We’re going to use the kmeans function, saving it to an object called k5. We just need to tell it our dataframe name, how many centers (k) we want, and we’ll use a sensible default for how many different configurations to try.

```
k5 <- kmeans(playersscaled, centers = 5, nstart = 25)
```

Let’s look at what we get.

```
k5
```

```
## K-means clustering with 5 clusters of sizes 242, 58, 495, 989, 863
##
## Cluster means:
##           MP          FG%         3P%          AST          TOV          PTS
## 1 -1.3602270 -1.98012908 -1.65325745 -0.9392522 -1.1172744 -1.1417360
## 2 -1.3785814  3.22934178  3.92207534 -0.9251265 -1.1485250 -1.1077564
## 3  1.2279089  0.26624900  0.17613521  1.5255303  1.5159849  1.4306790
## 4  0.5451701  0.14083415  0.12854961  0.1392254  0.2539075  0.3248641
## 5 -0.8549890  0.02411493 -0.04833668 -0.7090093 -0.7700257 -0.7982928
##
## Clustering vector:
## [1] 3 4 4 4 4 5 5 5 4 4 3 4 5 5 5 3 4 5 5 1 1 3 4 4 4 4 4 5 2 4 4 3 4 4 5 1
## [38] 3 3 4 4 5 4 5 5 3 4 4 4 5 5 4 3 4 4 4 4 5 5 1 4 3 5 4 4 5 5 5 5 4 4 3 4
## [75] 4 5 1 5 1 3 4 4 4 5 5 5 3 4 4 4 4 5 1 3 4 4 5 5 5 3 4 5 4 5 5 5 5 2 4 4
## [112] 5 5 5 4 5 1 3 4 4 4 5 5 5 3 4 4 4 5 1 1 3 4 4 5 5 5 3 4 4 4 4 5 5 3 3 4 4 4
## [149] 4 1 3 4 4 4 5 5 5 1 3 4 4 4 4 3 5 1 3 3 3 4 4 4 4 5 5 4 4 4 4 5 5 5 5 5 5 3
## [186] 4 1 1 3 4 4 4 5 5 1 5 3 3 4 4 4 5 5 5 2 1 3 4 4 3 4 4 5 5 5 1 3 4 5 5 5 1 1
## [223] 3 3 4 4 4 4 4 4 1 3 4 4 3 5 3 4 3 4 4 5 5 5 4 4 4 4 4 4 5 1 3 3 3 4 4 4 5 5 1
```

```

## [260] 1 3 3 4 4 4 5 5 5 1 3 3 4 4 4 5 5 5 3 4 4 4 1 3 3 4 4 5 1 4 3 3 4 3 4 5 5
## [297] 3 4 4 4 5 5 5 5 5 3 4 4 5 1 5 1 1 3 3 4 4 5 5 3 3 4 4 4 4 5 5 4 4 4 4 5 5 5
## [334] 3 3 4 4 4 4 5 3 3 4 5 5 5 1 3 4 4 4 4 5 1 3 3 4 5 4 5 5 1 1 3 4 4 4 4 4 5 5
## [371] 5 3 3 4 4 5 5 5 5 5 3 3 4 3 4 5 5 5 1 3 3 4 5 5 3 3 4 4 5 1 3 3 5 5 5
## [408] 3 3 4 5 5 1 3 4 4 4 4 5 5 2 3 4 4 4 4 5 5 5 3 3 4 5 5 4 3 4 4 4 4 5 1 3
## [445] 3 4 4 5 5 5 3 3 4 4 5 1 3 4 4 4 5 5 5 4 3 4 4 4 4 5 5 5 3 4 4 4 5 5 5 1
## [482] 3 3 4 4 5 1 1 3 3 4 5 5 4 5 3 3 4 4 4 5 1 1 4 4 4 3 5 5 3 4 4 4 4 5 1 3
## [519] 4 4 4 5 2 5 3 3 4 4 4 5 1 2 3 3 4 3 5 5 5 1 3 4 4 4 4 4 5 5 5 5 3 3 4 5 5
## [556] 5 3 3 4 4 4 5 5 1 4 4 4 4 5 1 1 1 4 4 4 4 5 5 5 1 3 3 3 5 5 5 5 1 1 3 3 4 4
## [593] 5 4 3 3 4 4 4 4 5 5 2 3 4 4 5 1 1 1 3 4 4 4 4 4 5 5 1 1 3 3 4 4 5 3 4 4 4 5
## [630] 5 1 3 4 3 4 5 5 1 3 4 4 4 5 2 1 3 4 4 4 4 5 4 4 4 4 4 4 4 5 5 5 3 4 3 5 5 5
## [667] 3 4 4 3 5 5 5 1 2 3 4 4 5 5 5 3 3 4 4 4 4 5 5 5 1 3 4 4 4 4 3 5 5 4 4 4 5 1 1
## [704] 1 4 3 4 4 4 4 2 5 1 4 3 3 4 3 5 1 3 3 4 4 5 5 1 5 4 3 4 5 5 4 4 3 4 4 4 5 4
## [741] 5 5 3 3 4 5 3 3 4 5 5 5 5 2 1 3 4 4 4 5 5 3 4 4 4 4 4 5 5 3 3 3 4 4 5 2 4 4 4
## [778] 5 5 5 5 3 4 4 4 4 5 5 1 1 4 4 4 4 4 2 5 2 5 5 4 4 3 4 4 4 5 3 4 4 4 5 5 5
## [815] 5 4 3 3 4 4 1 5 1 1 3 4 4 4 4 5 1 3 4 4 4 5 5 5 5 4 4 4 4 4 4 5 1 1 3 4 4 4
## [852] 5 5 4 5 3 4 4 4 5 3 4 5 5 5 3 4 4 4 5 5 5 3 3 4 4 5 5 5 2 3 4 4 4 1 4 3 4
## [889] 5 4 5 5 5 2 3 4 4 1 5 5 1 4 3 4 4 5 4 4 3 4 5 5 1 1 1 3 3 4 4 5 1 1 3 4 5
## [926] 5 5 4 4 4 4 5 4 5 5 5 3 3 4 3 5 5 5 5 5 3 4 4 4 4 5 5 3 4 4 4 5 5 5 3 4 4
## [963] 4 4 5 5 3 4 4 5 5 5 1 1 3 3 4 4 4 4 5 5 5 1 3 3 4 5 5 5 5 5 3 3 4 5 4 1 3
## [1000] 4 4 5 5 5 1 3 3 4 4 5 5 5 5 5 3 4 4 4 4 4 4 5 5 1 4 3 4 4 4 4 4 5 2 1 3 3 4
## [1037] 4 1 5 3 4 4 4 5 1 5 4 3 4 5 4 5 4 1 4 3 4 4 4 4 5 1 3 4 3 3 4 5 5 5 1 3 4 3 5
## [1074] 5 1 1 3 4 3 4 5 5 5 5 5 5 3 4 4 4 5 5 5 5 3 3 3 4 5 5 5 5 1 3 4 4 4 5 4 5
## [1111] 5 2 3 3 5 4 5 5 5 5 3 3 4 4 4 5 4 5 5 5 1 4 4 4 4 5 4 4 5 5 2 5 1 3 4 4 4 5
## [1148] 5 3 3 3 4 5 1 3 3 4 4 5 1 4 3 4 4 4 5 5 5 3 3 4 4 4 4 5 1 5 3 4 4 5 5 2 1 5
## [1185] 1 3 4 3 5 1 5 1 3 4 4 5 1 5 1 5 1 3 4 4 4 4 5 5 5 1 5 3 4 3 4 4 5 3 3 4 4 5
## [1222] 3 3 4 5 5 4 3 4 4 4 5 5 3 3 4 4 4 4 1 5 2 3 4 4 4 5 5 1 4 4 4 4 4 5 4 5 5 3 3
## [1259] 4 5 5 1 1 5 3 3 3 4 5 5 1 3 4 4 4 4 4 1 3 4 4 5 1 1 3 3 4 4 4 5 2 5 4 3 3 4
## [1296] 4 5 1 3 4 4 4 5 2 1 3 4 5 5 5 5 1 1 3 4 4 4 4 4 5 5 4 4 4 4 4 5 3 4 4 4 4 4 5
## [1333] 5 1 5 3 3 4 4 4 4 5 5 5 4 4 4 3 4 5 4 5 4 3 4 5 4 4 4 5 5 4 3 5 4 5 5 2 1 3 3 4
## [1370] 4 4 5 1 1 3 4 4 4 4 5 1 2 2 4 4 4 4 5 4 5 1 3 4 4 4 5 5 5 5 1 3 3 4 4 4 5 5 3 4
## [1407] 3 5 5 2 3 4 4 4 4 5 5 1 1 4 4 4 4 4 4 5 1 2 1 3 5 5 5 5 5 1 3 3 4 4 4 4 5 5 2 1
## [1444] 1 1 4 5 5 5 3 4 4 4 5 5 1 4 4 4 4 4 4 5 3 4 4 4 4 4 5 5 1 5 3 3 5 5 5 5 5
## [1481] 3 4 4 5 5 5 1 3 3 4 5 5 5 5 3 3 4 4 4 5 5 3 3 3 4 4 5 2 1 4 3 4 4 4 5 5 5 4
## [1518] 4 4 4 5 5 5 5 4 4 4 4 4 5 5 5 3 3 4 4 4 5 1 1 3 3 3 4 4 4 5 3 4 4 4 5 5 5 5 4 4
## [1555] 4 4 4 3 4 1 4 3 3 5 3 4 4 4 4 5 5 5 3 4 5 5 5 5 1 3 4 3 4 5 5 5 5 2 3 4 4
## [1592] 4 5 3 4 3 5 1 5 5 4 4 4 4 4 4 4 5 1 1 4 4 4 4 4 4 4 5 5 1 3 4 3 4 5 5 5 4 5 5 2
## [1629] 1 1 1 3 3 3 4 4 5 3 4 4 4 4 4 5 1 4 4 4 4 4 4 5 5 1 3 3 4 5 4 5 5 2 1 3 4 4 4
## [1666] 5 1 1 4 4 4 5 2 3 3 4 4 4 4 5 5 5 4 3 4 4 4 4 5 5 5 3 4 4 4 4 5 1 5 3 4 4 4 5 4
## [1703] 5 5 5 3 3 4 4 4 5 3 4 4 4 4 4 5 5 5 3 4 4 4 5 5 5 5 2 3 3 4 4 4 4 4 5 5 5 4 3 4 3 4
## [1740] 3 4 4 4 4 4 5 1 3 4 4 4 5 5 1 2 3 3 4 4 4 5 5 1 3 4 4 4 4 4 1 4 3 4 4 4 5 5 3 4 4 5
## [1777] 5 5 5 1 4 3 3 4 5 5 5 3 4 4 4 2 5 5 5 1 3 3 3 4 4 4 4 5 2 3 3 4 4 4 5 3 4 4 4 4
## [1814] 5 4 5 1 5 2 1 4 3 3 5 5 5 3 3 4 4 5 5 1 1 3 3 4 4 4 5 5 5 3 3 3 3 4 1 3 4 3
## [1851] 4 4 5 5 5 5 4 4 4 3 3 5 2 3 4 3 5 5 5 5 1 3 4 4 4 4 5 1 4 3 3 5 4 5 2 1 3 4
## [1888] 3 5 5 5 1 4 4 4 4 4 4 4 5 1 5 4 4 4 4 5 4 5 5 3 4 4 4 5 5 4 3 4 4 4 5 5 5 2 3 4 4
## [1925] 4 4 5 5 5 3 4 3 4 5 5 1 1 3 4 4 4 5 5 5 4 3 4 4 5 4 5 3 4 4 5 5 5 1 4 4 4 4 4 5

```

```

## [1962] 3 4 4 4 4 4 5 2 2 3 3 3 4 4 5 5 1 1 1 3 4 4 3 4 5 5 5 5 5 4 4 5 4 5 5 1 3
## [1999] 3 4 4 5 5 5 3 4 4 3 5 4 5 5 5 3 3 3 4 4 1 3 4 4 4 4 5 5 3 4 4 4 1 1 2 2 3
## [2036] 4 5 4 1 3 3 3 4 5 1 3 4 4 4 4 4 5 5 2 1 3 4 4 5 5 1 3 4 5 5 5 1 1 5 3 3 4 4
## [2073] 5 5 5 5 3 3 3 5 5 5 1 3 3 4 4 4 5 5 5 1 3 4 4 3 5 5 5 1 3 3 4 4 4 5 5 2 1
## [2110] 4 4 4 4 4 5 5 5 3 3 3 4 5 5 5 5 4 4 5 5 5 1 4 3 4 4 4 4 5 5 3 4 3 4 4 5 1
## [2147] 3 4 4 4 5 5 1 5 3 3 4 4 5 5 5 3 3 4 4 4 5 5 2 3 4 4 4 5 5 2 1 2 3 3 4 4 5 5
## [2184] 3 4 4 4 5 5 5 1 4 3 4 4 4 4 5 5 2 5 4 4 3 4 4 4 5 5 5 5 4 4 4 3 4 5 5 1 3
## [2221] 4 4 4 5 5 5 3 4 4 5 5 3 3 4 4 5 5 5 5 3 3 4 4 4 4 5 5 1 5 3 4 4 4 5 5 5
## [2258] 1 4 4 3 4 4 5 5 3 4 3 4 4 5 5 5 4 4 4 4 4 5 4 5 1 5 3 3 3 4 5 5 5 4 3 4 4 5
## [2295] 5 4 4 4 4 4 5 5 5 2 3 4 4 4 5 5 1 3 4 4 4 5 5 4 4 4 4 5 4 5 5 3 4 4 4 5 5
## [2332] 4 3 4 4 4 4 5 5 5 3 4 3 4 5 5 5 2 3 3 3 5 1 3 3 4 5 5 2 1 4 3 4 4 4 4 4
## [2369] 4 4 4 1 1 3 3 5 5 3 3 4 4 5 1 3 4 4 5 5 5 3 3 4 4 4 1 5 3 4 5 5 5 1 1 3
## [2406] 3 4 4 4 5 5 5 2 3 4 4 4 1 5 5 4 4 3 4 4 4 5 5 5 4 3 4 4 4 4 5 5 5 1 3 3 4 5
## [2443] 5 3 3 4 4 4 5 5 2 3 4 4 4 5 5 5 3 4 4 4 4 4 5 2 3 3 4 4 5 5 5 5 3 4 3 4 4
## [2480] 5 1 1 3 3 4 4 5 5 1 3 3 4 3 4 4 5 5 5 3 3 5 5 5 1 5 3 3 3 4 4 4 5 5 1 5 1 3 4
## [2517] 4 4 4 3 3 4 4 4 4 1 1 5 5 4 5 5 5 5 5 1 3 4 4 4 4 5 5 5 1 3 4 4 4 4 5 5 3 4 4
## [2554] 5 5 5 5 3 4 4 5 5 5 1 1 1 3 4 4 4 5 5 5 1 1 3 3 4 4 5 5 5 1 3 3 4 4 4 4 5 1 5
## [2591] 3 3 4 4 5 5 3 4 4 4 4 5 5 5 5 3 3 3 5 5 5 5 1 4 4 3 4 1 5 1 1 3 3 4 5
## [2628] 1 1 3 4 4 4 5 5 5 5 3 4 4 4 5 5 5 5
##
## Within cluster sum of squares by cluster:
## [1] 319.4278 250.5515 1094.8283 1341.0495 1430.7781
## (between_SS / total_SS = 72.1 %)
##
## Available components:
##
## [1] "cluster"      "centers"       "totss"         "withinss"      "tot.withinss"
## [6] "betweenss"    "size"          "iter"          "ifault"

```

Interpreting this output, the very first thing you need to know is that **the cluster numbers are meaningless**. They aren't ranks. They aren't anything. After you have taken that on board, look at the cluster sizes at the top. Clusters 1 and 2 are pretty large compared to others. That's notable. Then we can look at the cluster means. For reference, 0 is going to be average. So group 1 is below average on minutes played. Group 2 is slightly above, group 5 is well above.

So which group is Cam Mack in? Well, first we have to put our data back together again. In K5, there is a list of cluster assignments in the same order we put them in, but recall we have no names. So we need to re-combine them with our original data. We can do that with the following:

```
playercluster <- data.frame(playersselected, k5$cluster)
```

Now we have a dataframe called playercluster that has our player names and what cluster they are in. The fastest way to find Cam Mack is to double click on the playercluster table in the environment and use the search in the top right of the table. Because this is based on some random selections of points to start

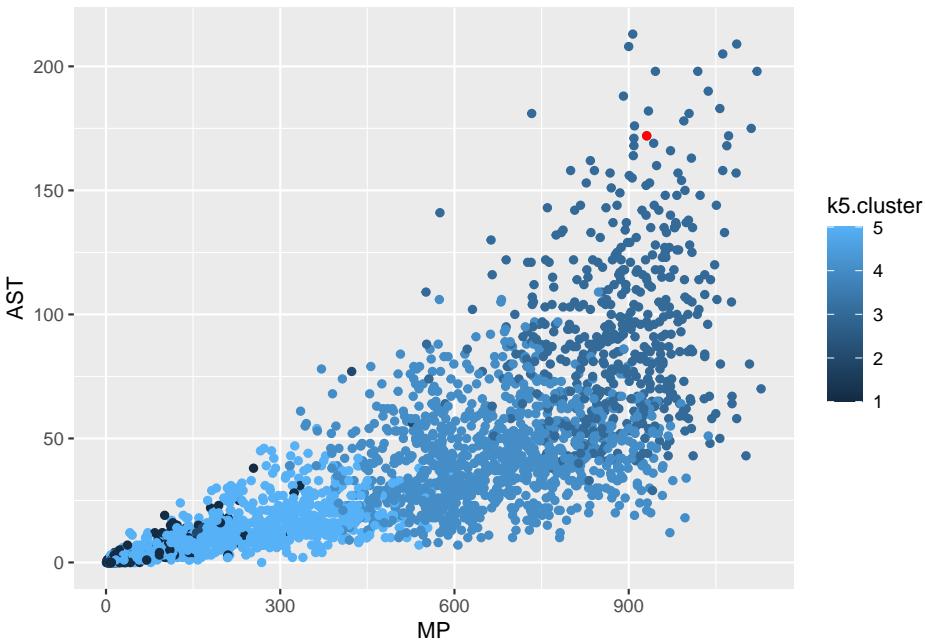
the groupings, these may change from person to person, but Mack is in Group 1 in my data.

We now have a dataset and can plot it like anything else. Let's get Cam Mack and then plot him against the rest of college basketball on assists versus minutes played.

```
cm <- playercluster %>% filter(Player == "Cam Mack")  
  
cm  
  
##      Player           Team Pos MP FG. X3P. AST TOV PTS k5.cluster  
## 1 Cam Mack Nebraska Cornhuskers   G 931 0.39 0.339 172 71 324            3
```

So Cam's in cluster 3, which if you look at our clusters, puts him in the cluster with all above average metrics. What does that look like? We know Cam was an assist machine, so where do group 5 people grade out on assists?

```
ggplot() +  
  geom_point(data=playercluster, aes(x=MP, y=AST, color=k5.cluster)) +  
  geom_point(data=cm, aes(x=MP, y=AST), color="red")
```



Not bad, not bad. But who are Cam Mack's peers? If we look at the numbers in Group 3, there's 495 of them. So let's limit them to just Big Ten guards. Unfortunately, my scraper didn't quite work and in the place of Conference is the coach's name. So I'm going to have to do this the hard way and make a list of Big Ten teams and filter on that. Then I'll sort by minutes played.

```

big10 <- c("Nebraska Cornhuskers", "Iowa Hawkeyes", "Minnesota Golden Gophers", "Illinoi
playercluster %>% filter(k5.cluster == 3) %>% filter(Team %in% big10) %>% arrange(desc

##           Player             Team Pos   MP    FG.   X3P.  AST TOV PTS
## 1      Marcus Carr Minnesota Golden Gophers  G 1004 0.377 0.341 181  76 419
## 2      Anthony Cowan Maryland Terrapins   G  965 0.379 0.331 133  59 454
## 3       Cam Mack Nebraska Cornhuskers  G  931 0.390 0.339 172  71 324
## 4   Eric Hunter Jr.  Purdue Boilermakers  G  907 0.419 0.373  77  57 297
## 5     Zavier Simpson Michigan Wolverines  G  907 0.472 0.354 213  85 351
## 6      Ayo Dosunmu Illinois Fighting Illini  G  891 0.483 0.295  86  72 443
## 7     D'Mitrik Trice Wisconsin Badgers   G  883 0.397 0.390 116  50 289
## 8    Cassius Winston Michigan State Spartans  G  868 0.432 0.409 157  85 497
## 9      CJ Walker Ohio State Buckeyes   G  795 0.420 0.338  94  47 227
## 10     Pat Spencer Northwestern Wildcats  G  787 0.458 0.250 103  63 286
## 11      Geo Baker Rutgers Scarlet Knights  G  736 0.389 0.261  91  45 273
##   k5.cluster
## 1      3
## 2      3
## 3      3
## 4      3
## 5      3
## 6      3
## 7      3
## 8      3
## 9      3
## 10     3
## 11     3

```

So there are the 11 guards most like Cam Mack in the Big Ten. Safe to say, these are the 11 best guards in the conference.

### 13.1 Advanced metrics

How much does this change if we change the metrics? I used pretty standard box score metrics above. What if we did it using Player Efficiency Rating, True Shooting Percentage, Point Production, Assist Percentage, Win Shares Per 40 Minutes and Box Plus Minus (you can get definitions of all of them by hovering over the stats on Nebraksa's stats page).

We'll repeat the process. Filter out players who don't play, players with stats missing, and just focus on those stats listed above.

```

playersadvanced <- players %>%
  filter(MP>0) %>%
  filter(Pos == "G") %>%

```

```
select(Player, Team, Pos, PER, `TS%`, PProd, `AST%`, `WS/40`, BPM) %>%
na.omit()
```

Now to scale them.

```
playersadvscaled <- playersadvanced %>%
  select(PER, `TS%`, PProd, `AST%`, `WS/40`, BPM) %>%
  mutate_all(scale) %>%
  na.omit()
```

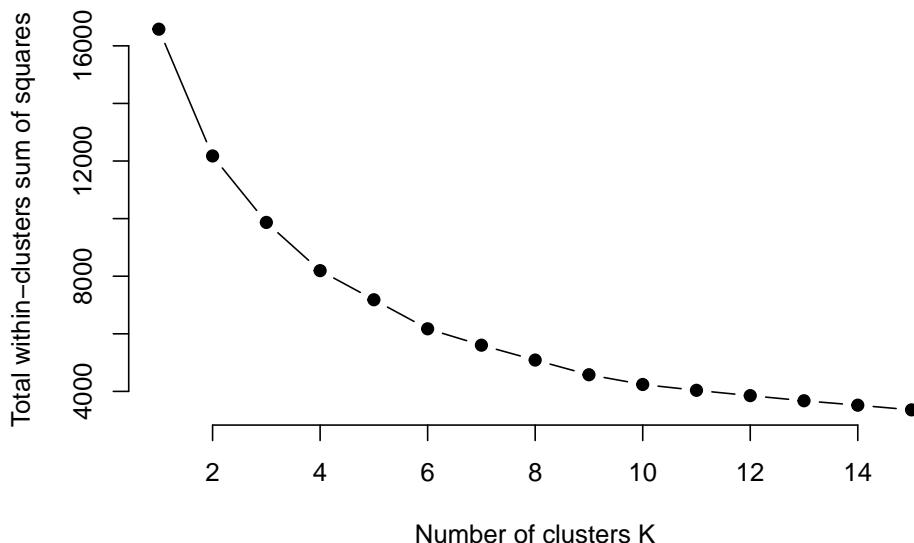
Let's find the optimal number of clusters.

```
# function to compute total within-cluster sum of square
wss <- function(k) {
  kmeans(playersadvscaled, k, nstart = 10 )$tot.withinss
}

# Compute and plot wss for k = 1 to k = 15
k.values <- 1:15

# extract wss for 2-15 clusters
wss_values <- map_dbl(k.values, wss)

## Warning: did not converge in 10 iterations
plot(k.values, wss_values,
      type="b", pch = 19, frame = FALSE,
      xlab="Number of clusters K",
      ylab="Total within-clusters sum of squares")
```



Looks like 5 again.

```
advk5 <- kmeans(playersadvscaled, centers = 5, nstart = 25)
```

What do we have here?

```
advk5
```

```
## K-means clustering with 5 clusters of sizes 1477, 1030, 166, 3, 88
##
## Cluster means:
##          PER         TS%       PProd        AST%      WS/40       BPM
## 1 -0.1853781 -0.06873982 -0.5343488 -0.2696473 -0.1184126 -0.1622670
## 2  0.4087876  0.24526817  1.0477637  0.4472182  0.3563642  0.4575732
## 3 -2.1693836 -2.42045625 -1.1410523 -0.7345409 -2.2297633 -2.2703594
## 4 14.5515907  4.92713793 -1.1318648 -1.3849128 16.0721265  8.9433106
## 5  1.9228983  2.68087293 -1.1040127  0.7241224  1.4745973  1.3456549
##
## Clustering vector:
## [1] 2 1 1 1 1 1 1 1 2 2 2 1 2 1 5 2 2 1 1 1 1 3 2 1 1 1 2 1 1 1 5 2 2 2 1
## [38] 1 3 2 2 2 2 1 1 1 1 2 2 2 2 1 1 5 3 2 2 2 1 2 2 1 1 5 2 2 2 1 1 1 1 3 2 2 1
## [75] 1 1 2 2 2 1 2 1 3 1 3 2 1 2 1 1 1 1 2 2 2 1 1 3 2 2 1 1 1 1 3 2 2 1 1 1 1 3 2 2 1
## [112] 1 1 1 1 1 2 2 2 1 1 1 1 3 2 2 2 1 1 1 1 2 2 2 1 1 5 5 3 2 1 1 1 1 1 2 2 2
## [149] 1 1 1 1 3 2 2 2 1 2 1 1 2 2 1 1 1 1 5 1 3 2 2 2 1 2 2 2 3 2 2 2 2 1 1 1 1
## [186] 2 2 2 1 1 1 1 1 1 1 2 2 1 3 2 2 2 2 1 1 1 1 2 2 1 1 1 1 1 5 4 1 2 2 1 2 2 1
## [223] 1 1 1 1 2 2 2 1 1 1 3 3 2 2 2 2 1 1 1 1 2 1 2 2 1 2 2 2 1 1 1 1 5 1 3 2 2 2
## [260] 1 1 1 3 2 2 2 2 1 1 1 1 5 3 3 2 2 1 1 1 1 1 3 2 2 2 1 1 1 1 2 1 1 1 2 2 2
## [297] 1 5 1 2 2 2 1 1 3 2 2 2 2 2 1 1 1 2 2 1 2 1 1 1 1 2 1 1 1 1 1 1 2 2 2 1 1 1 2 2 1
## [334] 1 1 1 2 2 1 1 1 1 1 2 2 1 1 1 1 1 2 2 1 1 1 1 1 2 2 2 1 1 1 1 3 2 2 2 1 1 1
## [371] 1 2 2 2 1 1 1 1 3 3 2 2 2 2 2 1 1 2 2 2 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 1 1 1
## [408] 1 3 2 2 1 1 1 1 2 2 2 1 1 5 3 2 2 1 1 3 1 3 2 2 1 1 1 1 1 1 2 2 2 1 1 1 1 1 1 1 1
## [445] 5 2 2 2 2 2 1 1 1 2 2 2 1 1 1 2 2 2 2 1 1 3 3 2 2 2 1 1 1 1 1 2 2 2 1 1 1 1 2 2 2 1
## [482] 1 3 2 2 1 1 1 1 1 2 2 2 2 1 1 1 1 2 2 1 1 1 1 1 1 3 2 2 1 1 1 1 1 3 2 2 2
## [519] 1 1 2 1 1 2 2 2 2 1 1 1 1 1 2 1 1 2 1 1 1 3 3 2 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1
## [556] 2 2 2 2 1 1 1 1 5 2 2 2 2 1 1 1 3 2 2 2 2 2 1 1 1 1 1 2 2 2 1 1 5 2 2 1 1 1
## [593] 1 1 1 2 1 2 1 1 1 1 3 2 2 1 1 1 1 3 2 2 2 1 1 1 1 3 3 2 2 2 1 1 1 1 2 2 2
## [630] 1 1 1 1 1 1 2 2 2 1 1 3 3 2 2 2 2 1 1 1 1 1 1 2 2 2 1 1 2 2 2 1 1 1 1 3 2
## [667] 2 2 1 1 1 3 2 2 1 1 1 5 3 1 2 2 2 2 2 1 1 2 1 2 2 2 1 1 3 2 2 2 1 1 1 1 2 2 2
## [704] 2 1 2 1 1 1 1 4 2 1 1 1 1 2 2 1 1 1 1 1 1 3 2 2 2 1 2 1 1 1 1 2 2 2 1 1 1 2 2 2 1
## [741] 1 3 1 2 2 2 2 1 2 5 1 3 2 2 2 2 2 1 1 3 2 2 1 1 1 1 1 1 2 2 2 1 1 1 3 2 2 2
## [778] 2 2 1 1 1 1 2 2 2 1 2 2 1 1 1 1 1 1 2 2 1 1 1 3 2 2 2 1 2 2 1 1 1 5 3 2 2 2 1 1
## [815] 1 1 5 2 2 1 1 1 1 1 1 2 1 1 2 1 1 1 1 3 2 2 1 1 2 1 1 1 1 1 5 3 2 2 2 1 1
## [852] 1 1 2 2 2 1 1 1 1 2 2 2 2 1 1 1 1 1 2 2 1 1 1 1 3 2 2 2 1 1 1 1 1 1 2 2 2
## [889] 2 1 2 5 3 1 2 1 2 1 2 1 2 2 2 2 1 2 2 1 1 1 1 1 2 2 1 2 1 1 1 2 2 2 1 1 1 2 2 2 1
## [926] 1 1 5 2 2 1 1 5 3 2 2 1 1 1 1 1 1 5 2 2 1 1 1 1 3 2 2 2 2 1 2 1 1 1 1 1 2 2 2 2 1
## [963] 1 1 2 2 2 2 1 1 1 2 2 1 1 2 1 1 1 2 2 1 1 1 1 1 2 2 2 2 1 1 1 1 1 1 1 1 2 2 2 2
## [1000] 1 1 2 2 2 1 1 1 1 2 2 2 1 1 1 1 2 2 1 1 1 1 1 3 2 2 2 2 1 1 1 1 1 1 3 2 2 2 2
```

```
## [1037] 1 1 1 1 1 2 2 2 1 1 3 2 2 1 1 1 1 3 2 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 5 1 2
## [1074] 2 2 1 1 1 1 5 1 2 2 2 1 3 1 2 2 2 1 1 1 1 1 1 1 2 1 1 1 1 1 3 1 2 1 1 1 1 3
## [1111] 2 2 2 2 1 1 1 1 3 2 2 2 1 1 1 3 2 2 2 2 1 1 1 1 1 1 1 5 2 1 1 1 1 1 2 2 2
## [1148] 1 1 1 1 1 2 2 1 1 1 1 1 1 1 1 2 2 1 1 1 1 1 1 1 2 2 2 2 1 1 2 1 1 1 1 1 2 2 1
## [1185] 1 1 1 1 5 5 5 3 2 1 2 2 1 1 2 2 2 2 1 3 2 2 2 2 1 1 3 2 2 2 2 1 1 1 1 1 2 2 2
## [1222] 2 2 2 1 1 1 2 2 2 1 1 5 1 5 3 2 2 2 2 1 1 1 1 1 2 2 2 1 1 1 3 1 3 2 2 1 2 1 1
## [1259] 1 3 1 2 2 2 2 1 3 2 2 1 1 1 2 2 1 1 1 2 2 2 1 1 1 1 1 2 2 1 1 1 1 1 5 2 1 1
## [1296] 1 1 1 1 2 2 2 1 1 1 1 2 2 2 1 1 1 3 1 2 2 2 1 1 1 3 2 2 2 2 2 3 2 2 2 1
## [1333] 3 3 2 2 2 1 1 1 5 1 2 2 2 2 1 1 1 2 2 2 2 1 5 3 5 2 1 1 2 5 1 1 1 2 2 2 2
## [1370] 1 1 5 2 2 2 1 2 1 3 2 2 1 2 1 1 1 1 5 2 2 2 1 1 1 1 1 2 1 2 2 2 1 1 2 2 2
## [1407] 1 1 1 1 1 2 2 1 1 1 1 1 5 3 2 2 2 2 1 1 1 3 2 2 2 2 1 1 5 1 1 1 1 1 1 1 1
## [1444] 1 3 2 2 2 1 2 1 1 1 2 2 1 1 1 2 2 2 1 1 5 2 1 1 1 1 1 3 3 3 1 2 2 2 1 1
## [1481] 1 5 1 3 2 1 1 1 1 1 3 2 2 1 2 1 1 1 1 3 3 1 1 1 1 1 2 2 2 2 1 1 1 3 3 2
## [1518] 2 2 2 1 1 1 2 2 2 2 1 1 1 1 1 2 2 1 1 1 1 1 2 2 1 1 1 1 1 2 2 1 1 1 1 1 1
## [1555] 2 2 1 1 1 2 2 1 2 2 1 5 1 2 2 2 1 1 1 2 2 2 1 1 1 1 1 2 2 2 2 1 1 1 1 1
## [1592] 1 2 2 1 1 1 3 3 2 2 2 1 1 1 2 2 2 1 1 1 1 1 2 2 2 1 1 2 1 3 2 2 2 1 1 2 2 2
## [1629] 1 1 1 1 2 1 1 1 1 1 1 2 2 2 1 1 1 1 5 2 1 2 2 1 2 2 2 1 1 1 1 2 1 2 1
## [1666] 1 1 1 3 1 2 2 1 1 1 2 1 1 1 2 2 2 1 1 1 1 1 1 5 1 1 3 2 2 2 1 1 1 1 2 2 2
## [1703] 1 1 1 3 2 2 2 2 1 1 1 3 2 2 1 1 1 1 1 5 3 2 2 2 1 1 1 3 3 2 2 2 1 5 3 2
## [1740] 2 2 1 2 1 1 1 1 3 2 2 2 1 1 1 1 1 5 2 2 2 2 1 1 1 2 2 2 1 1 2 1 1 1 2 2 2 1
## [1777] 1 1 2 2 1 1 1 1 1 1 3 2 2 1 1 1 5 1 1 5 2 2 2 1 2 1 1 1 2 2 2 2 2 1 1 1
## [1814] 1 1 2 1 2 1 1 1 1 2 2 2 2 1 1 1 1 3 2 2 1 1 1 1 2 2 1 1 1 1 2 2 1 1 1 1 1
## [1851] 2 3 2 2 2 1 1 1 1 2 2 1 1 1 1 1 1 5 3 2 2 2 2 2 1 1 5 2 2 2 1 1 2 1 1 1
## [1888] 1 1 5 1 2 5 1 2 2 2 1 1 1 1 2 2 1 2 1 1 1 3 2 2 1 1 1 1 5 5 2 2 2 2 3 3 2
## [1925] 2 2 2 1 1 5 1 1 1 2 2 2 2 2 2 1 3 2 2 2 1 1 1 1 5 1 2 2 1 2 1 1 3 3 2 2 2 1
## [1962] 1 1 1 3 2 2 2 1 1 1 3 2 2 1 2 2 1 1 1 1 2 1 2 1 1 1 1 2 2 2 1 1 1 2 2 2
## [1999] 2 1 1 5 5 1 3 2 2 1 1 1 2 1 1 1 2 2 2 1 1 1 1 3 2 2 2 1 1 1 2 2 1 1 1 1 5
## [2036] 2 2 1 1 1 1 1 2 2 1 1 2 2 1 1 1 1 1 5 1 1 2 2 2 1 1 1 1 1 3 3 2 2 1 1 2
## [2073] 1 1 1 1 1 1 2 2 1 1 1 3 2 2 1 1 1 1 2 2 2 1 1 2 1 1 1 2 2 2 2 3 2 1
## [2110] 1 1 1 1 1 3 2 2 1 1 1 5 5 2 2 1 2 1 2 2 1 1 3 2 2 2 2 1 1 1 1 5 3 2 1
## [2147] 1 1 1 3 2 2 1 1 1 3 1 3 2 2 1 1 1 1 1 2 2 2 1 1 1 3 2 1 1 1 1 1 1 1 1
## [2184] 2 2 2 2 1 1 5 5 5 3 2 2 2 2 1 1 1 1 5 3 2 2 1 1 1 1 1 5 2 2 2 2 1 1 1 1 2
## [2221] 1 1 1 5 1 1 2 2 2 1 1 1 1 1 1 2 2 2 1 1 1 3 3 2 1 1 1 1 1 1 1 1 1 2 2 2 1 1
## [2258] 1 1 2 2 2 1 1 1 2 5 2 2 1 1 1 1 1 4 2 2 2 1 1 1 1 5 2 2 2 2 1 2 1 1 2 2 2
## [2295] 1 1 1 1 1 5 1 2 2 2 1 1 1 1 1 2 2 2 2 1 1 1 3 2 2 2 1 1 1 1 1 2 2 1 1 2 2 1 1
## [2332] 1 2 2 1 1 1 2 1 1 1 2 2 2 2 1 1 1 3 1 2 2 2 2 1 1 1 1 2 2 2 2 1 1 1 1 2 2 2
## [2369] 1 1 1 1 1 1 1 1 1 1 1 3 2 2 2 1 1 1 1 1 2 2 1 1 1 1 2 2 2 2 1 1 1 1 5 2
## [2406] 1 2 1 1 1 3 2 2 2 1 1 1 2 2 1 2 1 1 1 1 2 2 1 1 1 1 2 2 2 2 1 1 1 1 1 5 2
## [2443] 2 2 1 1 1 1 5 2 2 2 1 3 2 2 1 1 1 5 1 2 2 2 1 1 2 1 1 2 1 1 1 2 2 2 1 1 2 2
## [2480] 2 2 1 1 3 3 2 2 2 1 1 1 2 2 1 2 2 1 5 2 2 1 1 1 1 3 2 2 2 2 2 1 1 1 5 3
## [2517] 2 2 2 1 1 1 1 2 2 2 1 1 2 1 1 1 5 2 2 1 1 1 1 1 3 3 2 2 1 1 1 3 3 2 2 2
## [2554] 1 1 1 5 5 2 1 1 1 1 1 2 1 1 2 1 1 1 1 2 2 2 1 1 1 1 1 2 2 2 1 2 2 3 3
## [2591] 2 2 1 1 1 1 1 2 2 2 2 1 1 1 5 2 2 1 1 1 1 1 3 2 2 2 1 1 1 1 1 3 2 2 2 2 2
## [2628] 1 3 2 2 2 2 1 1 1 3 1 1 1 1 1 1 1 3 2 2 2 1 1 1 1 3 2 2 2 2 1 1 1 3 2 2 2 2 1 2 2
## [2665] 2 1 1 1 1 1 2 2 2 1 1 1 3 1 3 2 2 2 1 1 1 3 3 2 2 2 1 1 1 1 5 3 2 2 1 1 1
## [2702] 1 1 1 2 2 2 1 1 1 2 2 1 1 1 1 1 1 1 2 2 2 1 1 1 3 3 2 1 2 1 1 1 3 3
```

```
## [2739] 2 2 2 1 1 5 3 2 2 2 1 1 1 5 1 5 1 2 2 2 1 1 1 1 1
##
## Within cluster sum of squares by cluster:
## [1] 2608.03836 1795.99837 1233.54705 29.44812 1412.27741
## (between_SS / total_SS = 57.3 %)
##
## Available components:
##
## [1] "cluster"      "centers"       "totss"         "withinss"      "tot.withinss"
## [6] "betweenss"    "size"          "iter"          "ifault"
```

