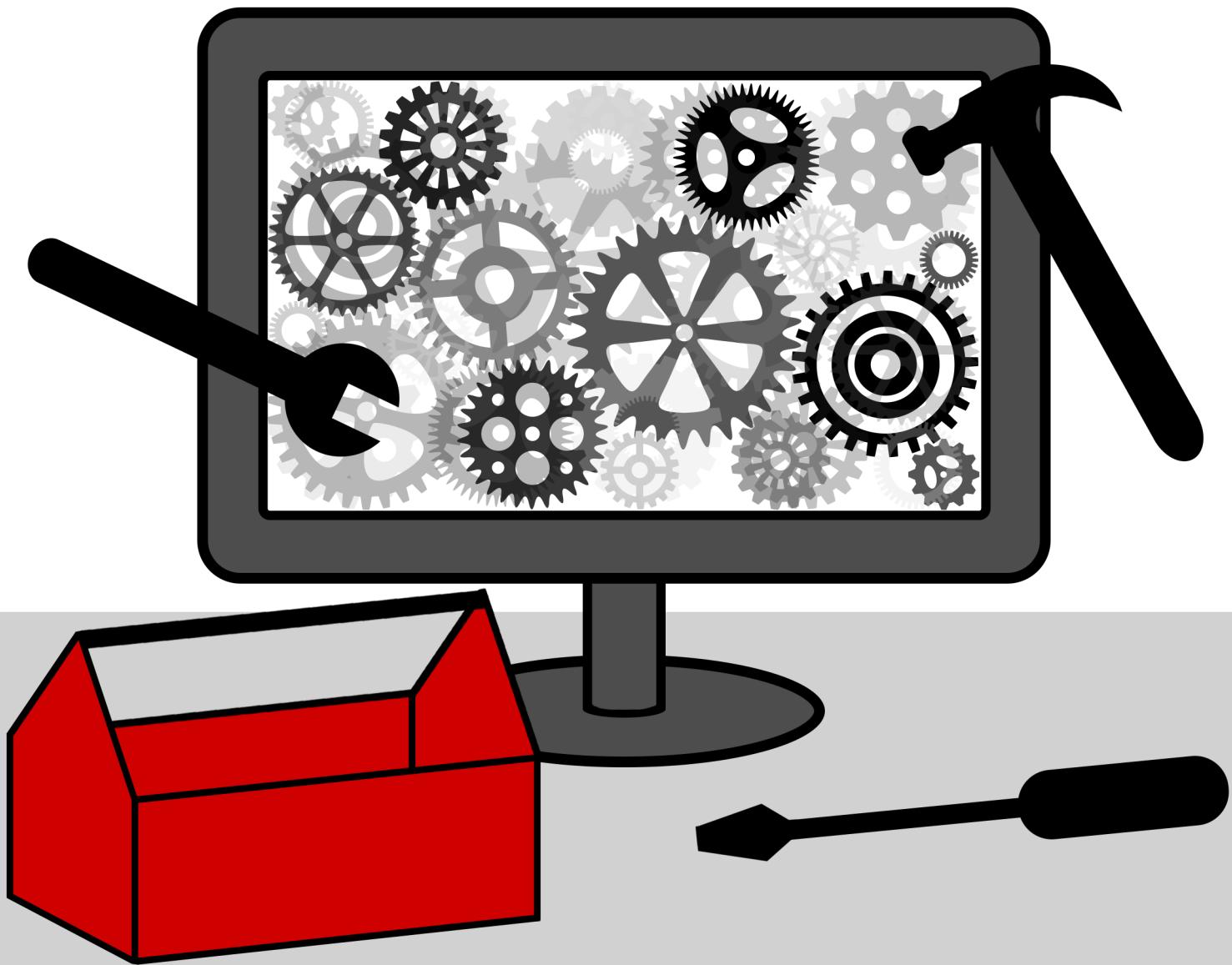


Mastering Software Development in R



Roger D Peng

Sean Kross

Brooke Anderson

Mastering Software Development in R

Roger D. Peng, Sean Kross and Brooke Anderson

This book is for sale at <http://leanpub.com/msdr>

This version was published on 2016-12-09



This is a [Leanpub](#) book. Leanpub empowers authors and publishers with the Lean Publishing process. [Lean Publishing](#) is the act of publishing an in-progress ebook using lightweight tools and many iterations to get reader feedback, pivot until you have the right book and build traction once you do.

© 2016 Roger D. Peng, Sean Kross and Brooke Anderson

Also By These Authors

Books by [Roger D. Peng](#)

[R Programming for Data Science](#)

[The Art of Data Science](#)

[Exploratory Data Analysis with R](#)

[Executive Data Science](#)

[Report Writing for Data Science in R](#)

[Conversations On Data Science](#)

Books by [Sean Kross](#)

[Developing Data Products in R](#)

Contents

Introduction	i
Setup	i
1. The R Programming Environment	1
1.1 Crash Course on R Syntax	1
1.2 The Importance of Tidy Data	12
1.3 Reading Tabular Data with the <code>readr</code> Package	14
1.4 Reading Web-Based Data	17
1.5 Basic Data Manipulation	24
1.6 Working with Dates, Times, Time Zones	42
1.7 Text Processing and Regular Expressions	52
1.8 The Role of Physical Memory	63
1.9 Working with Large Datasets	67
1.10 Diagnosing Problems	72
2. Advanced R Programming	74
2.1 Control Structures	74
2.2 Functions	78
2.3 Functional Programming	87
2.4 Expressions & Environments	99
2.5 Error Handling and Generation	105
2.6 Debugging	111
2.7 Profiling and Benchmarking	117
2.8 Non-standard Evaluation	127
2.9 Object Oriented Programming	128
2.10 Gaining Your ‘tidyverse’ Citizenship	140
3. Building R Packages	142
3.1 Before You Start	142
3.2 R Packages	143
3.3 The <code>devtools</code> Package	147
3.4 Documentation	150
3.5 Data Within a Package	160
3.6 Software Testing Framework for R Packages	162
3.7 Passing CRAN checks	165
3.8 Open Source Licensing	168

CONTENTS

3.9 Version Control and GitHub	171
3.10 Software Design and Philosophy	180
3.11 Continuous Integration	183
3.12 Cross Platform Development	185
4. Building Data Visualization Tools	192
4.1 Basic Plotting With ggplot2	192
4.2 Customizing ggplot2 Plots	217
4.3 Mapping	248
4.4 htmlWidgets	298
4.5 The grid Package	336
4.6 Building a New Theme	341
4.7 Building New Graphical Elements	349
About the Authors	360

Introduction



NOTE: This book is under active development.

This book is designed to be used in conjunction with the course sequence *Mastering Software Development in R*, available on Coursera. The book covers R software development for building data science tools. As the field of data science evolves, it has become clear that software development skills are essential for producing useful data science results and products. You will obtain rigorous training in the R language, including the skills for handling complex data, building R packages and developing custom data visualizations. You will learn modern software development practices to build tools that are highly reusable, modular, and suitable for use in a team-based environment or a community of developers.

Setup

This book makes use of the following R packages, which should be installed to take full advantage of the examples.

```
bookdown  
choroplethr  
choroplethrMaps  
data.table  
datasets  
devtools  
dlnm  
dplyr  
faraway  
GGally  
ggmap  
ggplot2  
ggthemes  
ghit  
GISTools  
grid  
gridExtra  
hexbin  
httr  
knitr  
leaflet  
lubridate  
magrittr
```

```
maps  
methods  
microbenchmark  
package  
pander  
plotly  
profvis  
pryr  
purrr  
rappdirs  
raster  
RColorBrewer  
readr  
rmarkdown  
sp  
stats  
stringr  
testthat  
tidyR  
tidyverse  
tigris  
titanic  
viridisLite
```

You can install all of these packages with the following code:

```
install.packages(c("bookdown", "choroplethr", "choroplethrMaps",  
"data.table", "datasets", "devtools", "dlnm", "dplyr", "faraway",  
"GGally", "ggmap", "ggplot2", "ggthemes", "ghit", "GISTools",  
"grid", "gridExtra", "hexbin", "httr", "knitr", "leaflet",  
"lubridate", "magrittr", "maps", "methods", "microbenchmark",  
"package", "pander", "plotly", "profvis", "pryr", "purrr",  
"rappdirs", "raster", "RColorBrewer", "readr", "rmarkdown", "sp",  
"stats", "stringr", "testthat", "tidyR", "tidyverse", "tigris",  
"titanic", "viridisLite"))
```

1. The R Programming Environment

This chapter provides a rigorous introduction to the R programming language, with a particular focus on using R for software development in a data science setting. Whether you are part of a data science team or working individually within a community of developers, this chapter will give you the knowledge of R needed to make useful contributions in those settings.

As the first chapter in this book, the chapter provides the essential foundation of R needed for the following chapters. We cover basic R concepts and language fundamentals, key concepts like tidy data and related “tidyverse” tools, processing and manipulation of complex and large datasets, handling textual data, and basic data science tasks. Upon finishing this chapter, you will have fluency at the R console and will be able to create tidy datasets from a wide range of possible data sources.

The learning objectives for this chapter are to:

- Develop fluency in using R at the console
- Execute basic arithmetic operations
- Subset and index R objects
- Remove missing values from an R object
- Modify object attributes and metadata
- Describe differences in different R classes and data types
- Read tabular data into R and read in web data via web scraping tools and APIs
- Define tidy data and to transform non-tidy data into tidy data
- Manipulate and transform a variety of data types, including dates, times, and text data
- Describe how memory is used in R sessions to store R objects
- Read and manipulate large datasets
- Describe how to diagnose programming problems and to look up answers from the web or forums

1.1 Crash Course on R Syntax

Note: Some of the material in this section is taken from [R Programming for Data Science](#).

The learning objectives for this section are to:

- Develop fluency in using R at the console
- Execute basic arithmetic operations
- Subset and index R objects
- Remove missing values from an R object

- Modify object attributes and metadata
- Describe differences in different R classes and data types

At the R prompt we type expressions. The `<-` symbol (*gets arrow*) is the assignment operator.

```
x <- 1
print(x)
[1] 1
x
[1] 1
msg <- "hello"
```

The grammar of the language determines whether an expression is complete or not.

```
x <- ## Incomplete expression
```

The `#` character indicates a comment. Anything to the right of the `#` (including the `#` itself) is ignored. This is the only comment character in R. Unlike some other languages, R does not support multi-line comments or comment blocks.

Evaluation

When a complete expression is entered at the prompt, it is evaluated and the result of the evaluated expression is returned. The result may be *auto-printed*.

```
x <- 5 ## nothing printed
x      ## auto-printing occurs
[1] 5
print(x) ## explicit printing
[1] 5
```

The `[1]` shown in the output indicates that `x` is a vector and 5 is its first element.

Typically with interactive work, we do not explicitly print objects with the `print` function; it is much easier to just auto-print them by typing the name of the object and hitting return/enter. However, when writing scripts, functions, or longer programs, there is sometimes a need to explicitly print objects because auto-printing does not work in those settings.

When an R vector is printed you will notice that an index for the vector is printed in square brackets `[]` on the side. For example, see this integer sequence of length 20.

```
x <- 11:30
x
[1] 11 12 13 14 15 16 17 18 19 20 21 22
[13] 23 24 25 26 27 28 29 30
```

The numbers in the square brackets are not part of the vector itself, they are merely part of the *printed output*.

With R, it's important that one understand that there is a difference between the actual R object and the manner in which that R object is printed to the console. Often, the printed output may have additional bells and whistles to make the output more friendly to the users. However, these bells and whistles are not inherently part of the object.

Note that the `:` operator is used to create integer sequences.

R Objects

R has five basic or “atomic” classes of objects:

- character
- numeric (real numbers)
- integer
- complex
- logical (True/False)

The most basic type of R object is a vector. Empty vectors can be created with the `vector()` function. There is really only one rule about vectors in R, which is: **A vector can only contain objects of the same class.**

But of course, like any good rule, there is an exception, which is a *list*, which we will get to a bit later. A list is represented as a vector but can contain objects of different classes. Indeed, that's usually why we use them.

There is also a class for “raw” objects, but they are not commonly used directly in data analysis and we won't cover them here.

Numbers

Numbers in R are generally treated as numeric objects (i.e. double precision real numbers). This means that even if you see a number like “1” or “2” in R, which you might think of as integers, they are likely represented behind the scenes as numeric objects (so something like “1.00” or “2.00”). This isn't important most of the time...except when it is.

If you explicitly want an integer, you need to specify the `L` suffix. So entering `1` in R gives you a numeric object; entering `1L` explicitly gives you an integer object.

There is also a special number `Inf` which represents infinity. This allows us to represent entities like `1 / 0`. This way, `Inf` can be used in ordinary calculations; e.g. `1 / Inf` is 0.

The value `NaN` represents an undefined value (“not a number”); e.g. `0 / 0`; `NaN` can also be thought of as a missing value (more on that later)

Creating Vectors

[Watch a video of this section](#)

The `c()` function can be used to create vectors of objects by concatenating things together.

```
x <- c(0.5, 0.6)      ## numeric
x <- c(TRUE, FALSE)   ## logical
x <- c(T, F)          ## logical
x <- c("a", "b", "c") ## character
x <- 9:29              ## integer
x <- c(1+0i, 2+4i)    ## complex
```

Note that in the above example, `T` and `F` are short-hand ways to specify `TRUE` and `FALSE`. However, in general one should try to use the explicit `TRUE` and `FALSE` values when indicating logical values. The `T` and `F` values are primarily there for when you're feeling lazy.

You can also use the `vector()` function to initialize vectors.

```
x <- vector("numeric", length = 10)
x
[1] 0 0 0 0 0 0 0 0 0 0
```

Mixing Objects

There are occasions when different classes of R objects get mixed together. Sometimes this happens by accident but it can also happen on purpose. So what happens with the following code?

```
y <- c(1.7, "a")    ## character
y <- c(TRUE, 2)     ## numeric
y <- c("a", TRUE)   ## character
```

In each case above, we are mixing objects of two different classes in a vector. But remember that the only rule about vectors says this is not allowed. When different objects are mixed in a vector, *coercion* occurs so that every element in the vector is of the same class.

In the example above, we see the effect of *implicit coercion*. What R tries to do is find a way to represent all of the objects in the vector in a reasonable fashion. Sometimes this does exactly what you want and...sometimes not. For example, combining a numeric object with a character object will create a character vector, because numbers can usually be easily represented as strings.

Explicit Coercion

Objects can be explicitly coerced from one class to another using the `as.*` functions, if available.

```
x <- 0:6
class(x)
[1] "integer"
as.numeric(x)
[1] 0 1 2 3 4 5 6
as.logical(x)
[1] FALSE TRUE TRUE TRUE TRUE TRUE
as.character(x)
[1] "0" "1" "2" "3" "4" "5" "6"
```

Sometimes, R can't figure out how to coerce an object and this can result in `NAs` being produced.

```
x <- c("a", "b", "c")
as.numeric(x)
Warning: NAs introduced by coercion
[1] NA NA NA
as.logical(x)
[1] NA NA NA
as.complex(x)
Warning: NAs introduced by coercion
[1] NA NA NA
```

When nonsensical coercion takes place, you will usually get a warning from R.

Matrices

Matrices are vectors with a *dimension* attribute. The dimension attribute is itself an integer vector of length 2 (number of rows, number of columns)

```
m <- matrix(nrow = 2, ncol = 3)
m
[,1] [,2] [,3]
[1,] NA NA NA
[2,] NA NA NA
dim(m)
[1] 2 3
attributes(m)
$dim
[1] 2 3
```

Matrices are constructed *column-wise*, so entries can be thought of starting in the “upper left” corner and running down the columns.

```
m <- matrix(1:6, nrow = 2, ncol = 3)
m
 [,1] [,2] [,3]
[1,]    1    3    5
[2,]    2    4    6
```

Matrices can also be created directly from vectors by adding a dimension attribute.

```
m <- 1:10
m
[1] 1 2 3 4 5 6 7 8 9 10
dim(m) <- c(2, 5)
m
 [,1] [,2] [,3] [,4] [,5]
[1,]    1    3    5    7    9
[2,]    2    4    6    8   10
```

Matrices can be created by *column-binding* or *row-binding* with the `cbind()` and `rbind()` functions.

```
x <- 1:3
y <- 10:12
cbind(x, y)
  x  y
[1,] 1 10
[2,] 2 11
[3,] 3 12
rbind(x, y)
 [,1] [,2] [,3]
x     1     2     3
y    10    11    12
```

Lists

Lists are a special type of vector that can contain elements of different classes. Lists are a very important data type in R and you should get to know them well. Lists, in combination with the various “apply” functions discussed later, make for a powerful combination.

Lists can be explicitly created using the `list()` function, which takes an arbitrary number of arguments.

```
x <- list(1, "a", TRUE, 1 + 4i)
x
[[1]]
[1] 1

[[2]]
[1] "a"

[[3]]
[1] TRUE

[[4]]
[1] 1+4i
```

We can also create an empty list of a prespecified length with the `vector()` function

```
x <- vector("list", length = 5)
x
[[1]]
NULL

[[2]]
NULL

[[3]]
NULL

[[4]]
NULL

[[5]]
NULL
```

Factors

Factors are used to represent categorical data and can be unordered or ordered. One can think of a factor as an integer vector where each integer has a *label*. Factors are important in statistical modeling and are treated specially by modelling functions like `lm()` and `glm()`.

Using factors with labels is *better* than using integers because factors are self-describing. Having a variable that has values “Male” and “Female” is better than a variable that has values 1 and 2.

Factor objects can be created with the `factor()` function.

```
x <- factor(c("yes", "yes", "no", "yes", "no"))
x
[1] yes yes no yes no
Levels: no yes
table(x)
x
no yes
2 3
## See the underlying representation of factor
unclass(x)
[1] 2 2 1 2 1
attr(,"levels")
[1] "no"  "yes"
```

Often factors will be automatically created for you when you read a dataset in using a function like `read.table()`. Those functions often, as a default, create factors when they encounter data that look like characters or strings.

The order of the levels of a factor can be set using the `levels` argument to `factor()`. This can be important in linear modelling because the first level is used as the baseline level. This feature can also be used to customize order in plots that include factors, since by default factors are plotted in the order of their levels.

```
x <- factor(c("yes", "yes", "no", "yes", "no"))
x ## Levels are put in alphabetical order
[1] yes yes no yes no
Levels: no yes
x <- factor(c("yes", "yes", "no", "yes", "no"),
            levels = c("yes", "no"))
x
[1] yes yes no yes no
Levels: yes no
```

Missing Values

Missing values are denoted by `NA` or `NaN` for undefined mathematical operations.

- `is.na()` is used to test objects if they are `NA`
- `is.nan()` is used to test for `NaN`
- `NA` values have a class also, so there are integer `NA`, character `NA`, etc.
- A `NaN` value is also `NA` but the converse is not true

```
## Create a vector with NAs in it
x <- c(1, 2, NA, 10, 3)
## Return a logical vector indicating which elements are NA
is.na(x)
[1] FALSE FALSE TRUE FALSE FALSE
## Return a logical vector indicating which elements are NaN
is.nan(x)
[1] FALSE FALSE FALSE FALSE FALSE

## Now create a vector with both NA and NaN values
x <- c(1, 2, NaN, NA, 4)
is.na(x)
[1] FALSE FALSE TRUE TRUE FALSE
is.nan(x)
[1] FALSE FALSE TRUE FALSE FALSE
```

Data Frames

Data frames are used to store tabular data in R. They are an important type of object in R and are used in a variety of statistical modeling applications. Hadley Wickham's package `dplyr` has an optimized set of functions designed to work efficiently with data frames, and `ggplot2` plotting functions work best with data stored in data frames.

Data frames are represented as a special type of list where every element of the list has to have the same length. Each element of the list can be thought of as a column and the length of each element of the list is the number of rows.

Unlike matrices, data frames can store different classes of objects in each column. Matrices must have every element be the same class (e.g. all integers or all numeric).

In addition to column names, indicating the names of the variables or predictors, data frames have a special attribute called `row.names` which indicate information about each row of the data frame.

Data frames are usually created by reading in a dataset using the `read.table()` or `read.csv()`. However, data frames can also be created explicitly with the `data.frame()` function or they can be coerced from other types of objects like lists.

Data frames can be converted to a matrix by calling `data.matrix()`. While it might seem that the `as.matrix()` function should be used to coerce a data frame to a matrix, almost always, what you want is the result of `data.matrix()`.

```
x <- data.frame(foo = 1:4, bar = c(T, T, F, F))
x
  foo    bar
1 1  TRUE
2 2  TRUE
3 3 FALSE
4 4 FALSE
nrow(x)
[1] 4
ncol(x)
[1] 2
```

Names

R objects can have names, which is very useful for writing readable code and self-describing objects. Here is an example of assigning names to an integer vector.

```
x <- 1:3
names(x)
NULL
names(x) <- c("New York", "Seattle", "Los Angeles")
x
  New York      Seattle Los Angeles
  1           2           3
names(x)
[1] "New York"     "Seattle"      "Los Angeles"
```

Lists can also have names, which is often very useful.

```
x <- list("Los Angeles" = 1, Boston = 2, London = 3)
x
$`Los Angeles`
[1] 1

$Boston
[1] 2

$London
[1] 3
names(x)
[1] "Los Angeles" "Boston"      "London"
```

Matrices can have both column and row names.

```
m <- matrix(1:4, nrow = 2, ncol = 2)
dimnames(m) <- list(c("a", "b"), c("c", "d"))
m
  c d
a 1 3
b 2 4
```

Column names and row names can be set separately using the `colnames()` and `rownames()` functions.

```
colnames(m) <- c("h", "f")
rownames(m) <- c("x", "z")
m
  h f
x 1 3
z 2 4
```

Note that for data frames, there is a separate function for setting the row names, the `row.names()` function. Also, data frames do not have column names, they just have names (like lists). So to set the column names of a data frame just use the `names()` function. Yes, I know its confusing. Here's a quick summary:

Object	Set column names	Set row names
data frame	<code>names()</code>	<code>row.names()</code>
matrix	<code>colnames()</code>	<code>rownames()</code>

Attributes

In general, R objects can have attributes, which are like metadata for the object. These metadata can be very useful in that they help to describe the object. For example, column names on a data frame help to tell us what data are contained in each of the columns. Some examples of R object attributes are

- `names`, `dimnames`
- dimensions (e.g. matrices, arrays)
- class (e.g. `integer`, `numeric`)
- `length`
- other user-defined attributes/metadata

Attributes of an object (if any) can be accessed using the `attributes()` function. Not all R objects contain attributes, in which case the `attributes()` function returns `NULL`.

Summary

There are a variety of different builtin-data types in R. In this section we have reviewed the following

- atomic classes: numeric, logical, character, integer, complex
- vectors, lists
- factors
- missing values
- data frames and matrices

All R objects can have attributes that help to describe what is in the object. Perhaps the most useful attributes are names, such as column and row names in a data frame, or simply names in a vector or list. Attributes like dimensions are also important as they can modify the behavior of objects, like turning a vector into a matrix.

1.2 The Importance of Tidy Data

The learning objectives for this section are to:

- Define tidy data and to transform non-tidy data into tidy data

One unifying concept of this book is the notion of **tidy data**. As defined by Hadley Wickham in his 2014 paper published in the *Journal of Statistical Software*, a **tidy dataset** has the following properties:

1. Each variable forms a column.
2. Each observation forms a row.
3. Each type of observational unit forms a table.

The purpose of defining tidy data is to highlight the fact that *most data do not start out life as tidy*. In fact, much of the work of data analysis may involve simply making the data tidy (at least this has been our experience). Once a dataset is tidy, it can be used as input into a variety of other functions that may transform, model, or visualize the data.

As a quick example, consider the following data illustrating death rates in Virginia in 1940 in a classic table format:

	Rural Male	Rural Female	Urban Male	Urban Female
50-54	11.7	8.7	15.4	8.4
55-59	18.1	11.7	24.3	13.6
60-64	26.9	20.3	37.0	19.3
65-69	41.0	30.9	54.6	35.1
70-74	66.0	54.3	71.1	50.0

While this format is canonical and is useful for quickly observing the relationship between multiple variables, it is not tidy. This format violates the tidy form because there are variables in both the rows and columns. In this case the variables are age category, gender, and urban-ness. Finally, the death rate itself, which is the fourth variable, is presented inside the table.

Converting this data to tidy format would give us

```

library(tidyr)
library(dplyr)

VAdeaths %>%
 tbl_df() %>%
  mutate(age = row.names(VAdeaths)) %>%
  gather(key, death_rate, -age) %>%
  separate(key, c("urban", "gender"), sep = " ") %>%
  mutate(age = factor(age), urban = factor(urban), gender = factor(gender))
# A tibble: 20 × 4
  age   urban gender death_rate
  <fctr> <fctr> <fctr>     <dbl>
1 50-54 Rural   Male    11.7
2 55-59 Rural   Male    18.1
3 60-64 Rural   Male    26.9
4 65-69 Rural   Male    41.0
5 70-74 Rural   Male    66.0
6 50-54 Rural   Female  8.7
7 55-59 Rural   Female  11.7
8 60-64 Rural   Female  20.3
9 65-69 Rural   Female  30.9
10 70-74 Rural   Female  54.3
11 50-54 Urban   Male    15.4
12 55-59 Urban   Male    24.3
13 60-64 Urban   Male    37.0
14 65-69 Urban   Male    54.6
15 70-74 Urban   Male    71.1
16 50-54 Urban   Female  8.4
17 55-59 Urban   Female  13.6
18 60-64 Urban   Female  19.3
19 65-69 Urban   Female  35.1
20 70-74 Urban   Female  50.0

```

The “Tidyverse”

There are a number of R packages that take advantage of the tidy data form and can be used to do interesting things with data. Many (but not all) of these packages are written by Hadley Wickham and the collection of packages is sometimes referred to as the “tidyverse” because of their dependence on and presumption of tidy data. “Tidyverse” packages include

- **ggplot2**: a plotting system based on the grammar of graphics
- **magrittr**: defines the `%>%` operator for chaining functions together in a series of operations on data
- **dplyr**: a suite of (fast) functions for working with data frames
- **tidyverse**: easily tidy data with `spread()` and `gather()` functions

We will be using these packages extensively in this book.

The “tidyverse” package can be used to install all of the packages in the tidyverse at once. For example, instead of starting an R script with this:

```
library(dplyr)
library(tidyr)
library(readr)
library(ggplot2)
```

You can start with this:

```
library(tidyverse)
```

1.3 Reading Tabular Data with the `readr` Package

The learning objectives for this section are to:

- Read tabular data into R and read in web data via web scraping tools and APIs

The `readr` package is the primary means by which we will read tabular data, most notably, comma-separated-value (CSV) files. The `readr` package has a few functions in it for reading and writing tabular data—we will focus on the `read_csv` function. The `readr` package is available on CRAN and the code for the package is maintained on GitHub.

The importance of the `read_csv` function is perhaps better understood from an historical perspective. R's built in `read.csv` function similarly reads CSV files, but the `read_csv` function in `readr` builds on that by removing some of the quirks and “gotchas” of `read.csv` as well as dramatically optimizing the speed with which it can read data into R. The `read_csv` function also adds some nice user-oriented features like a progress meter and a compact method for specifying column types.

The only required argument to `read_csv` is a character string specifying the path to the file to read. A typical call to `read_csv` will look as follows.

```
library(readr)
teams <- read_csv("data/team_standings.csv")
Parsed with column specification:
cols(
  Standing = col_integer(),
  Team = col_character()
)
teams
# A tibble: 32 × 2
  Standing      Team
  <int>      <chr>
1       1      Spain
2       2 Netherlands
3       3    Germany
4       4   Uruguay
5       5 Argentina
6       6     Brazil
```

```

7      7      Ghana
8      8      Paraguay
9      9      Japan
10     10     Chile
# ... with 22 more rows

```

By default, `read_csv` will open a CSV file and read it in line-by-line. It will also (by default), read in the first few rows of the table in order to figure out the type of each column (i.e. integer, character, etc.). In the code example above, you can see that `read_csv` has correctly assigned an integer class to the “Standing” variable in the input data and a character class to the “Team” variable. From the `read_csv` help page:

If [the argument for `col_types` is] ‘NULL’, all column types will be imputed from the first 1000 rows on the input. This is convenient (and fast), but not robust. If the imputation fails, you’ll need to supply the correct types yourself.

You can also specify the type of each column with the `col_types` argument. In general, it’s a good idea to specify the column types explicitly. This rules out any possible guessing errors on the part of `read_csv`. Also, specifying the column types explicitly provides a useful safety check in case anything about the dataset should change without you knowing about it.

```
teams <- read_csv("data/team_standings.csv", col_types = "cc")
```

Note that the `col_types` argument accepts a compact representation. Here “cc” indicates that the first column is `character` and the second column is `character` (there are only two columns). Using the `col_types` argument is useful because often it is not easy to automatically figure out the type of a column by looking at a few rows (especially if a column has many missing values).