Looks like this time, cluster 1 is all below average and cluster 5 is mostly above. Which cluster is Cam Mack in?

```
playeradvcluster <- data.frame(playersadvanced, advk5$cluster)
```

```
cmandv <- playeradvcluster %>% filter(Player == "Cam Mack")
```

```
cmandv
```

```
##           Player              Team Pos PER   TS. PProd AST. WS.40 BPM
## 1 Cam Mack Nebraska Cornhuskers G 15.9 0.481 382 36.4 0.081 3.9
##   advk5.cluster
## 1                 2
```

Cluster 2 on my dataset. So in this season, we can say he's in a big group of players who are all above average on these advanced metrics.

Now who are his Big Ten peers?

```
playeradvcluster %>%
  filter(advk5.cluster == 2) %>%
  filter(Team %in% big10) %>%
  arrange(desc(PProd))
```

	Player	Team	Pos	PER	TS.	PProd	AST.	WS.40	BPM
## 1	Cassius Winston	Michigan State Spartans	G	22.5	0.567	495	37.2		
## 2	Marcus Carr	Minnesota Golden Gophers	G	18.1	0.491	468	37.0		
## 3	Anthony Cowan	Maryland Terrapins	G	21.1	0.551	466	29.2		
## 4	Zavier Simpson	Michigan Wolverines	G	19.7	0.535	436	43.6		
## 5	Ayo Dosunmu	Illinois Fighting Illini	G	19.6	0.551	411	21.1		
## 6	Cam Mack	Nebraska Cornhuskers	G	15.9	0.481	382	36.4		
## 7	Joe Wieskamp	Iowa Hawkeyes	G	19.7	0.566	368	9.3		
## 8	Ron Harper Jr.	Rutgers Scarlet Knights	G	20.4	0.541	318	7.6		
## 9	D'Mitrik Trice	Wisconsin Badgers	G	16.2	0.531	312	27.5		
## 10	Haanif Cheatham	Nebraska Cornhuskers	G	15.6	0.551	308	9.1		
## 11	Andres Feliz	Illinois Fighting Illini	G	20.1	0.542	306	22.1		
## 12	Pat Spencer	Northwestern Wildcats	G	15.9	0.527	299	28.6		
## 13	Eric Hunter Jr.	Purdue Boilermakers	G	14.0	0.524	297	17.0		

## 14	Dachon Burke	Nebraska Cornhuskers	G 14.3 0.461	294 9.7
## 15	Myreon Jones	Penn State Nittany Lions	G 21.0 0.589	288 20.0
## 16	Eli Brooks	Michigan Wolverines	G 13.8 0.521	283 11.2
## 17	Geo Baker	Rutgers Scarlet Knights	G 16.5 0.483	282 23.6
## 18	Aaron Wiggins	Maryland Terrapins	G 14.5 0.496	282 11.0
## 19	Gabe Kalscheur	Minnesota Golden Gophers	G 10.5 0.491	273 9.2
## 20	Aljami Durham	Indiana Hoosiers	G 14.3 0.573	262 18.3
## 21	Darryl Morsell	Maryland Terrapins	G 13.4 0.518	252 17.0
## 22	Trent Frazier	Illinois Fighting Illini	G 11.5 0.495	251 12.0
## 23	Brad Davison	Wisconsin Badgers	G 14.3 0.546	250 11.9
## 24	Jahaad Proctor	Purdue Boilermakers	G 15.5 0.506	249 13.9
## 25	Devonte Green	Indiana Hoosiers	G 16.1 0.507	249 19.6
## 26	CJ Walker	Ohio State Buckeyes	G 13.9 0.534	247 23.3
## 27	Duane Washington Jr.	Ohio State Buckeyes	G 16.6 0.544	244 13.4
## 28	Franz Wagner	Michigan Wolverines	G 15.5 0.552	242 5.6
## 29	Myles Dread	Penn State Nittany Lions	G 13.7 0.498	237 15.6
## 30	Eric Ayala	Maryland Terrapins	G 10.3 0.470	237 17.8
## 31	Izaiah Brockington	Penn State Nittany Lions	G 15.5 0.530	232 9.4
## 32	Jacob Young	Rutgers Scarlet Knights	G 12.3 0.455	227 17.4
## 33	Boo Buie	Northwestern Wildcats	G 12.7 0.475	223 22.8
## 34	Sasha Stefanovic	Purdue Boilermakers	G 13.9 0.549	221 11.3
## 35	Thorir Thorbjarnarson	Nebraska Cornhuskers	G 14.1 0.606	218 8.8
## 36	CJ Fredrick	Iowa Hawkeyes	G 17.2 0.649	218 16.9
## 37	Connor McCaffery	Iowa Hawkeyes	G 12.5 0.491	214 20.5
## 38	Caleb McConnell	Rutgers Scarlet Knights	G 14.2 0.493	213 13.3
## 39	Alan Griffin	Illinois Fighting Illini	G 27.6 0.652	209 6.1
## 40	D.J. Carton	Ohio State Buckeyes	G 17.8 0.593	208 27.0
## 41	Joe Toussaint	Iowa Hawkeyes	G 13.4 0.466	206 26.9
## 42	Brevin Pritzl	Wisconsin Badgers	G 13.5 0.545	205 5.5
## 43	Payton Willis	Minnesota Golden Gophers	G 13.4 0.517	193 13.3
## 44	Kobe King	Wisconsin Badgers	G 14.4 0.513	188 12.2
## 45	David DeJulius	Michigan Wolverines	G 13.3 0.522	187 12.9
## 46	Nojel Eastern	Purdue Boilermakers	G 11.5 0.434	182 19.3
## 47	Luther Muhammad	Ohio State Buckeyes	G 12.2 0.557	181 10.9
## 48	Rob Phinisee	Indiana Hoosiers	G 11.5 0.480	175 25.8
## 49	Jamari Wheeler	Penn State Nittany Lions	G 10.4 0.564	147 19.6
## 50	Paul Mulcahy	Rutgers Scarlet Knights	G 13.0 0.605	128 18.7
## 51	Foster Loyer	Michigan State Spartans	G 16.5 0.662	86 22.7
## 52	Tommy Luce	Purdue Boilermakers	G 14.0 0.385	12 43.0
## 53	Tino Malnati	Northwestern Wildcats	G 24.3 0.500	3 40.4
##	WS.40	BPM advk5.cluster		
## 1	0.203	8.0	2	
## 2	0.146	6.7	2	
## 3	0.204	9.4	2	
## 4	0.157	7.3	2	
## 5	0.161	6.2	2	

```
## 6 0.081 3.9 2
## 7 0.149 7.3 2
## 8 0.186 7.4 2
## 9 0.133 6.8 2
## 10 0.091 2.5 2
## 11 0.185 7.3 2
## 12 0.076 3.0 2
## 13 0.120 5.9 2
## 14 0.049 1.1 2
## 15 0.183 8.4 2
## 16 0.107 5.8 2
## 17 0.139 6.9 2
## 18 0.123 6.2 2
## 19 0.079 3.8 2
## 20 0.122 4.6 2
## 21 0.115 5.5 2
## 22 0.130 5.2 2
## 23 0.137 6.7 2
## 24 0.139 5.0 2
## 25 0.114 4.9 2
## 26 0.142 6.9 2
## 27 0.151 4.4 2
## 28 0.125 7.0 2
## 29 0.129 6.5 2
## 30 0.089 3.3 2
## 31 0.123 3.4 2
## 32 0.071 1.8 2
## 33 0.039 0.2 2
## 34 0.133 6.8 2
## 35 0.097 4.3 2
## 36 0.133 6.7 2
## 37 0.109 6.5 2
## 38 0.121 5.1 2
## 39 0.252 10.6 2
## 40 0.152 6.6 2
## 41 0.071 2.5 2
## 42 0.123 5.7 2
## 43 0.120 6.0 2
## 44 0.100 3.3 2
## 45 0.113 4.5 2
## 46 0.085 4.5 2
## 47 0.129 6.4 2
## 48 0.080 3.3 2
## 49 0.095 5.8 2
## 50 0.130 5.7 2
## 51 0.183 5.0 2
```

```
## 52 0.102 2.2      2
## 53 0.161 7.8      2
```

Sorting on Points Produced, Cam Mack is sixth out of the 53 guards in the Big Ten who land in Cluster 2. Seems advanced metrics take a little bit of the shine off of Cam. But then, so does leaving the program after one suspension-riddled season.



# Chapter 14

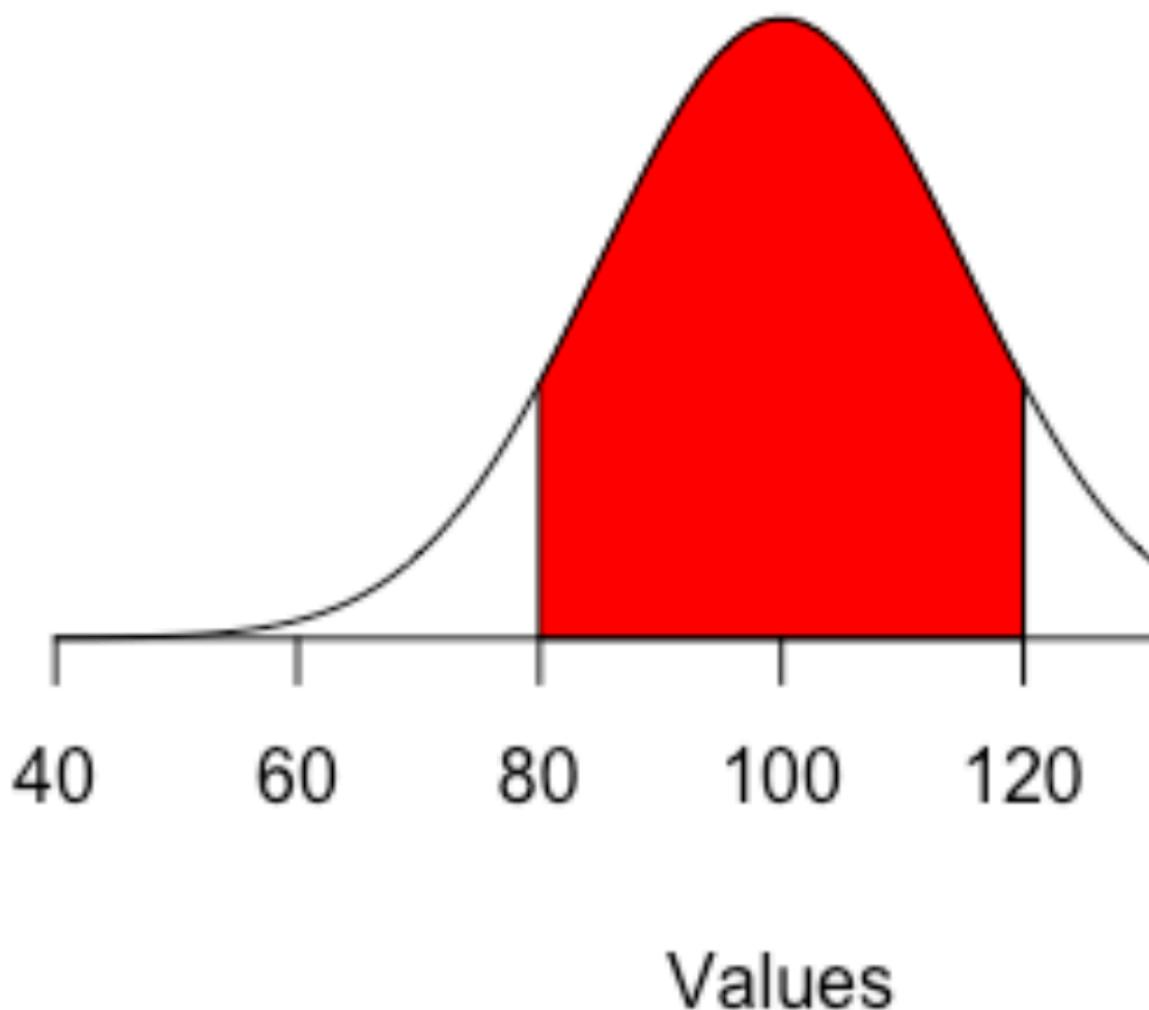
## Simulations

In the 2017-2018 season, James Palmer Jr. took 139 three point attempts and made 43 of them for a .309 shooting percentage. A few weeks into the next season, he was 7 for 39 – a paltry .179.

Is something wrong or is this just bad luck?

Luck is something that comes up a lot in sports. Is a team unlucky? Or a player? One way we can get to this, we can get to that is by simulating things based on their typical percentages. Simulations work by choosing random values within a range based on a distribution. The most common distribution is the normal or binomial distribution. The normal distribution is where the most cases appear around the mean, 66 percent of cases are within one standard deviation from the mean, and the further away from the mean you get, the more rare things become.

## Normal Distribution



Let's simulate 39 three point attempts 1000 times with his season long shooting percentage and see if this could just be random chance or something else.

We do this using a base R function called `rbinom` or binomial distribution. So what that means is there's a normally distributed chance that James Palmer Jr. is going to shoot above and below his career three point shooting percentage.

If we randomly assign values in that distribution 1000 times, how many times will it come up 7, like this example?

```
set.seed(1234)

simulations <- rbinom(n = 1000, size = 39, prob = .309)

table(simulations)

## simulations
##   3   4   5   6   7   8   9   10  11  12  13  14  15  16  17  18  19  20  21  22
##   1   4   5  12  35  44  76 117 134 135 135  99  71  53  37  21  15   2   3   1
```

How do we read this? The first row and the second row form a pair. The top row is the number of shots made. The number immediately under it is the number of simulations where that occurred.

simulations												
3	4	5	6	7	8	9	10	11	12	13	14	15
1	4	5	12	35	44	76	117	134	135	135	99	71
18	19	20	21	22								
21	15	2	3	1								

So what we see is given his season long shooting percentage, it's not out of the realm of randomness that with just 39 attempts for Palmer, he's only hit only 7. In 1000 simulations, it comes up 35 times. Is he below where he should be? Yes. Will he likely improve and soon? Unless something is very wrong, yes. And indeed, by the end of the season, he finished with a .313 shooting percentage from 3 point range. So we can say he was just unlucky.

## 14.1 Cold streaks

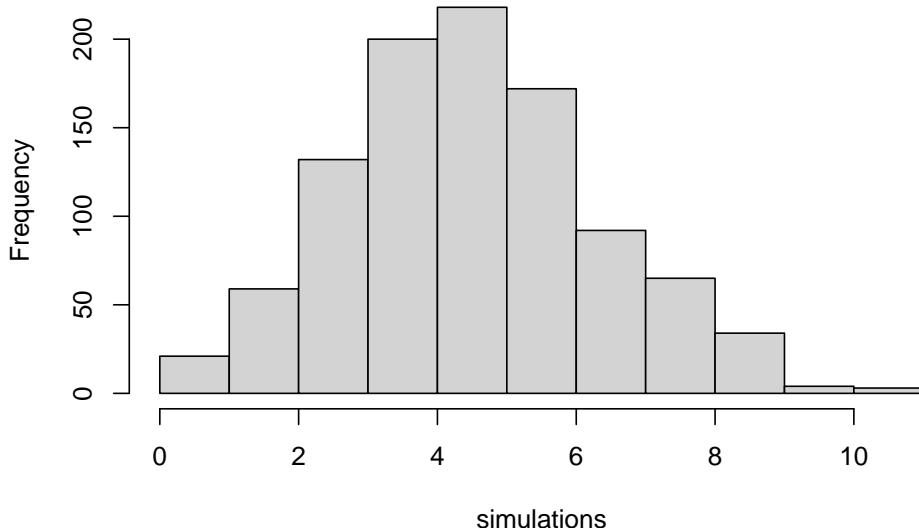
During the Western Illinois game in the 2018-2019 season, the team, shooting .329 on the season from behind the arc, went 0-15 in the second half. How strange is that?

```
set.seed(1234)

simulations <- rbinom(n = 1000, size = 15, prob = .329)
```

```
hist(simulations)
```

**Histogram of simulations**



```
table(simulations)
```

```
## simulations
##   0   1   2   3   4   5   6   7   8   9   10  11
##   5  16  59 132 200 218 172  92  65  34   4   3
```

Short answer: Really weird. If you simulate 15 threes 1000 times, sometimes you'll see them miss all of them, but only a few times – five times, in this case. Most of the time, the team won't go 0-15 even once. So going ice cold is not totally out of the realm of random chance, but it's highly unlikely.

# Chapter 15

## Intro to ggplot

With `ggplot2`, we dive into the world of programmatic data visualization. The `ggplot2` library implements something called the grammar of graphics. The main concepts are:

- aesthetics - which in this case means the data which we are going to plot
- geometries - which means the shape the data is going to take
- scales - which means any transformations we might make on the data
- facets - which means how we might graph many elements of the same dataset in the same space
- layers - which means how we might lay multiple geometries over top of each other to reveal new information.

Hadley Wickham, who is behind all of the libraries we have used in this course to date, wrote about his layered grammar of graphics in this 2009 paper that is worth your time to read.

Here are some `ggplot2` resources you'll want to keep handy:

- The `ggplot` documentation.
- The `ggplot` cookbook

Let's dive in using data we've already seen before – football attendance. This workflow will represent a clear picture of what your work in this class will be like for much of the rest of the semester. One way to think of this workflow is that your R Notebook is now your digital sketchbook, where you will try different types of visualizations to find ones that work. Then, you will either write the code that adds necessary and required parts to finish it, or you'll export your work into a program like Illustrator to finish the work.

To begin, we'll use data we've seen before: college football attendance.

Now load the tidyverse.

```
library(tidyverse)
```

And the data.

```
attendance <- read_csv('data/attendance.csv')

## 
## -- Column specification -----
## cols(
##   Institution = col_character(),
##   Conference = col_character(),
##   `2013` = col_double(),
##   `2014` = col_double(),
##   `2015` = col_double(),
##   `2016` = col_double(),
##   `2017` = col_double(),
##   `2018` = col_double()
## )
```

First, let's get a top 10 list by announced attendance in the most recent season we have data. We'll use the same tricks we used in the filtering assignment.

```
attendance %>%
  arrange(desc(`2018`)) %>%
  top_n(10) %>%
  select(Institution, `2018`)

## Selecting by 2018

## # A tibble: 10 x 2
##   Institution `2018`
##   <chr>        <dbl>
## 1 Michigan     775156
## 2 Penn St.    738396
## 3 Ohio St.    713630
## 4 Alabama      710931
## 5 LSU          705733
## 6 Texas A&M  698908
## 7 Tennessee    650887
## 8 Georgia       649222
## 9 Nebraska     623240
## 10 Oklahoma    607146
```

That looks good, so let's save it to a new data frame and use that data frame instead going forward.

```
top10 <- attendance %>%
  arrange(desc(`2018`)) %>%
  top_n(10) %>%
```

```
select(Institution, `2018`)
## Selecting by 2018
```

## 15.1 The bar chart

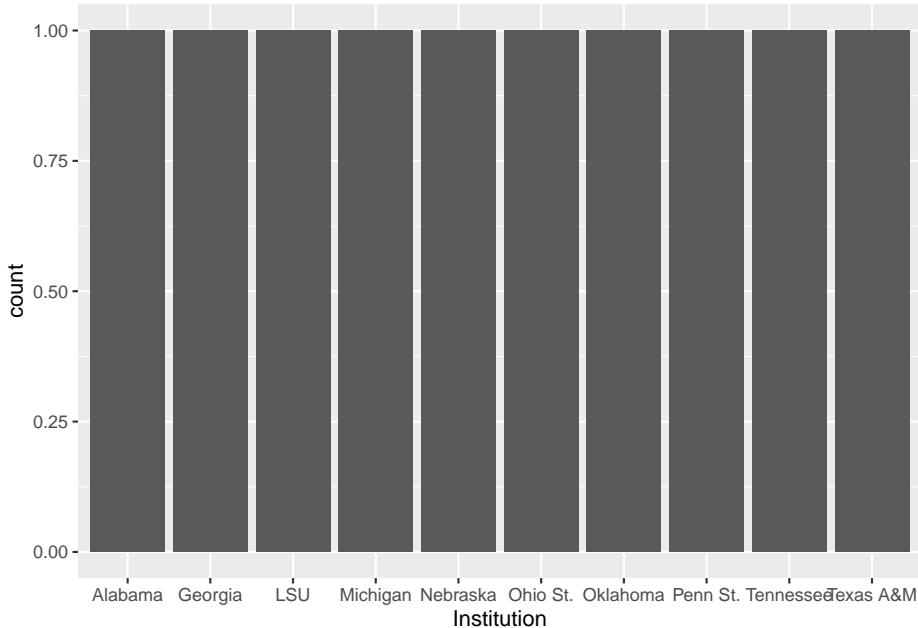
The easiest thing we can do is create a simple bar chart of our data. **Bar charts show magnitude. They invite you to compare how much more or less one thing is compared to others.**

We could, for instance, create a bar chart of the total attendance. To do that, we simply tell `ggplot2` what our dataset is, what element of the data we want to make the bar chart out of (which is the aesthetic), and the geometry type (which is the geom). It looks like this:

```
ggplot() + geom_bar(data=top10, aes(x=Institution))
```

Note: `top10` is our data, `aes` means aesthetics, `x=Institution` explicitly tells `ggplot2` that our x value – our horizontal value – is the `Institution` field from the data, and then we add on the `geom_bar()` as the geometry. And what do we get when we run that?

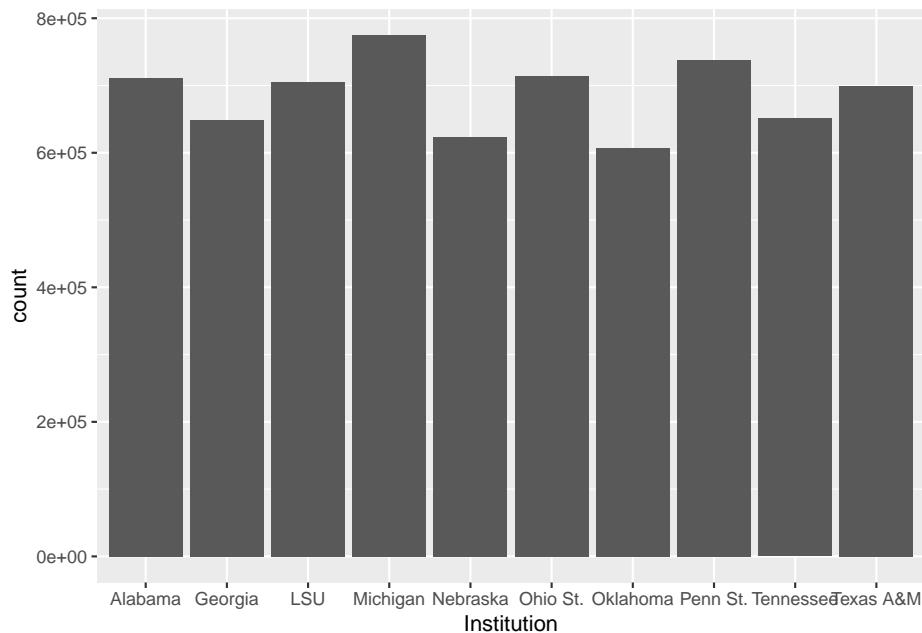
```
ggplot(top10, aes(x=Institution)) + geom_bar()
```



We get ... weirdness. We expected to see bars of different sizes, but we get all with a count of 1. What gives? Well, this is the default behavior. What we have here is something called a histogram, where `ggplot2` helpfully counted up

the number of times the Institution appears and counted them up. Since we only have one record per Institution, the count is always 1. How do we fix this? By adding `weight` to our aesthetic.

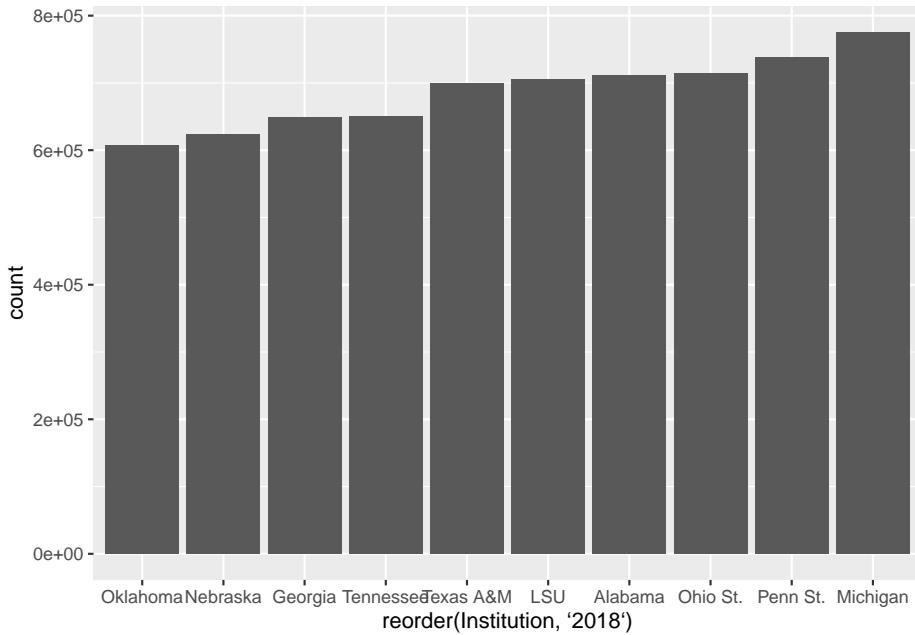
```
ggplot() +
  geom_bar(data=top10, aes(x=Institution, weight=~2018))
```



Closer. But ... what order is that in? And what happened to our count numbers on the left? Why are they in scientific notation?