The `read_csv` function will also read compressed files automatically. There is no need to decompress the file first or use the `gzfile` connection function. The following call reads a gzip-compressed CSV file containing download logs from the RStudio CRAN mirror.

```

logs <- read_csv("data/2016-07-19.csv.gz", n_max = 10)
Parsed with column specification:
cols(
  date = col_date(format = ""),
  time = col_time(format = ""),
  size = col_integer(),
  r_version = col_character(),
  r_arch = col_character(),
  r_os = col_character(),
  package = col_character(),
  version = col_character(),
  country = col_character(),
  ip_id = col_integer()
)

```

Note that the message (“Parse with column specification ...”) printed after the call indicates that `read_csv` may have had some difficulty identifying the type of each column. This can be solved by using the `col_types` argument.

```
logs <- read_csv("data/2016-07-20.csv.gz", col_types = "ccicccccci", n_max = 10)
logs
# A tibble: 10 × 10
  date     time    size r_version r_arch      r_os    package
  <chr>    <chr>   <int>    <chr>    <chr>    <chr>    <chr>
1 2016-07-20 06:04:55 144723    3.3.1    i386    mingw32  gtools
2 2016-07-20 06:04:51 2049711   3.3.0    i386    mingw32  rmarkdown
3 2016-07-20 06:04:35 26252     <NA>     <NA>     <NA>    R.methodsS3
4 2016-07-20 06:04:34 556091    2.13.1   x86_64   mingw32  tibble
5 2016-07-20 06:03:46 313363    2.13.1   x86_64   mingw32  iterators
6 2016-07-20 06:03:47 378892    3.3.1    x86_64   mingw32  foreach
7 2016-07-20 06:04:46 41228     3.3.1    x86_64   linux-gnu moments
8 2016-07-20 06:04:34 403177    <NA>     <NA>     <NA>    R.oo
9 2016-07-20 06:04:53 525       3.3.0    x86_64   linux-gnu rgl
10 2016-07-20 06:04:29 755720    3.2.5    x86_64   mingw32 geosphere
# ... with 3 more variables: version <chr>, country <chr>, ip_id <int>
```

You can specify the column type in a more detailed fashion by using the various `col_*` functions. For example, in the log data above, the first column is actually a date, so it might make more sense to read it in as a Date variable. If we wanted to just read in that first column, we could do

```
logdates <- read_csv("data/2016-07-20.csv.gz",
                      col_types = cols_only(date = col_date()),
                      n_max = 10)
logdates
# A tibble: 10 × 1
  date
  <date>
1 2016-07-20
2 2016-07-20
3 2016-07-20
4 2016-07-20
5 2016-07-20
6 2016-07-20
7 2016-07-20
8 2016-07-20
9 2016-07-20
10 2016-07-20
```

Now the `date` column is stored as a `Date` object which can be used for relevant date-related computations (for example, see the `lubridate` package).

The `read_csv` function has a `progress` option that defaults to TRUE. This option provides a nice progress meter while the CSV file is being read. However, if you are using `read_csv` in a function, or perhaps embedding it in a loop, it's probably best to set `progress = FALSE`.

The `readr` package includes a variety of functions in the `read_*` family that allow you to read in data from different formats of flat files. The following table gives a guide to several functions in the `read_*` family.

<code>readr</code> function	Use
<code>read_csv</code>	Reads comma-separated file
<code>read_csv2</code>	Reads semicolon-separated file
<code>read_tsv</code>	Reads tab-separated file
<code>read_delim</code>	General function for reading delimited files
<code>read_fwf</code>	Reads fixed width files
<code>read_log</code>	Reads log files

1.4 Reading Web-Based Data

The learning objectives for this section are to:

- Read in web data via web scraping tools and APIs

Not only can you read in data locally stored on your computer, with R it is also fairly easy to read in data stored on the web.

Flat files online

The simplest way to do this is if the data is available online as a flat file (see note below). For example, the “Extended Best Tracks” for the North Atlantic are hurricane tracks that include both the best estimate of the central location of each storm and also gives estimates of how far winds of certain speeds extended from the storm’s center in four quadrants of the storm (northeast, northwest, southeast, southwest) at each measurement point. You can see this file online [here](#).

How can you tell if you’ve found a flat file online? Here are a couple of clues:

- It will not have any formatting. Instead, it will look online as if you opened a file in a text editor on your own computer.
- It will often have a web address that ends with a typical flat file extension (“`.csv`”, “`.txt`”, or “`.fwf`”, for example).

Here are a couple of examples of flat files available online:

- Population mean county centers for Colorado counties, from the US Census
- Weather in Fort Collins, Colorado, in 2015, from Weather Underground

If you copy and paste the web address for this file, you'll see that the url for this example hurricane data file is non-secure (starts with `http:`) and that it ends with a typical flat file extension (`.txt`, in this case). You can read this file into your R session using the same `readr` function that you would use to read it in if the file were stored on your computer.

First, you can create an R object with the filepath to the file. In the case of online files, that's the url. To fit the long web address comfortably in an R script window, you can use the `paste0` function to paste pieces of the web address together:

```
ext_tracks_file <- paste0("http://rammb.cira.colostate.edu/research/",
                           "tropical_cyclones/tc_extended_best_track_dataset/",
                           "data/ebtrk_atlc_1988_2015.txt")
```

Next, since this web-based file is a fixed width file, you'll need to define the width of each column, so that R will know where to split between columns. You can then use the `read_fwf` function from the `readr` package to read the file into your R session. This data, like a lot of weather data, uses the string `"-99"` for missing data, and you can specify that missing value character with the `na` argument in `read_fwf`. Also, the online file does not include column names, so you'll have to use the [data documentation file](#) for the dataset to determine and set those yourself.

```
library(readr)

# Create a vector of the width of each column
ext_tracks_widths <- c(7, 10, 2, 2, 3, 5, 5, 6, 4, 5, 4, 4, 5, 3, 4, 3, 3, 3,
                       4, 3, 3, 3, 4, 3, 3, 3, 2, 6, 1)

# Create a vector of column names, based on the online documentation for this data
ext_tracks_colnames <- c("storm_id", "storm_name", "month", "day",
                        "hour", "year", "latitude", "longitude",
                        "max_wind", "min_pressure", "rad_max_wind",
                        "eye_diameter", "pressure_1", "pressure_2",
                        paste("radius_34", c("ne", "se", "sw", "nw"), sep = "_"),
                        paste("radius_50", c("ne", "se", "sw", "nw"), sep = "_"),
                        paste("radius_64", c("ne", "se", "sw", "nw"), sep = "_"),
                        "storm_type", "distance_to_land", "final")

# Read the file in from its url
ext_tracks <- read_fwf(ext_tracks_file,
                        fwf_widths(ext_tracks_widths, ext_tracks_colnames),
                        na = "-99")
```

```
ext_tracks[1:3, 1:9]
# A tibble: 3 × 9
  storm_id storm_name month day hour year latitude longitude max_wind
  <chr>     <chr>   <chr> <chr> <chr> <int>    <dbl>    <dbl>    <int>
1 AL0188    ALBERTO    08    05    18  1988     32.0    77.5     20
2 AL0188    ALBERTO    08    06    00  1988     32.8    76.2     20
3 AL0188    ALBERTO    08    06    06  1988     34.0    75.2     20
```

For some fixed width files, you may be able to save the trouble of counting column widths by using the `fwf_empty` function in the `readr` package. This function guesses the widths of columns based on the positions of empty columns. However, the example hurricane dataset we are using here is a bit too messy for this- in some cases, there are values from different columns that are not separated by white space. Just as it is typically safer for you to specify column types yourself, rather than relying on R to correctly guess them, it is also safer when reading in a fixed width file to specify column widths yourself.

You can use some `dplyr` functions to check out the dataset once it's in R (there will be much more about `dplyr` in the next section). For example, the following call prints a sample of four rows of data from Hurricane Katrina, with, for each row, the date and time, maximum wind speed, minimum pressure, and the radius of maximum winds of the storm for that observation:

```
library(dplyr)

ext_tracks %>%
  filter(storm_name == "KATRINA") %>%
  select(month, day, hour, max_wind, min_pressure, rad_max_wind) %>%
  sample_n(4)
# A tibble: 4 × 6
  month   day   hour max_wind min_pressure rad_max_wind
  <chr> <chr> <chr>   <int>      <int>        <int>
1 10     29    12       30      1000         40
2 10     28    18       30      1001         75
3 08     29    06      125      913         20
4 08     24    18       40      1003         55
```

With the functions in the `readr` package, you can also read in flat files from secure urls (ones that starts with `https:`). (This is not true with the `read.table` family of functions from base R.) One example where it is common to find flat files on secure sites is on GitHub. If you find a file with a flat file extension in a GitHub repository, you can usually click on it and then choose to view the “Raw” version of the file, and get to the flat file version of the file.

For example, the CDC Epidemic Prediction Initiative has a GitHub repository with data on Zika cases, including files on cases in Brazil. When we wrote this, the most current file was available [here](#), with the raw version (i.e., a flat file) available [by clicking the “Raw” button on the top right of the first site](#).

```

zika_file <- paste0("https://raw.githubusercontent.com/cdcepi/zika/master/",
                      "Brazil/COES_Microcephaly/data/COES_Microcephaly-2016-06-25.csv")
zika_brazil <- read_csv(zika_file)

zika_brazil %>%
  select(location, value, unit)
# A tibble: 210 × 3
  location     value   unit
  <chr>       <int>  <chr>
1 Brazil-Acre      2 cases
2 Brazil-Alagoas    75 cases
3 Brazil-Amapa       7 cases
4 Brazil-Amazonas     8 cases
5 Brazil-Bahia      263 cases
6 Brazil-Ceara      124 cases
7 Brazil-Distrito_Federal 5 cases
8 Brazil-Espirito_Santo 13 cases
9 Brazil-Goias        14 cases
10 Brazil-Maranhao     131 cases
# ... with 200 more rows

```

Requesting data through a web API

Web APIs are growing in popularity as a way to access open data from government agencies, companies, and other organizations. “API” stands for “Application Program Interface”; an API provides the rules for software applications to interact. In the case of open data APIs, they provide the rules you need to know to write R code to request and pull data from the organization’s web server into your R session. Usually, some of the computational burden of querying and subsetting the data is taken on by the source’s server, to create the subset of requested data to pass to your computer. In practice, this means you can often pull the subset of data you want from a very large available dataset without having to download the full dataset and load it locally into your R session.

As an overview, the basic steps for accessing and using data from a web API when working in R are:

- Figure out the API rules for HTTP requests
- Write R code to create a request in the proper format
- Send the request using GET or POST HTTP methods
- Once you get back data from the request, parse it into an easier-to-use format if necessary

To get the data from an API, you should first read the organization’s API documentation. An organization will post details on what data is available through their API(s), as well as how to set up HTTP requests to get that data— to request the data through the API, you will typically need to send the organization’s web server an HTTP request using a GET or POST method. The API documentation details will typically show an example GET or POST request

for the API, including the base URL to use and the possible query parameters that can be used to customize the dataset request.

For example, the National Aeronautics and Space Administration (NASA) has an API for pulling [the Astronomy Picture of the Day](#). In [their API documentation](#), they specify that the base URL for the API request should be “<https://api.nasa.gov/planetary/apod>” and that you can include parameters to specify the date of the daily picture you want, whether to pull a high-resolution version of the picture, and a NOAA API key you have requested from NOAA.

Many organizations will require you to get an API key and use this key in each of your API requests. This key allows the organization to control API access, including enforcing rate limits per user. API rate limits restrict how often you can request data (e.g., an hourly limit of 1,000 requests per user for NASA APIs).

API keys should be kept private, so if you are writing code that includes an API key, be very careful not to include the actual key in any code made public (including any code in public GitHub repositories). One way to do this is to save the value of your key in a file named `.Renviron` in your home directory. This file should be a plain text file and must end in a blank line. Once you’ve saved your API key to a global variable in that file (e.g., with a line added to the `.Renviron` file like `NOAA_API_KEY="abdafjsiopnab038"`), you can assign the key value to an R object in an R session using the `sys.getenv` function (e.g., `noaa_api_key <- sys.getenv("NOAA_API_KEY")`), and then use this object (`noaa_api_key`) anywhere you would otherwise have used the character string with your API key.



To find more R packages for accessing and exploring open data, check out the [Open Data CRAN task view](#). You can also browse through the [ROpenSci packages](#), all of which have GitHub repositories where you can further explore how each package works. ROpenSci is an organization with the mission to create open software tools for science. If you create your own package to access data relevant to scientific research through an API, consider submitting it for peer-review through ROpenSci.

The `riem` package, developed by Maelle Salmon and an ROpenSci package, is an excellent and straightforward example of how you can use R to pull open data through a web API. This package allows you to pull weather data from airports around the world directly from the [Iowa Environmental Mesonet](#). To show you how to pull data into R through an API, in this section we will walk you through code in the `riem` package or code based closely on code in the package.

To get a certain set of weather data from the Iowa Environmental Mesonet, you can send an HTTP request specifying a base URL, “<https://mesonet.agron.iastate.edu/cgi-bin/request/asos.py/>”, as well as some parameters describing the subset of dataset you want (e.g., date ranges, weather variables, output format). Once you know the rules for the names and possible values of these parameters (more on that below), you can submit an HTTP GET request using the `GET` function from the `httr` package.

When you are making an HTTP request using the `GET` or `POST` functions from the `httr` package, you can include the key-value pairs for any query parameters as a list object in the `query` argument of the function. For example, suppose you want to get wind speed in miles per

hour (`data = "sped"`) for Denver, CO, (`station = "DEN"`) for the month of June 2016 (`year1 = "2016"`, `month1 = "6"`, etc.) in Denver's local time zone (`tz = "America/Denver"`) and in a comma-separated file (`format = "comma"`). To get this weather dataset, you can run:

```
library(httr)
meso_url <- "https://mesonet.agron.iastate.edu/cgi-bin/request/asos.py/"
denver <- GET(url = meso_url,
              query = list(station = "DEN",
                           data = "sped",
                           year1 = "2016",
                           month1 = "6",
                           day1 = "1",
                           year2 = "2016",
                           month2 = "6",
                           day2 = "30",
                           tz = "America/Denver",
                           format = "comma")) %>%
content() %>%
read_csv(skip = 5, na = "M")

denver %>% slice(1:3)
# A tibble: 3 × 3
  station      valid    sped
  <chr>     <dttm> <dbl>
1 DEN 2016-06-01 01:00:00 6.9
2 DEN 2016-06-01 01:05:00 6.9
3 DEN 2016-06-01 01:10:00 6.9
```

The `content` call in this code extracts the content from the response to the HTTP request sent by the `GET` function. The Iowa Environmental Mesonet API offers the option to return the requested data in a comma-separated file (`format = "comma"` in the `GET` request), so here `content` and `read_csv` are used to extract and read in that csv file. Usually, data will be returned in a JSON format instead. We include more details later in this section on parsing data returned in a JSON format.

The only tricky part of this process is figuring out the available parameter names (e.g., `station`) and possible values for each (e.g., "DEN" for Denver). Currently, the details you can send in an HTTP request through Iowa Environmental Mesonet's API include:

- A four-character weather station identifier (`station`)
- The weather variables (e.g., temperature, wind speed) to include (`data`)
- Starting and ending dates describing the range for which you'd like to pull data (`year1`, `month1`, `day1`, `year2`, `month2`, `day2`)
- The time zone to use for date-times for the weather observations (`tz`)

- Different formatting options (e.g., delimiter to use in the resulting data file [`format`], whether to include longitude and latitude)

Typically, these parameter names and possible values are explained in the API documentation. In some cases, however, the documentation will be limited. In that case, you may be able to figure out possible values, especially if the API specifies a GET rather than POST method, by playing around with the website's point-and-click interface and then looking at the url for the resulting data pages. For example, if you look at the [Iowa Environmental Mesonet's page for accessing this data](#), you'll notice that the point-and-click web interface allows you the options in the list above, and if you click through to access a dataset using this interface, the web address of the data page includes these parameter names and values.

The `riem` package implements all these ideas in three very clean and straightforward functions. You can explore the code behind this package and see how these ideas can be incorporated into a small R package, in the `/R` directory of the [package's GitHub page](#).

R packages already exist for many open data APIs. If an R package already exists for an API, you can use functions from that package directly, rather than writing your own code using the API protocols and `httr` functions. Other examples of existing R packages to interact with open data APIs include:

- `twitteR`: Twitter
- `rnoaa`: National Oceanic and Atmospheric Administration
- `Quandl`: Quandl (financial data)
- `RGoogleAnalytics`: Google Analytics
- `censusr`, `acs`: United States Census
- `WDI`, `wbstats`: World Bank
- `GuardianR`, `rdian`: The Guardian Media Group
- `blsAPI`: Bureau of Labor Statistics
- `rtimes`: New York Times
- `dataRetrieval`, `waterData`: United States Geological Survey

If an R package doesn't exist for an open API and you'd like to write your own package, find out more about writing API packages with [this vignette for the httr package](#). This document includes advice on error handling within R code that accesses data through an open API.

Scraping web data

You can also use R to pull and clean web-based data that is not accessible through a web API or as an online flat file. In this case, the strategy will often be to pull in the full web page file (often in HTML or XML) and then parse or clean it within R.

The `rvest` package is a good entry point for handling more complex collection and cleaning of web-based data. This package includes functions, for example, that allow you to select certain elements from the code for a web page (e.g., using the `html_node` and `xml_node`

functions), to parse tables in an HTML document into R data frames (`html_table`), and to parse, fill out, and submit HTML forms (`html_form`, `set_values`, `submit_form`). Further details on web scraping with R are beyond the scope of this course, but if you’re interested, you can find out more through the [rvest GitHub README](#).

Parsing JSON, XML, or HTML data

Often, data collected from the web, including the data returned from an open API or obtained by scraping a web page, will be in JSON, XML, or HTML format. To use data in a JSON, XML, or HTML format in R, you need to parse the file from its current format and convert it into an R object more useful for analysis.

Typically, JSON-, XML-, or HTML-formatted data is parsed into a list in R, since list objects allow for a lot of flexibility in the structure of the data. However, if the data is structured appropriately, you can often parse data into another type of object (a data frame, for example, if the data fits well into a two-dimensional format of rows and columns). If the data structure of the data that you are pulling in is complex but consistent across different observations, you may alternatively want to create a custom object type to parse the data into.

There are a number of packages for parsing data from these formats, including `jsonlite` and `xmll2`. To find out more about parsing data from typical web formats, and for more on working with web-based documents and data, see the [CRAN task view for Web Technologies and Services](#)

1.5 Basic Data Manipulation

The learning objectives for this section are to:

- Transform non-tidy data into tidy data
- Manipulate and transform a variety of data types, including dates, times, and text data

The two packages `dplyr` and `tidyverse`, both “tidyverse” packages, allow you to quickly and fairly easily clean up your data. These packages are not very old, and so much of the example R code you might find in books or online might not use the functions we use in examples in this section (although this is quickly changing for new books and for online examples). Further, there are many people who are used to using R base functionality to clean up their data, and some of them still do not use these packages much when cleaning data. We think, however, that `dplyr` is easier for people new to R to learn than learning how to clean up data using base R functions, and we also think it produces code that is much easier to read, which is useful in maintaining and sharing code.

For many of the examples in this section, we will use the `ext_tracks hurricane` dataset we input from a url as an example in a previous section of this book. If you need to load a version of that data, we have also saved it locally, so you can create an R object with the example data for this section by running:

```

ext_tracks_file <- "data/ebtrk_atlc_1988_2015.txt"
ext_tracks_widths <- c(7, 10, 2, 2, 3, 5, 5, 6, 4, 5, 4, 4, 5, 3, 4, 3, 3, 3,
                      4, 3, 3, 3, 4, 3, 3, 3, 2, 6, 1)
ext_tracks_colnames <- c("storm_id", "storm_name", "month", "day",
                        "hour", "year", "latitude", "longitude",
                        "max_wind", "min_pressure", "rad_max_wind",
                        "eye_diameter", "pressure_1", "pressure_2",
                        paste("radius_34", c("ne", "se", "sw", "nw"), sep = "_"),
                        paste("radius_50", c("ne", "se", "sw", "nw"), sep = "_"),
                        paste("radius_64", c("ne", "se", "sw", "nw"), sep = "_"),
                        "storm_type", "distance_to_land", "final")
ext_tracks <- read_fwf(ext_tracks_file,
                        fwf_widths(ext_tracks_widths, ext_tracks_colnames),
                        na = "-99")

```

Piping

The `dplyr` and `tidyverse` functions are often used in conjunction with piping, which is done with the `%>%` function from the `magrittr` package. Piping can be done with many R functions, but is especially common with `dplyr` and `tidyverse` functions. The concept is straightforward—the pipe passes the data frame output that results from the function right before the pipe to input it as the first argument of the function right after the pipe.

Here is a generic view of how this works in code, for a pseudo-function named `function` that inputs a data frame as its first argument:

```

# Without piping
function(dataframe, argument_2, argument_3)

# With piping
dataframe %>%
  function(argument_2, argument_3)

```

For example, without piping, if you wanted to see the time, date, and maximum winds for Katrina from the first three rows of the `ext_tracks` hurricane data, you could run:

```

katrina <- filter(ext_tracks, storm_name == "KATRINA")
katrina_reduced <- select(katrina, month, day, hour, max_wind)
head(katrina_reduced, 3)
# A tibble: 3 × 4
  month   day   hour max_wind
  <chr> <chr> <chr>    <int>
1     10    28     18      30
2     10    29     00      30
3     10    29     06      30

```

In this code, you are creating new R objects at each step, which makes the code cluttered and also requires copying the data frame several times into memory. As an alternative, you could just wrap one function inside another:

```
head(select(filter(ext_tracks, storm_name == "KATRINA"),
            month, day, hour, max_wind), 3)
# A tibble: 3 × 4
  month   day   hour max_wind
  <chr> <chr> <chr>    <int>
1 10     28     18      30
2 10     29     00      30
3 10     29     06      30
```

This avoids re-assigning the data frame at each step, but quickly becomes ungainly, and it's easy to put arguments in the wrong layer of parentheses. Piping avoids these problems, since at each step you can send the output from the last function into the next function as that next function's first argument:

```
ext_tracks %>%
  filter(storm_name == "KATRINA") %>%
  select(month, day, hour, max_wind) %>%
  head(3)
# A tibble: 3 × 4
  month   day   hour max_wind
  <chr> <chr> <chr>    <int>
1 10     28     18      30
2 10     29     00      30
3 10     29     06      30
```

Summarizing data

The `dplyr` and `tidyverse` packages have numerous functions (sometimes referred to as “verbs”) for cleaning up data. We’ll start with the functions to summarize data.

The primary of these is `summarize`, which inputs a data frame and creates a new data frame with the requested summaries. In conjunction with `summarize`, you can use other functions from `dplyr` (e.g., `n`, which counts the number of observations in a given column) to create this summary. You can also use R functions from other packages or base R functions to create the summary.

For example, say we want a summary of the number of observations in the `ext_tracks` hurricane dataset, as well as the highest measured maximum windspeed (given by the column `max_wind` in the dataset) in any of the storms, and the lowest minimum pressure (`min_pressure`). To create this summary, you can run:

```
ext_tracks %>%
  summarize(n_obs = n(),
            worst_wind = max(max_wind),
            worst_pressure = min(min_pressure))
# A tibble: 1 × 3
  n_obs worst_wind worst_pressure
  <int>     <int>          <int>
1 11824        160             0
```

This summary provides particularly useful information for this example data, because it gives an unrealistic value for minimum pressure (0 hPa). This shows that this dataset will need some cleaning. The highest wind speed observed for any of the storms, 160 knots, is more reasonable.

You can also use `summarize` with functions you've written yourself, which gives you a lot of power in summarizing data in interesting ways. As a simple example, if you wanted to present the maximum wind speed in the summary above using miles per hour rather than knots, you could write a function to perform the conversion, and then use that function within the `summarize` call:

```
knots_to_mph <- function(knots){
  mph <- 1.152 * knots
}

ext_tracks %>%
  summarize(n_obs = n(),
            worst_wind = knots_to_mph(max(max_wind)),
            worst_pressure = min(min_pressure))
# A tibble: 1 × 3
  n_obs worst_wind worst_pressure
  <int>     <dbl>          <int>
1 11824      184.32             0
```

So far, we've only used `summarize` to create a single-line summary of the data frame. In other words, the summary functions are applied across the entire dataset, to return a single value for each summary statistic. However, often you might want summaries stratified by a certain grouping characteristic of the data. For the hurricane data, for example, you might want to get the worst wind and worst pressure by storm, rather than across all storms.

You can do this by grouping your data frame by one of its column variables, using the function `group_by`, and then using `summarize`. The `group_by` function does not make a visible change to a data frame, although you can see, if you print out a grouped data frame, that the new grouping variable will be listed under "Groups" at the top of a print-out:

```
ext_tracks %>%
  group_by(storm_name, year) %>%
  head()
Source: local data frame [6 x 29]
Groups: storm_name, year [1]

  storm_id storm_name month day hour year latitude longitude max_wind
  <chr>     <chr> <chr> <chr> <chr> <int>    <dbl>    <dbl>    <int>
1 AL0188   ALBERTO    08   05   18  1988     32.0    77.5     20
2 AL0188   ALBERTO    08   06   00  1988     32.8    76.2     20
3 AL0188   ALBERTO    08   06   06  1988     34.0    75.2     20
4 AL0188   ALBERTO    08   06   12  1988     35.2    74.6     25
5 AL0188   ALBERTO    08   06   18  1988     37.0    73.5     25
6 AL0188   ALBERTO    08   07   00  1988     38.7    72.4     25
# ... with 20 more variables: min_pressure <int>, rad_max_wind <int>,
# eye_diameter <int>, pressure_1 <int>, pressure_2 <int>,
# radius_34_ne <int>, radius_34_se <int>, radius_34_sw <int>,
# radius_34_nw <int>, radius_50_ne <int>, radius_50_se <int>,
# radius_50_sw <int>, radius_50_nw <int>, radius_64_ne <int>,
# radius_64_se <int>, radius_64_sw <int>, radius_64_nw <int>,
# storm_type <chr>, distance_to_land <int>, final <chr>
```

As a note, since hurricane storm names repeat at regular intervals until they are retired, to get a separate summary for each unique storm, this example requires grouping by both `storm_name` and `year`.

Even though applying the `group_by` function does not cause a noticeable change to the data frame itself, you'll notice the difference in grouped and ungrouped data frames when you use `summarize` on the data frame. If a data frame is grouped, all summaries are calculated and given separately for each unique value of the grouping variable:

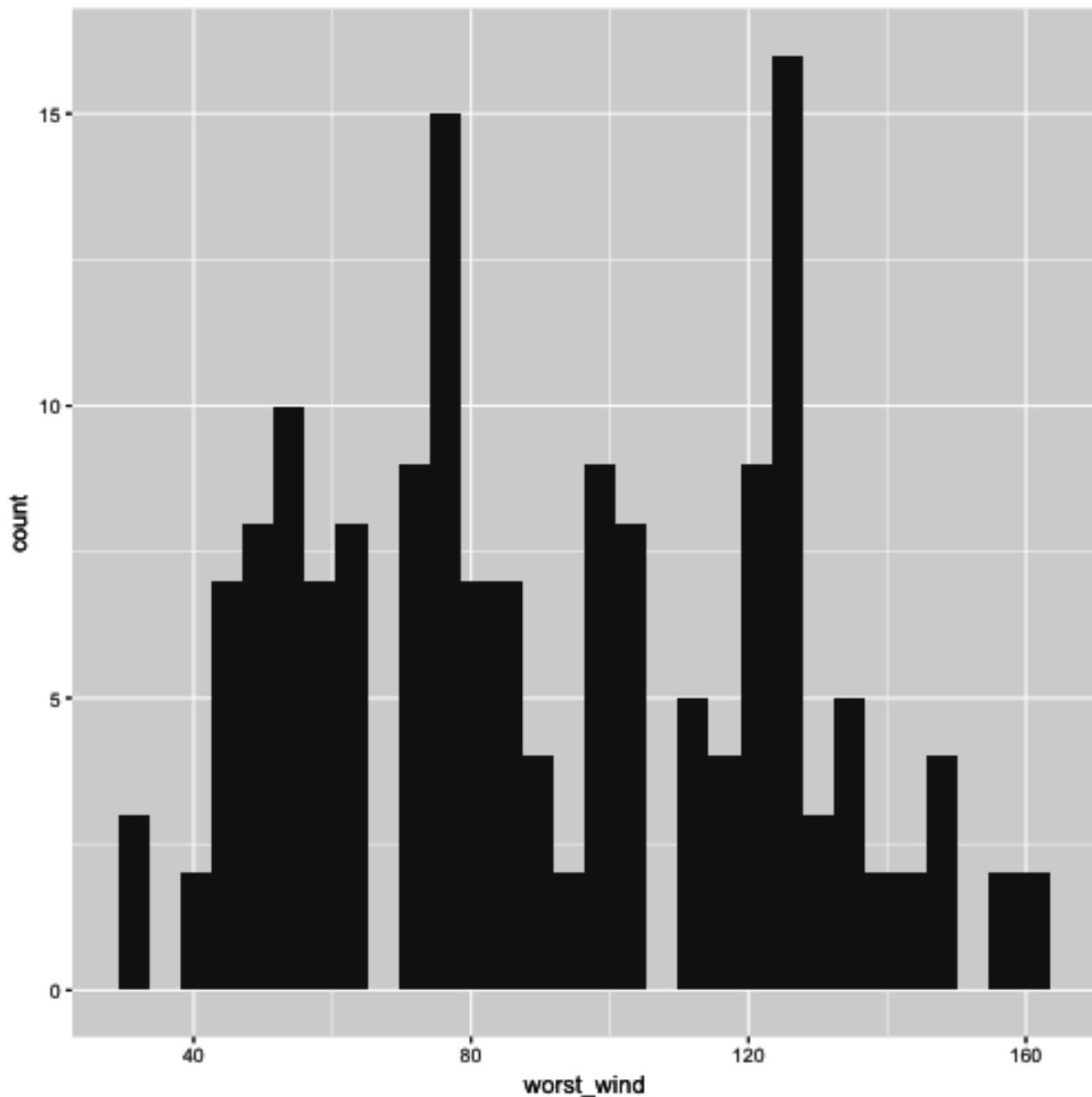
```
ext_tracks %>%
  group_by(storm_name, year) %>%
  summarize(n_obs = n(),
            worst_wind = max(max_wind),
            worst_pressure = min(min_pressure))
Source: local data frame [378 x 5]
Groups: storm_name [?]

  storm_name year n_obs worst_wind worst_pressure
  <chr>     <int> <int>      <int>        <int>
1 ALBERTO    1988    13       35        1002
2 ALBERTO    1994    31       55        993
3 ALBERTO    2000    87      110        950
4 ALBERTO    2006    37       60        969
5 ALBERTO    2012    20       50        995
6 ALEX        1998    26       45        1002
7 ALEX        2004    25      105        957
8 ALEX        2010    30       90        948
9 ALLISON    1989    28       45        999
```

```
10    ALLISON 1995    33        65        982
# ... with 368 more rows
```

This grouping / summarizing combination can be very useful for quickly plotting interesting summaries of a dataset. For example, to plot a histogram of maximum wind speed observed for each storm (Figure @ref(fig:windhistogram)), you could run:

```
library(ggplot2)
ext_tracks %>%
  group_by(storm_name) %>%
  summarize(worst_wind = max(max_wind)) %>%
  ggplot(aes(x = worst_wind)) + geom_histogram()
```



Histogram of the maximum wind speed observed during a storm for all Atlantic basin tropical storms, 1988–2015.