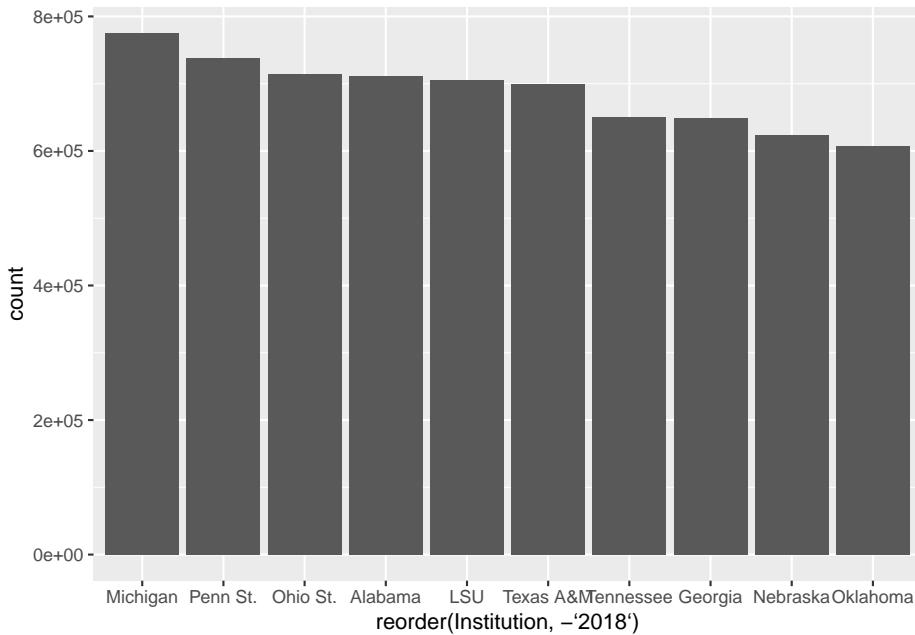
Let's deal with the ordering first. `ggplot2`'s default behavior is to sort the data by the x axis variable. So it's in alphabetical order. To change that, we have to `reorder` it. With `reorder`, we first have to tell `ggplot` what we are reordering, and then we have to tell it HOW we are reordering it. So it's `reorder(FIELD, SORTFIELD)`.

```
ggplot() + geom_bar(data=top10, aes(x=reorder(Institution, ~2018), weight=~2018))
```



Better. We can argue about if the right order is smallest to largest or largest to smallest. But this gets us close. By the way, to sort it largest to smallest, put a negative sign in front of the sort field.

```
ggplot() + geom_bar(data=top10, aes(x=reorder(Institution, -`2018`), weight=`2018`))
```



## 15.2 Scales

To fix the axis labels, we need try one of the other main elements of the `ggplot2` library, which is transform a scale. More often than not, that means doing something like putting it on a logarithmic scale or some other kind of transformation. In this case, we're just changing how it's represented. The default in `ggplot2` for large values is to express them as scientific notation. Rarely ever is that useful in our line of work. So we have to transform them into human readable numbers.

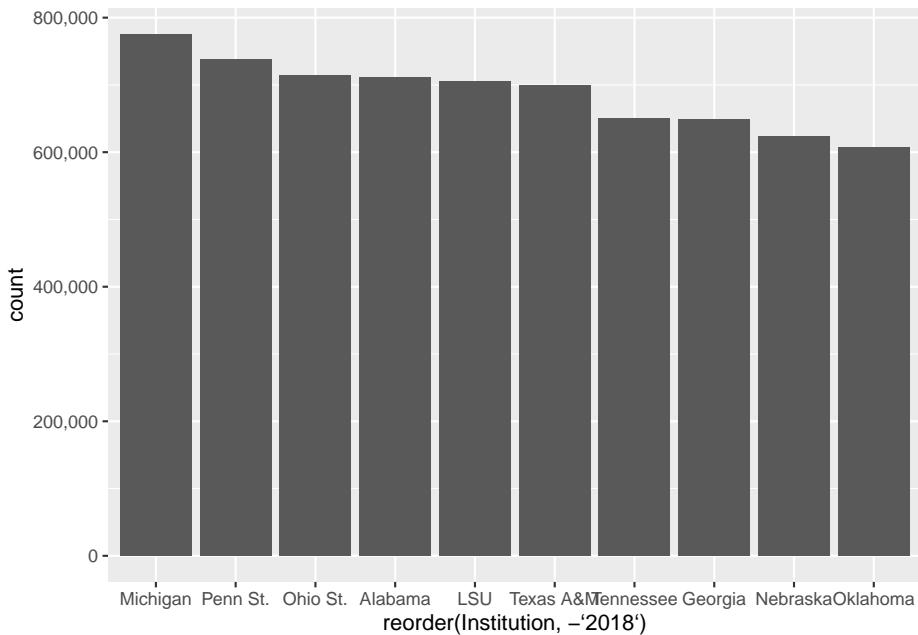
The easiest way to do this is to use a library called `scales` and it's already installed.

```
library(scales)
```

```
##  
## Attaching package: 'scales'  
  
## The following object is masked from 'package:purrr':  
##  
##     discard  
  
## The following object is masked from 'package:readr':  
##  
##     col_factor
```

To alter the scale, we add a piece to our plot with `+` and we tell it which scale is getting altered and what kind of data it is. In our case, our Y axis is what is needing to be altered, and it's continuous data (meaning it can be any number between x and y, vs discrete data which are categorical). So we need to add `scale_y_continuous` and the information we want to pass it is to alter the labels with a function called `comma`.

```
ggplot() +  
  geom_bar(data=top10, aes(x=reorder(Instiution, -`2018`), weight=`2018`)) +  
  scale_y_continuous(labels=comma)
```

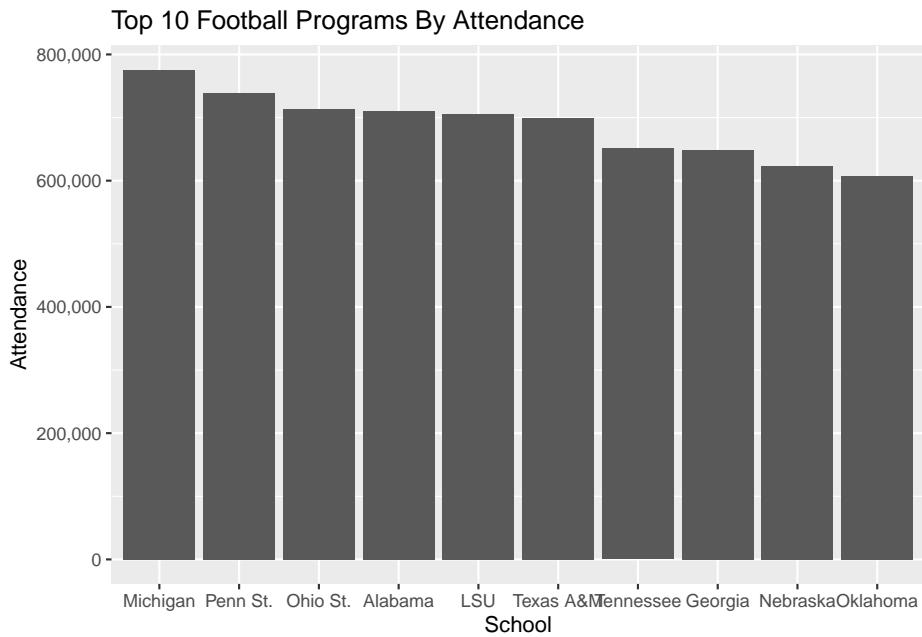


Better.

## 15.3 Styling

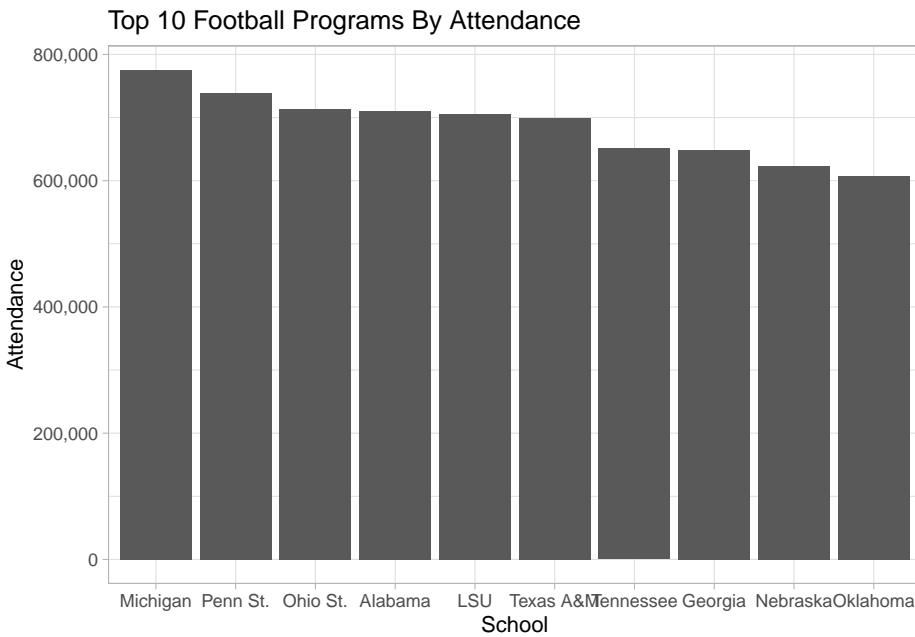
We are going to spend a lot more time on styling, but let's add some simple labels to this with a new bit called `labs` which is short for labels.

```
ggplot() +
  geom_bar(data=top10, aes(x=reorder(Institution, -`2018`), weight=`2018`)) +
  scale_y_continuous(labels=comma) +
  labs(
    title="Top 10 Football Programs By Attendance",
    x="School",
    y="Attendance"
  )
```



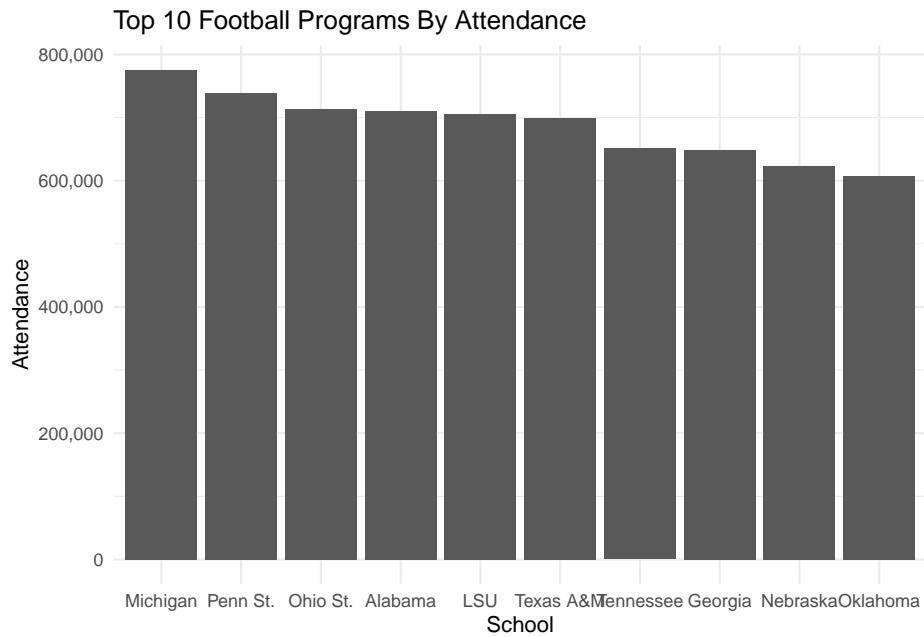
The library has lots and lots of ways to alter the styling – we can programmatically control nearly every part of the look and feel of the chart. One simple way is to apply themes in the library already. We do that the same way we've done other things – we add them. Here's the light theme.

```
ggplot() +
  geom_bar(data=top10, aes(x=reorder(Institution, -`2018`), weight=`2018`)) +
  scale_y_continuous(labels=comma) +
  labs(
    title="Top 10 Football Programs By Attendance",
    x="School",
    y="Attendance") +
  theme_light()
```



Or the minimal theme:

```
ggplot() +  
  geom_bar(data=top10, aes(x=reorder(Institution, -`2018`), weight=`2018`)) +  
  scale_y_continuous(labels=comma) +  
  labs(  
    title="Top 10 Football Programs By Attendance",  
    x="School",  
    y="Attendance") +  
  theme_minimal()
```

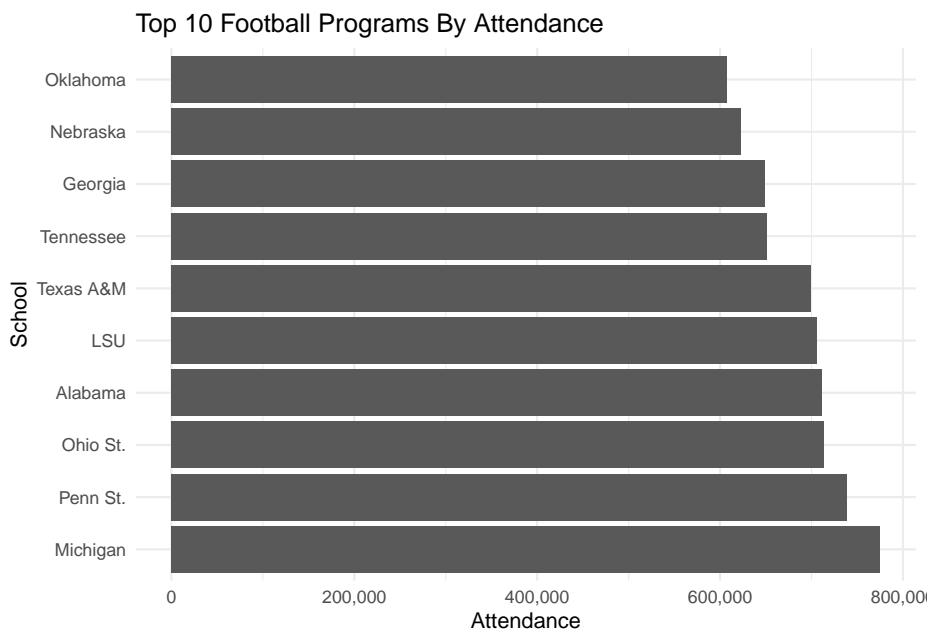


Later on, we'll write our own themes. For now, the built in ones will get us closer to something that looks good.

## 15.4 One last trick: coord flip

Sometimes, we don't want vertical bars. Maybe we think this would look better horizontal. How do we do that? By adding `coord_flip()` to our code. It does what it says – it inverts the coordinates of the figures.

```
ggplot() +
  geom_bar(data=top10, aes(x=reorder(Instiution, -`2018`), weight=`2018`)) +
  scale_y_continuous(labels=comma) +
  labs(
    title="Top 10 Football Programs By Attendance",
    x="School",
    y="Attendance") +
  theme_minimal() +
  coord_flip()
```





# Chapter 16

## Stacked bar charts

One of the elements of data visualization excellence is **inviting comparison**. Often that comes in showing **what proportion a thing is in relation to the whole thing**. With bar charts, we're showing magnitude of the whole thing. If we have information about the parts of the whole, we **can stack them on top of each other to compare them, showing both the whole and the components**. And it's a simple change to what we've already done.

We're going to use a dataset of college football games from this season.

Load the tidyverse.

```
library(tidyverse)
```

And the data.

```
football <- read_csv("data/footballlogs20.csv")
```

```
##  
## -- Column specification -----  
## cols(  
##   .default = col_double(),  
##   Date = col_date(format = ""),  
##   HomeAway = col_character(),  
##   Opponent = col_character(),  
##   Result = col_character(),  
##   TeamFull = col_character(),  
##   TeamURL = col_character(),  
##   Outcome = col_character(),  
##   Team = col_character(),  
##   Conference = col_character()  
## )  
## i Use `spec()` for the full column specifications.
```

What we have here is every game in college football this season. The question we want to answer is this: Who had the most prolific offenses in the Big Ten? And how did they get there?

So to make this chart, we have to just add one thing to a bar chart like we did in the previous chapter. However, it's not that simple.

We have game data, and we need season data. To get that, we need to do some group by and sum work. And since we're only interested in the Big Ten, we have some filtering to do too. For this, we're going to measure offensive production by rushing yards and passing yards. So if we have all the games a team played, and the rushing and passing yards for each of those games, what we need to do to get the season totals is just add them up.

```
football %>%
  group_by(Conference, Team) %>%
  summarise(
    SeasonRushingYards = sum(RushingYds),
    SeasonPassingYards = sum(PassingYds),
    ) %>% filter(Conference == "Big Ten Conference")

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

## # A tibble: 14 x 4
## # Groups:   Conference [1]
##   Conference       Team   SeasonRushingYards   SeasonPassingYards
##   <chr>           <chr>            <dbl>              <dbl>
## 1 Big Ten Conference Illinois          1569               1223
## 2 Big Ten Conference Indiana          726                1806
## 3 Big Ten Conference Iowa            1368               1581
## 4 Big Ten Conference Maryland        722                1320
## 5 Big Ten Conference Michigan        786                1502
## 6 Big Ten Conference Michigan State 635                1672
## 7 Big Ten Conference Minnesota      1343               1394
## 8 Big Ten Conference Nebraska       1611               1521
## 9 Big Ten Conference Northwestern   1299               1490
## 10 Big Ten Conference Ohio State    1654               1521
## 11 Big Ten Conference Penn State    1569               2304
## 12 Big Ten Conference Purdue       489                1854
## 13 Big Ten Conference Rutgers      1259               1786
## 14 Big Ten Conference Wisconsin    1030               1123
```

By looking at this, we can see we got what we needed. We have 14 teams and numbers that look like season totals for yards. Save that to a new dataframe.

```
football %>%
  group_by(Conference, Team) %>%
  summarise(
    SeasonRushingYards = sum(RushingYds),
```

```

SeasonPassingYards = sum(PassingYds),
) %>% filter(Conference == "Big Ten Conference") -> yards

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

```

Now, the problem we have is that ggplot wants long data and this data is wide. So we need to use `tidyverse` to make it long, just like we did in the transforming data chapter.

```

yards %>%
  pivot_longer(
    cols=starts_with("Season"),
    names_to="Type",
    values_to="Yards")

## # A tibble: 28 x 4
## # Groups:   Conference [1]
##   Conference     Team     Type        Yards
##   <chr>       <chr>    <chr>      <dbl>
## 1 Big Ten Conference Illinois SeasonRushingYards 1569
## 2 Big Ten Conference Illinois SeasonPassingYards 1223
## 3 Big Ten Conference Indiana SeasonRushingYards 726
## 4 Big Ten Conference Indiana SeasonPassingYards 1806
## 5 Big Ten Conference Iowa   SeasonRushingYards 1368
## 6 Big Ten Conference Iowa   SeasonPassingYards 1581
## 7 Big Ten Conference Maryland SeasonRushingYards 722
## 8 Big Ten Conference Maryland SeasonPassingYards 1320
## 9 Big Ten Conference Michigan SeasonRushingYards 786
## 10 Big Ten Conference Michigan SeasonPassingYards 1502
## # ... with 18 more rows

```

What you can see now is that we have two rows for each team: One for rushing yards, one for passing yards. This is what ggplot needs. Save it to a new dataframe.

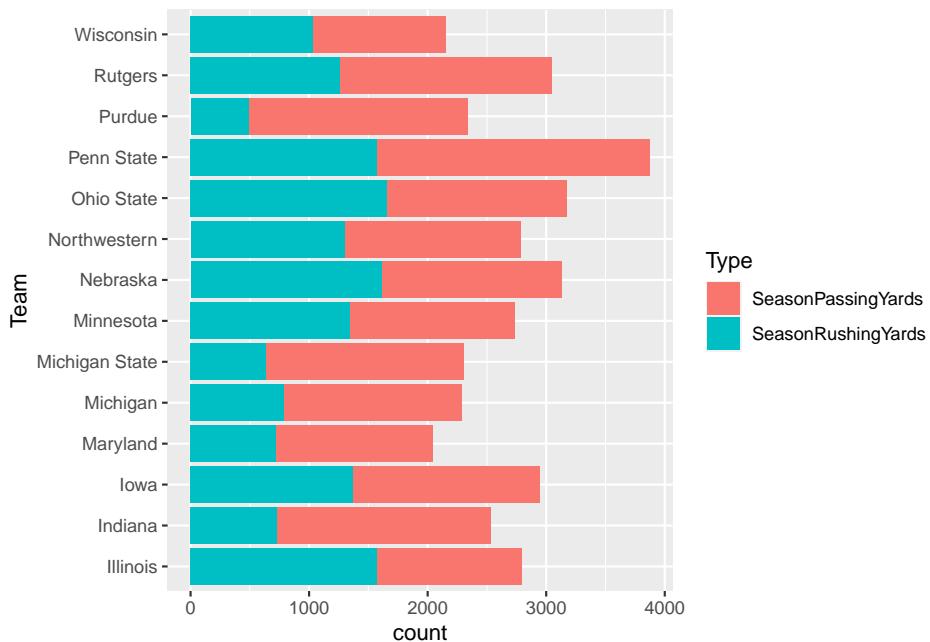
```

yards %>%
  pivot_longer(
    cols=starts_with("Season"),
    names_to="Type",
    values_to="Yards") -> yardswide

```

Building on what we learned in the last chapter, we know we can turn this into a bar chart with an x value, a weight and a `geom_bar`. What we are going to add is a `fill`. The `fill` will stack bars on each other based on which element it is. In this case, we can fill the bar by Type, which means it will stack the number of rushing yards on top of passing yards and we can see how they compare.

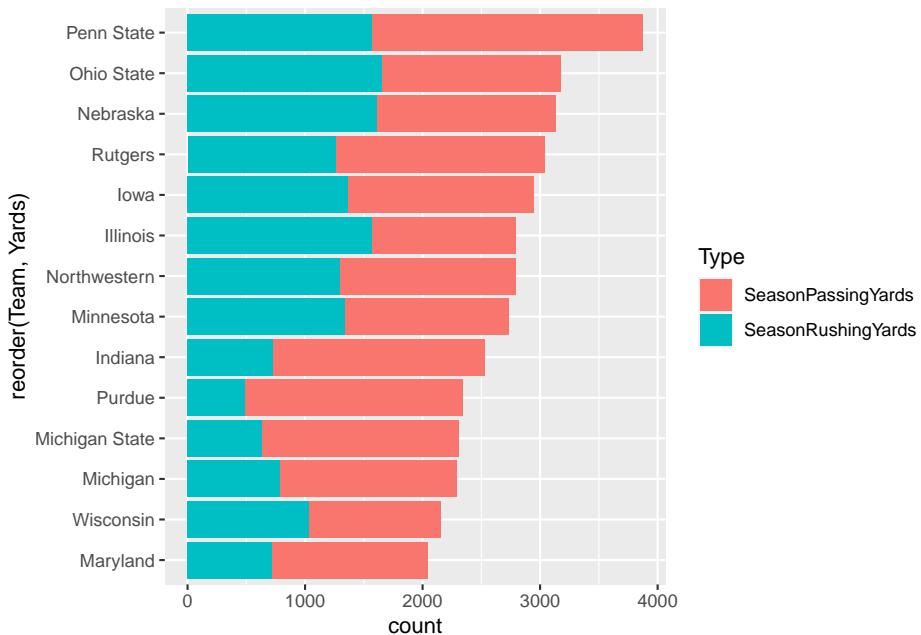
```
ggplot(yardswide, aes(x=Team, weight=Yards, fill=Type)) +
  geom_bar() +
  coord_flip()
```



What's the problem with this chart?

There's a couple of things, one of which we'll deal with now: The ordering is alphabetical (from the bottom up). So let's `reorder` the teams by Yards.

```
ggplot(yardswide, aes(x=reorder(Team, Yards), weight=Yards, fill=Type)) +
  geom_bar() +
  coord_flip()
```



And just like that ... Penn State comes out on top? Huh. And look who is third.

What else is the problem here? Hint: there was a global pandemic going on.



# Chapter 17

## Circular bar plots

At the 27:36 mark in the Half Court Podcast, former Omaha World Herald Writer Chris Heady said “November basketball doesn’t matter, but it shows you where you are.”

It’s a tempting phrase to believe, especially a day after Nebraska lost the first game of the Fred Hoiberg era at home to a baseball school, UC Riverside. And it wasn’t close. The Huskers, because of a total roster turnover, were a complete mystery before the game. And what happened during it wasn’t pretty, so there was a little soul searching going on in Lincoln.

But does November basketball really not matter?

Let’s look, using a new form of chart called a circular bar plot. It’s a chart type that combines several forms we’ve used before: bar charts to show magnitude, stacked bar charts to show proportion, but we’re going to add bending the chart around a circle to add some visual interestingness to it. We’re also going to use time as an x-axis value to make a not subtle circle of time reference – a common technique with circular bar charts.

We’ll use a dataset of every college basketball game last season.

Load your libraries.

```
library(tidyverse)
library(lubridate)
```

And load your data.

```
logs <- read_csv("data/logs20.csv")

## 
## -- Column specification ----- 
## cols(
```

```

##   .default = col_double(),
##   Date = col_date(format = ""),
##   HomeAway = col_character(),
##   Opponent = col_character(),
##   W_L = col_character(),
##   Blank = col_logical(),
##   Team = col_character(),
##   Conference = col_character(),
##   season = col_character()
## )
## i Use `spec()` for the full column specifications.

```

## 17.1 Does November basketball matter?

So let's test the notion of November Basketball Doesn't Matter. What matters in basketball? Let's start simple: Wins.

Sports Reference's win columns are weird, so we need to scan through them and find W and L and we'll give them numbers using `case_when`. I'm also going to filter out tournament basketball.

```

winlosslogs <- logs %>% mutate(winloss = case_when(
  grepl("W", W_L) ~ 1,
  grepl("L", W_L) ~ 0)
)

```

Now we can group by date and conference and sum up the wins. How many wins by day does each conference get?

```

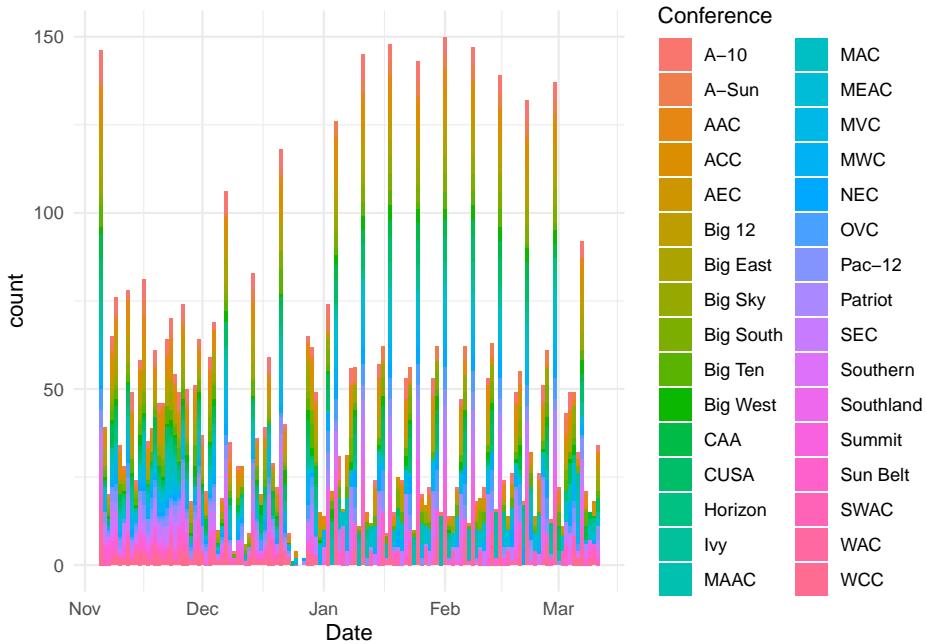
dates <- winlosslogs %>% group_by(Date, Conference) %>% summarise(wins = sum(winloss))

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

```

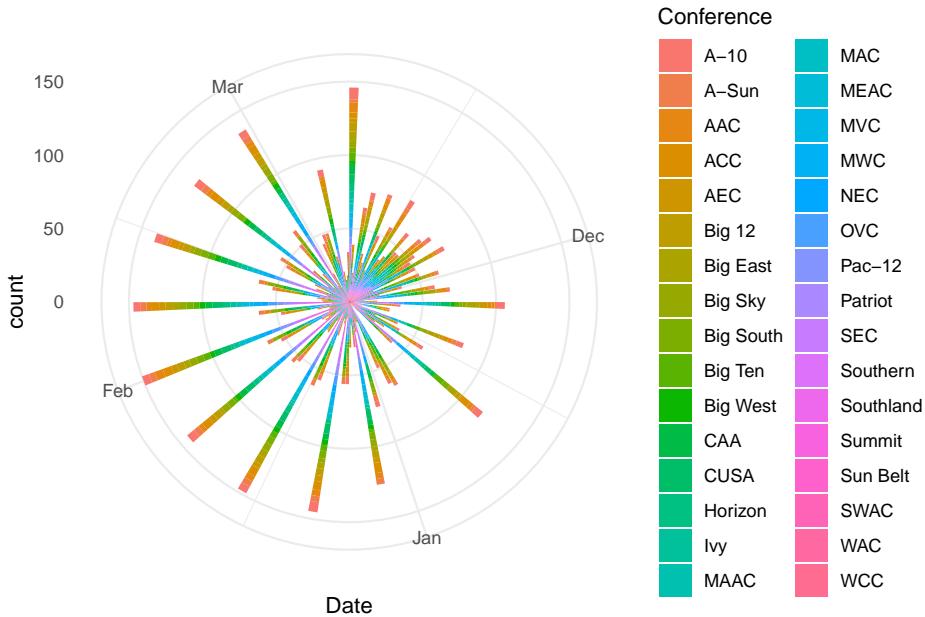
Earlier, we did stacked bar charts. We have what we need to do that now.

```
ggplot() + geom_bar(data=dates, aes(x=Date, weight=wins, fill=Conference)) + theme_minimal()
```



Eek. This is already looking not great. But to make it a circular bar chart, we add `coord_polar()` to our chart.

```
ggplot() + geom_bar(data=dates, aes(x=Date, weight=wins, fill=Conference)) + theme_minimal() + co
```



Based on that, the day is probably too thin a slice, and there's way too many

conferences in college basketball. Let's group this by months and filter out all but the power five conferences.

```
p5 <- c("SEC", "Big Ten", "Pac-12", "Big 12", "ACC")
```

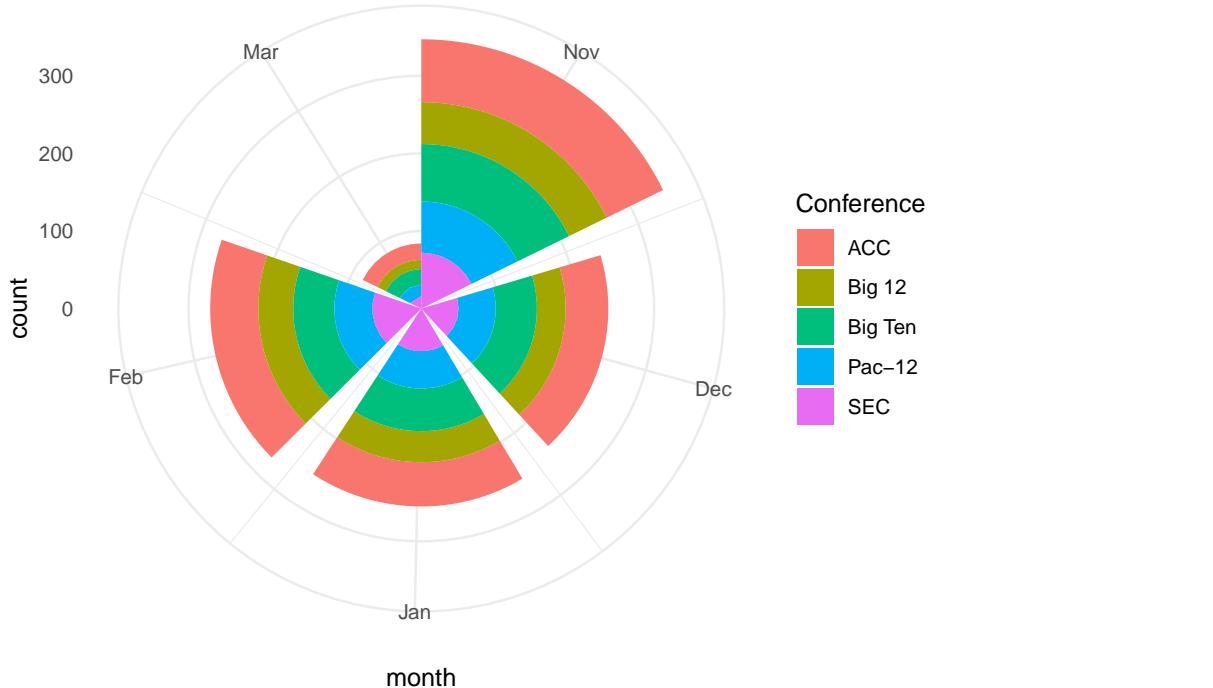
To get months, we're going to use a function in the library `lubridate` called `floor_date`, which combined with `mutate` will give us a field of just months.

```
wins <- winlosslogs %>% mutate(month = floor_date(Date, unit="months")) %>% group_by(mo
```

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

Now we can use `wins` to make our circular bar chart of wins by month in the Power Five.

```
ggplot() + geom_bar(data=wins, aes(x=month, weight=wins, fill=Conference)) + theme_min
```

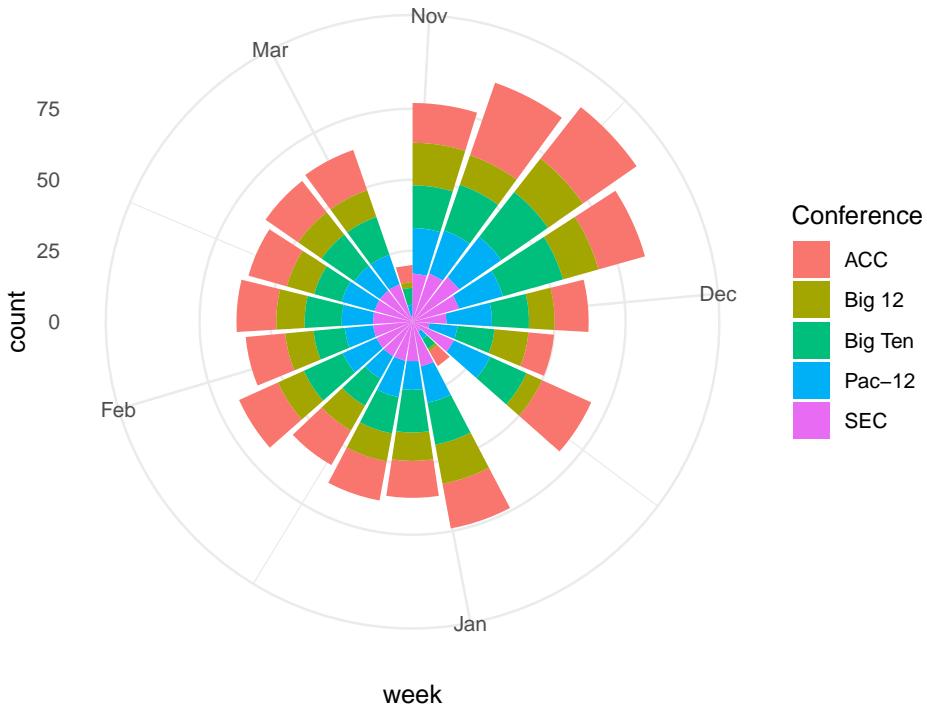


Yikes. That looks a lot like a broken pie chart. So months are too thick of a slice. Let's use weeks in our floor date to see what that gives us.

```
wins <- winlosslogs %>% mutate(week = floor_date(Date, unit="weeks")) %>% group_by(wee
```

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

```
ggplot() + geom_bar(data=wins, aes(x=week, weight=wins, fill=Conference)) + theme_min
```



That looks better. But what does it say? Does November basketball matter? What this is saying is ... yeah, it kinda does. The reason? Lots of wins get piled up in November and December, during non-conference play. So if you are a team with NCAA tournament dreams, you need to win games in November to make sure your tournament resume is where it needs to be come March. Does an individual win or loss matter? Probably not. But your record in November does.