We will show a few basic examples of plotting using `ggplot2` functions in this chapter of the book. We will cover plotting much more thoroughly in a later section of the specialization.

From Figure @ref(fig:windhistogram), we can see that only two storms had maximum wind speeds at or above 160 knots (we'll check this later with some other `dplyr` functions).



You cannot make changes to a variable that is being used to group a data frame. If you try, you will get the error `Error: cannot modify grouping variable`. If you get this error, use the `ungroup` function to remove grouping within a data frame, and then you will be able to mutate any of the variables in the data.

Selecting and filtering data

When cleaning up data, you will need to be able to create subsets of the data, by selecting certain columns or filtering down to certain rows. These actions can be done using the `dplyr` functions `select` and `filter`.

The `select` function subsets certain columns of a data frame. The most basic way to use `select` is select certain columns by specifying their full column names. For example, to select the storm name, date, time, latitude, longitude, and maximum wind speed from the `ext_tracks` dataset, you can run:

```
ext_tracks %>%
  select(storm_name, month, day, hour, year, latitude, longitude, max_wind)
# A tibble: 11,824 × 8
  storm_name month   day   hour   year latitude longitude max_wind
    <chr>     <chr> <chr> <chr> <int>     <dbl>     <dbl>     <int>
1 ALBERTO     08     05     18  1988     32.0      77.5      20
2 ALBERTO     08     06     00  1988     32.8      76.2      20
3 ALBERTO     08     06     06  1988     34.0      75.2      20
4 ALBERTO     08     06     12  1988     35.2      74.6      25
5 ALBERTO     08     06     18  1988     37.0      73.5      25
6 ALBERTO     08     07     00  1988     38.7      72.4      25
7 ALBERTO     08     07     06  1988     40.0      70.8      30
8 ALBERTO     08     07     12  1988     41.5      69.0      35
9 ALBERTO     08     07     18  1988     43.0      67.5      35
10 ALBERTO    08     08     00  1988     45.0      65.5      35
# ... with 11,814 more rows
```

There are several functions you can use with `select` that give you more flexibility, and so allow you to select columns without specifying the full names of each column. For example, the `starts_with` function can be used within a `select` function to pick out all the columns that start with a certain text string. As an example of using `starts_with` in conjunction with `select`, in the `ext_tracks` hurricane data, there are a number of columns that say how far from the storm center winds of certain speeds extend. Tropical storms often have asymmetrical wind fields, so these wind radii are given for each quadrant of the storm (northeast, southeast, northwest, and southeast of the storm's center). All of the columns with the radius to which winds of 34 knots or more extend start with "radius_34". To get a dataset with storm names, location, and radii of winds of 34 knots, you could run:

```
ext_tracks %>%
  select(storm_name, latitude, longitude, starts_with("radius_34"))
# A tibble: 11,824 × 7
  storm_name latitude longitude radius_34_ne radius_34_se radius_34_sw
    <chr>     <dbl>     <dbl>      <int>      <int>      <int>
1 ALBERTO     32.0      77.5        0          0          0
2 ALBERTO     32.8      76.2        0          0          0
3 ALBERTO     34.0      75.2        0          0          0
4 ALBERTO     35.2      74.6        0          0          0
5 ALBERTO     37.0      73.5        0          0          0
6 ALBERTO     38.7      72.4        0          0          0
7 ALBERTO     40.0      70.8        0          0          0
8 ALBERTO     41.5      69.0       100        100         50
9 ALBERTO     43.0      67.5       100        100         50
10 ALBERTO    45.0      65.5       NA         NA         NA
# ... with 11,814 more rows, and 1 more variables: radius_34_nw <int>
```

Other functions that can be used with `select` in a similar way include:

- `ends_with`: Select all columns that end with a certain string (for example, `select(ext_tracks, ends_with("ne"))`) to get all the wind radii for the northeast quadrant of a storm for the hurricane example data)
- `contains`: Select all columns that include a certain string (`select(ext_tracks, contains("34"))`) to get all wind radii for 34-knot winds)
- `matches`: Select all columns that match a certain relative expression (`select(ext_tracks, matches("_[0-9][0-9]_"))`) to get all columns where the column name includes two numbers between two underscores, a pattern that matches all of the wind radii columns)

While `select` picks out certain columns of the data frame, `filter` picks out certain rows. With `filter`, you can specify certain conditions using R's logical operators, and the function will return rows that meet those conditions.

R's logical operators include:

Operator	Meaning	Example
<code>==</code>	Equals	<code>storm_name == KATRINA</code>
<code>!=</code>	Does not equal	<code>min_pressure != 0</code>
<code>></code>	Greater than	<code>latitude > 25</code>
<code>>=</code>	Greater than or equal to	<code>max_wind >= 160</code>
<code><</code>	Less than	<code>min_pressure < 900</code>
<code><=</code>	Less than or equal to	<code>distance_to_land <= 0</code>
<code>%in%</code>	Included in	<code>storm_name %in% c("KATRINA", "ANDREW")</code>
<code>is.na()</code>	Is a missing value	<code>is.na(radius_34_ne)</code>

If you are ever unsure of how to write a logical statement, but know how to write its opposite, you can use the `!` operator to negate the whole statement. For example, if you

wanted to get all storms *except* those named “KATRINA” and “ANDREW”, you could use `!(storm_name %in% c("KATRINA", "ANDREW"))`. A common use of this is to identify observations with non-missing data (e.g., `!(is.na(radius_34_ne))`).

A logical statement, run by itself on a vector, will return a vector of the same length with `TRUE` every time the condition is met and `FALSE` every time it is not.

```
head(ext_tracks$hour)
[1] "18" "00" "06" "12" "18" "00"
head(ext_tracks$hour == "00")
[1] FALSE TRUE FALSE FALSE FALSE TRUE
```

When you use a logical statement within `filter`, it will return just the rows where the logical statement is true:

```
ext_tracks %>%
  select(storm_name, hour, max_wind) %>%
  head(9)
# A tibble: 9 × 3
  storm_name hour max_wind
  <chr>     <chr>    <int>
1 ALBERTO    18        20
2 ALBERTO    00        20
3 ALBERTO    06        20
4 ALBERTO    12        25
5 ALBERTO    18        25
6 ALBERTO    00        25
7 ALBERTO    06        30
8 ALBERTO    12        35
9 ALBERTO    18        35

ext_tracks %>%
  select(storm_name, hour, max_wind) %>%
  filter(hour == "00") %>%
  head(3)
# A tibble: 3 × 3
  storm_name hour max_wind
  <chr>     <chr>    <int>
1 ALBERTO    00        20
2 ALBERTO    00        25
3 ALBERTO    00        35
```

Filtering can also be done after summarizing data. For example, to determine which storms had maximum wind speed equal to or above 160 knots, run:

```
ext_tracks %>%
  group_by(storm_name, year) %>%
  summarize(worst_wind = max(max_wind)) %>%
  filter(worst_wind >= 160)
Source: local data frame [2 x 3]
Groups: storm_name [2]

  storm_name   year worst_wind
  <chr>     <int>      <int>
1 GILBERT    1988        160
2 WILMA      2005        160
```

If you would like to string several logical conditions together and select rows where all or any of the conditions are true, you can use the “and” (&) or “or” (|) operators. For example, to pull out observations for Hurricane Andrew when it was at or above Category 5 strength (137 knots or higher), you could run:

```
ext_tracks %>%
  select(storm_name, month, day, hour, latitude, longitude, max_wind) %>%
  filter(storm_name == "ANDREW" & max_wind >= 137)
# A tibble: 2 × 7
  storm_name month   day   hour latitude longitude max_wind
  <chr>     <chr> <chr> <dbl>    <dbl>      <int>
1 ANDREW     08     23     12     25.4      74.2      145
2 ANDREW     08     23     18     25.4      75.8      150
```



Some common errors that come up when using logical operators in R are:

- If you want to check that two things are equal, make sure you use double equal signs (==), not a single one. At best, a single equals sign won’t work; in some cases, it will cause a variable to be re-assigned (= can be used for assignment, just like <-).
- If you are trying to check if one thing is equal to one of several things, use %in% rather than ==. For example, if you want to filter to rows of `ext_tracks` with storm names of “KATRINA” and “ANDREW”, you need to use `storm_name %in% c("KATRINA", "ANDREW")`, not `storm_name == c("KATRINA", "ANDREW")`.
- If you want to identify observations with missing values (or without missing values), you must use the `is.na` function, not == or !=. For example, `is.na(radius_34_ne)` will work, but `radius_34_ne == NA` will not.

Adding, changing, or renaming columns

The `mutate` function in `dplyr` can be used to add new columns to a data frame or change existing columns in the data frame. As an example, I’ll use the `worldcup` dataset from the package `faraway`, which statistics from the 2010 World Cup. To load this example data frame, you can run:

```
library(faraway)
data(worldcup)
```

This dataset has observations by player, including the player's team, position, amount of time played in this World Cup, and number of shots, passes, tackles, and saves. This dataset is currently not tidy, as it has one of the variables (players' names) as rownames, rather than as a column of the data frame. You can use the `mutate` function to move the player names to its own column:

```
worldcup <- worldcup %>%
  mutate(player_name = rownames(worldcup))

worldcup %>% slice(1:3)
#> # A tibble: 3 × 9
#>   Team    Position  Time  Shots  Passes  Tackles  Saves player_name
#>   <fctr>   <fctr> <int> <int>   <int>   <int>   <chr>
#> 1 Algeria Midfielder    16     0      6      0      0     Abdoun
#> 2 Japan   Midfielder   351     0    101     14      0       Abe
#> 3 France   Defender    180     0     91      6      0     Abidal
```

You can also use `mutate` in coordination with `group_by` to create new columns that give summaries within certain windows of the data. For example, the following code will add a column with the average number of shots for a player's position added as a new column. While this code is summarizing the original data to generate the values in this column, `mutate` will add these repeated summary values to the original dataset by group, rather than returning a data frame with a single row for each of the grouping variables (try replacing `mutate` with `summarize` in this code to make sure you understand the difference).

```
worldcup <- worldcup %>%
  group_by(Position) %>%
  mutate(ave_shots = mean(Shots)) %>%
  ungroup()

worldcup %>% slice(1:3)
#> # A tibble: 3 × 9
#>   Team    Position  Time  Shots  Passes  Tackles  Saves player_name
#>   <fctr>   <fctr> <int> <int>   <int>   <int>   <chr>
#> 1 Algeria Midfielder    16     0      6      0      0     Abdoun
#> 2 Japan   Midfielder   351     0    101     14      0       Abe
#> 3 France   Defender    180     0     91      6      0     Abidal
#> # ... with 1 more variables: ave_shots <dbl>
```

If there is a column that you want to rename, but not change, you can use the `rename` function. For example:

```

worldcup %>%
  rename(Name = player_name) %>%
  slice(1:3)
# A tibble: 3 × 9
  Team    Position   Time Shots Passes Tackles Saves     Name ave_shots
  <fctr>      <fctr> <int> <int>  <int>  <int> <int> <chr>      <dbl>
1 Algeria Midfielder    16     0      6      0      0 Abdoun  2.394737
2 Japan   Midfielder   351     0    101     14      0 Abe    2.394737
3 France  Defender     180     0     91      6      0 Abidal 1.164894

```

Spreading and gathering data

The `tidyverse` package includes functions to transfer a data frame between *long* and *wide*. Wide format data tends to have different attributes or variables describing an observation placed in separate columns. Long format data tends to have different attributes encoded as levels of a single variable, followed by another column that contains the values of the observation at those different levels.

In the section on tidy data, we showed an example that used `gather` to convert data into a tidy format. The data is first in an untidy format:

```
data("VADeaths")
head(VADeaths)
```

	Rural	Male	Rural	Female	Urban	Male	Urban	Female
50-54	11.7		8.7		15.4		8.4	
55-59	18.1		11.7		24.3		13.6	
60-64	26.9		20.3		37.0		19.3	
65-69	41.0		30.9		54.6		35.1	
70-74	66.0		54.3		71.1		50.0	

After changing the age categories from row names to a variable (which can be done with the `mutate` function), the key problem with the tidyness of the data is that the variables of urban / rural and male / female are not in their own columns, but rather are embedded in the structure of the columns. To fix this, you can use the `gather` function to gather values spread across several columns into a single column, with the column names gathered into a “key” column. When gathering, exclude any columns that you don’t want “gathered” (`age` in this case) by including the column names with a the minus sign in the `gather` function. For example:

```

data("VADeaths")
library(tidyverse)

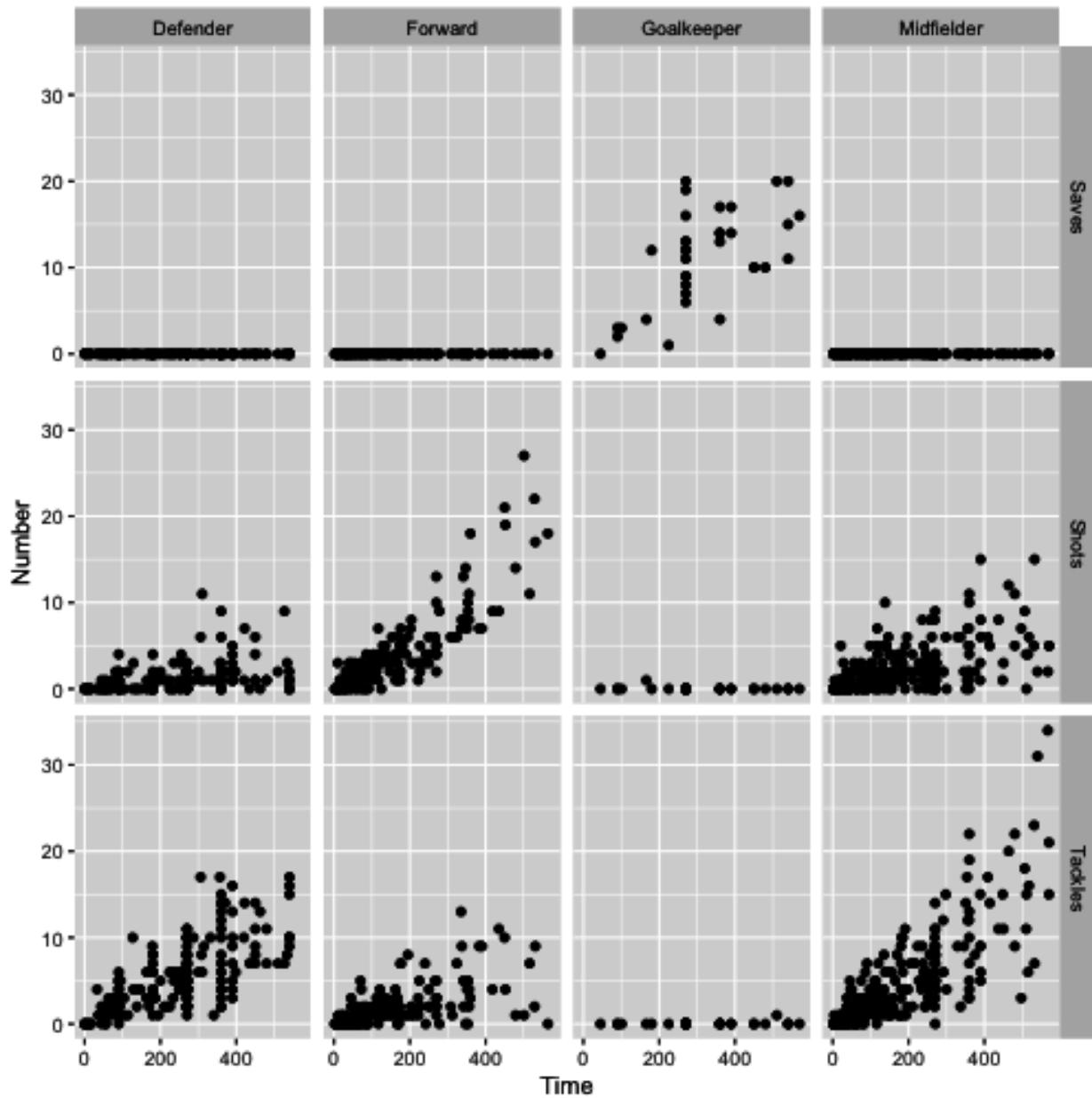
# Move age from row names into a column
VADeaths <- VADeaths %>%
 tbl_df() %>%
  mutate(age = row.names(VADeaths))
VADeaths
# A tibble: 5 × 5
`Rural Male` `Rural Female` `Urban Male` `Urban Female`   age
<dbl>        <dbl>       <dbl>        <dbl> <chr>
1    11.7        8.7        15.4        8.4 50-54
2    18.1       11.7        24.3       13.6 55-59
3    26.9       20.3        37.0       19.3 60-64
4    41.0       30.9        54.6       35.1 65-69
5    66.0       54.3        71.1       50.0 70-74

# Gather everything but age to tidy data
VADeaths %>%
  gather(key = key, value = death_rate, -age)
# A tibble: 20 × 3
  age      key death_rate
  <chr>    <chr>     <dbl>
1 50-54 Rural Male    11.7
2 55-59 Rural Male    18.1
3 60-64 Rural Male    26.9
4 65-69 Rural Male    41.0
5 70-74 Rural Male    66.0
6 50-54 Rural Female  8.7
7 55-59 Rural Female  11.7
8 60-64 Rural Female  20.3
9 65-69 Rural Female  30.9
10 70-74 Rural Female 54.3
11 50-54 Urban Male   15.4
12 55-59 Urban Male   24.3
13 60-64 Urban Male   37.0
14 65-69 Urban Male   54.6
15 70-74 Urban Male   71.1
16 50-54 Urban Female 8.4
17 55-59 Urban Female 13.6
18 60-64 Urban Female 19.3
19 65-69 Urban Female 35.1
20 70-74 Urban Female 50.0