## 17.2 Does it show you where you are?

So here is the problem we have:

1. We have data for every game. In the past, we were able to calculate the team wins and losses because the way the data records them is the Team is the main team, and they win or lose. The Opponent is recorded, but the Opponent has the mirror image of this game as well, where they are the Team. So essentially every game is in here twice – one for each team that plays in the game.
2. We need to attach the Opponent's winning percentage to each game so we can decide if it's a quality win for Team.
3. The Team name is not an exact copy of the Team name. So we can't join them using it.

So what we have to do is invert the process that we've done before. We need to group by the Opponent – because the names will be consistent then – and we need to invert the wins and losses. A win in the W\_L column is a win for the Team. That means each loss in the W\_L column is a WIN for the Opponent.

Once we invert, the data looks very similar to what we've done before. One other thing: I noticed there's some tournament games in here, so the filter at the end strips them out.

```
oppwinlosslogs <- logs %>% mutate(winloss = case_when(
  grepl("W", W_L) ~ 0,
  grepl("L", W_L) ~ 1)
) %>% filter(Date < "2020-03-19")
```

So now we have a dataframe called oppwinlosslogs that has an inverted winloss column. So now we can group by the Opponent and sum the wins and it will tell us how many games the Opponent won. We can also count the wins and get a winning percentage.

```
oppwinlosslogs %>% group_by(Opponent) %>% summarise(games=n(), wins=sum(winloss)) %>%
```

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

Now we have a dataframe of 659 opponent winning records. Wait, what? There's 353 teams in major college basketball, so why 659? If you look through it, there's a bunch of teams playing lower level teams. Given that they are lower level, they're likely cannon fodder and will lose the game, and we're going to filter them out in a minute.

Now we can join the opponent winning percentage to our winlosslogs data so we can answer our question about quality wins.

```
winlosslogs <- logs %>% mutate(winloss = case_when(
  grepl("W", W_L) ~ 1,
  grepl("L", W_L) ~ 0)
) %>% filter(Date < "2020-03-19")

winlosslogs %>% left_join(opprecord, by=("Opponent")) -> winswithopppct
```

Now that we have a table called winswithopppct, we can filter out teams non power 5 teams and teams that won less than 60 percent of their games and run the same calculations in the book.

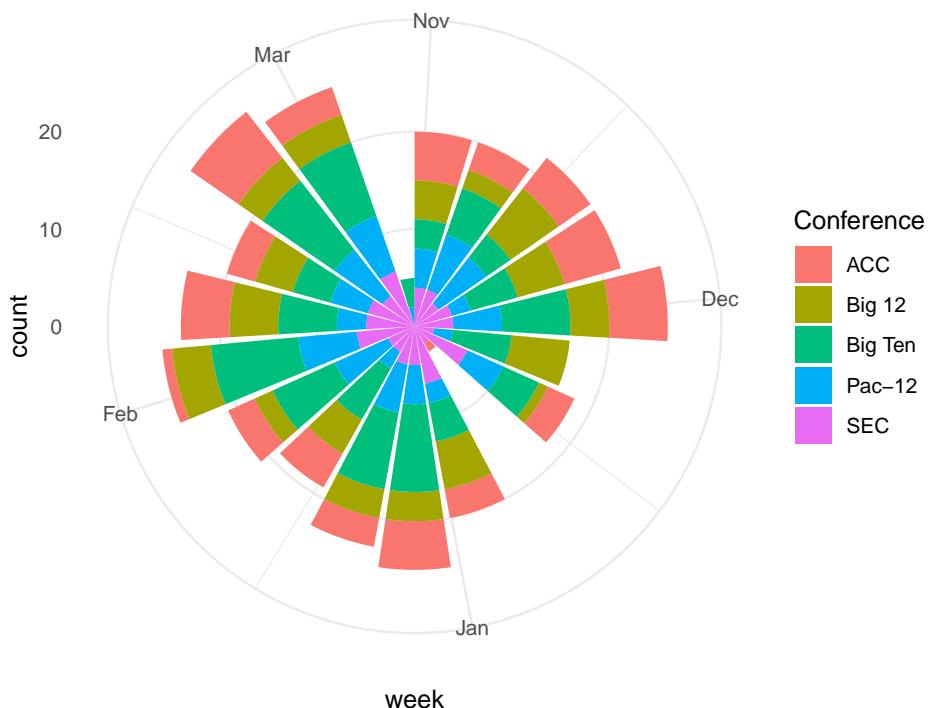
```
p5 <- c("SEC", "Big Ten", "Pac-12", "Big 12", "ACC")

winswithopppct %>% filter(winpct > .6) %>% mutate(week = floor_date(Date, unit="weeks"))

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

Now with our dataframe called qualitywins, we can chart it again.

```
ggplot() + geom_bar(data=qualitywins, aes(x=week, weight=wins, fill=Conference)) + theme_minimal()
```



Look at this chart and compare it to the first one.



# Chapter 18

## Waffle charts

Pie charts are the devil. They should be an instant F in any data visualization class. The problem? How carefully can you evaluate angles and area? Unless they are blindingly obvious and only a few categories, not well. If you've got 25 categories, how can you tell the difference between 7 and 9 percent? You can't.

So let's introduce a better way: The Waffle Chart. Some call it a square pie chart. I personally hate that. Waffles it is.

**A waffle chart is designed to show you parts of the whole – proportionality.** How many yards on offense come from rushing or passing. How many singles, doubles, triples and home runs make up a teams hits. How many shots a basketball team takes are two pointers versus three pointers.

First, install the library in the console. We want a newer version of the `waffle` library than is in CRAN – where you normally get libraries from – so copy and paste this into your console:

```
install.packages("waffle")
```

Now load it:

```
library(waffle)
```

### 18.1 Waffles two ways: Part 1

Let's look at the debacle that was Nebraska vs. Ohio State this past fall in college football. Here's the box score, which we'll use for this part of the walkthrough.

Maybe the easiest way to do waffle charts, at least at first, is to make vectors of your data and plug them in. To make a vector, we use the `c` or concatenate function.

So let's look at offense. Rushing vs passing.

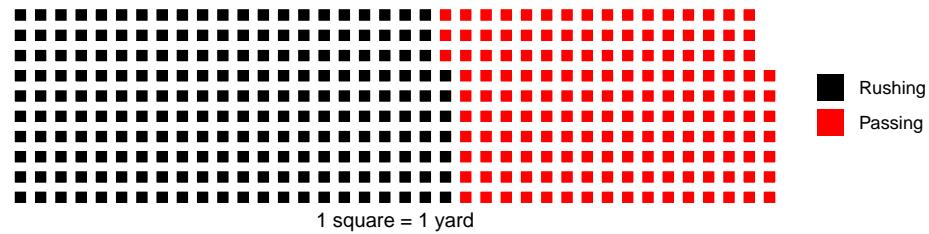
```
nu <- c("Rushing"=217, "Passing"=160)
oh <- c("Rushing"=215, "Passing"=276)
```

So what does the breakdown of the night look like?

The waffle library can break this down in a way that's easier on the eyes than a pie chart. We call the library, add the data, specify the number of rows, give it a title and an x value label, and to clean up a quirk of the library, we've got to specify colors.

```
waffle(
  nu,
  rows = 10,
  title="Nebraska's offense",
  xlab="1 square = 1 yard",
  colors = c("black", "red")
)
```

Nebraska's offense

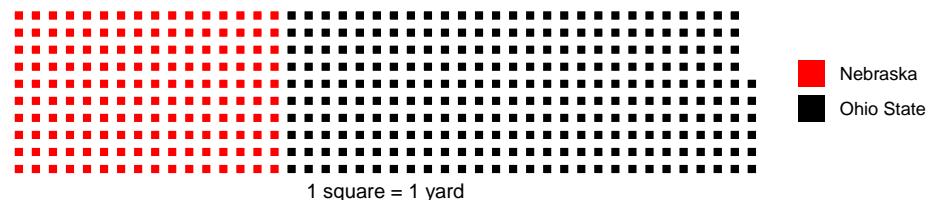


Or, we could make this two teams in the same chart.

```
passing <- c("Nebraska"=160, "Ohio State"=276)

waffle(
  passing,
  rows = 10,
  title="Nebraska vs Ohio State: passing",
  xlab="1 square = 1 yard",
  colors = c("red", "black")
)
```

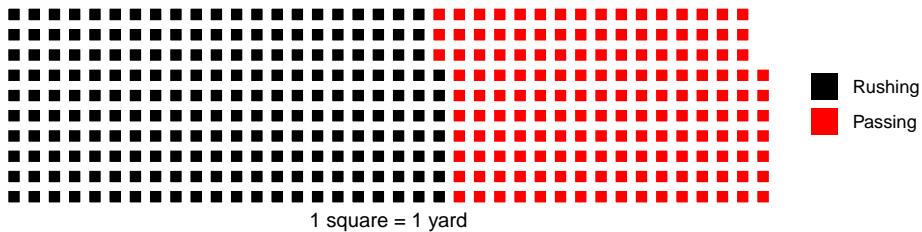
Nebraska vs Ohio State: passing



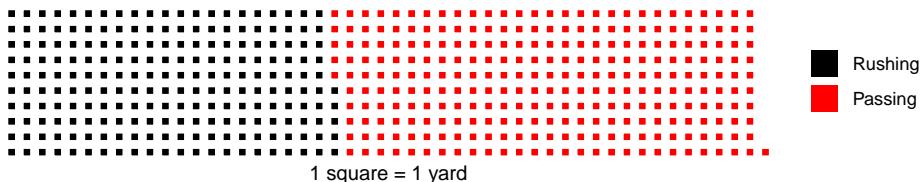
So what does it look like if we compare the two teams using the two vectors in the same chart? To do that – and I am not making this up – you have to create a waffle iron. Get it? Waffle charts? Iron?

```
iron(
  waffle(nu,
    rows = 10,
    title="Nebraska's offense",
    xlab="1 square = 1 yard",
    colors = c("black", "red")
  ),
  waffle(oh,
    rows = 10,
    title="Ohio State's offense",
    xlab="1 square = 1 yard",
    colors = c("black", "red")
  )
)
```

**Nebraska's offense**



**Ohio State's offense**



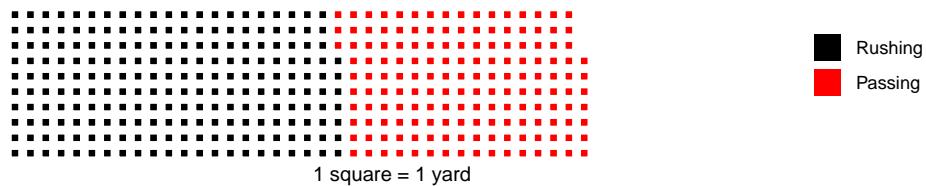
What do you notice about this chart? Notice how the squares aren't the same size? Well, Ohio State out-gained Nebraska by a long way. So the squares aren't the same size because the numbers aren't the same. We can fix that by adding an unnamed padding number so the number of yards add up to the same thing. Let's make the total for everyone be 491, Ohio State's total yards of offense. So to do that, we need to add a padding of 114 to Nebraska. REMEMBER: Don't name it or it'll show up in the legend.

```
nu <- c("Rushing"=217, "Passing"=160, 114)
oh <- c("Rushing"=215, "Passing"=276)
```

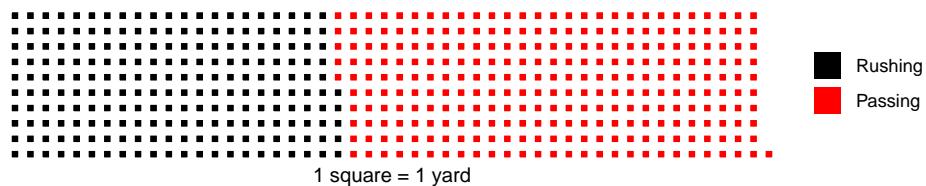
Now, in our waffle iron, if we don't give that padding a color, we'll get an error. So we need to make it white. Which, given our white background, means it will disappear.

```
iron(
  waffle(nu,
    rows = 10,
    title="Nebraska's offense",
    xlab="1 square = 1 yard",
    colors = c("black", "red", "white")
  ),
  waffle(oh,
    rows = 10,
    title="Ohio State's offense",
    xlab="1 square = 1 yard",
    colors = c("black", "red", "white")
  )
)
```

**Nebraska's offense**



**Ohio State's offense**



One last thing we can do is change the 1 square = 1 yard bit – which makes the squares really small in this case – by dividing our vector. Remember what you learned in Swirl about math on vectors?

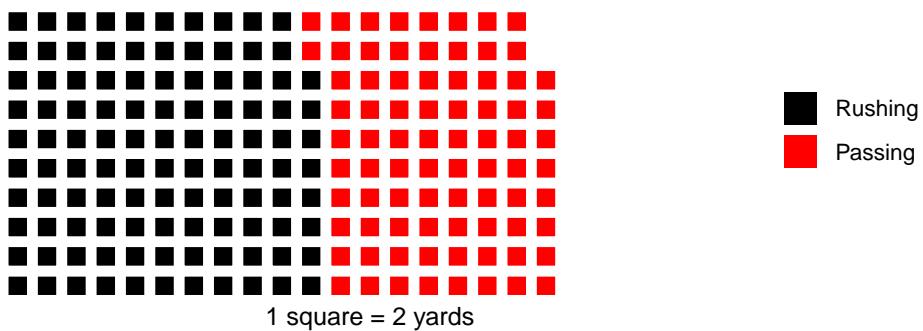
```
iron(
  waffle(nu/2,
    rows = 10,
    title="Nebraska's offense",
    xlab="1 square = 2 yards",
    colors = c("black", "red", "white")
  ),
  waffle(oh/2,
```

```

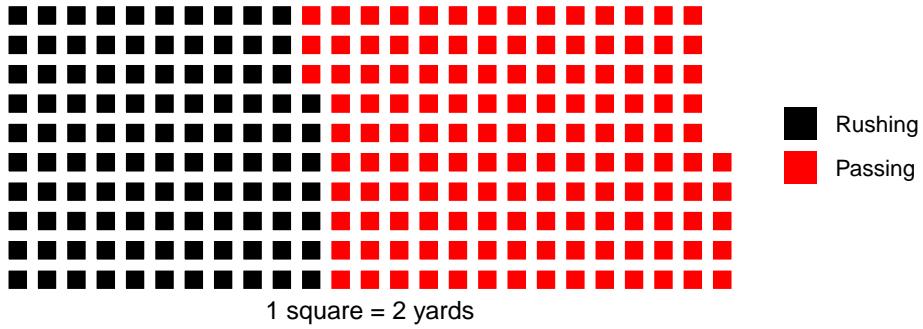
rows = 10,
title="Ohio State's offense",
xlab="1 square = 2 yards",
colors = c("black", "red", "white")
)
)

```

### Nebraska's offense



### Ohio State's offense



News flash: Ohio State beat Nebraska.

## 18.2 Waffles two ways: Part 2

For this part, we want a newer version of the `waffle` library than is in CRAN – where you normally get libraries from.

**WARNING:** This didn't work in a variety of environments, so it may not work on yours.

Copy and paste this into your console:

```
install.packages("waffle", repos = "https://cinc.rud.is")
```

At first, this way might seem harder than doing it the way we just walked through, but the benefits will come later when it's far, far easier to style this chart, where the previous charts are harder.

We have the log of every game in college football – you can get it here – and we can find this game with some simple filtering.

```
fblogs <- read_csv("data/footballlogs19.csv")

##
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Date = col_date(format = ""),
##   HomeAway = col_character(),
##   Opponent = col_character(),
##   Result = col_character(),
##   TeamFull = col_character(),
##   TeamURL = col_character(),
##   Outcome = col_character(),
##   Team = col_character(),
##   Conference = col_character()
## )
## i Use `spec()` for the full column specifications.

fblogs %>% filter(Team == "Nebraska" & Opponent == "Ohio State")

## # A tibble: 1 x 54
##   Game Date      HomeAway Opponent Result PassingCmp PassingAtt PassingPct
##   <dbl> <date>    <chr>     <chr>     <chr>     <dbl>      <dbl>      <dbl>
## 1      5 2019-09-28 <NA>       Ohio St~ L (7~        8          17        47.1
## # ... with 46 more variables: PassingYds <dbl>, PassingTD <dbl>,
## #   RushingAtt <dbl>, RushingYds <dbl>, RushingAvg <dbl>, RushingTD <dbl>,
## #   OffensivePlays <dbl>, OffensiveYards <dbl>, OffenseAvg <dbl>,
## #   FirstDownPass <dbl>, FirstDownRush <dbl>, FirstDownPen <dbl>,
## #   FirstDownTotal <dbl>, Penalties <dbl>, PenaltyYds <dbl>, Fumbles <dbl>,
## #   Interceptions <dbl>, TotalTurnovers <dbl>, TeamFull <chr>, TeamURL <chr>,
## #   Outcome <chr>, TeamScore <dbl>, OpponentScore <dbl>, DefPassingCmp <dbl>,
## #   DefPassingAtt <dbl>, DefPassingPct <dbl>, DefPassingYds <dbl>,
## #   DefPassingTD <dbl>, DefRushingAtt <dbl>, DefRushingYds <dbl>,
## #   DefRushingAvg <dbl>, DefRushingTD <dbl>, DefPlays <dbl>, DefYards <dbl>,
## #   DefAvg <dbl>, DefFirstDownPass <dbl>, DefFirstDownRush <dbl>,
## #   DefFirstDownPen <dbl>, DefFirstDownTotal <dbl>, DefPenalties <dbl>,
## #   DefPenaltyYds <dbl>, DefFumbles <dbl>, DefInterceptions <dbl>,
## #   DefTotalTurnovers <dbl>, Team <chr>, Conference <chr>
```

That's the game. So now we need to make this long data – same as we did with the stacked bar charts – and we'll focus on total yards.

```
fblogs %>%
  filter(Team == "Nebraska" & Opponent == "Ohio State") %>%
  select(Team, OffensiveYards, DefYards) %>%
  pivot_longer(
    cols=c("OffensiveYards", "DefYards"),
    names_to="Type",
    values_to="Yards"
  )

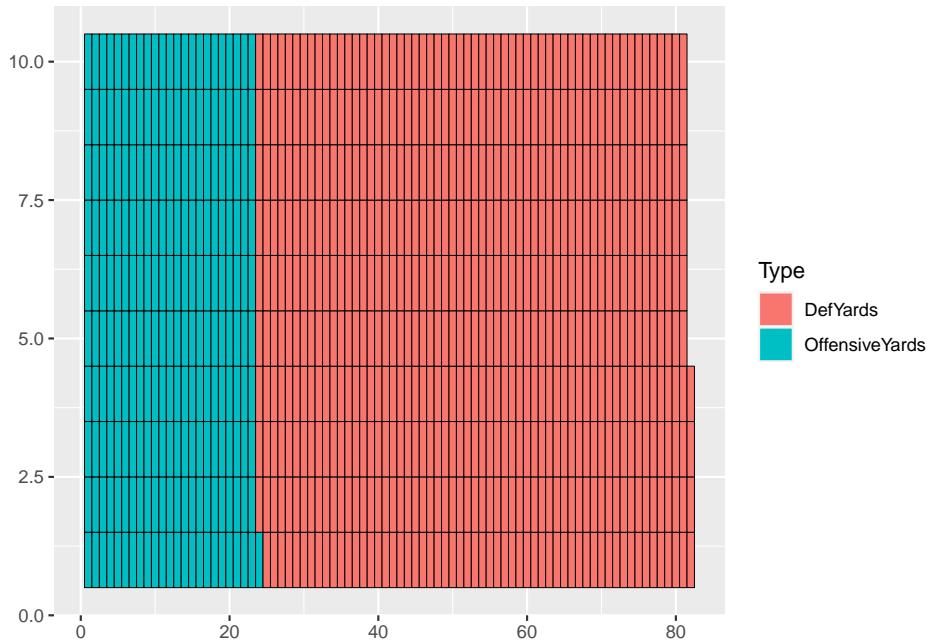
## # A tibble: 2 x 3
##   Team     Type       Yards
##   <chr>   <chr>     <dbl>
## 1 Nebraska OffensiveYards  231
## 2 Nebraska DefYards      583
```

That does what we want, so let's save that to a new dataframe.

```
nuoh <- fbLogs %>%
  filter(Team == "Nebraska" & Opponent == "Ohio State") %>%
  select(Team, OffensiveYards, DefYards) %>%
  pivot_longer(
    cols=c("OffensiveYards", "DefYards"),
    names_to="Type",
    values_to="Yards"
  )
```

Now we can use a new geom – `geom_waffle` – that the `waffle` library has added to `ggplot`. The `geom_waffle` takes two required inputs: `fill` and `value`, but otherwise, it looks the same as previous things we've done.

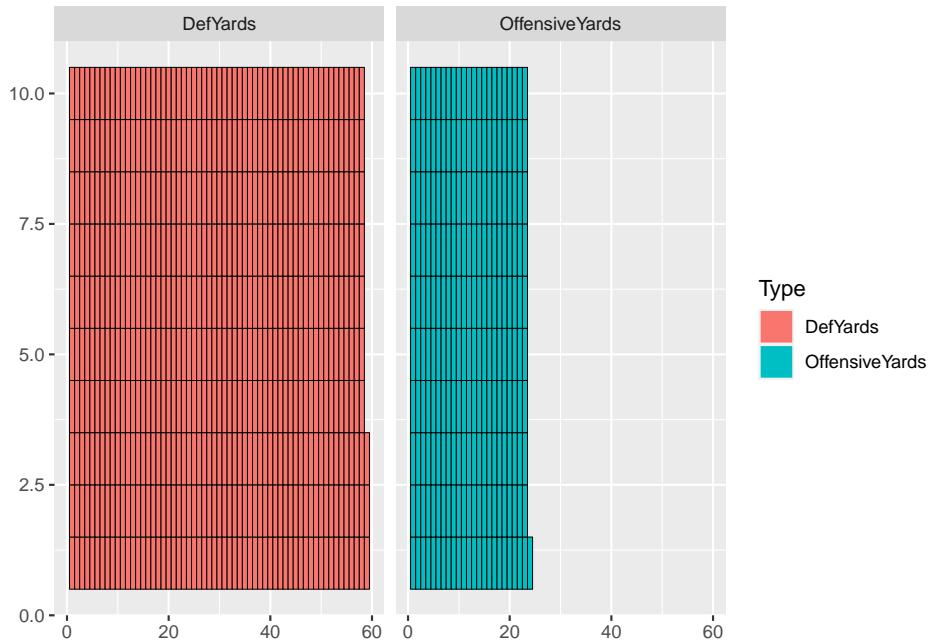
```
ggplot() + geom_waffle(
  data=nuoh,
  aes(fill=Type, values=Yards))
```



First, we can see that going this route changes the boxes to narrow rectangles. That's to fill the space given.

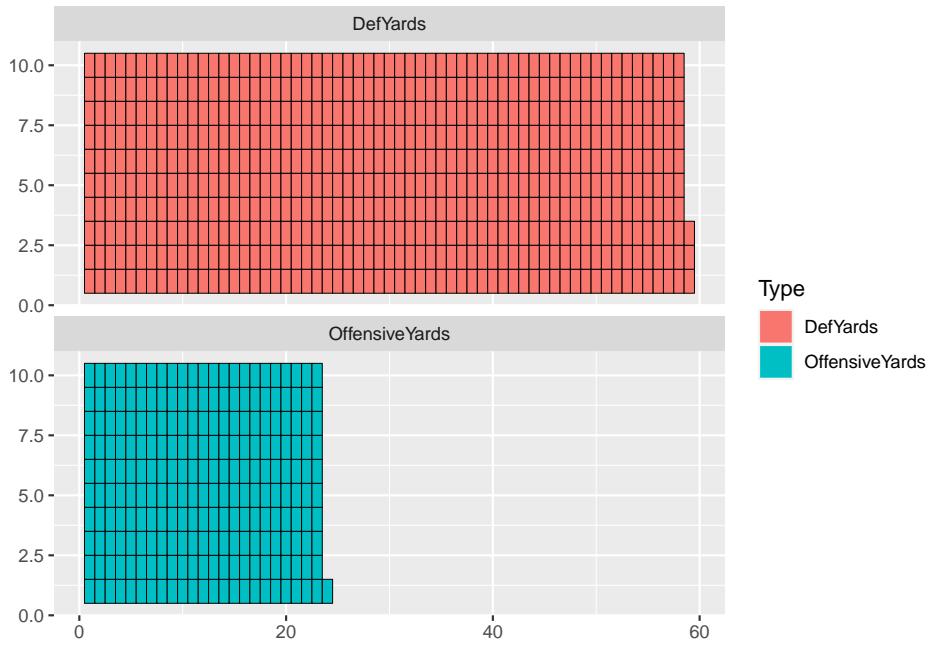
If we want to split them, we can use something called a `facet_wrap` which we will spend a whole class on later, so don't worry about this now. Just know we can split it by Type this way.

```
ggplot() + geom_waffle(
  data=nuoh,
  aes(fill=Type, values=Yards)
) + facet_wrap(~Type)
```



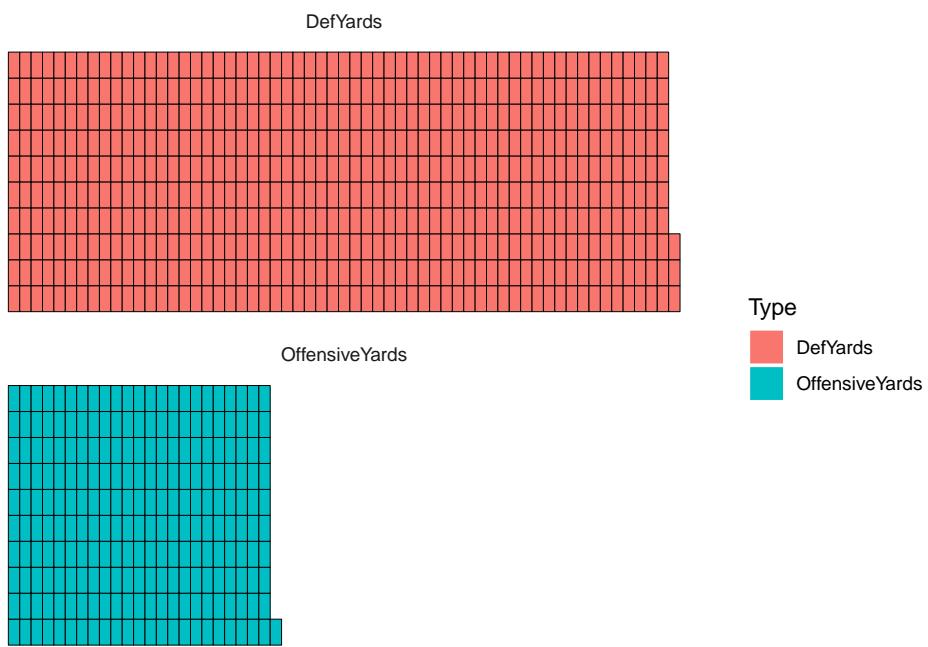
Now we can stack the two charts so they aren't side by side using `ncol` inside the `facet_wrap`.

```
ggplot() + geom_waffle(  
  data=nuoh,  
  aes(fill=Type, values=Yards)  
) + facet_wrap(~Type, ncol=1)
```



Now it's just a matter of formatting, labeling and general cleanup. We'll focus on that later as well, but here's a quick way to get started, which we did in the bar chart chapter.

```
ggplot() + geom_waffle(
  data=nuoh,
  aes(fill=Type, values=Yards)
) +
  facet_wrap(~Type, ncol=1) +
  theme_minimal() +
  theme_enhance_waffle()
```





# Chapter 19

## Line charts

So far, we've talked about bar charts – stacked or otherwise – are good for showing relative size of a thing compared to another thing. Stacked Bars and Waffle charts are good at showing proportions of a whole.

**Line charts are good for showing change over time.**

Let's look at how we can answer this question: Why was Nebraska terrible at basketball last season?

We'll need the logs of every game in college basketball for this.

Let's start getting all that we need. We can use the tidyverse shortcut.

```
library(tidyverse)
```

And now load the data.

```
logs <- read_csv("data/logs20.csv")  
  
##  
## -- Column specification -----  
## cols(  
##   .default = col_double(),  
##   Date = col_date(format = ""),  
##   HomeAway = col_character(),  
##   Opponent = col_character(),  
##   W_L = col_character(),  
##   Blank = col_logical(),  
##   Team = col_character(),  
##   Conference = col_character(),  
##   season = col_character()  
## )  
## i Use `spec()` for the full column specifications.
```

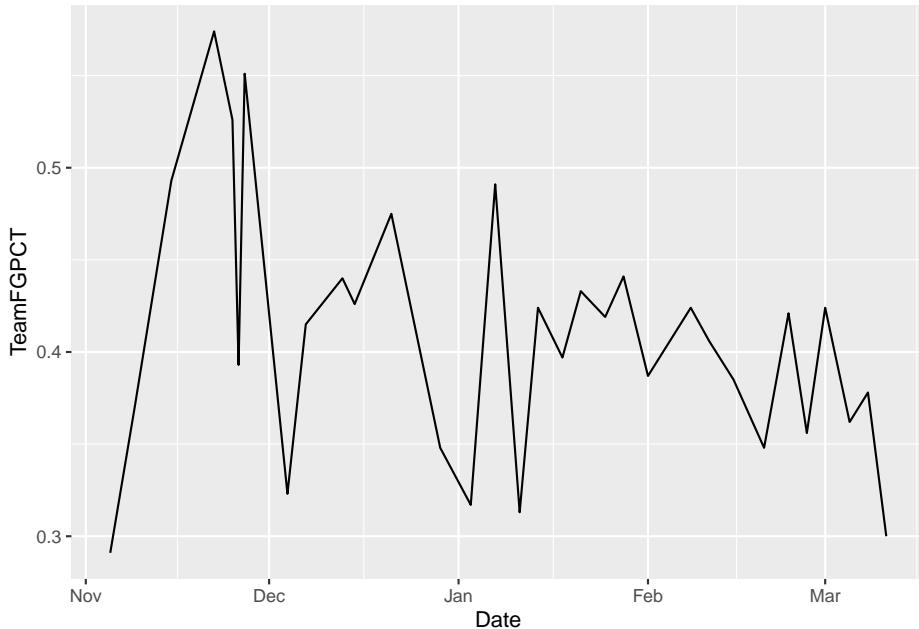
This data has every game from every team in it, so we need to use filtering to limit it, because we just want to look at Nebraska. If you don't remember, flip back to chapter 6.