```

Even if your data is in a tidy format, `gather` is occasionally useful for pulling data together to take advantage of faceting, or plotting separate plots based on a grouping variable. For example, if you'd like to plot the relationship between the time a player played in the World Cup and his number of saves, tackles, and shots, with a separate graph for each position (Figure @ref(fig:facetworldcup)), you can use `gather` to pull all the numbers of saves, tackles, and shots into a single column (`Number`) and then use faceting to plot them as separate graphs:

```
library(tidyr)
library(ggplot2)
worldcup %>%
  select(Position, Time, Shots, Tackles, Saves) %>%
  gather(Type, Number, -Position, -Time) %>%
  ggplot(aes(x = Time, y = Number)) +
  geom_point() +
  facet_grid(Type ~ Position)
```



Example of a faceted plot created by taking advantage of the `gather` function to pull together data.

The `spread` function is less commonly needed to tidy data. It can, however, be useful for creating summary tables. For example, if you wanted to print a table of the average number and range of passes by position for the top four teams in this World Cup (Spain, Netherlands, Uruguay, and Germany), you could run:

```
library(knitr)

# Summarize the data to create the summary statistics you want
wc_table <- worldcup %>%
  filter(Team %in% c("Spain", "Netherlands", "Uruguay", "Germany")) %>%
  select(Team, Position, Passes) %>%
  group_by(Team, Position) %>%
  summarize(ave_passes = mean(Passes),
            min_passes = min(Passes),
            max_passes = max(Passes),
            pass_summary = paste0(round(ave_passes), " (",
                                  min_passes, ", ",
                                  max_passes, ")"))
  select(Team, Position, pass_summary)
# What the data looks like before using `spread`
wc_table
Source: local data frame [16 x 3]
Groups: Team [4]

      Team    Position   pass_summary
      <fctr>    <fctr>     <chr>
1  Germany   Defender  190 (44, 360)
2  Germany     Forward   90 (5, 217)
3  Germany Goalkeeper   99 (99, 99)
4  Germany Midfielder  177 (6, 423)
5 Netherlands   Defender  182 (30, 271)
6 Netherlands     Forward   97 (12, 248)
7 Netherlands Goalkeeper 149 (149, 149)
8 Netherlands Midfielder 170 (22, 307)
9    Spain   Defender  213 (1, 402)
10   Spain     Forward   77 (12, 169)
11   Spain Goalkeeper   67 (67, 67)
12   Spain Midfielder 212 (16, 563)
13 Uruguay   Defender   83 (22, 141)
14 Uruguay     Forward  100 (5, 202)
15 Uruguay Goalkeeper  75 (75, 75)
16 Uruguay Midfielder 100 (1, 252)

# Use spread to create a prettier format for a table
wc_table %>%
  spread(Position, pass_summary) %>%
  kable()
```

Team	Defender	Forward	Goalkeeper	Midfielder
Germany	190 (44, 360)	90 (5, 217)	99 (99, 99)	177 (6, 423)
Netherlands	182 (30, 271)	97 (12, 248)	149 (149, 149)	170 (22, 307)
Spain	213 (1, 402)	77 (12, 169)	67 (67, 67)	212 (16, 563)
Uruguay	83 (22, 141)	100 (5, 202)	75 (75, 75)	100 (1, 252)

Notice in this example how `spread` has been used at the very end of the code sequence to convert the summarized data into a shape that offers a better tabular presentation for a report. In the `spread` call, you first specify the name of the column to use for the new column names (`Position` in this example) and then specify the column to use for the cell values (`pass_summary` here).

In this code, I've used the `kable` function from the `knitr` package to create the summary table in a table format, rather than as basic R output. This function is very useful for formatting basic tables in R markdown documents. For more complex tables, check out the `pander` and `xtable` packages.

Merging datasets

Often, you will have data in two separate datasets that you'd like to combine based on a common variable or variables. For example, for the World Cup example data we've been using, it would be interesting to add in a column with the final standing of each player's team. We've included data with that information in a file called "team_standings.csv", which can be read into the R object `team_standings` with the call:

```
team_standings <- read_csv("data/team_standings.csv")
team_standings %>% slice(1:3)
# A tibble: 3 × 2
  Standing     Team
    <int>   <chr>
1       1     Spain
2       2 Netherlands
3       3 Germany
```

This data frame has one observation per team, and the team names are consistent with the team names in the `worlcup` data frame.

You can use the different functions from the `*_join` family to merge this team standing data with the player statistics in the `worlcup` data frame. Once you've done that, you can use other data cleaning tools from `dplyr` to quickly pull and explore interesting parts of the dataset. The main arguments for the `*_join` functions are the object names of the two data frames to join and `by`, which specifies which variables to use to match up observations from the two dataframes.

There are several functions in the `*_join` family. These functions all merge together two data frames; they differ in how they handle observations that exist in one but not both data frames. Here are the four functions from this family that you will likely use the most often:

Function	What it includes in merged data frame
<code>left_join</code>	Includes all observations in the left data frame, whether or not there is a match in the right data frame
<code>right_join</code>	Includes all observations in the right data frame, whether or not there is a match in the left data frame
<code>inner_join</code>	Includes only observations that are in both data frames
<code>full_join</code>	Includes all observations from both data frames

In this table, the “left” data frame refers to the first data frame input in the `*_join` call, while the “right” data frame refers to the second data frame input into the function. For example, in the call

```
left_join(world_cup, team_standings, by = "Team")
```

the `world_cup` data frame is the “left” data frame and the `team_standings` data frame is the “right” data frame. Therefore, using `left_join` would include all rows from `world_cup`, whether or not the player had a team listed in `team_standings`, while `right_join` would include all the rows from `team_standings`, whether or not there were any players from that team in `world_cup`.



Remember that if you are using piping, the first data frame (“left” for these functions) is by default the dataframe created by the code right before the pipe. When you merge data frames as a step in piped code, therefore, the “left” data frame is the one piped into the function while the “right” data frame is the one stated in the `*_join` function call.

As an example of merging, say you want to create a table of the top 5 players by shots on goal, as well as the final standing for each of these player’s teams, using the `worldcup` and `team_standings` data. You can do this by running:

```
data(worldcup)
worldcup %>%
  mutate(Name = rownames(worldcup),
        Team = as.character(Team)) %>%
  select(Name, Position, Shots, Team) %>%
  arrange(desc(Shots)) %>%
  slice(1:5) %>%
  left_join(team_standings, by = "Team") %>% # Merge in team standings
  rename("Team Standing" = Standing) %>%
  kable()
```

Name	Position	Shots	Team	Team Standing
212	Forward	27	Ghana	7
560	Forward	22	Spain	1
370	Forward	21	Argentina	5
514	Forward	19	Uruguay	4
174	Forward	18	Uruguay	4

In addition to the merging in this code, there are a few other interesting things to point out:

- The code uses the `as.character` function within a `mutate` call to change the team name from a factor to a character in the `worldcup` data frame. When merging two data frames, it's safest if the column you're using to merge has the same class in each data frame. The "Team" column is a character class in the `team_standings` data frame but a factor class in the `worldcup` data frame, so this call converts that column to a character class in `worldcup`. The `left_join` function will still perform a merge if you don't include this call, but it will throw a warning that it is coercing the column in `worldcup` to a character vector. It's generally safer to do this yourself explicitly.
- It uses the `select` function both to remove columns we're not interested in and also to put the columns we want to keep in the order we'd like for the final table.
- It uses `arrange` followed by `slice` to pull out the top 5 players and order them by number of shots.
- For one of the column names, we want to use "Team Standing" rather than the current column name of "Standing". This code uses `rename` at the very end to make this change right before creating the table. You can also use the `col.names` argument in the `kable` function to customize all the column names in the final table, but this `rename` call is a quick fix since we just want to change one column name.

1.6 Working with Dates, Times, Time Zones

The learning objectives for this section are to:

- Transform non-tidy data into tidy data
- Manipulate and transform a variety of data types, including dates, times, and text data

R has special object classes for dates and date-times. It is often worthwhile to convert a column in a data frame to one of these special object types, because you can do some very useful things with date or date-time objects, including pull out the month or day of the week from the observations in the object, or determine the time difference between two values.

Many of the examples in this section use the `ext_tracks` object loaded earlier in the book. If you need to reload that, you can use the following code to do so:

```

ext_tracks_file <- "data/ebtrk_atlc_1988_2015.txt"
ext_tracks_widths <- c(7, 10, 2, 2, 3, 5, 5, 6, 4, 5, 4, 4, 5, 3, 4, 3, 3, 3,
                      4, 3, 3, 3, 4, 3, 3, 3, 2, 6, 1)
ext_tracks_colnames <- c("storm_id", "storm_name", "month", "day",
                        "hour", "year", "latitude", "longitude",
                        "max_wind", "min_pressure", "rad_max_wind",
                        "eye_diameter", "pressure_1", "pressure_2",
                        paste("radius_34", c("ne", "se", "sw", "nw"), sep = "_"),
                        paste("radius_50", c("ne", "se", "sw", "nw"), sep = "_"),
                        paste("radius_64", c("ne", "se", "sw", "nw"), sep = "_"),
                        "storm_type", "distance_to_land", "final")
ext_tracks <- read_fwf(ext_tracks_file,
                        fwf_widths(ext_tracks_widths, ext_tracks_colnames),
                        na = ".99")

```

Converting to a date or date-time class

The `lubridate` package (another package from the “tidyverse”) has some excellent functions for working with dates in R. First, this package includes functions to transform objects into date or date-time classes. For example, the `ymd_hm` function (along with other functions in the same family: `ymd`, `ymd_h`, and `ymd_hms`) can be used to convert a vector from character class to R’s data and datetime classes, `POSIXlt` and `POSIXct`, respectively.

Functions in this family can be used to parse character strings into dates, regardless of how the date is formatted, as long as the date is in the order: year, month, day (and, for time values, hour, minute). For example:

```

library(lubridate)

ymd("2006-03-12")
[1] "2006-03-12"
ymd("'06 March 12")
[1] "2006-03-12"
ymd_hm("06/3/12 6:30 pm")
[1] "2006-03-12 18:30:00 UTC"

```

The following code shows how to use the `ymd_h` function to transform the date and time information in a subset of the hurricane example data called `andrew_tracks` (the storm tracks for Hurricane Andrew) to a date-time class (`POSIXct`). This code also uses the `unite` function from the `tidyverse` package to join together date components that were originally in separate columns before applying `ymd_h`.

```

library(dplyr)
library(tidyr)

andrew_tracks <- ext_tracks %>%
  filter(storm_name == "ANDREW" & year == "1992") %>%
  select(year, month, day, hour, max_wind, min_pressure) %>%
  unite(datetime, year, month, day, hour) %>%
  mutate(datetime = ymd_h(datetime))

head(andrew_tracks, 3)
# A tibble: 3 × 3
  datetime max_wind min_pressure
  <dttm>     <int>      <int>
1 1992-08-16 18:00:00      25        1010
2 1992-08-17 00:00:00      30        1009
3 1992-08-17 06:00:00      30        1008
class(andrew_tracks$datetime)
[1] "POSIXct" "POSIXt"

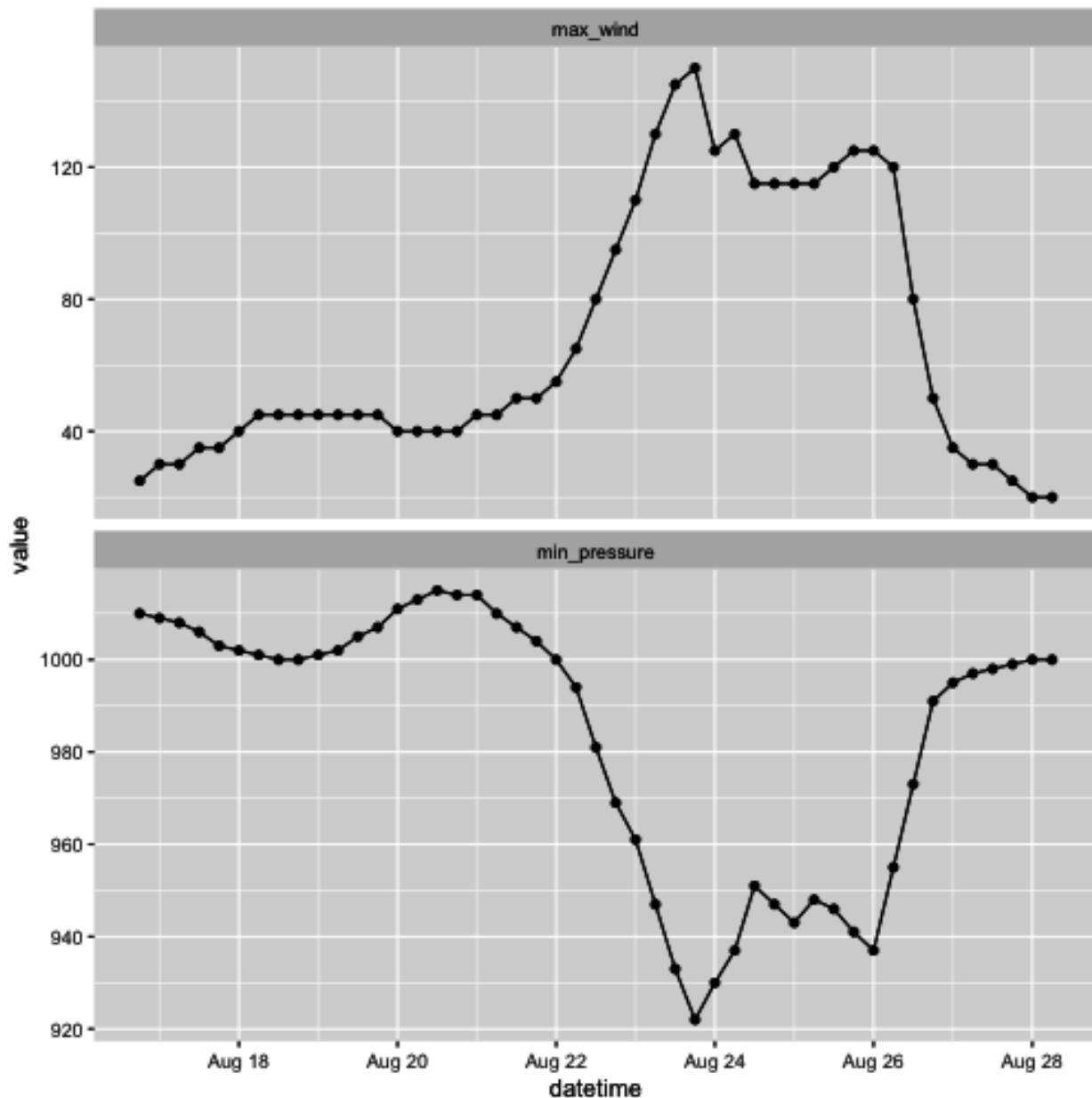
```

Now that the `datetime` variable in this dataset has been converted to a date-time class, the variable becomes much more useful. For example, if you plot a time series using `datetime`, `ggplot2` can recognize that this object is a date-time and will make sensible axis labels. The following code plots maximum wind speed and minimum air pressure at different observation times for Hurricane Andrew (Figure @ref(fig:andrewwind))—check the axis labels to see how they've been formatted. Note that this code uses `gather` from the `tidyr` package to enable easy facetting, to create separate plots for wind speed and air pressure.

```

andrew_tracks %>%
  gather(measure, value, -datetime) %>%
  ggplot(aes(x = datetime, y = value)) +
  geom_point() + geom_line() +
  facet_wrap(~ measure, ncol = 1, scales = "free_y")

```



Example of how variables in a date-time class can be parsed for sensible axis labels.

Pulling out date and time elements

Once an object is in a date or date-time class (`POSIXlt` or `POSIXct`, respectively), there are other functions in the `lubridate` package you can use to pull certain elements out of it. For example, you can use the functions `year`, `months`, `mday`, `wday`, `yday`, `weekdays`, `hour`, `minute`, and `second` to pull the year, month, month day, etc., of the date. The following code uses the `datetime` variable in the Hurricane Andrew track data to add new columns for the year, month, weekday, year day, and hour of each observation:

```
andrew_tracks %>%
  select(datetime) %>%
  mutate(year = year(datetime),
         month = months(datetime),
         weekday = weekdays(datetime),
         yday = yday(datetime),
         hour = hour(datetime)) %>%
  slice(1:3)
# A tibble: 3 × 6
  datetime   year month weekday  yday   hour
  <dttm>   <dbl> <chr>   <chr> <dbl> <int>
1 1992-08-16 18:00:00 1992 August Sunday    229     18
2 1992-08-17 00:00:00 1992 August Monday     230      0
3 1992-08-17 06:00:00 1992 August Monday     230      6
```

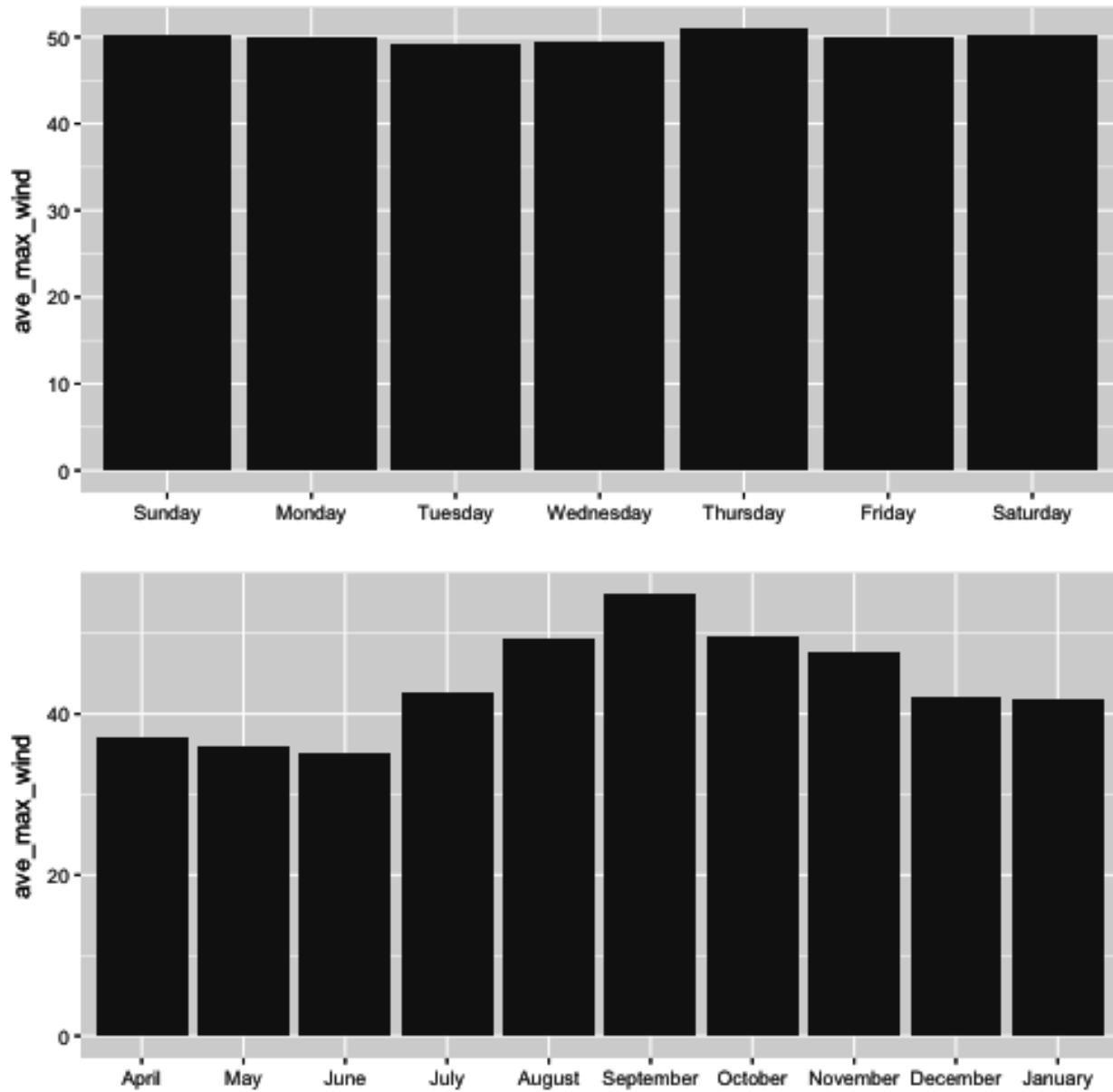
This functionality makes it easy to look at patterns in the `max_wind` value by different time groupings, like weekday and month. For example, the following code puts together some of the `dplyr` and `tidyverse` data cleaning tools and `ggplot2` plotting functions with these `lubridate` functions to look at the average value of `max_wind` storm observations by day of the week and by month (Figure @ref(fig:stormbytimegroups)).

```
check_tracks <- ext_tracks %>%
  select(month, day, hour, year, max_wind) %>%
  unite(datetime, year, month, day, hour) %>%
  mutate(datetime = ymd_h(datetime),
         weekday = weekdays(datetime),
         weekday = factor(weekday, levels = c("Sunday", "Monday",
                                              "Tuesday", "Wednesday",
                                              "Thursday", "Friday",
                                              "Saturday")),
         month = months(datetime),
         month = factor(month, levels = c("April", "May", "June",
                                           "July", "August", "September",
                                           "October", "November",
                                           "December", "January")))

check_weekdays <- check_tracks %>%
  group_by(weekday) %>%
  summarize(ave_max_wind = mean(max_wind)) %>%
  rename(grouping = weekday)
check_months <- check_tracks %>%
  group_by(month) %>%
  summarize(ave_max_wind = mean(max_wind)) %>%
  rename(grouping = month)

a <- ggplot(check_weekdays, aes(x = grouping, y = ave_max_wind)) +
  geom_bar(stat = "identity") + xlab("")
b <- a %+% check_months

library(gridExtra)
grid.arrange(a, b, ncol = 1)
```



Example of using `lubridate` functions to explore data with a date variable by different time groupings

Based on Figure @ref(fig:stormbytimegroups), there's little pattern in storm intensity by day of the week, but there is a pattern by month, with the highest average wind speed measurements in observations in September and neighboring months (and no storm observations in February or March).

There are a few other interesting things to note about this code:

- To get the weekday and month values in the right order, the code uses the `factor` function in conjunction with the `levels` option, to control the order in which R sets

the factor levels. By specifying the order we want to use with `levels`, the plot prints out using this order, rather than alphabetical order (try the code without the `factor` calls for month and weekday and compare the resulting graphs to the ones shown here).

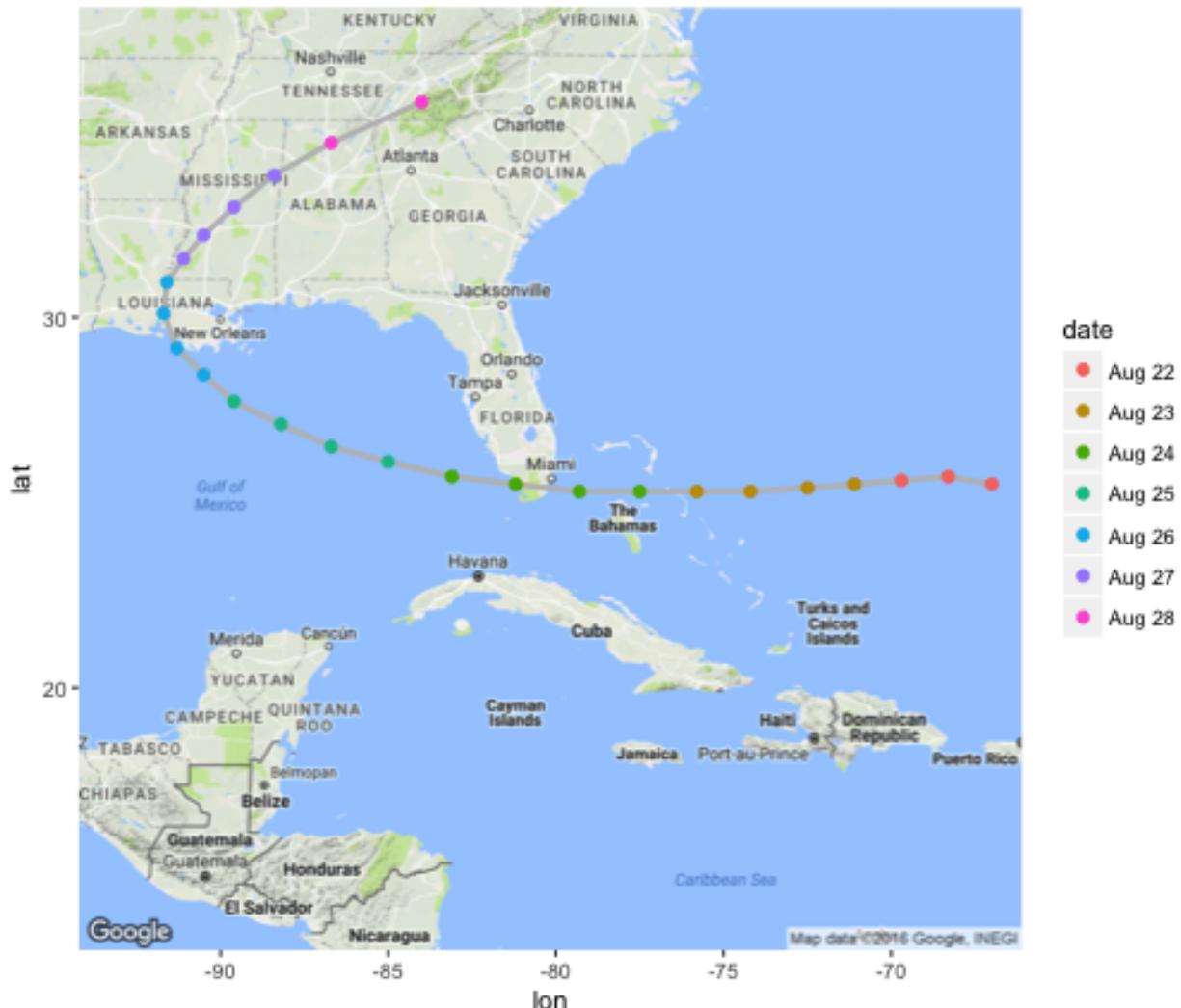
- The `grid.arrange` function, from the `gridExtra` package, allows you to arrange different `ggplot` objects in the same plot area. Here, I've used it to put the bar charts for weekday (a) and for month (b) together in one column (`ncol = 1`).
- If you ever have `ggplot` code that you would like to re-use for a new plot with a different data frame, you can save a lot of copying and pasting by using the `%+%` function. This function takes a `ggplot` object (`a` in this case, which is the bar chart by weekday) and substitutes a different data frame (`check_months`) for the original one (`check_weekdays`), but otherwise maintains all code. Note that we used `rename` to give the x-variable the same name in both datasets so we could take advantage of the `%+%` function.