```
nu <- logs %>% filter(Team == "Nebraska Cornhuskers")
```

Because this data has just Nebraska data in it, the dates are formatted correctly, and the data is long data (instead of wide), we have what we need to make line charts.

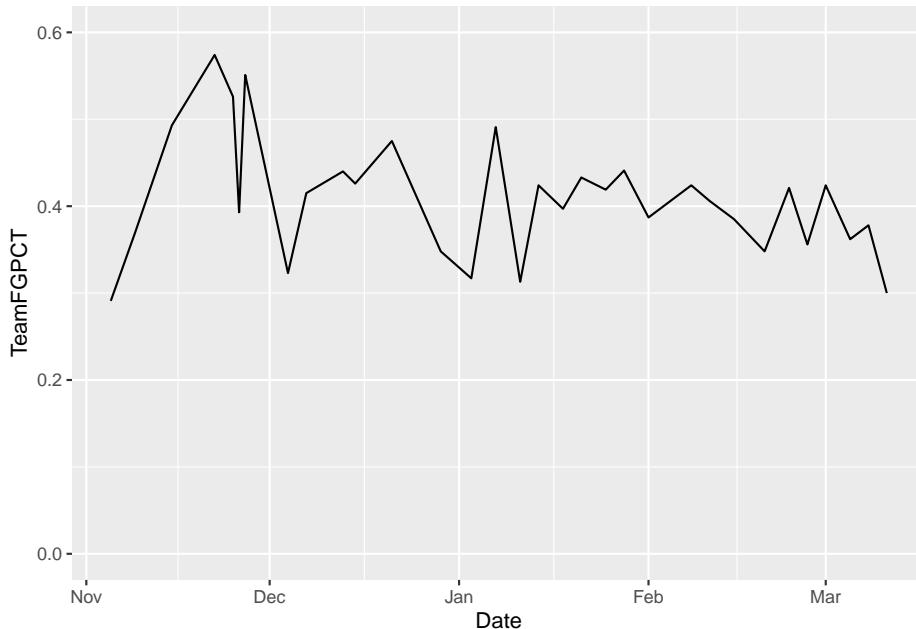
Line charts, unlike bar charts, do have a y-axis. So in our ggplot step, we have to define what our x and y axes are. In this case, the x axis is our Date – the most common x axis in line charts is going to be a date of some variety – and y in this case is up to us. We've seen from previous walkthroughs that how well a team shoots the ball has a lot to do with how well a team does in a season, so let's chart that.

```
ggplot() + geom_line(data=nu, aes(x=Date, y=TeamFGPCT))
```



See a problem here? Note the Y axis doesn't start with zero. That makes this look worse than it is (and that February swoon is pretty bad). To make the axis what you want, you can use `scale_x_continuous` or `scale_y_continuous` and pass in a list with the bottom and top value you want. You do that like this:

```
ggplot() +
  geom_line(data=nu, aes(x=Date, y=TeamFGPCT)) +
  scale_y_continuous(limits = c(0, .6))
```



Note also that our X axis labels are automated. It knows it's a date and it just labels it by month.

## 19.1 This is too simple.

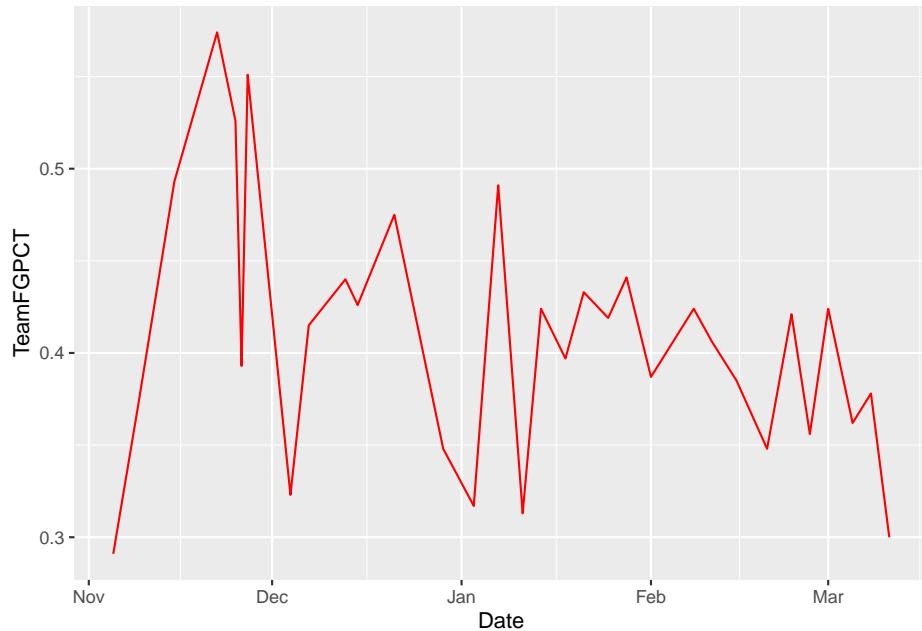
With datasets, we want to invite comparison. So let's answer the question visually. Let's put two lines on the same chart. How does Nebraska compare to Michigan State and Purdue, the eventual regular season co-champions?

```
msu <- logs %>% filter(Team == "Michigan State Spartans")
```

In this case, because we have two different datasets, we're going to put everything in the geom instead of the ggplot step. We also have to explicitly state what dataset we're using by saying `data=` in the geom step.

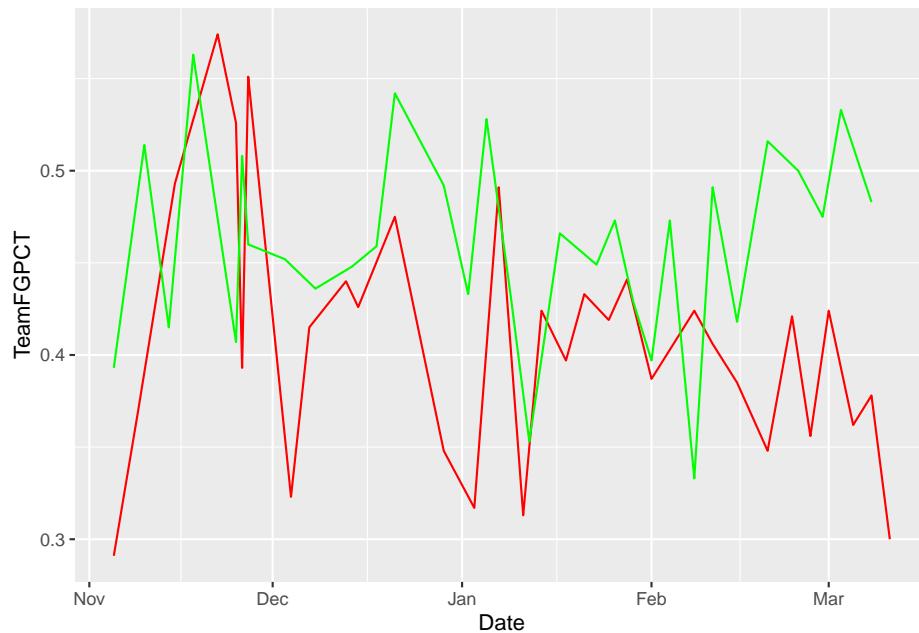
First, let's chart Nebraska. Read carefully. First we set the data. Then we set our aesthetic. Unlike bars, we need an X and a Y variable. In this case, our X is the date of the game, Y is the thing we want the lines to move with. In this case, the Team Field Goal Percentage – TeamFGPCT.

```
ggplot() + geom_line(data=nu, aes(x=Date, y=TeamFGPCT), color="red")
```



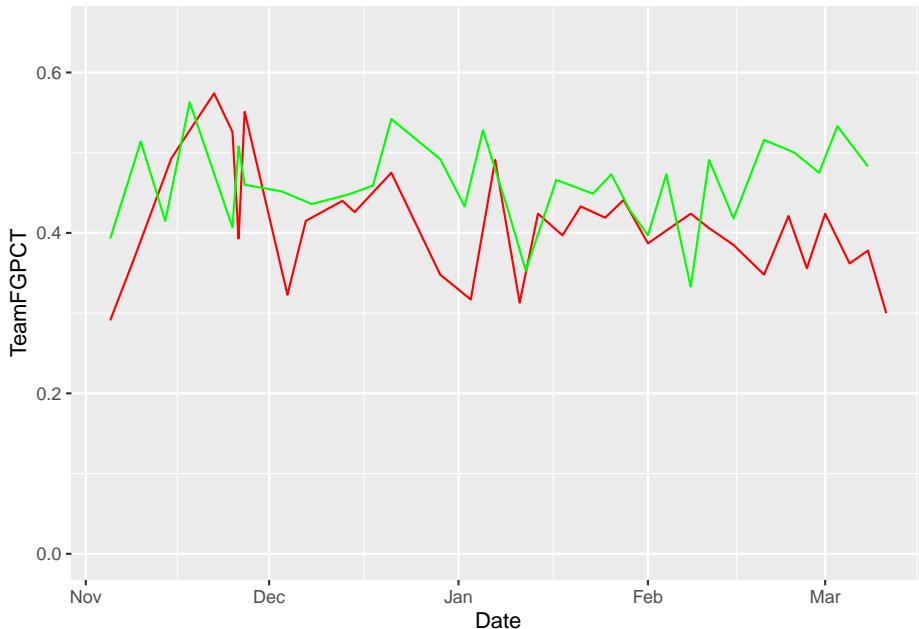
Now, by using `+`, we can add Michigan State to it. REMEMBER COPY AND PASTE IS A THING. Nothing changes except what data you are using.

```
ggplot() +  
  geom_line(data=nu, aes(x=Date, y=TeamFGPCT), color="red") +  
  geom_line(data=msu, aes(x=Date, y=TeamFGPCT), color="green")
```



Let's flatten our lines out by zeroing the Y axis.

```
ggplot() +  
  geom_line(data=nu, aes(x=Date, y=TeamFGPCT), color="red") +  
  geom_line(data=msu, aes(x=Date, y=TeamFGPCT), color="green") +  
  scale_y_continuous(limits = c(0, .65))
```



So visually speaking, the difference between Nebraska and Michigan State's season is that Michigan State stayed mostly on an even keel, and Nebraska went on a two month swoon.

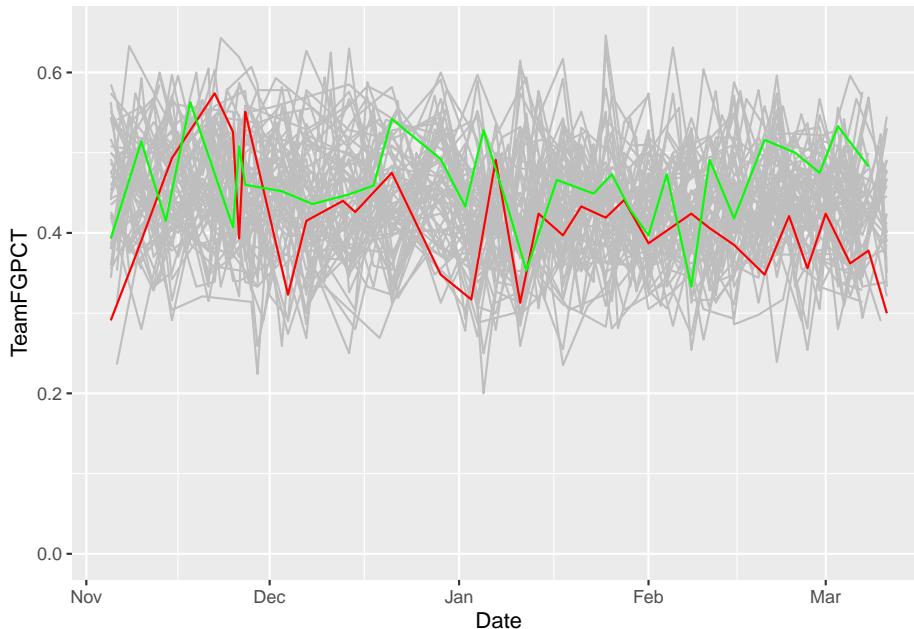
## 19.2 But what if I wanted to add a lot of lines.

Fine. How about all Power Five Schools? This data for example purposes. You don't have to do it.

```
powerfive <- c("SEC", "Big Ten", "Pac-12", "Big 12", "ACC")
p5conf <- logs %>% filter(Conference %in% powerfive)
```

I can keep layering on layers all day if I want. And if my dataset has more than one team in it, I need to use the `group` command. And, the layering comes in order – so if you're going to layer a bunch of lines with a smaller group of lines, you want the bunch on the bottom. So to do that, your code stacks from the bottom. The first geom in the code gets rendered first. The second gets layered on top of that. The third gets layered on that and so on.

```
ggplot() +
  geom_line(data=p5conf, aes(x=Date, y=TeamFGPCT, group=Team), color="grey") +
  geom_line(data=nu, aes(x=Date, y=TeamFGPCT), color="red") +
  geom_line(data=msu, aes(x=Date, y=TeamFGPCT), color="green") +
  scale_y_continuous(limits = c(0, .65))
```

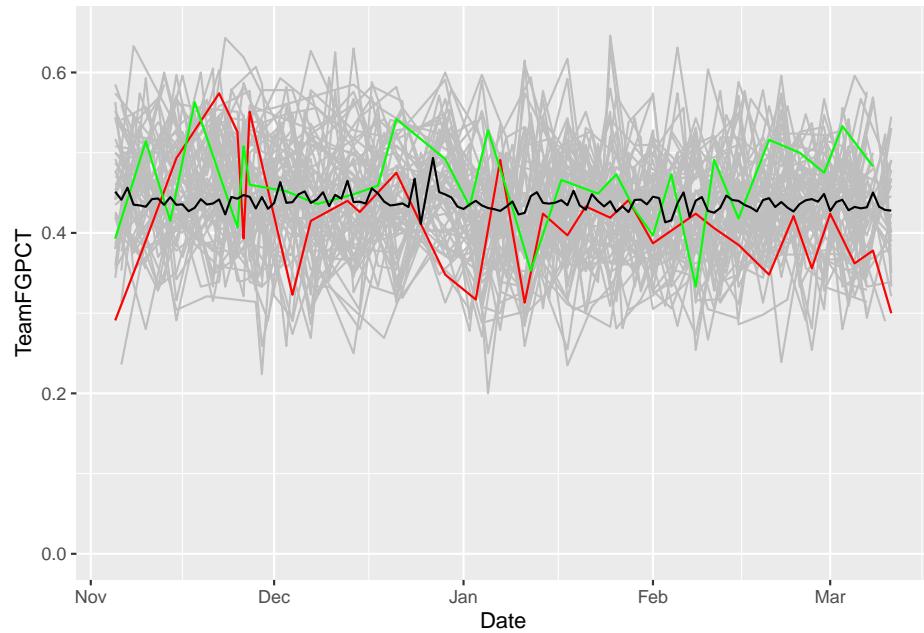


What do we see here? How has Nebraska and Michigan State's season evolved against all the rest of the teams in college basketball?

But how does that compare to the average? We can add that pretty easily by creating a new dataframe with it and add another geom\_line.

```
average <- logs %>% group_by(Date) %>% summarise(mean_shooting=mean(TeamFGPCT))
```

```
## `summarise()` ungrouping output (override with `.`groups` argument)
ggplot() +
  geom_line(data=p5conf, aes(x=Date, y=TeamFGPCT, group=Team), color="grey") +
  geom_line(data=nu, aes(x=Date, y=TeamFGPCT), color="red") +
  geom_line(data=msu, aes(x=Date, y=TeamFGPCT), color="green") +
  geom_line(data=average, aes(x=Date, y=mean_shooting), color="black") +
  scale_y_continuous(limits = c(0, .65))
```



# Chapter 20

## Step charts

Step charts are **a method of showing progress** toward something. They combine showing change over time – **cumulative change over time** – with magnitude. They're good at inviting comparison.

There's great examples out there. First is the Washignton Post looking at Lebron passing Jordan's career point total. Another is John Burn-Murdoch's work at the Financial Times (which is paywalled) about soccer stars. Here's an example of his work outside the paywall.

To replicate this, we need cumulative data – data that is the running total of data at a given point. So think of it this way – Nebraska scores 50 points in a basketball game and then 50 more the next, their cumulative total at two games is 100 points.

Step charts can be used for all kinds of things – showing how a player's career has evolved over time, how a team fares over a season, or franchise history. Let's walk through an example.

Let's look at Fred Hoiberg's first team.

We'll need the tidyverse.

```
library(tidyverse)
```

And we need to load our logs data we just downloaded.

```
logs <- read_csv("data/logs20.csv")
```

```
##  
## -- Column specification -----  
## cols(  
##   .default = col_double(),  
##   Date = col_date(format = ""),
```

```

##   HomeAway = col_character(),
##   Opponent = col_character(),
##   W_L = col_character(),
##   Blank = col_logical(),
##   Team = col_character(),
##   Conference = col_character(),
##   season = col_character()
## )
## i Use `spec()` for the full column specifications.

```

Here we're going to look at the scoring differential of teams. If you score more than your opponent, you win. So it stands to reason that if you score a lot more than your opponent over the course of a season, you should be very good, right? Let's see.

The first thing we're going to do is calculate that differential. Then, we'll group it by the team. After that, we're going to summarize using a new function called `cumsum` or cumulative sum – the sum for each game as we go forward. So game 1's `cumsum` is the differential of that game. Game 2's `cumsum` is Game 1 + Game 2. Game 3 is Game 1 + 2 + 3 and so on.

```

difflogs <- logs %>%
  mutate(Differential = TeamScore - OpponentScore) %>%
  group_by(Team) %>%
  mutate(CumDiff = cumsum(Differential))

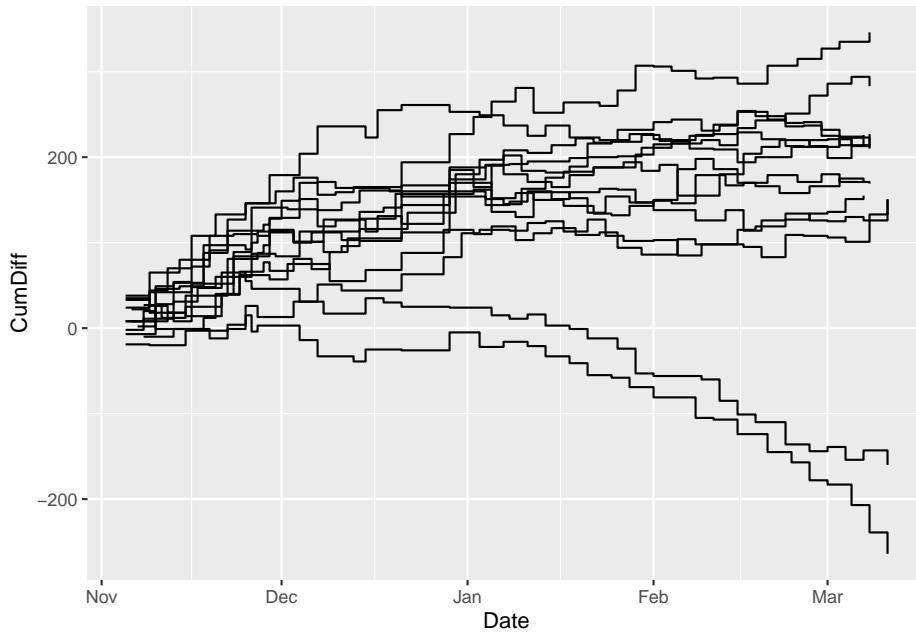
```

Now that we have the cumulative sum for each, let's filter it down to just Big Ten teams.

```
bigdiff <- difflogs %>% filter(Conference == "Big Ten")
```

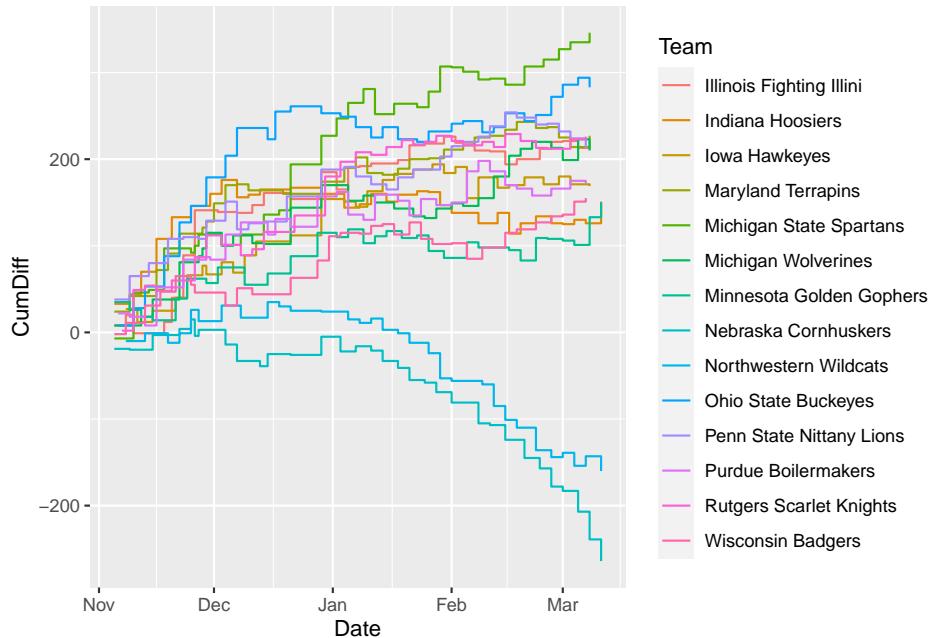
The step chart is its own geom, so we can employ it just like we have the others. It works almost exactly the same as a line chart, but it uses the cumulative sum instead of a regular value and, as the name implies, creates a step like shape to the line instead of a curve.

```
ggplot() + geom_step(data=bigdiff, aes(x=Date, y=CumDiff, group=Team))
```



Let's try a different element of the aesthetic: color, but this time inside the aesthetic. Last time, we did the color outside. When you put it inside, you pass it a column name and ggplot will color each line based on what thing that is, and it will create a legend that labels each line that thing.

```
ggplot() + geom_step(data=bigdiff, aes(x=Date, y=CumDiff, group=Team, color=Team))
```



From this, we can see two teams in the Big Ten had negative point differentials last season – Nebraska and Northwestern. But which is which?

Let's look at the top team plus Nebraska.

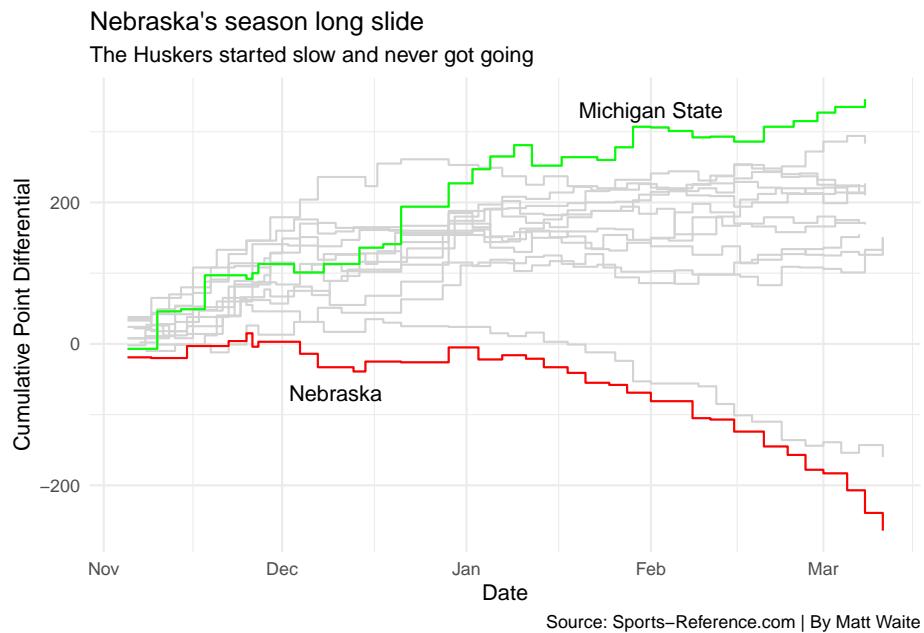
```
nu <- bigdiff %>% filter(Team == "Nebraska Cornhuskers")
ms <- bigdiff %>% filter(Team == "Michigan State Spartans")
```

Let's introduce a couple of new things here. First, note when I take the color OUT of the aesthetic, the legend disappears.

The second thing I'm going to add is the annotation layer. In this case, I am adding a text annotation layer, and I can specify where by adding in a x and a y value where I want to put it. This takes some finesse. After that, I'm going to add labels and a theme.

```
ggplot() +
  geom_step(data=bigdiff, aes(x=Date, y=CumDiff, group=Team), color="light grey") +
  geom_step(data=nu, aes(x=Date, y=CumDiff, group=Team), color="red") +
  geom_step(data=ms, aes(x=Date, y=CumDiff, group=Team), color="green") +
  annotate("text", x=(as.Date("2019-12-10")), y=-70, label="Nebraska") +
  annotate("text", x=(as.Date("2020-02-01")), y=330, label="Michigan State") +
  labs(
    x="Date",
    y="Cumulative Point Differential",
    title="Nebraska's season long slide",
    subtitle="The Huskers started slow and never got going",
```

```
caption="Source: Sports-Reference.com | By Matt Waite") +  
theme_minimal()
```





# Chapter 21

## Ridge charts

Ridgeplots are useful for when you want to show how different groupings compare with a large number of datapoints. So let's look at how we do this, and in the process, we learn about ggplot extensions. The extensions add functionality to ggplot, which doesn't out of the box have ridgeplots (sometimes called joyplots).

In the console, run this: `install.packages("ggridges")`

Now we can add those libraries.

```
library(tidyverse)
library(ggridges)
```

So for this, let's look at every basketball game played since the 2014-15 season.

We load that like this.

```
logs <- read_csv("data/logs1519.csv")

## Warning: Missing column names filled in: 'X1' [1]
##
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Date = col_date(format = ""),
##   HomeAway = col_character(),
##   Opponent = col_character(),
##   W_L = col_character(),
##   Blank = col_logical(),
##   Team = col_character(),
##   Conference = col_character(),
##   season = col_character()
```

```
## )
## i Use `spec()` for the full column specifications.
```

So I want to group teams by wins. Wins are the only thing that matter – ask Tim Miles. So our data has a column called W\_L that lists if the team won or lost. The problem is it doesn’t just say W or L. If the game went to overtime, it lists that. That complicates counting wins. And, with ridgeplots, I want to be able to separate EVERY GAME by how many wins the team had over a SEASON. So I’ve got some work to do.

First, here’s a trick to find a string of text and make that. It’s called `grepl` and the basic syntax is `grepl` for this string in this field and then do what I tell you. In this case, we’re going to create a new field called `winloss` look for W or L (and ignore any OT notation) and give wins a 1 and losses a 0.

```
winlosslogs <- logs %>% mutate(winloss = case_when(
  grepl("W", W_L) ~ 1,
  grepl("L", W_L) ~ 0)
)
```

Now I’m going to add up all the winlosses for each team, which should give me the number of wins for each team.

```
winlosslogs %>% group_by(Team, Conference, season) %>% summarise(TeamWins = sum(winloss))
```

```
## `summarise()` regrouping output by 'Team', 'Conference' (override with `groups` argument)
```

Now that I have season win totals, I can join that data back to my log data so each game has the total number of wins in each season.

```
winlosslogs %>% left_join(teamseasonwins, by=c("Team", "Conference", "season")) -> wintotallogs
```

Now I can use that same `case_when` logic to create some groupings. So I want to group teams together by how many wins they had over the season. For no good reason, I started with more than 25 wins, then did groups of 5 down to 10 wins. If you had fewer than 10 wins, God help your program.

The way to create a new field based on groupings like that is to use `case_when`, which says, basically, when This Thing Is True, Do This. So in our case, we’re going to pass a couple of logical statements that when they are both true, our data gets labeled how we want to label it. So we `mutate` a field called `grouping` and then use `case_when`.

```
wintotallogs %>% mutate(grouping = case_when(
  TeamWins > 25 ~ "More than 25 wins",
  TeamWins >= 20 & TeamWins <=25 ~ "20-25 wins",
  TeamWins >= 15 & TeamWins <=19 ~ "15-19 wins",
  TeamWins >= 10 & TeamWins <=14 ~ "10-14 wins",
  TeamWins < 10 ~ "Less than 10 wins")
) -> wintotalgroupinglogs
```

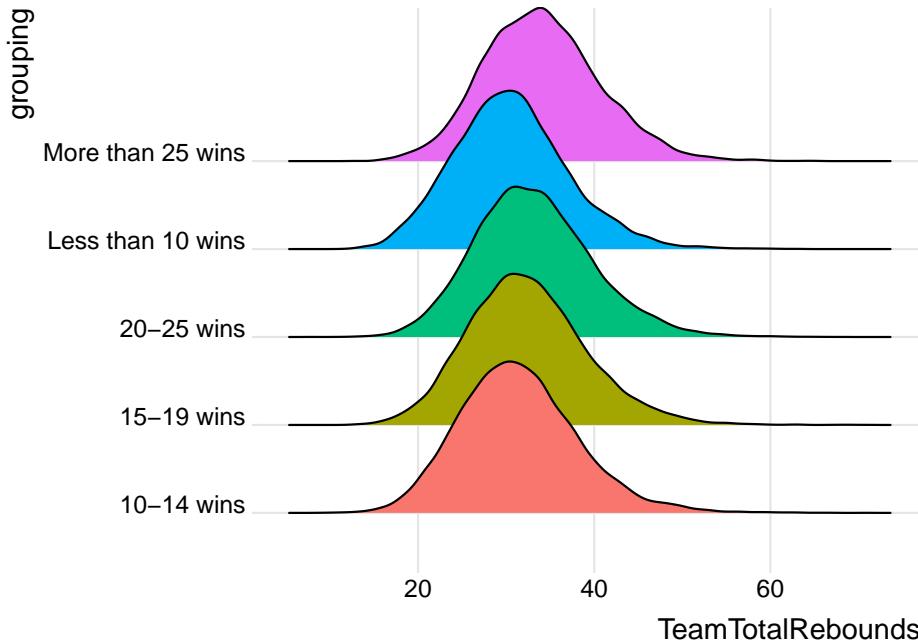
So my `wintotalgroupinglogs` table has each game, with a field that gives the total number of wins each team had that season and labeling each game with one of five groupings. I could use `dplyr` to do `group_by` on those five groups and find some things out about them, but ridgeplots do that visually.

Let's look at the differences in rebounds by those five groups. Do teams that win more than 25 games rebound better than teams that win fewer games?

The answer might surprise you.

```
ggplot(wintotalgroupinglogs, aes(x = TeamTotalRebounds, y = grouping, fill = grouping)) +
  geom_density_ridges() +
  theme_ridges() +
  theme(legend.position = "none")
```

```
## Picking joint bandwidth of 0.88
## Warning: Removed 2 rows containing non-finite values (stat_density_ridges).
```



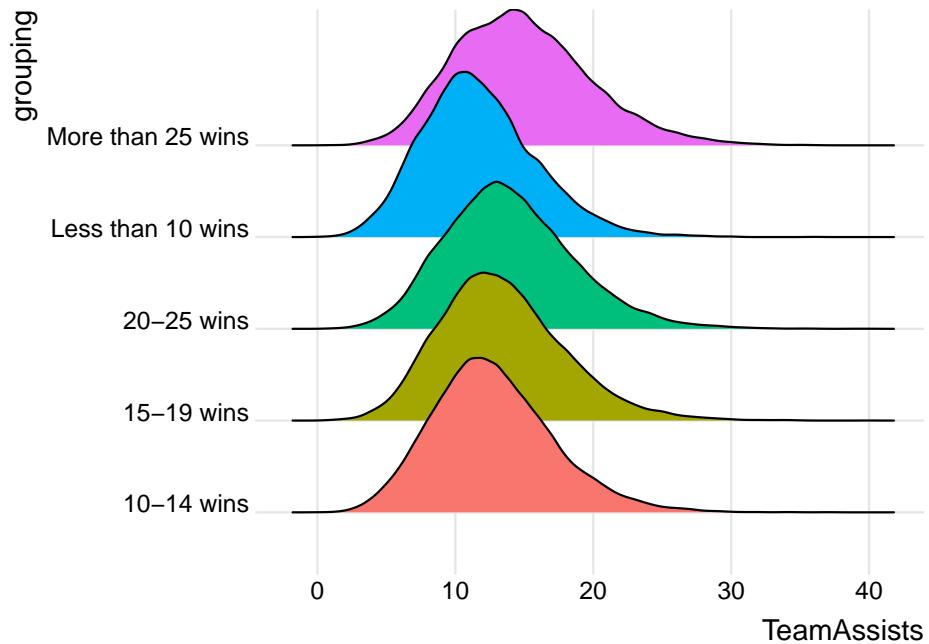
Answer? Not really. Game to game, maybe. Over five seasons? The differences are indistinguishable.

How about assists?

```
ggplot(wintotalgroupinglogs, aes(x = TeamAssists, y = grouping, fill = grouping)) +
  geom_density_ridges() +
  theme_ridges() +
  theme(legend.position = "none")
```

```
## Picking joint bandwidth of 0.601
```

```
## Warning: Removed 2 rows containing non-finite values (stat_density_ridges).
```

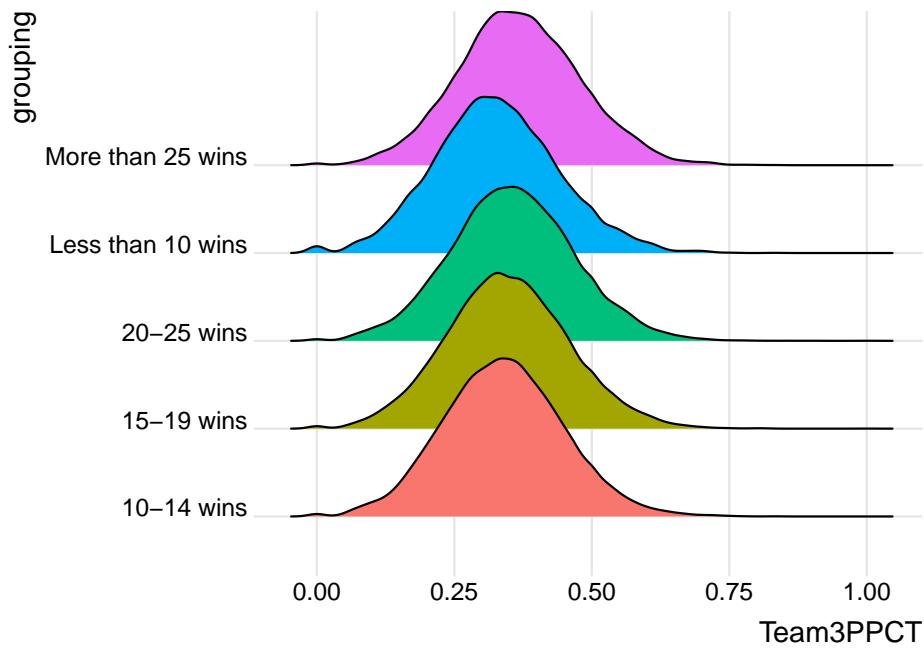


There's a little better, especially between top and bottom.

```
ggplot(wintotalgroupinglogs, aes(x = Team3PPCT, y = grouping, fill = grouping)) +
  geom_density_ridges() +
  theme_ridges() +
  theme(legend.position = "none")
```

```
## Picking joint bandwidth of 0.0156
```

```
## Warning: Removed 2 rows containing non-finite values (stat_density_ridges).
```

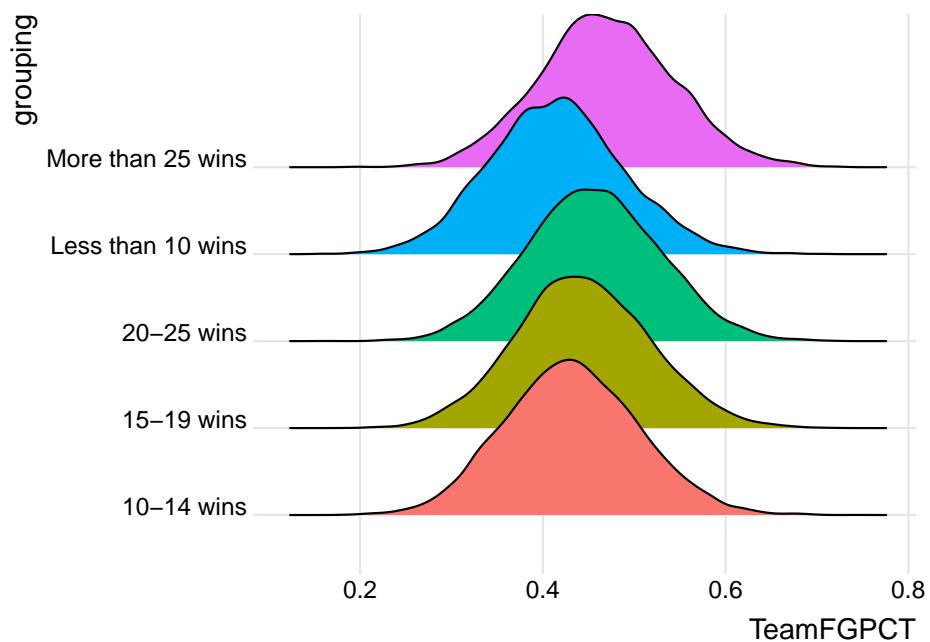


If you've been paying attention this semester, you know what's coming next.

```
ggplot(wintotalgroupinglogs, aes(x = TeamFGPCT, y = grouping, fill = grouping)) +
  geom_density_ridges() +
  theme_ridges() +
  theme(legend.position = "none")
```

```
## Picking joint bandwidth of 0.0102
```

```
## Warning: Removed 2 rows containing non-finite values (stat_density_ridges).
```



## Chapter 22

# Dumbbell and lollipop charts

Second to my love of waffle charts because I'm always hungry, dumbbell charts are an excellently named way of **showing the difference between two things on a number line** – a start and a finish, for instance. Or the difference between two related things. Say, turnovers and assists.

Lollipop charts – another excellent name – are a variation on bar charts. They do a good job of showing magnitude and difference between things.

To use both of them, you need to add a new library:

```
install.packages("ggalt")
```

Let's give it a whirl.

```
library(tidyverse)
library(ggalt)
```

### 22.1 Dumbbell plots

For this, let's use college football game logs.

And load it.

```
logs <- read_csv("data/footballlogs20.csv")
```

```
## 
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Date = col_date(format = ""),
##   ...)
```

```

##   HomeAway = col_character(),
##   Opponent = col_character(),
##   Result = col_character(),
##   TeamFull = col_character(),
##   TeamURL = col_character(),
##   Outcome = col_character(),
##   Team = col_character(),
##   Conference = col_character()
## )
## i Use `spec()` for the full column specifications.

```

For the first example, let's look at the difference between a team's giveaways – turnovers lost – versus takeaways, turnovers gained. To get this, we're going to add up all offensive turnovers and defensive turnovers for a team in a season and take a look at where they come out. To make this readable, I'm going to focus on the Big Ten.

```

turnovers <- logs %>%
  group_by(Team, Conference) %>%
  summarise(
    Giveaways = sum(TotalTurnovers),
    Takeaways = sum(DefTotalTurnovers)) %>%
  filter(Conference == "Big Ten Conference")

```

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

Now, the way that the `geom_dumbbell` works is pretty simple when viewed through what we've done before. There's just some tweaks.

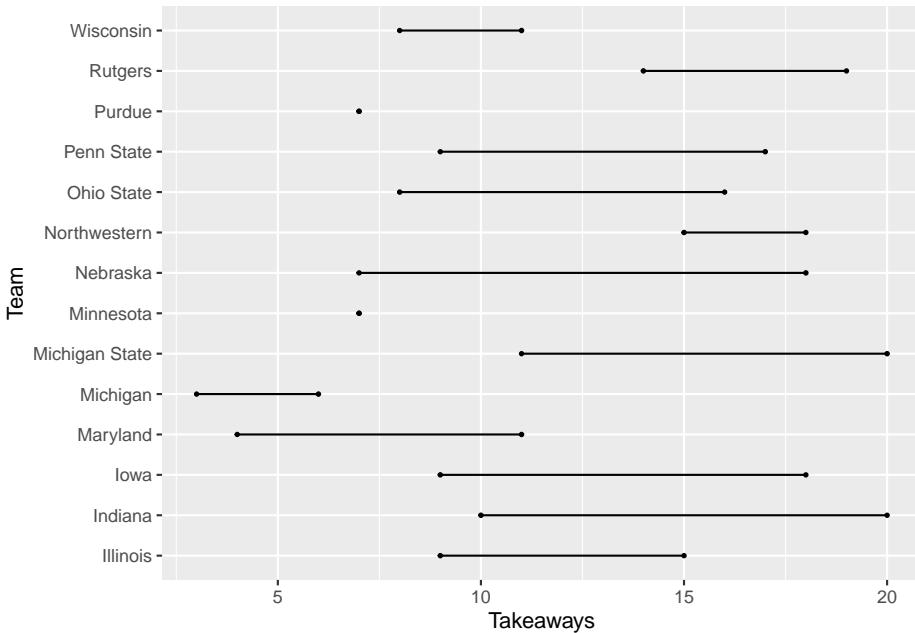
First: We start with the y axis. The reason is we want our dumbbells going left and right, so the label is going to be on the y axis.

Second: Our x is actually two things: x and xend. What you put in there will decide where on the line the dot appears.

```

ggplot() +
  geom_dumbbell(
    data=turnovers,
    aes(y=Team, x=Takeaways, xend=Giveaways)
  )

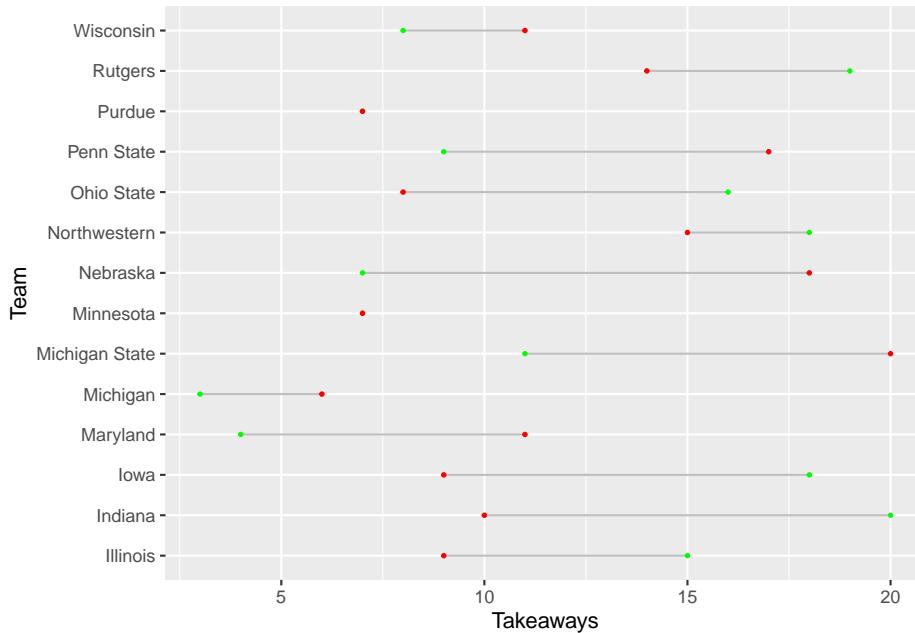
```



Well, that's a chart alright, but what dot is the giveaways and what are the takeaways? To fix this, we'll add colors.

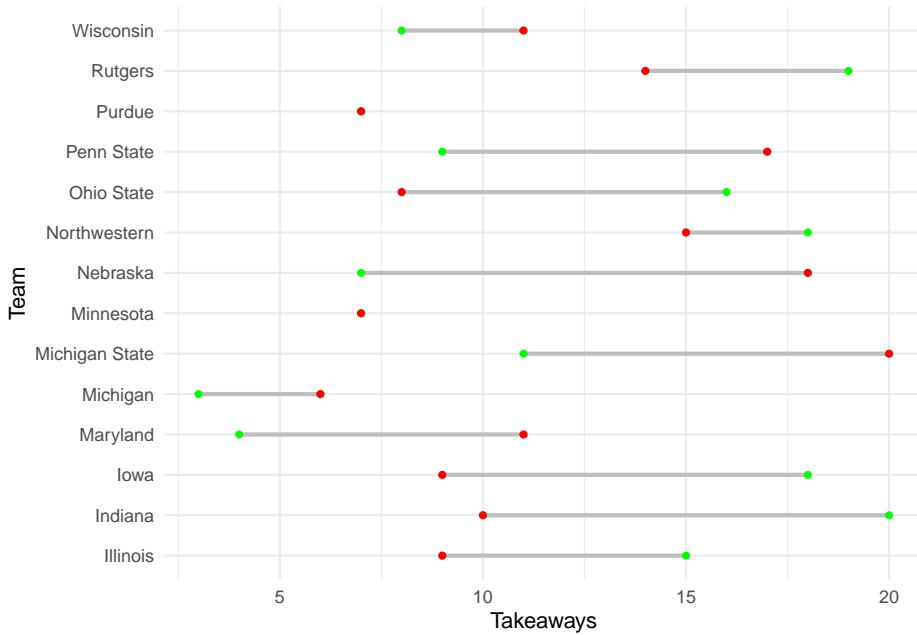
So our choice of colors here is important. We want giveaways to be seen as bad and takeaways to be seen as good. So let's try red for giveaways and green for takeaways. To make this work, we'll need to do three things: first, use the English spelling of color, so `colour`. The, uh, `colour` is the bar between the dots, the `x_colour` is the color of the x value dot and the `xend_colour` is the color of the xend dot. So in our setup, takeaways are x, they're good, so they're green.

```
ggplot() +
  geom_dumbbell(
    data=turnovers,
    aes(y=Team, x=Takeaways, xend=Giveaways),
    colour = "grey",
    colour_x = "green",
    colour_xend = "red")
```



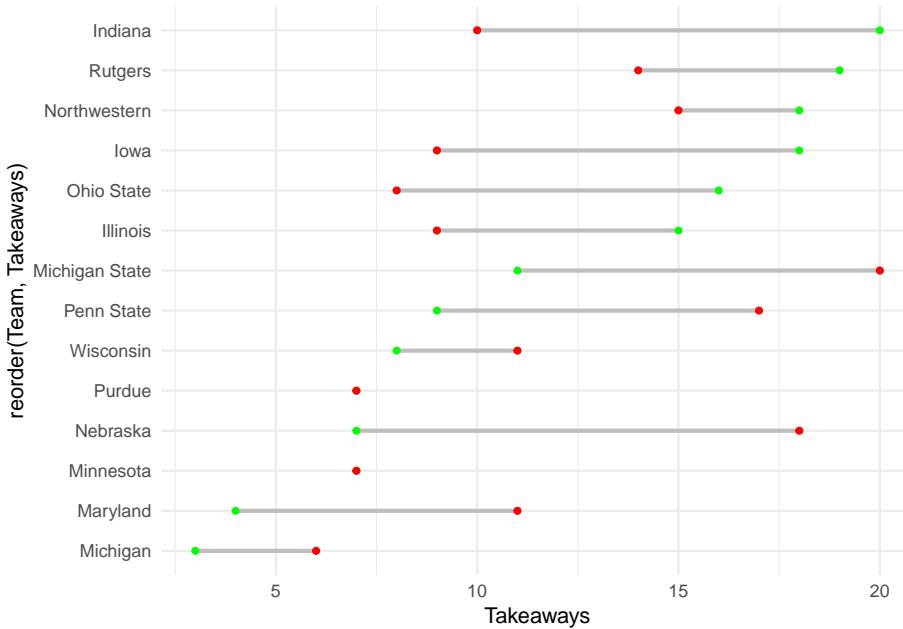
Better. Let's make two more tweaks. First, let's make the whole thing bigger with a `size` element. And let's add `theme_minimal` to clean out some cruft.

```
ggplot() +
  geom_dumbbell(
    data=turnovers,
    aes(y=Team, x=Takeaways, xend=Giveaways),
    size = 1,
    color = "grey",
    colour_x = "green",
    colour_xend = "red") +
  theme_minimal()
```



And now we have a chart that tells a story – got green on the right? That's good. A long distance between green and red? Better. But what if we sort it by good turnovers?

```
ggplot() +
  geom_dumbbell(
    data=turnovers,
    aes(y=reorder(Team, Takeaways), x=Takeaways, xend=Giveaways),
    size = 1,
    color = "grey",
    colour_x = "green",
    colour_xend = "red") +
  theme_minimal()
```

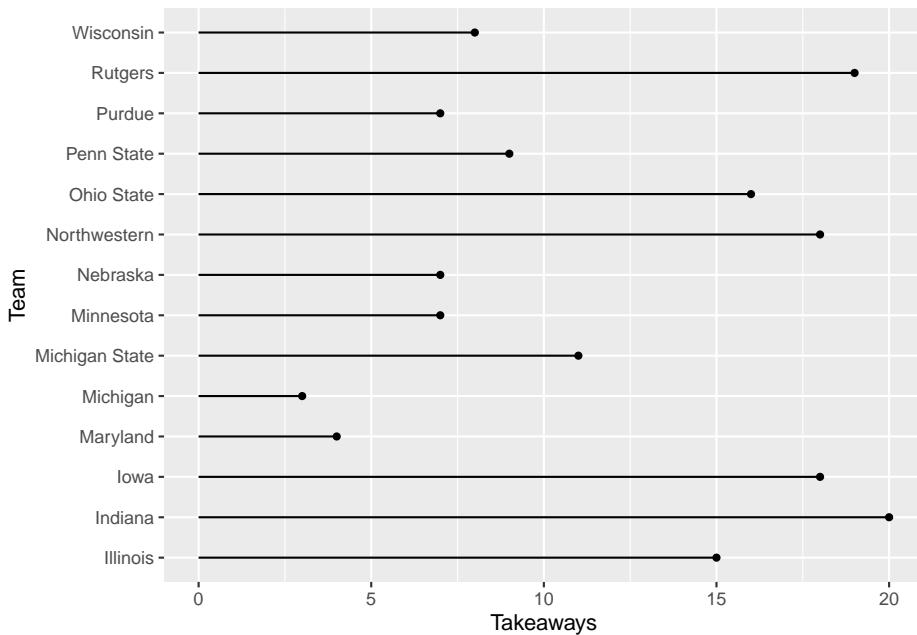


Believe it or not, Indiana had the most takeaways in the Big Ten last season. Don't sleep on the Fighting Basketball School.

## 22.2 Lollipop charts

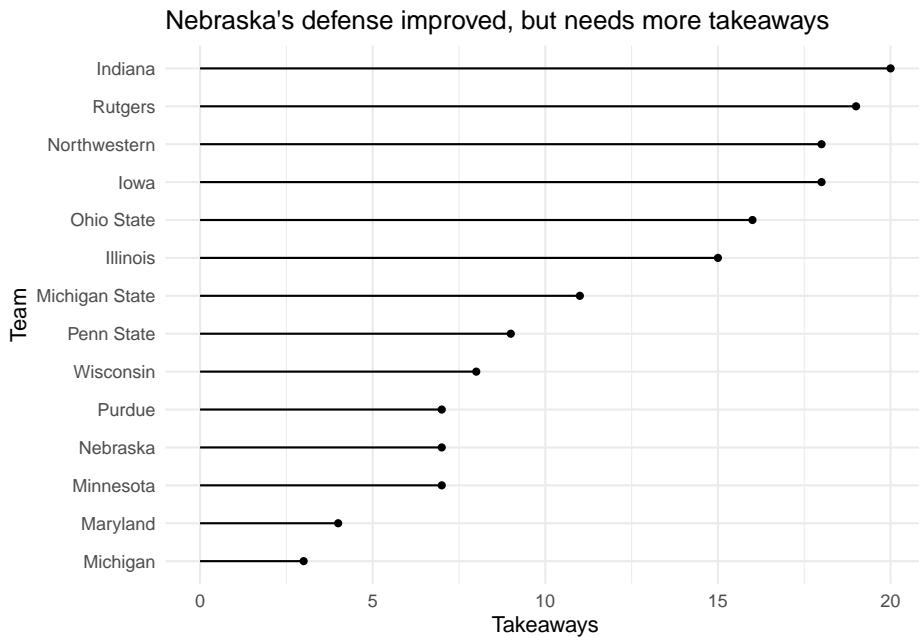
Sticking with takeaways, lollipops are similar to bar charts in that they show magnitude. And like dumbbells, they are similar in that we start with a y – the traditional lollipop chart is on its side – and we only need one x. The only additional thing we need to add is that we need to tell it that it is a horizontal chart.

```
ggplot() +
  geom_lollipop(
    data=turnovers,
    aes(y=Team, x=Takeaways),
    horizontal = TRUE
  )
```



We can do better than this with a simple theme\_minimal and some better labels.

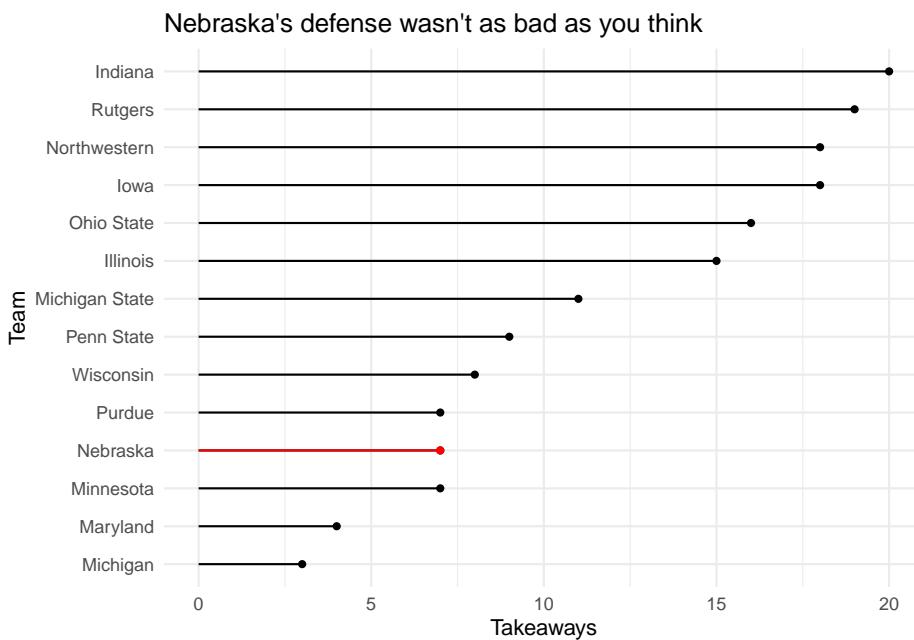
```
ggplot() +
  geom_lollipop(
    data=turnovers,
    aes(y=reorder(Team, Takeaways), x=Takeaways),
    horizontal = TRUE
  ) + theme_minimal() +
  labs(title = "Nebraska's defense improved, but needs more takeaways", y="Team")
```



How about some layering?

```
nu <- turnovers %>% filter(Team == "Nebraska")

ggplot() +
  geom_lollipop(
    data=turnovers,
    aes(y=reorder(Team, Takeaways), x=Takeaways),
    horizontal = TRUE
  ) +
  geom_lollipop(
    data=nu,
    aes(y=Team, x=Takeaways),
    horizontal = TRUE,
    color = "red"
  ) +
  theme_minimal() +
  labs(title = "Nebraska's defense wasn't as bad as you think", y="Team")
```



The headline says it all.



# Chapter 23

## Scatterplots

On the Monday, Sept. 21, 2020 edition of the Pick Six Podcast, Omaha World Herald reporter Sam McKewon talked a little about the Nebraska mens basketball team. Specifically the conversation was about a new roster release, and how the second year of Fred Hoiberg ball was going to look very different, starting with the heights of the players. After a near complete roster turnover, the players on the team now were nearly all taller than 6'4", and one of the shorter ones is penciled in as the starting point guard.

Why is that important? One reason, McKewon posited, is that teams made a lot of three point shots on Nebraska. In fact, Nebraska finished dead last in the conference in three points shots made against them. McKewon chalked this up to bad perimeter defense, and that Nebraska needed to improve it. Being taller – or more specifically having the longer arms that go with being taller – will help with that, McKewon said.

Better perimeter defense, better team.

The question before you is this: is that true? Does keeping a lid on your opponent's ability to score three pointers mean more wins?

This is what we're going to start to answer today. And we'll do it with scatterplots and regressions. Scatterplots are very good at showing **relationships between two numbers**.

First, we need libraries and every college basketball game last year.

Load the tidyverse.

```
library(tidyverse)
```

And the data.

```
logs <- read_csv("data/logs20.csv")

##
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Date = col_date(format = ""),
##   HomeAway = col_character(),
##   Opponent = col_character(),
##   W_L = col_character(),
##   Blank = col_logical(),
##   Team = col_character(),
##   Conference = col_character(),
##   season = col_character()
## )
## i Use `spec()` for the full column specifications.
```

To do this, we need all teams and their season stats. How much, team to team, does a thing matter? That's the question you're going to answer.

In our case, we want to know how much do three point shots made influence wins? How much difference can we explain in wins by knowing how many threes the other team made against you? We're going to total up the number of threes each team allowed and their season wins in one swoop.

To do this, we need to use conditional logic – `case_when` in this case – to determine if the team won or lost the game. In this case, we'll create a new column called `winloss`. Case when statements can be read like this: When This is True, Do This. This bit of code – which you can use in a lot of contexts in this class – uses the `grepl` function to look for the letter W in the `W_L` column and, if it finds it, makes `winloss` 1. If it finds an L, it makes it 0. Sum your `winloss` column and you have your season win total. The reason we have to use `grepl` to find W or L is because Sports Reference will record overtime wins differently than regular wins. Same with losses.

```
winlosslogs <- logs %>%
  mutate(
    winloss = case_when(
      grepl("W", W_L) ~ 1,
      grepl("L", W_L) ~ 0
    )
  )
```

Now we can get a dataframe together that gives us the total wins for each team, and the total three point shots made. We'll call that new dataframe `threedef`.

```
threedef <- winlosslogs %>%
  group_by(Team) %>%
  summarise(
```

```

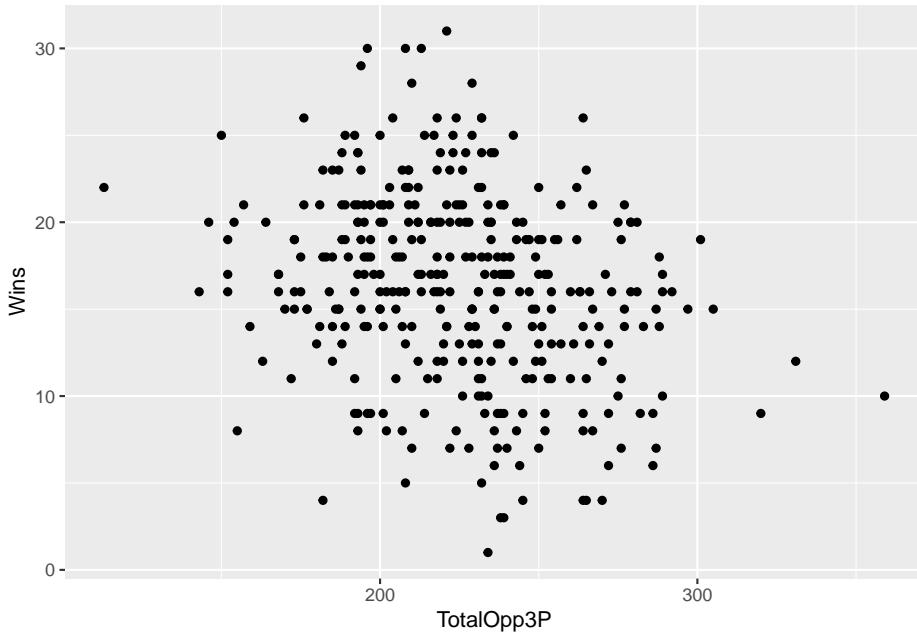
Wins = sum(winloss),
TotalOpp3P = sum(Opponent3P)
)

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

```

Now let's look at the scatterplot. With a scatterplot, we put what predicts the thing on the X axis, and the thing being predicted on the Y axis. In this case, X is our three pointers given up, y is our wins.

```
ggplot() + geom_point(data=threedef, aes(x=TotalOpp3P, y=Wins))
```



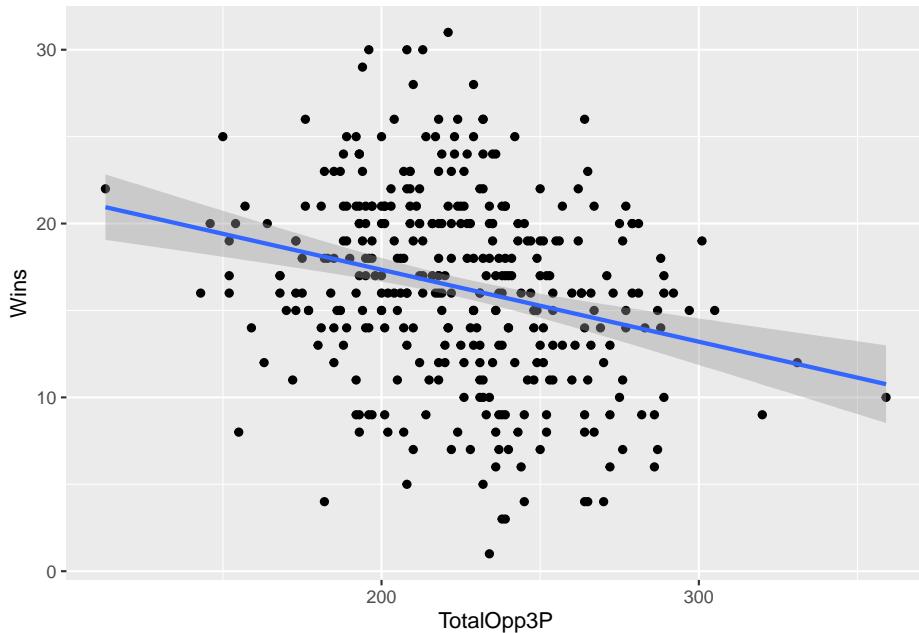
Let's talk about this. This seems kind of random, but clustered around the middle and maybe sloping down to the right. That would mean the more threes you give up, the less you win. And that makes intuitive sense. But can we get a better sense of this? Yes, by adding another geom – `geom_smooth`. It's identical to our `geom_point`, but we add a method to the end, which in this case we're using the linear method or `lm`.

```

ggplot() +
  geom_point(data=threedef, aes(x=TotalOpp3P, y=Wins)) +
  geom_smooth(data=threedef, aes(x=TotalOpp3P, y=Wins), method="lm")

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

```



So it does slope down to the right like we expect, but this still doesn't look good to me. It's very spread out. Which is a clue that you should be asking a question here: how strong of a relationship is this? How much can threes given up explain wins? Can we put some numbers to this?

Of course we can. We can apply a linear model to this – remember Chapter 9? We're going to create an object called `fit`, and then we're going to put into that object a linear model – `lm` – and the way to read this is “wins are predicted by opponent threes”. Then we just want the summary of that model.

```
fit <- lm(Wins ~ TotalOpp3P, data = threedef)
summary(fit)

##
## Call:
## lm(formula = Wins ~ TotalOpp3P, data = threedef)
##
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -14.9345  -3.4593   0.3006   3.6036  14.5274 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 25.619712   1.859947 13.774 < 2e-16 ***
## TotalOpp3P -0.041390   0.008184 -5.057 6.87e-07 ***
## ---
```

```

## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.313 on 351 degrees of freedom
## Multiple R-squared:  0.06792,   Adjusted R-squared:  0.06526
## F-statistic: 25.58 on 1 and 351 DF,  p-value: 6.869e-07

```

Remember from Chapter 9: There's just a few things you really need.

The first thing: R-squared. In this case, the Adjusted R-squared value is 0.06526, which we can interpret as shooting percentage predicts about 6.5 percent of the variance in wins. Which sounds not great.

Second: The P-value. We want anything less than .05. If it's above .05, the change between them is not statistically significant – it's probably explained by random chance. In our case, we have 6.869e-07, which is to say 6.869 with 7 zeros in front of it, or .00000006869. Is that less than .05? Yes. Yes it is. So this is not random. Again, we would expect this, so it's a good logic test.

Third: The coefficient. In this case, the coefficient for TeamOpp3P is -0.041390. What this model predicts, given that and the intercept of 25.619712, is this: Every team starts with about 26 wins. For every 100 three pointers the other team makes, you lose 4.139 games off that total. So if you give up 100 threes in a season, you'll be a 20 win team. Give up 200, you're a 17 win team, and so on. How am I doing that? Remember your algebra and  $y = mx + b$ . In this case,  $y$  is the wins,  $m$  is the coefficient,  $x$  is the number of threes given up and  $b$  is the intercept.

Let's use Nebraska as an example. They had 276 threes scored on them in the last season.

$y = -0.041390 * 276 + 25.619712$  or 14.19 wins.

How many wins did Nebraska have? 7.

What does that mean? It means that as disappointing a season as it was, Nebraska UNDERPERFORMED according to this model. But our R-squared is only 6.5 percent. Put another way: 93.5 percent of the difference in wins between teams is predicted by something else.

Where is Nebraska on the plot? We know we can use layering for that.

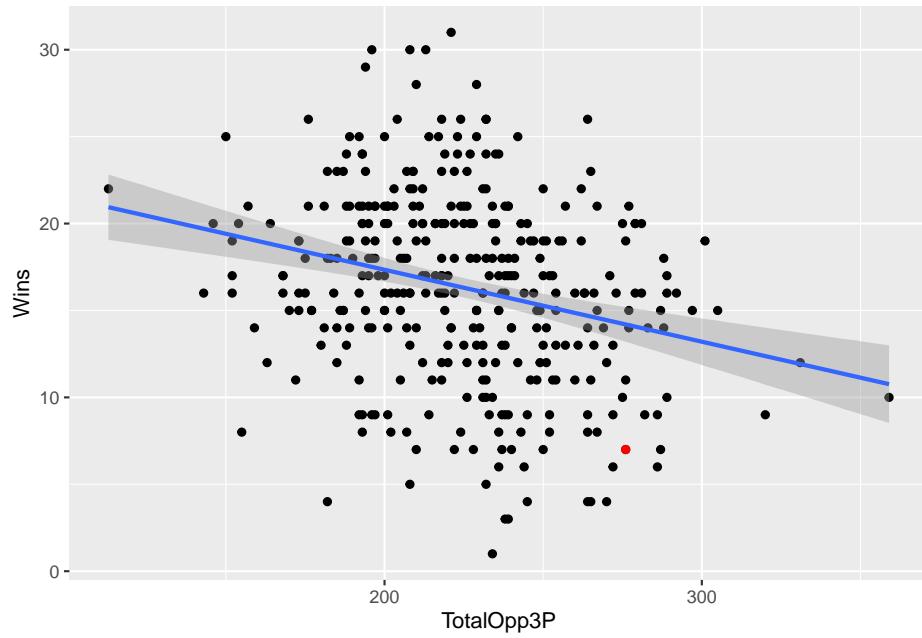
```

nu <- threedef %>% filter(Team == "Nebraska Cornhuskers")

ggplot() +
  geom_point(data=threedef, aes(x=TotalOpp3P, y=Wins)) +
  geom_smooth(data=threedef, aes(x=TotalOpp3P, y=Wins), method="lm") +
  geom_point(data=nu, aes(x=TotalOpp3P, y=Wins), color="red")

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

```



# Chapter 24

## Bubble charts

Here is the real talk: Bubble charts are hard. The reason they are hard is not because of the code, or the complexity or anything like that. They're a scatterplot with magnitude added – the size of the dot in the scatterplot has meaning. The hard part is seeing when a bubble chart works and when it doesn't.

If you want to see it work spectacularly well, watch a semi-famous Ted Talk by Hans Rosling from 2006 where bubble charts were the centerpiece. It's worth watching. It'll change your perspective on the world. No seriously. It will.

And since then, people have wanted bubble charts. And we're back to the original problem: They're hard. There's a finite set of circumstances where they work.

First, I'm going to show you an example of them not working to illustrate the point.

I'm going to load up my libraries.

```
library(tidyverse)
```

So for this example, I want to look at where Big Ten teams compare to the rest of college football last season. Is the Big Ten's reputation for tough games and defenses earned? Can we see patterns in good team vs bad teams?

I'm going to create a scatterplot with offensive yards per play on the X axis and defensive yards per play on the y axis. We can then divide the grid into four quadrants. Teams with high yards per offensive play and low defensive yards per play are teams with good offenses and good defenses. The opposite means bad defense, bad offense. Then, to drive the point home, I'm going to make the dot the size of the total wins on the season – the bubble in my bubble charts.

We'll use this season's college football games.

And load it.

```
logs <- read_csv("data/footballlogs20.csv")
```

```
## 
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Date = col_date(format = ""),
##   HomeAway = col_character(),
##   Opponent = col_character(),
##   Result = col_character(),
##   TeamFull = col_character(),
##   TeamURL = col_character(),
##   Outcome = col_character(),
##   Team = col_character(),
##   Conference = col_character()
## )
## i Use `spec()` for the full column specifications.
```

To do this, I've got some work to do. First, I need to mutate the outcomes of the games into 1s and 0s so I can add up the wins. We've done this before, so this won't be new to you, just adjusted slightly from basketball data.

```
winlosslogs <- logs %>%
  mutate(
    wins = case_when(
      grepl("W", Outcome) ~ 1,
      grepl("L", Outcome) ~ 0
    )
  )
```

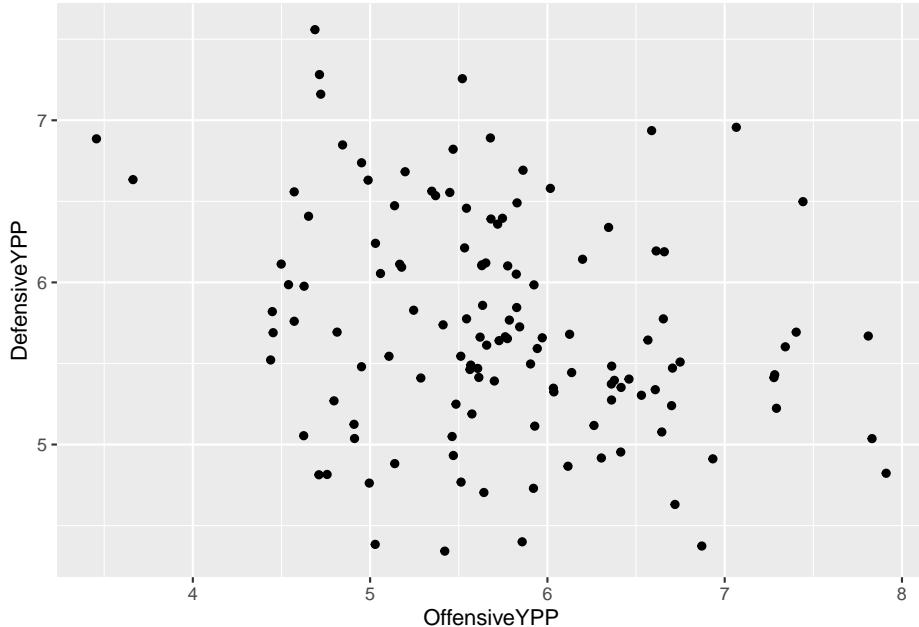
Now I have some more work to do. My football logs data has the yards per play of each game, and I could average those together and get something very close to what I'm going to do, but averaging each game's yards per play is not the same thing as calculating it, so we're going to calculate it.

```
winlosslogs %>%
  group_by(Team, Conference) %>%
  summarise(
    TotalPlays = sum(OffensivePlays),
    TotalYards = sum(OffensiveYards),
    DefensivePlays = sum(DefPlays),
    DefensiveYards = sum(DefYards),
    TotalWins = sum(wins)) %>%
  mutate(
    OffensiveYPP = TotalYards/TotalPlays,
    DefensiveYPP = DefensiveYards/DefensivePlays) -> ypp
```

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

A bubble chart is just a scatterplot with one additional element in the aesthetic – a size. Here's the scatterplot version.

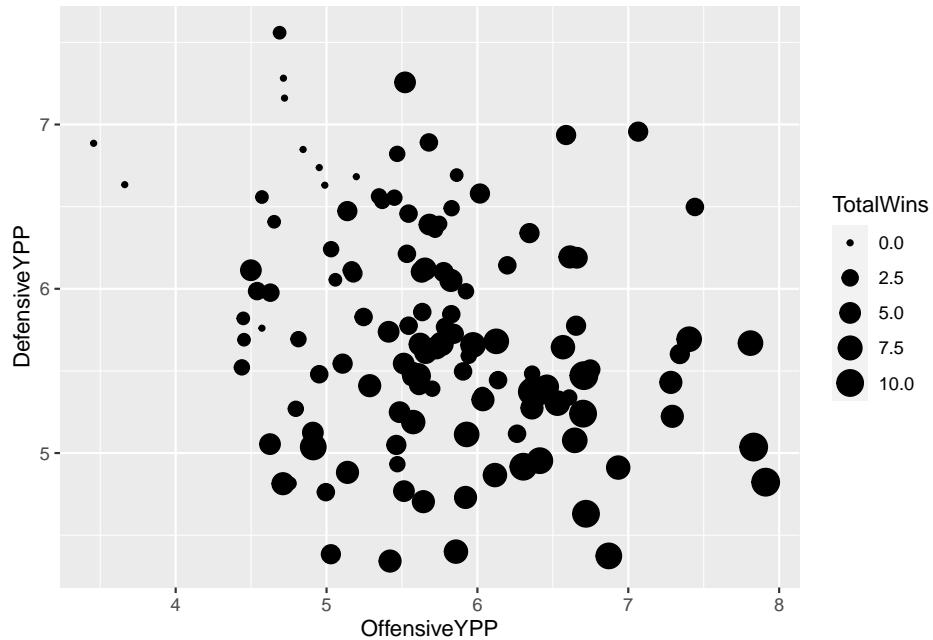
```
ggplot() + geom_point(data=yp, aes(x=OffensiveYPP, y=DefensiveYPP))
```



Looks kind of random, eh? In this case, that's not that bad because we're not claiming a relationship. We're saying the location on the chart has meaning. So, do teams on the bottom right – good offense, good defense – win more games?

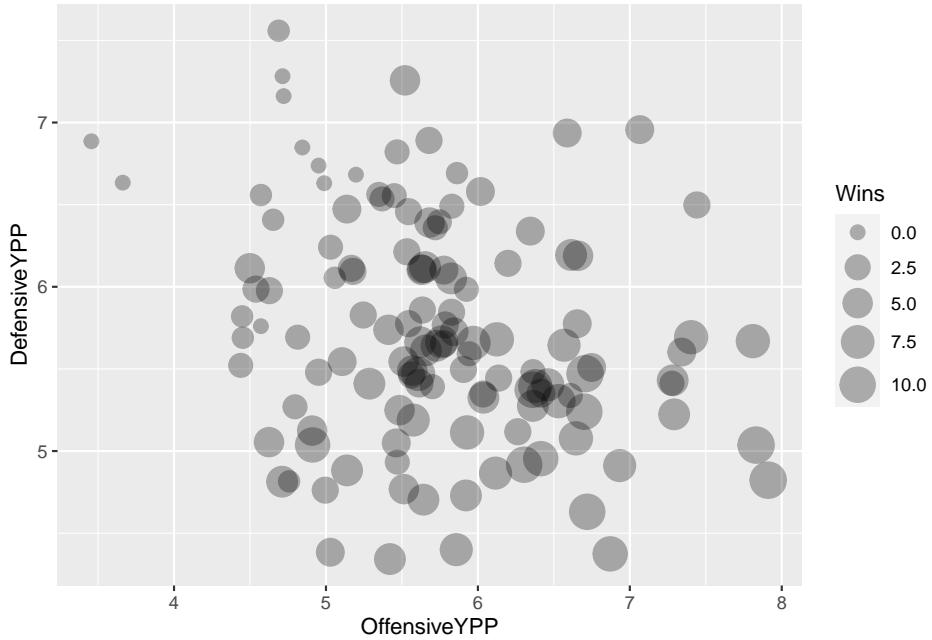
Let's add the size element.

```
ggplot() +
  geom_point(
    data=yp,
    aes(x=OffensiveYPP, y=DefensiveYPP, size=TotalWins)
  )
```



What does this chart tell you? We can see a general pattern that there are more big dots on the bottom right than the upper left. But we can make this more readable by adding an alpha element outside the aesthetic – alpha in this case is transparency – and we can manually change the size of the dots by adding `scale_size` and a `range`.

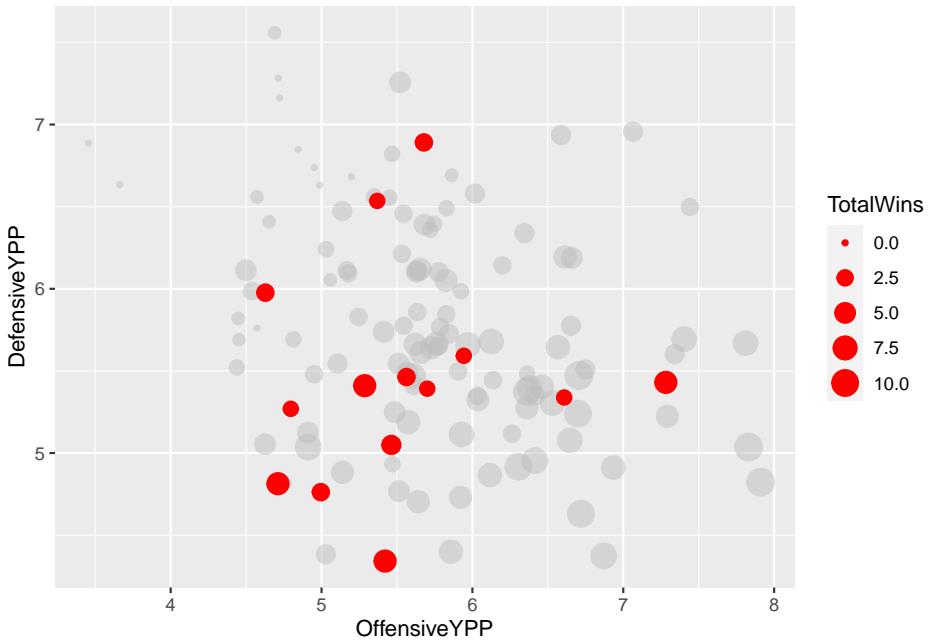
```
ggplot() +
  geom_point(
    data=yp,
    aes(x=OffensiveYPP, y=DefensiveYPP, size=TotalWins),
    alpha = .3) +
  scale_size(range = c(3, 8), name="Wins")
```



And by now, you now know to add in the Big Ten as a layer, I would hope.

```
bigten <- spp %>% filter(Conference == "Big Ten Conference")
```

```
ggplot() +
  geom_point(
    data=spp,
    aes(x=OffensiveYPP, y=DefensiveYPP, size=TotalWins),
    color="grey",
    alpha=.5) +
  geom_point(
    data=bigten,
    aes(x=OffensiveYPP, y=DefensiveYPP, size=TotalWins),
    color="red")
```



Let's add some things to this chart to help us out. First, let's add lines that show us the average of all teams for those two metrics. So first, we need to calculate those. Because I have grouped data, it's going to require me to ungroup it so I can get just the total average of those two numbers.

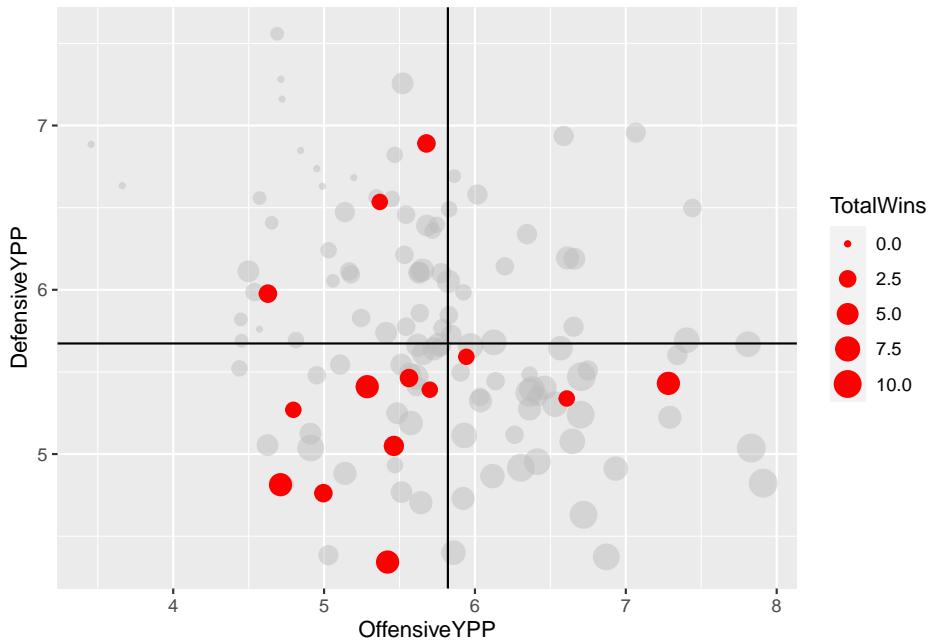
```
yppl %>%
  ungroup() %>%
  summarise(
    offense = mean(OffensiveYPP),
    defense = mean(DefensiveYPP)
  )

## # A tibble: 1 x 2
##   offense  defense
##     <dbl>    <dbl>
## 1      5.74     5.73
```

Now we can use those averages to add two more geoms – geom\_vline and geom\_hline, for vertical lines and horizontal lines.

```
ggplot() +
  geom_point(
    data=yppl,
    aes(x=OffensiveYPP, y=DefensiveYPP, size=TotalWins),
    color="grey",
    alpha=.5) +
  geom_point(
```

```
data=bigten,
aes(x=OffensiveYPP, y=DefensiveYPP, size=TotalWins),
color="red") +
geom_vline(xintercept = 5.820543) +
geom_hline(yintercept = 5.673263)
```



Now, let's add another new geom for us, using a new library called `ggrepel`, which will help us label the dots without overwriting other labels. So we'll have to install that in the console:

```
'install.packages("ggrepel")
library(ggrepel)
```

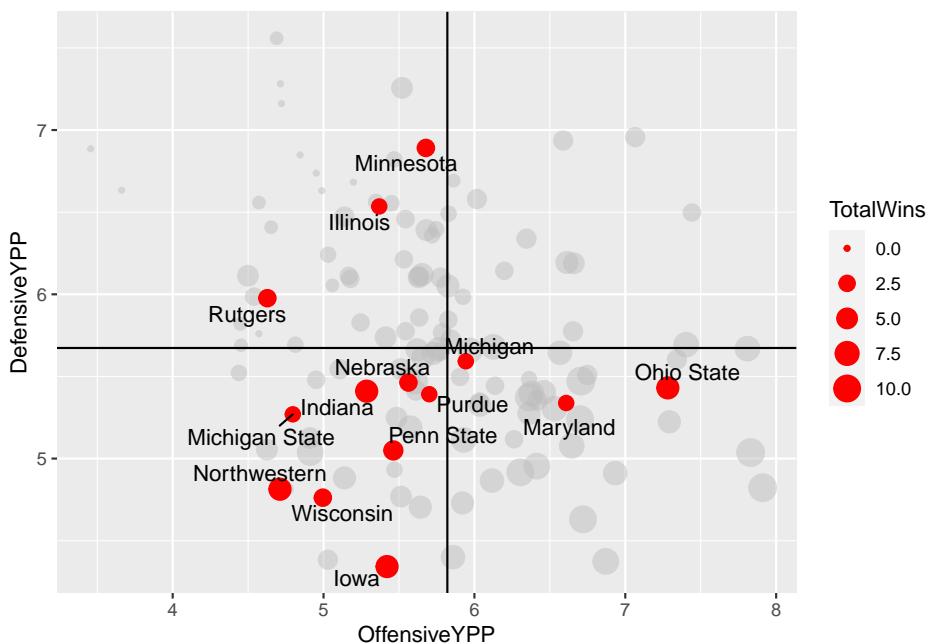
And with that, we can add labels to the dots. The `geom_text_repel` is pretty much the exact same thing as your Big Ten geom point, but instead of a size, you include a label.

```
ggplot() +
  geom_point(
    data=ypp,
    aes(x=OffensiveYPP, y=DefensiveYPP, size=TotalWins),
    color="grey",
    alpha=.5) +
  geom_point(
    data=bigten,
    aes(x=OffensiveYPP, y=DefensiveYPP, size=TotalWins),
```

```

    color="red") +
geom_vline(xintercept = 5.820543) +
geom_hline(yintercept = 5.673263) +
geom_text_repel(
  data=bigten,
  aes(x=OffensiveYPP, y=DefensiveYPP, label=Team)
)

```



Well, what do you know about that? Nebraska was ... really a mixed bag this season.

All that's left is some labels and some finishing touches.

```

ggsave("bubble.png")
ggplot() +
  geom_point(
    data=yppl,
    aes(x=OffensiveYPP, y=DefensiveYPP, size=TotalWins),
    color="grey",
    alpha=.5) +
  geom_point(
    data=bigten,
    aes(x=OffensiveYPP, y=DefensiveYPP, size=TotalWins),
    color="red") +
  geom_vline(xintercept = 5.820543) +
  geom_hline(yintercept = 5.673263) +
  geom_text_repel(

```

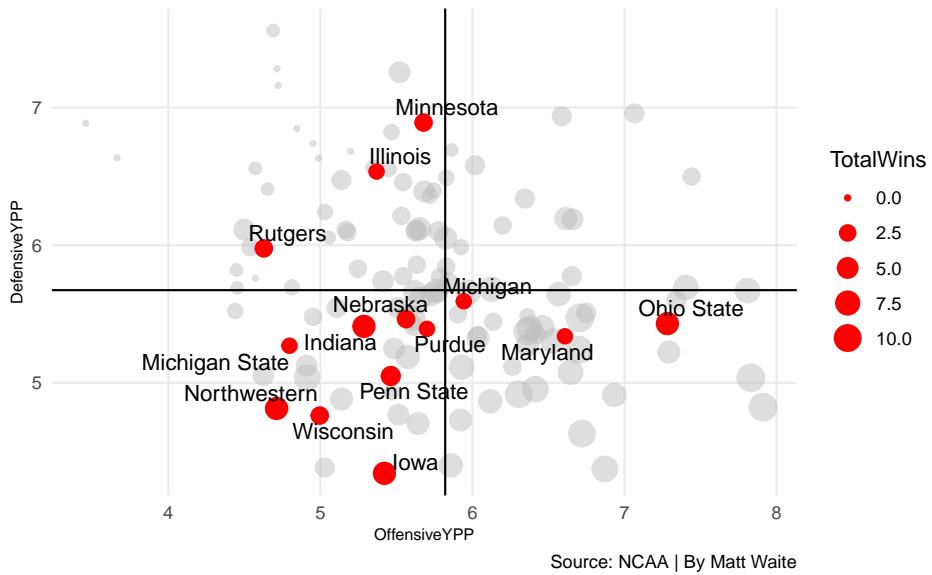
```

data=bigten,
aes(x=OffensiveYPP, y=DefensiveYPP, label=Team)
) +
labs(title="Nebraska: Average at football?", subtitle="The Huskers were kinda in the middle of theme",
plot.title = element_text(size = 16, face = "bold"),
axis.title = element_text(size = 8),
plot.subtitle = element_text(size=10),
panel.grid.minor = element_blank()
)

```

## Nebraska: Average at football?

The Huskers were kinda in the middle of everything, neither good nor bad at anything.





# Chapter 25

## Beeswarm plots

A beeswarm plot is sometimes called a column scatterplot. It's an effective way to show how individual things – teams, players, etc. – are distributed along a numberline. The column is a grouping – say positions in basketball – and the dots are players, and the dots cluster where the numbers are more common. So think of it like a histogram mixed with a scatterplot crossed with a bar chart.

An example will help.

```
First things first: Install ggbeeswarm with install.packages("ggbeeswarm")
```

Like ggalt and ggrepel, ggbeeswarm adds a couple new geoms to ggplot. We'll need to load it, the tidyverse and, for later, ggrepel.

```
library(tidyverse)
library(ggbeeswarm)
library(ggrepel)
```

Another bit of setup: we need to set the seed for the random number generator. The library “jitters” the dots in the beeswarm randomly. If we don't set the seed, we'll get different results each time. Setting the seed means we get the same look.

```
set.seed(1234)
```

So let's look at last year's basketball team as a group of shooters. The team was disappointing – we know that – but what kind of a problem is it going to be that we're returning basically no one from it?

First we'll load our player data.

```
players <- read_csv("data/players20.csv")
```

```
## 
## -- Column specification -----
```

```
## cols(
##   .default = col_double(),
##   Team = col_character(),
##   Player = col_character(),
##   Class = col_character(),
##   Pos = col_character(),
##   Height = col_character(),
##   Hometown = col_character(),
##   `High School` = col_character(),
##   Summary = col_character()
## )
## i Use `spec()` for the full column specifications.
```

We know this data has a lot of players who didn't play, so let's get rid of them.

```
activeplayers <- players %>% filter(MP>0)
```

Now let's ask what makes a good shooter? The best measure, in my book, is True Shooting Percentage. It's a combination of weighted field goal shooting – to account for three pointers – and free throws. Our data has TS%, but if we include *all* players, we'll have too many dots. So let's narrow it down. A decent tool for cutoffs? Field goal attempts. Let's get a quick look at them.

```
summary(activeplayers$FGA)
```

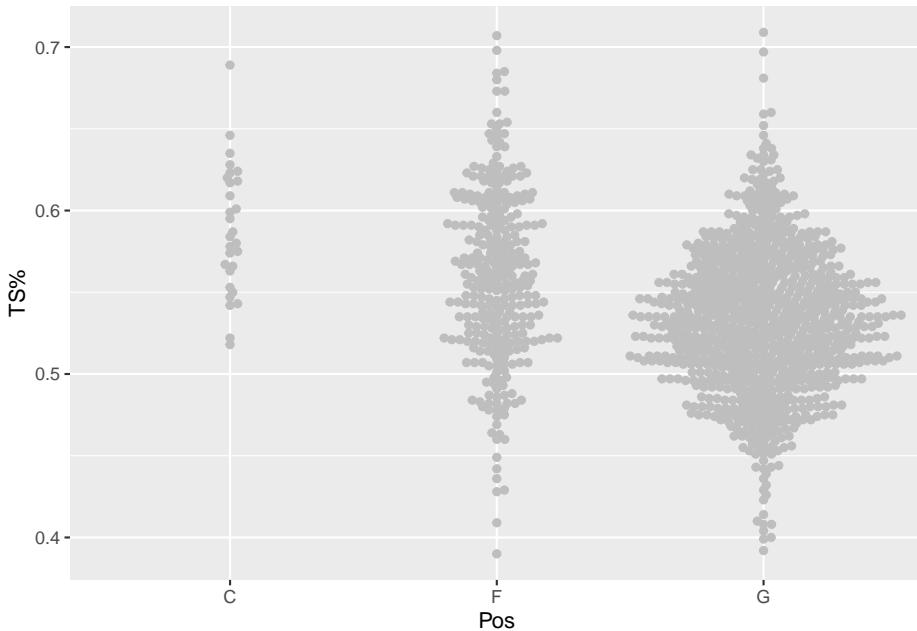
	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.00	23.25	96.00	121.90	197.00	611.00

The median number of shots is 96, but we only really care about good ones. So let's use 197 attempts as our cutoff.

```
shooters <- activeplayers %>% filter(FGA > 197)
```

Now we've got enough for a beeswarm plot. It works very much like you would expect – the group value is the x, the number is the y. We're going to beeswarm by position, and the dots will be true shooting percentage.

```
ggplot() + geom_beeswarm(data=shooters, aes(x=Pos, y=`TS%`), color="grey")
```



You can see that there's a lot fewer centers who have attempted more than 197 shots than guards, but then there's a lot more guards in college basketball than anything else. In the guards column, note that fat width of the swarm is between .5 and .55. So that means most guards who shoot more than 197 shots end up in that area. They're the average shooter at that level. You can see, some are better, some are worse.

So where are the Nebraska players in that mix?

We'll filter players on Nebraska who meet our criteria.

```
nu <- players %>%
  filter(Team == "Nebraska Cornhuskers") %>%
  filter(FGA>197) %>%
  arrange(desc(`TS%`))
```

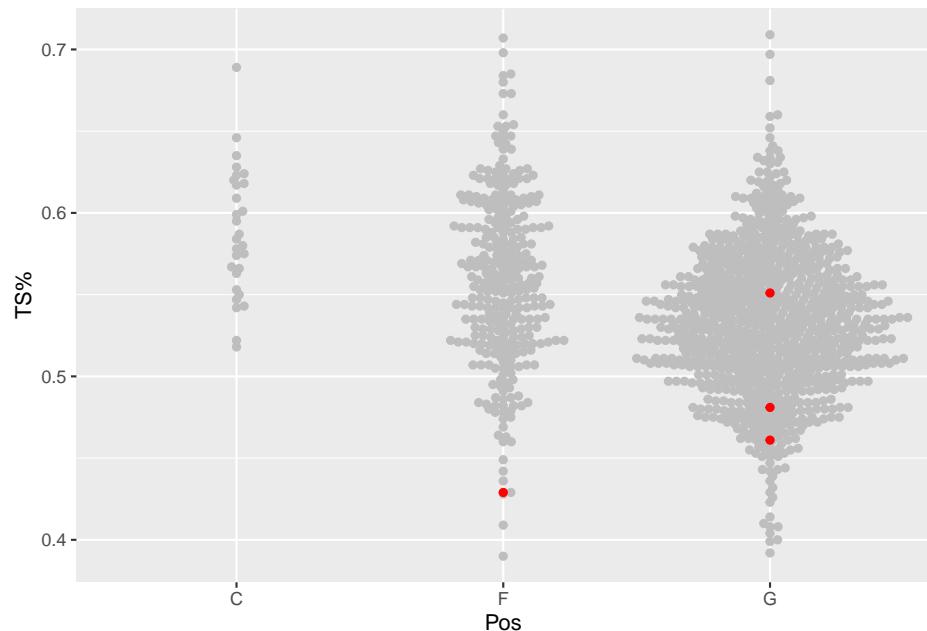
Four Cornhuskers took more than 197 shots. Number returning this season? Zero.

That's ... not a good start. Nothing coming back. But how good are they as true shooters?

When you add another beeswarm, we need to pass another element in – we need to tell it if we're grouping on the x value. Not sure why, but you'll get a warning if you don't.

```
ggplot() +
  geom_beeswarm(
```

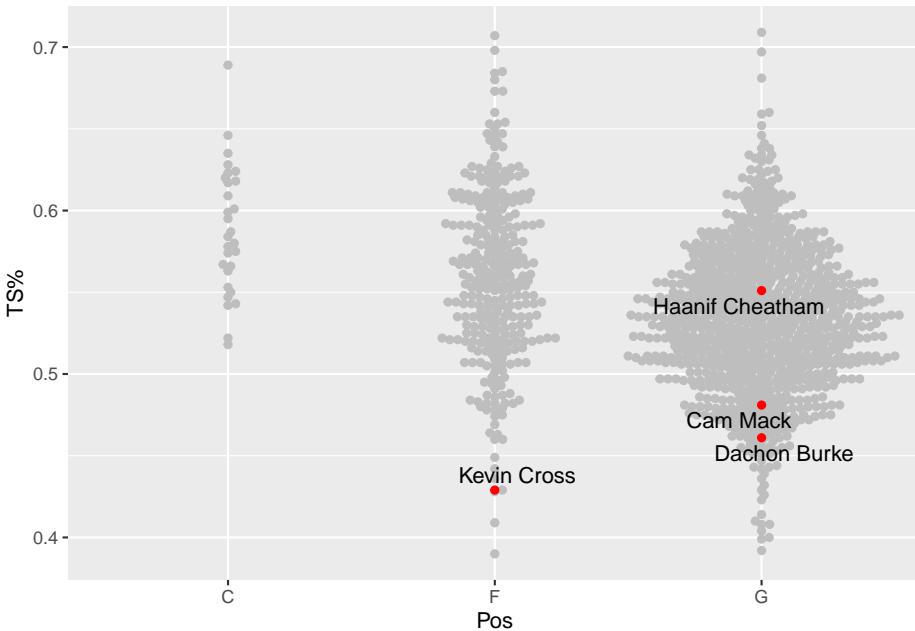
```
data=shooters,
groupOnX=TRUE,
aes(x=Pos, y=`TS%`), color="grey") +
geom_beeswarm(
  data=nu,
  groupOnX=TRUE,
  aes(x=Pos, y=`TS%`), color="red")
```



Ooof. Best we can muster is middle of the fat part. Who is that?

This is where we can use ggrepel. Let's add a text layer and label the dots.

```
ggplot() +
  geom_beeswarm(
    data=shooters,
    groupOnX=TRUE,
    aes(x=Pos, y=`TS%`), color="grey") +
  geom_beeswarm(
    data=nu,
    groupOnX=TRUE,
    aes(x=Pos, y=`TS%`), color="red") +
  geom_text_repel(
    data=nu,
    aes(x=Pos, y=`TS%`, label=Player))
```



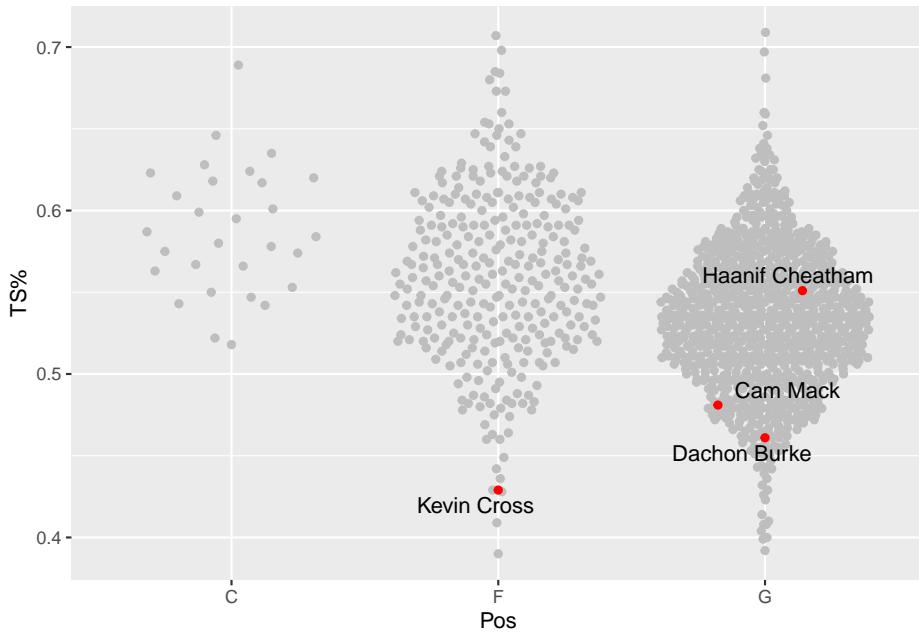
So Haanif Cheatham was our best shooter by true shooting percentage. The rest were below average shooters for that volume of shooting.

## 25.1 A few other options

The `ggbeeswarm` library has a couple of variations on the `geom_beeswarm` that may work better for your application. They are `geom_quasirandom` and `geom_jitter`.

There's not a lot to change from our example to see what they do.

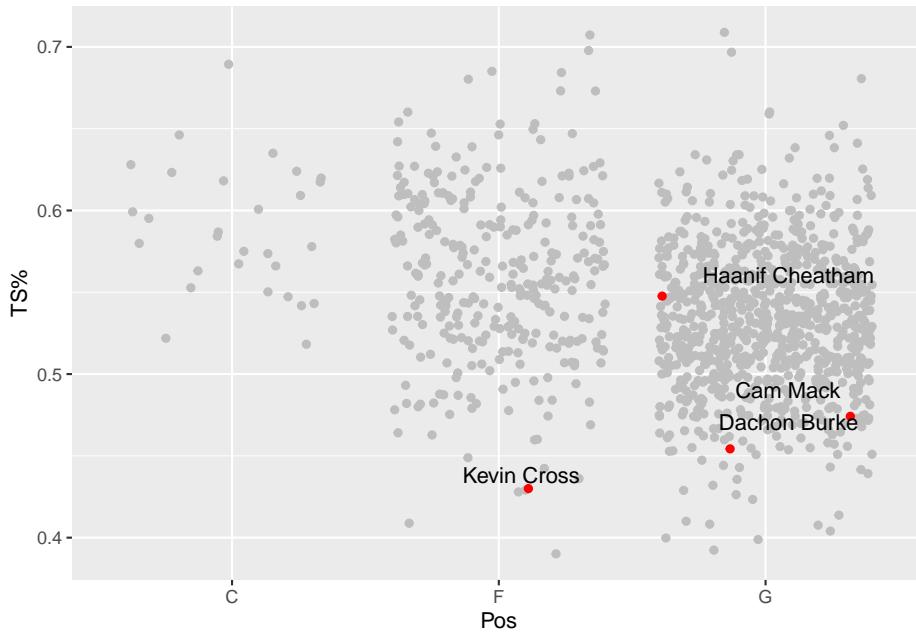
```
ggplot() +
  geom_quasirandom(
    data=shooters,
    groupOnX=TRUE,
    aes(x=Pos, y=~TS%~, color="grey") +
  geom_quasirandom(
    data=nu,
    groupOnX=TRUE,
    aes(x=Pos, y=~TS%~, color="red") +
  geom_text_repel(
    data=nu,
    aes(x=Pos, y=~TS%~, label=Player))
```



Quasirandom spreads out the dots you see in beeswarm using – you guessed it – quasirandom spacing.

For `geom_jitter`, we need to remove the `groupOnX` value. Why? No clue.

```
ggplot() +
  geom_jitter(
    data=shooters,
    aes(x=Pos, y=`TS%`), color="grey") +
  geom_jitter(
    data=nu,
    aes(x=Pos, y=`TS%`), color="red") +
  geom_text_repel(
    data=nu,
    aes(x=Pos, y=`TS%`, label=Player))
```



`geom_jitter` spreads out the dots evenly across the width of the column, randomly deciding where in the line of the true shooting percentage they appear.

Which one is right for you? You're going to have to experiment and decide. This is the art in the art and a science.



# Chapter 26

## Bump charts

The point of a bump chart is to show how the ranking of something changed over time – you could do this with the top 25 in football or basketball. I've seen it done with European soccer league standings over a season.

The requirements are that you have a row of data for a team, in that week, with their rank.

This is another extension to ggplot, and you'll install it the usual way:

```
install.packages("ggbump")
```

```
library(tidyverse)
library(ggbump)
```

Let's use last season's college football playoff rankings (this year wasn't done as of this writing):

```
rankings <- read_csv("data/cfbranking.csv")
```

```
## 
## -- Column specification -----
## cols(
##   Rank = col_double(),
##   Team = col_character(),
##   Record = col_character(),
##   Week = col_double(),
##   ShortTeam = col_character()
## )
```

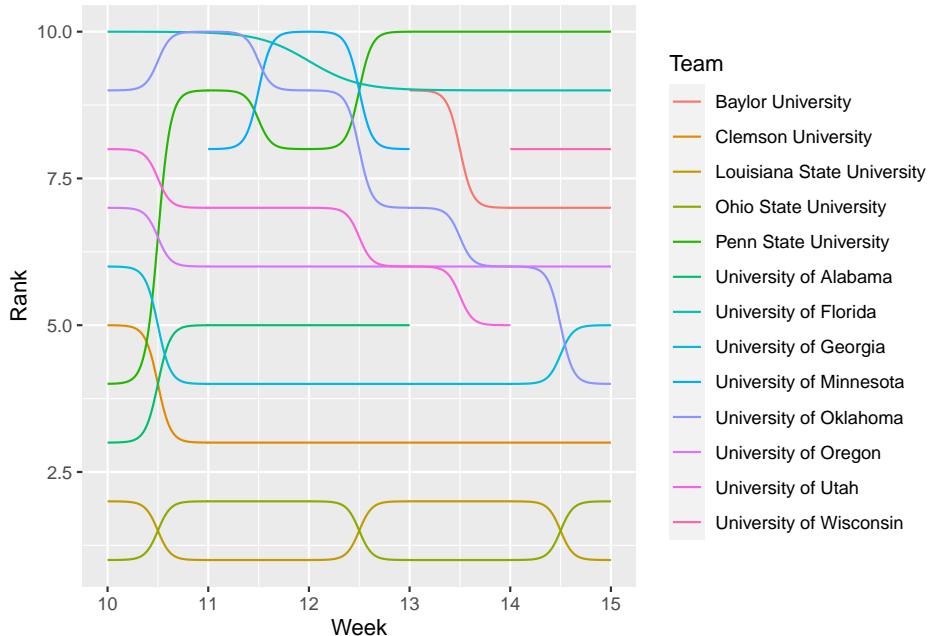
Given our requirements of a row of data for a team, in that week, with their rank, take a look at the data provided. We have 5 weeks of playoff rankings, so we should have five rows of LSU, and five rows of Ohio State. You can see the basic look of the data by using head()

```
head(rankings)
```

```
## # A tibble: 6 x 5
##   Rank Team           Record Week ShortTeam
##   <dbl> <chr>          <chr>  <dbl> <chr>
## 1     9 Baylor University 1-Oct    13 Baylor
## 2     7 Baylor University 1-Nov    14 Baylor
## 3     7 Baylor University 2-Nov    15 Baylor
## 4     5 Clemson University Sep-00 10 Clem.
## 5     3 Clemson University Oct-00 11 Clem.
## 6     3 Clemson University Nov-00 12 Clem.
```

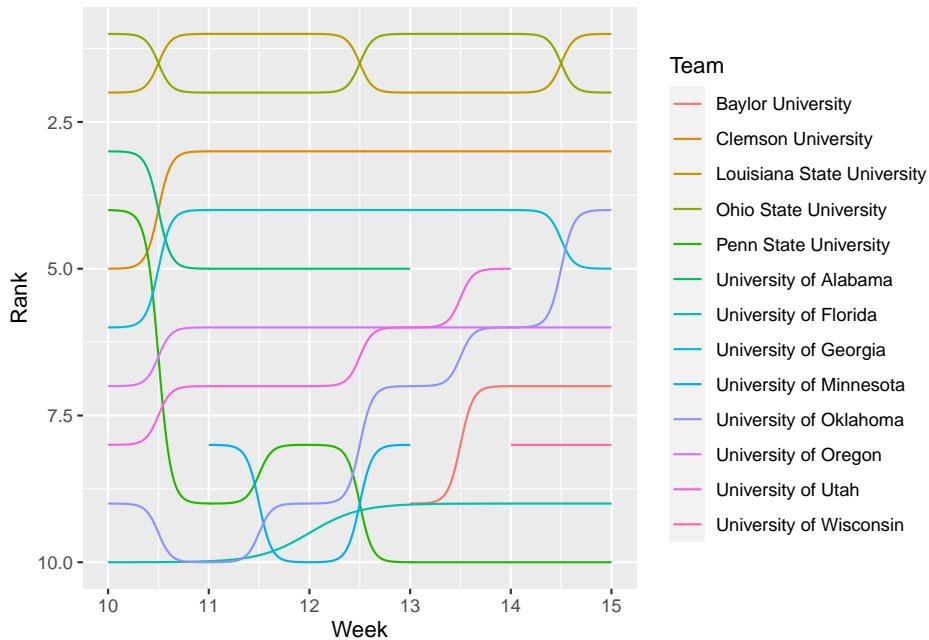
So Baylor was ranked in the 13th, 14th and 15th week, 9th, 7th and 7th, respectively. So our data is in the form we need it to be. Now we can make a bump chart. We'll start simple.

```
ggplot() + geom_bump(data=rankings, aes(x=Week, y=Rank, color=Team))
```



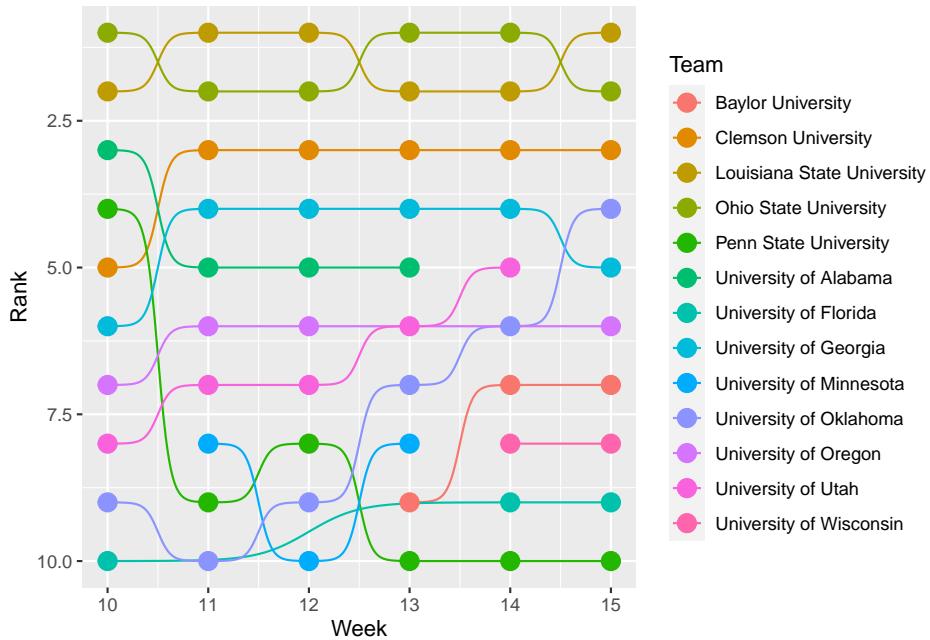
Well, it's a start. I'm immediately annoyed by the top teams being at the bottom. I learned a neat trick from ggbump that's been in ggplot all along – `scale_y_reverse()`

```
ggplot() + geom_bump(data=rankings, aes(x=Week, y=Rank, color=Team)) + scale_y_reverse()
```



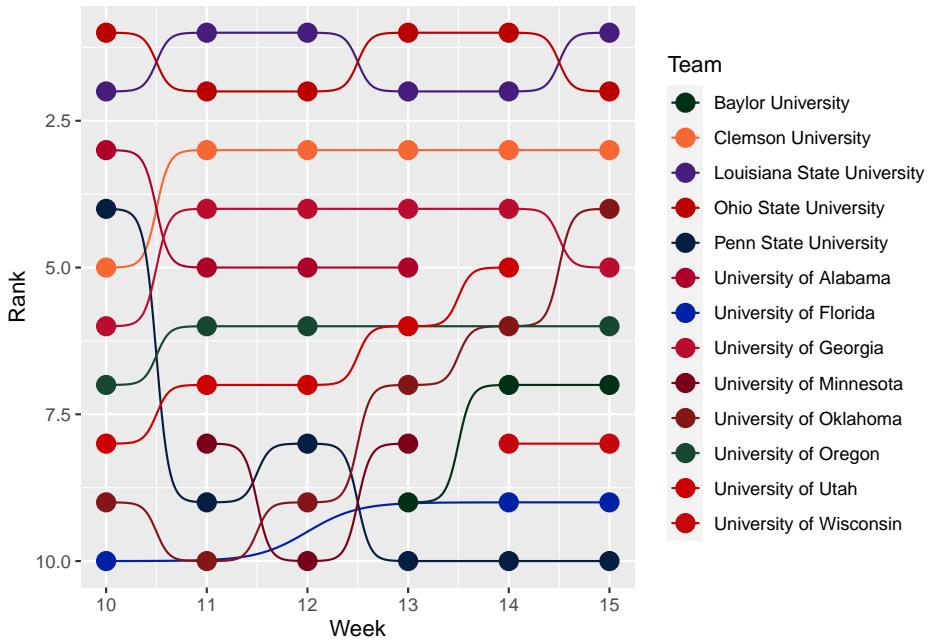
Better. But, still not great. Let's add a point at each week.

```
ggplot() +
  geom_bump(data=rankings, aes(x=Week, y=Rank, color=Team)) +
  geom_point(data=rankings, aes(x=Week, y=Rank, color=Team), size = 4) +
  scale_y_reverse()
```



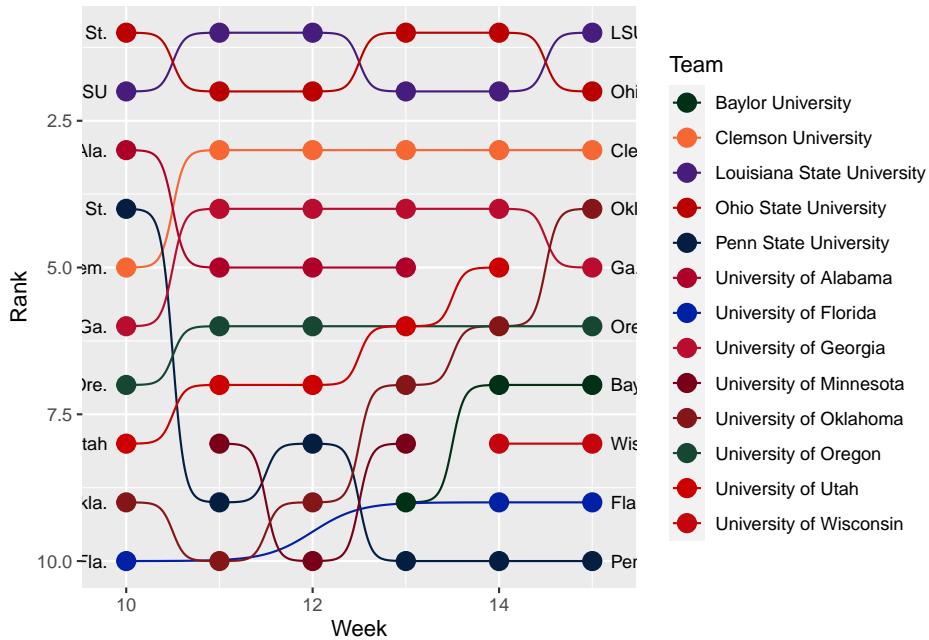
Another step. That makes it more subway map like. But the colors are all wrong. To fix this, we're going to use `scale_color_manual` and we're going to Google the hex codes for each team. The legend will tell you what order your `scale_color_manual` needs to be.

```
ggplot() +
  geom_bump(data=rankings, aes(x=Week, y=Rank, color=Team)) +
  geom_point(data=rankings, aes(x=Week, y=Rank, color=Team), size = 4) +
  scale_color_manual(values = c("#003015", "#F66733", "#461D7C", "#bb0000", "#041E42",
    "#0072BD", "#800080", "#808000", "#800000", "#808080", "#804000", "#800080")) +
  scale_y_reverse()
```



Another step. But the legend is annoying. And trying to find which red is Alabama vs Ohio State is hard. So what if we labeled each dot at the beginning and end? We can do that with some clever usage of `geom_text` and a little `dplyr` filtering inside the data step. We filter out the first and last weeks, then use `hjust` – horizontal justification – to move them left or right.

```
ggplot() +
  geom_bump(data=rankings, aes(x=Week, y=Rank, color=Team)) +
  geom_point(data=rankings, aes(x=Week, y=Rank, color=Team), size = 4) +
  geom_text(data = rankings %>% filter(Week == min(Week)), aes(x = Week - .2, y=Rank, label = Sh
  geom_text(data = rankings %>% filter(Week == max(Week)), aes(x = Week + .2, y=Rank, label = Sh
  scale_color_manual(values = c("#003015", "#F66733", "#461D7C", "#bb0000", "#041E42", "#AF002A", "
  scale_y_reverse()
```

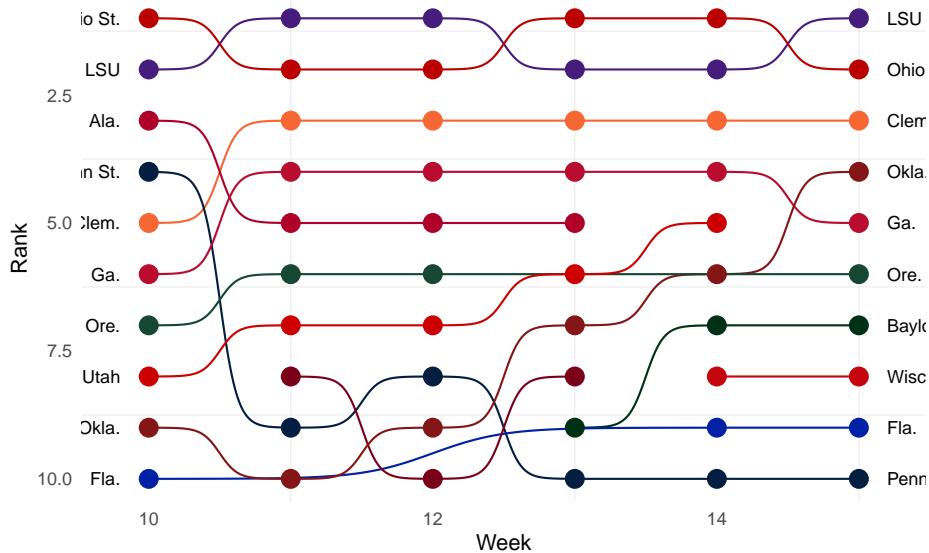


Better, but the legend is still there. We can drop it in a theme directive by saying `legend.position = "none"`. We'll also throw a `theme_minimal` on there to drop the default grey, and we'll add some better labeling.

```
ggplot() +
  geom_bump(data=rankings, aes(x=Week, y=Rank, color=Team)) +
  geom_point(data=rankings, aes(x=Week, y=Rank, color=Team), size = 4) +
  geom_text(data = rankings %>% filter(Week == min(Week)), aes(x = Week - .2, y=Rank, color=Team),
            hjust=1) +
  geom_text(data = rankings %>% filter(Week == max(Week)), aes(x = Week + .2, y=Rank, color=Team),
            hjust=-1) +
  labs(title="The race to the playoffs", subtitle="LSU and Ohio State were never out of contention",
       theme_minimal() +
       theme(
         legend.position = "none",
         panel.grid.major = element_blank()
       ) +
       scale_color_manual(values = c("#003015", "#F66733", "#461D7C", "#bb0000", "#041E42", "#0072BD"),
       scale_y_reverse())
```

### The race to the playoffs

LSU and Ohio State were never out of the top two spots. LSU deserved it.

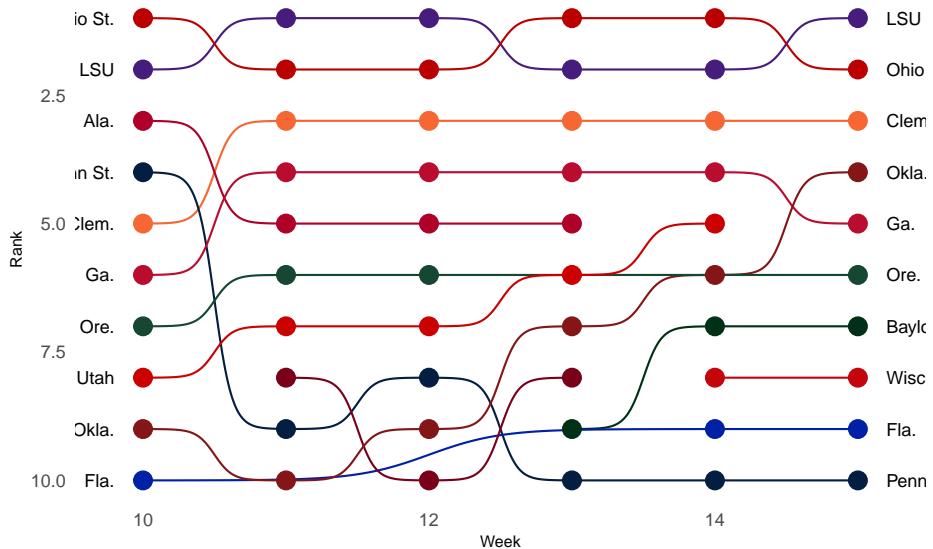


Now let's fix our text hierarchy.

```
ggplot() +
  geom_bump(data=rankings, aes(x=Week, y=Rank, color=Team)) +
  geom_point(data=rankings, aes(x=Week, y=Rank, color=Team), size = 4) +
  geom_text(data = rankings %>% filter(Week == min(Week)), aes(x = Week - .2, y=Rank, label = Shad),
            geom_text(data = rankings %>% filter(Week == max(Week)), aes(x = Week + .2, y=Rank, label = Shad),
            labs(title="The race to the playoffs", subtitle="LSU and Ohio State were never out of the top two spots",
            theme_minimal() +
            theme(
              legend.position = "none",
              panel.grid.major = element_blank(),
              plot.title = element_text(size = 16, face = "bold"),
              axis.title = element_text(size = 8),
              plot.subtitle = element_text(size=10),
              panel.grid.minor = element_blank()
            ) +
            scale_color_manual(values = c("#003015","#F66733", "#461D7C", "#bb0000", "#041E42", "#AF002A", "#002B36", "#333399", "#660033", "#990066"),
            scale_y_reverse()
```

## The race to the playoffs

LSU and Ohio State were never out of the top two spots. LSU deserved it.



And the last thing: anyone else annoyed at 7.5th place on the left? We can fix that too by specifying the breaks in `scale_y_reverse`. We can do that with the x axis as well, but since we haven't reversed it, we do that in `scale_x_continuous` with the same breaks. Also: forgot my source and credit line.

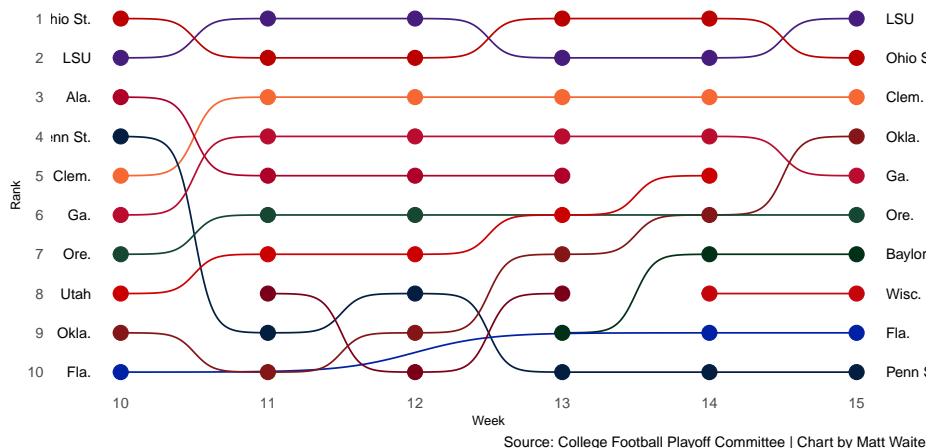
One last thing: Let's change the width of the chart to make Ohio State and Penn State fit. We can do that by adding `fig.width=X` in the `{r}` setup in your block. So something like this:

```
{r fig.width=8}
ggplot() +
  geom_bump(data=rankings, aes(x=Week, y=Rank, color=Team)) +
  geom_point(data=rankings, aes(x=Week, y=Rank, color=Team), size = 4) +
  geom_text(data = rankings %>% filter(Week == min(Week)), aes(x = Week - .2, y=Rank, color=Team),
  geom_text(data = rankings %>% filter(Week == max(Week)), aes(x = Week + .2, y=Rank, color=Team),
  labs(title="The race to the playoffs", subtitle="LSU and Ohio State were never out of the top two spots",
  theme_minimal() +
  theme(
    legend.position = "none",
    panel.grid.major = element_blank(),
    plot.title = element_text(size = 16, face = "bold"),
    axis.title = element_text(size = 8),
    plot.subtitle = element_text(size=10),
    panel.grid.minor = element_blank()
  ) +
  scale_color_manual(values = c("#003015","#F66733", "#461D7C", "#bb0000", "#041E42", "#800080"))
  
```

```
scale_x_continuous(breaks=c(10,11,12,13,14,15)) +
scale_y_reverse(breaks=c(1,2,3,4,5,6,7,8,9,10))
```

### The race to the playoffs

LSU and Ohio State were never out of the top two spots. LSU deserved it.



Source: College Football Playoff Committee | Chart by Matt Waite



# Chapter 27

## Tables

But not a table. A table with features.

Sometimes, the best way to show your data is with a table – simple rows and columns. It allows a reader to compare whatever they want to compare a little easier than a graph where you've chosen what to highlight. R has a neat package called `kableExtra`.

For this assignment, we're going to need a bunch of new libraries. Go over to the console and run these:

```
install.packages("kableExtra")
install.packages("formattable")
install.packages("htmltools")
install.packages("webshot")
webshot::install_phantomjs()
```

So what does all of these libraries do? Let's gather a few and use data of every game in the last 5 years.

Load libraries.

```
library(tidyverse)
library(kableExtra)
```

And the data.

```
logs <- read_csv("data/logs1520.csv")
```

```
##
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Date = col_date(format = ""),
##   ...)
```

```

##   HomeAway = col_character(),
##   Opponent = col_character(),
##   W_L = col_character(),
##   Blank = col_logical(),
##   Team = col_character(),
##   Conference = col_character(),
##   season = col_character()
## )
## i Use `spec()` for the full column specifications.

```

Let's ask this question: Which college basketball team saw the greatest increase in three point attempts last season as a percentage of shots? The simplest way to calculate that is by percent change.

We've got a little work to do, putting together ideas we've used before. What we need to end up with is some data that looks like this:

Team	2018–2019 season threes	2019–2020 season threes	pct change
------	-------------------------	-------------------------	------------

To get that, we'll need to do some filtering to get the right seasons, some grouping and summarizing to get the right number, some pivoting to get it organized correctly so we can mutate the percent change.

```

threechange <- logs %>%
  filter(season == "2018-2019" | season == "2019-2020") %>%
  group_by(Team, Conference, season) %>%
  summarise(Total3PA = sum(Team3PA)) %>%
  pivot_wider(names_from=season, values_from = Total3PA) %>%
  mutate(PercentChange = (`2019-2020` - `2018-2019`)/`2018-2019`) %>%
  arrange(desc(PercentChange)) %>%
  ungroup() %>%
  top_n(10) # just want a top 10 list

```

```
## `summarise()` regrouping output by 'Team', 'Conference' (override with `groups` argument)
```

```
## Selecting by PercentChange
```

We've output tables to the screen a thousand times in this class with `head`, but `kable` makes them look decent with very little code.

```
threechange %>% kable()
```

Team	Conference	2018-2019	2019-2020	PercentChange
Mississippi Valley State Delta Devils	SWAC	554	837	0.5108303
Valparaiso Crusaders	MVC	585	843	0.4410256
Ball State Cardinals	MAC	621	842	0.3558776
San Jose State Spartans	MWC	641	861	0.3432137
Alabama Crimson Tide	SEC	718	957	0.3328691
Minnesota Golden Gophers	Big Ten	603	762	0.2636816
Georgia Southern Eagles	Sun Belt	631	792	0.2551506
Tennessee Tech Golden Eagles	OVC	620	776	0.2516129
San Francisco Dons	WCC	728	899	0.2348901
McNeese State Cowboys	Southland	547	675	0.2340037

So there you have it. Mississippi Valley State changed their team so much they took 51 percent more threes last season from the season before. Where did Nebraska come out? Isn't Fred Ball supposed to be a lot of threes? We ranked 111th in college basketball in terms of change from the season before. Believe it or not, Nebraska took four fewer threes last season under Fred Ball than the last season of Tim Miles.

Kable has a mountain of customization options. The good news is that it works in a very familiar pattern. We'll start with default styling.

```
threechange %>%
  kable() %>%
  kable_styling()
```

Team	Conference	2018-2019	2019-2020	PercentChange
Mississippi Valley State Delta Devils	SWAC	554	837	0.5108303
Valparaiso Crusaders	MVC	585	843	0.4410256
Ball State Cardinals	MAC	621	842	0.3558776
San Jose State Spartans	MWC	641	861	0.3432137
Alabama Crimson Tide	SEC	718	957	0.3328691
Minnesota Golden Gophers	Big Ten	603	762	0.2636816
Georgia Southern Eagles	Sun Belt	631	792	0.2551506
Tennessee Tech Golden Eagles	OVC	620	776	0.2516129
San Francisco Dons	WCC	728	899	0.2348901
McNeese State Cowboys	Southland	547	675	0.2340037

Let's do more than the defaults, which you can see are pretty decent. Let's stripe every other row with a little bit of grey, and let's smush the width of the rows.

```
threechange %>%
  kable() %>%
  kable_styling(bootstrap_options = c("striped", "condensed"))
```

Team	Conference	2018-2019	2019-2020	PercentChange
Mississippi Valley State Delta Devils	SWAC	554	837	0.5108303
Valparaiso Crusaders	MVC	585	843	0.4410256
Ball State Cardinals	MAC	621	842	0.3558776
San Jose State Spartans	MWC	641	861	0.3432137
Alabama Crimson Tide	SEC	718	957	0.3328691
Minnesota Golden Gophers	Big Ten	603	762	0.2636816
Georgia Southern Eagles	Sun Belt	631	792	0.2551506
Tennessee Tech Golden Eagles	OVC	620	776	0.2516129
San Francisco Dons	WCC	728	899	0.2348901
McNeese State Cowboys	Southland	547	675	0.2340037

Throughout the semester, we've been using color and other signals to highlight things. Let's pretend we're doing a project on Minnesota. We can use `row_spec` to highlight things.

What `row_spec` is doing here is we're specifying which row – 6 – and making all the text on that row bold. We're making the color of the text white, because we're going to set the background to a color – in this case, the hex color for Minnesota gold.

```
threechange %>%
  kable() %>%
  kable_styling(bootstrap_options = c("striped", "condensed")) %>%
  row_spec(6, bold = T, color = "white", background = "#FBB93C")
```

Team	Conference	2018-2019	2019-2020	PercentChange
Mississippi Valley State Delta Devils	SWAC	554	837	0.5108303
Valparaiso Crusaders	MVC	585	843	0.4410256
Ball State Cardinals	MAC	621	842	0.3558776
San Jose State Spartans	MWC	641	861	0.3432137
Alabama Crimson Tide	SEC	718	957	0.3328691
<b>Minnesota Golden Gophers</b>	<b>Big Ten</b>	<b>603</b>	<b>762</b>	<b>0.2636816</b>
Georgia Southern Eagles	Sun Belt	631	792	0.2551506
Tennessee Tech Golden Eagles	OVC	620	776	0.2516129
San Francisco Dons	WCC	728	899	0.2348901
McNeese State Cowboys	Southland	547	675	0.2340037

There's also something called `column_spec` where we can change the styling on individual columns. What if we wanted to make all the team names bold?

```
threechange %>%
  kable() %>%
  kable_styling(bootstrap_options = c("striped", "condensed")) %>%
```