Working with time zones

The `lubridate` package also has functions for handling time zones. The hurricane tracks date-times are, as is true for a lot of weather data, in Coordinated Universal Time (UTC). This means that you can plot the storm track by date, but the dates will be based on UTC rather than local time near where the storm hit. Figure @ref(fig:andrewutc) shows the location of Hurricane Andrew by date as it neared and crossed the United States, based on date-time observations in UTC.

```
andrew_tracks <- ext_tracks %>%
  filter(storm_name == "ANDREW") %>%
  slice(23:47) %>%
  select(year, month, day, hour, latitude, longitude) %>%
  unite(datetime, year, month, day, hour) %>%
  mutate(datetime = ymd_h(datetime),
         date = format(datetime, "%b %d"))

library(ggmap)
miami <- get_map("miami", zoom = 5)
ggmap(miami) +
  geom_path(data = andrew_tracks, aes(x = -longitude, y = latitude),
            color = "gray", size = 1.1) +
  geom_point(data = andrew_tracks,
             aes(x = -longitude, y = latitude, color = date),
             size = 2)
```

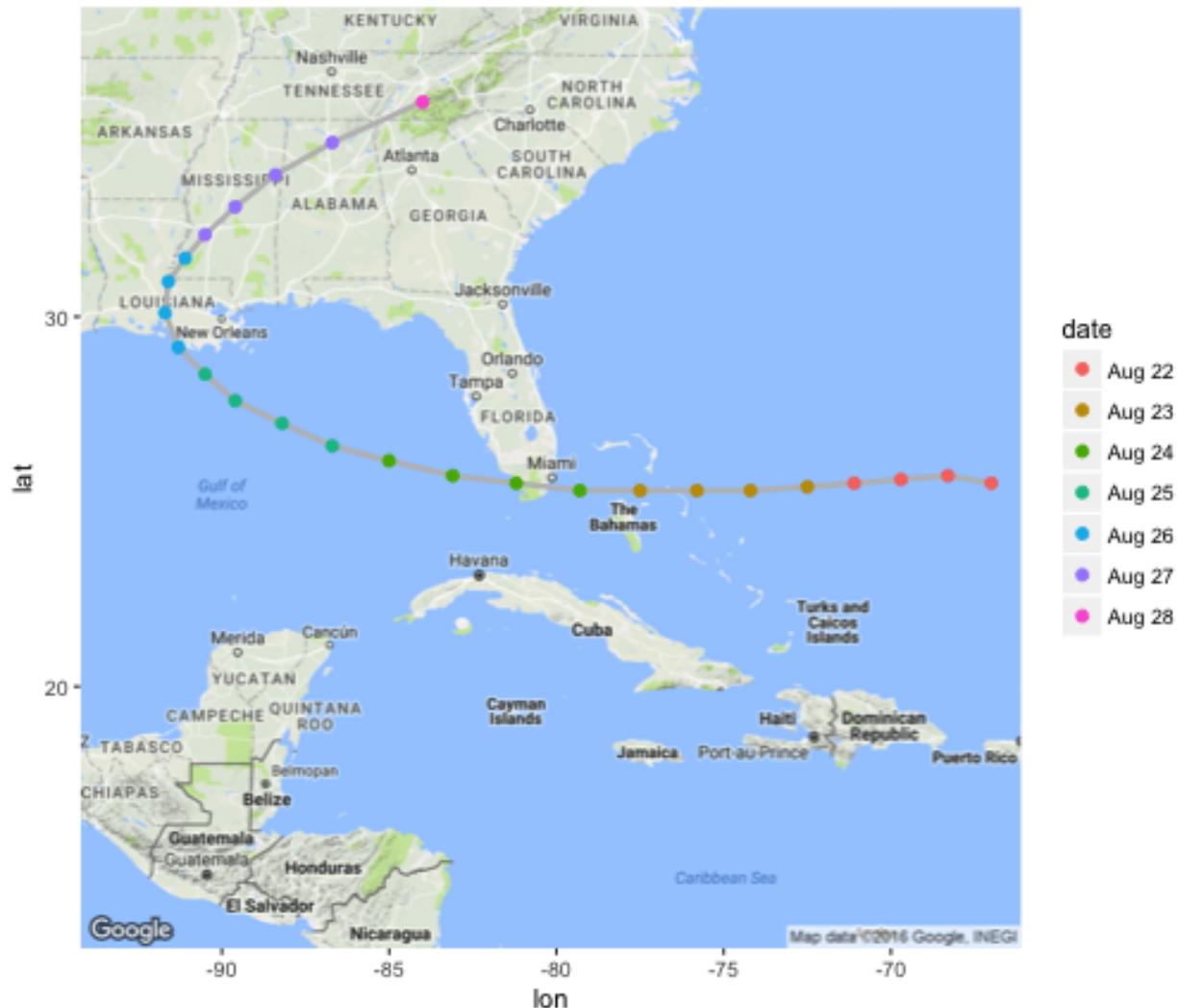


Hurricane Andrew tracks by date, based on UTC date times.

To create this plot using local time for Miami, FL, rather than UTC (Figure @ref(fig:andrewlocal)), you can use the `with_tz` function from `lubridate` to convert the `datetime` variable in the track data from UTC to local time. This function inputs a date-time object in the `POSIXct` class, as well as a character string with the time zone of the location for which you'd like to get local time, and returns the corresponding local time for that location.

```
andrew_tracks <- andrew_tracks %>%
  mutate(datetime = with_tz(datetime, tzzone = "America/New_York"),
         date = format(datetime, "%b %d"))

ggmap(miami) +
  geom_path(data = andrew_tracks, aes(x = -longitude, y = latitude),
            color = "gray", size = 1.1) +
  geom_point(data = andrew_tracks,
             aes(x = -longitude, y = latitude, color = date),
             size = 2)
```



Hurricane Andrew tracks by date, based on Miami, FL, local time.

With Figure @ref(fig:andrewlocal), it is clearer that Andrew made landfall in Florida on the morning of August 24 local time.



This section has only skimmed the surface of the date-time manipulations you can do with the `lubridate` package. For more on what this package can do, check out Garrett Grolemund and Hadley Wickham's article in the *Journal of Statistical Software* on the package—“[Dates and Times Made Easy with lubridate](#)”—or the current package vignette.

1.7 Text Processing and Regular Expressions

The learning objectives for this section are to:

- Transform non-tidy data into tidy data
- Manipulate and transform a variety of data types, including dates, times, and text data

Most common types of data are encoded in text, even if that text is representing numerical values, so being able to manipulate text as a software developer is essential. R provides several built-in tools for manipulating text, and there is a rich ecosystem of packages for R for text based analysis. First let's concentrate on some basic text manipulation functions.

Text Manipulation Functions in R

Text in R is represented as a string object, which looks like a phrase surrounded by quotation marks in the R console. For example "Hello!" and 'Strings are fun!' are both strings. You can tell whether an object is a string using the `is.character()` function. Strings are also known as characters in R.

You can combine several strings using the `paste()` function:

```
paste("Square", "Circle", "Triangle")
[1] "Square Circle Triangle"
```

By default the `paste()` function inserts a space between each word. You can insert a different string between each word by specifying the `sep` argument:

```
paste("Square", "Circle", "Triangle", sep = "+")
[1] "Square+Circle+Triangle"
```

A shortcut for combining all of the string arguments without any characters in between each of them is to use the `paste0()` function:

```
paste0("Square", "Circle", "Triangle")
[1] "SquareCircleTriangle"
```

You can also provide a vector of strings as an argument to `paste()`. For example:

```

shapes <- c("Square", "Circle", "Triangle")
paste("My favorite shape is a", shapes)
[1] "My favorite shape is a Square"    "My favorite shape is a Circle"
[3] "My favorite shape is a Triangle"

two_cities <- c("best", "worst")
paste("It was the", two_cities, "of times.")
[1] "It was the best of times."    "It was the worst of times."

```

As you can see, all of the possible string combinations are produced when you provide a vector of strings as an argument to `paste()`. You can also collapse all of the elements of a vector of strings into a single string by specifying the `collapse` argument:

```

paste(shapes, collapse = " ")
[1] "Square Circle Triangle"

```

Besides pasting strings together, there are a few other basic string manipulation functions you should be made aware of. The `nchar()` function counts the number of characters in a string:

```

nchar("Supercalifragilisticexpialidocious")
[1] 34

```

The `toupper()` and `tolower()` functions make strings all uppercase or lowercase respectively:

```

cases <- c("CAPS", "low", "Title")
tolower(cases)
[1] "caps"   "low"    "title"
toupper(cases)
[1] "CAPS"   "LOW"    "TITLE"

```

Regular Expressions

Now that we've covered the basics of string manipulation in R, let's discuss the more advanced topic of regular expressions. A regular expression is a string that defines a pattern that could be contained within another string. A regular expression can be used for searching for a string, searching within a string, or replacing one part of a string with another string. In this section I might refer to a regular expression as a regex, just know that they're the same thing.

Regular expressions use characters to define patterns of other characters. Although that approach may seem [problematic](#) at first, we'll discuss meta-characters (characters that describe other characters) and how you can use them to create powerful regular expressions.

One of the most basic functions in R that uses regular expressions is the `grep1()` function, which takes two arguments: a regular expression and a string to be searched. If the string contains the specified regular expression then `grep1()` will return `TRUE`, otherwise it will return `FALSE`. Let's take a look at one example:

```
regular_expression <- "a"
string_to_search <- "Maryland"

grep1(regular_expression, string_to_search)
[1] TRUE
```

In the example above we specify the regular expression "a" and store it in a variable called `regular_expression`. Remember that regular expressions are just strings! We also store the string "Maryland" in a variable called `string_to_search`. The regular expression "a" represents a single occurrence of the character "a". Since "a" is contained within "Maryland", `grep1()` returns the value `TRUE`. Let's try another simple example:

```
regular_expression <- "u"
string_to_search <- "Maryland"

grep1(regular_expression, string_to_search)
[1] FALSE
```

The regular expression "u" represents a single occurrence of the character "u", which is not a sub-string of "Maryland", therefore `grep1()` returns the value `FALSE`. Regular expressions can be much longer than single characters. You could for example search for smaller strings inside of a larger string:

```
grep1("land", "Maryland")
[1] TRUE
grep1("ryla", "Maryland")
[1] TRUE
grep1("Marly", "Maryland")
[1] FALSE
grep1("dany", "Maryland")
[1] FALSE
```

Since "land" and "ryla" are sub-strings of "Maryland", `grep1()` returns `TRUE`, however when a regular expression like "Marly" or "dany" is searched `grep1()` returns `FALSE` because neither are sub-strings of "Maryland".

There's a dataset that comes with R called `state.name` which is a vector of Strings, one for each state in the United States of America. We're going to use this vector in several of the following examples.

```
head(state.name)
[1] "Alabama"     "Alaska"       "Arizona"      "Arkansas"    "California"
[6] "Colorado"
```

Let's build a regular expression for identifying several strings in this vector, specifically a regular expression that will match names of states that both start and end with a vowel.

The state name could start and end with any vowel, so we won't be able to match exact sub-strings like in the previous examples. Thankfully we can use metacharacters to look for vowels and other parts of strings. The first metacharacter that we'll discuss is ". ". The metacharacter that only consists of a period represents any character other than a new line (we'll discuss new lines soon). Let's take a look at some examples using the peroid regex:

```
grep1(".", "Maryland")
[1] TRUE
grep1(".", "*&2[0+,%<@#~| ]")
[1] TRUE
grep1(".", "")
[1] FALSE
```

As you can see the period metacharacter is very liberal. This metacharacter is most useful when you don't care about a set of characters in a regular expression. For example:

```
grep1("a.b", c("aaa", "aab", "abb", "acadb"))
[1] FALSE TRUE TRUE TRUE
```

In the case above `grep1()` returns `TRUE` for all strings that contain an `a` followed by any other character followed by a `b`.

You can specify a regular expression that contains a certain number of characters or metacharacters using the enumeration metacharacters. The `+` metacharacter indicates that one or more of the preceding expression should be present and `*` indicates that zero or more of the preceding expression is present. Let's take a look at some examples using these metacharacters:

```
# Does "Maryland" contain one or more of "a" ?
grep1("a+", "Maryland")
[1] TRUE

# Does "Maryland" contain one or more of "x" ?
grep1("x+", "Maryland")
[1] FALSE

# Does "Maryland" contain zero or more of "x" ?
grep1("x*", "Maryland")
[1] TRUE
```

You can also specify exact numbers of expressions using curly brackets `{}`. For example "`a{5}`" specifies "a exactly five times," "`a{2,5}`" specifies "a between 2 and 5 times," and "`a{2,}`" specifies "a at least 2 times." Let's take a look at some examples:

```
# Does "Mississippi" contain exactly 2 adjacent "s" ?
grep1("s{2}", "Mississippi")
[1] TRUE

# This is equivalent to the expression above:
grep1("ss", "Mississippi")
[1] TRUE

# Does "Mississippi" contain between 1 and 3 adjacent "s" ?
grep1("s{2,3}", "Mississippi")
[1] TRUE

# Does "Mississippi" contain between 2 and 3 adjacent "i" ?
grep1("i{2,3}", "Mississippi")
[1] FALSE

# Does "Mississippi" contain between 2 adjacent "iss" ?
grep1("(iss){2}", "Mississippi")
[1] TRUE

# Does "Mississippi" contain between 2 adjacent "ss" ?
grep1("(ss){2}", "Mississippi")
[1] FALSE

# Does "Mississippi" contain the pattern of an "i" followed by
# 2 of any character, with that pattern repeated three times adjacently?
grep1("(i.{2}){3}", "Mississippi")
[1] TRUE
```

In the last three examples I used parentheses () to create a capturing group. A capturing group allows you to use quantifiers on other regular expressions. In the last example I first created the regex "i.{2}" which matches i followed by any two characters ("iss" or "ipp"). I then used a capture group to wrap that regex, and to specify exactly three adjacent occurrences of that regex.

You can specify sets of characters with regular expressions, some of which come built in, but you can build your own character sets too. First we'll discuss the built in character sets: words ("\w"), digits ("\d"), and whitespace characters ("\s"). Words specify any letter, digit, or a underscore, digits specify the digits 0 through 9, and whitespace specifies line breaks, tabs, or spaces. Each of these character sets have their own compliments: not words ("\W"), not digits ("\D"), and not whitespace characters ("\S"). Each specifies all of the characters not included in their corresponding character sets. Let's take a look at a few examples:

```
grep1("\\w", "abcdefghijklmnopqrstuvwxyz0123456789")
[1] TRUE

grep1("\\d", "0123456789")
[1] TRUE

# "\n" this regex for a new line and "\t" is the regex for a tab
grep1("\\s", "\n\t ")
[1] TRUE

grep1("\\d", "abcdefghijklmnopqrstuvwxyz")
[1] FALSE

grep1("\\D", "abcdefghijklmnopqrstuvwxyz")
[1] TRUE

grep1("\\w", "\n\t ")
[1] FALSE
```

You can also specify specific character sets using straight brackets []. For example a character set of just the vowels would look like: "[aeiou]". You can find the complement to a specific character by putting a carrot ^ after the first bracket. For example "[^aeiou]" matches all characters except the lowercase vowels. You can also specify ranges of characters using a hyphen - inside of the brackets. For example "[a-m]" matches all of the lowercase characters between a and m, while "[5-8]" matches any digit between 5 and 8 inclusive. Let's take a look at some examples using custom character sets:

```
grep1("[aeiou]", "rhythms")
[1] FALSE

grep1("[^aeiou]", "rhythms")
[1] TRUE

grep1("[a-m]", "xyz")
[1] FALSE

grep1("[a-m]", "ABC")
[1] FALSE

grep1("[a-zA-M]", "ABC")
[1] TRUE
```

You might be wondering how you can use regular expressions to match a particular punctuation mark since many punctuation marks are used as metacharacters! Putting two backslashes before a punctuation mark that is also a metacharacter indicates that you are looking for the symbol and not the metacharacter meaning. For example "\\." indicates you are trying to match a period in a string. Let's take a look at a few examples:

```
grep1("\\+", "tragedy + time = humor")
[1] TRUE
```

```
grep1("\\.", "http://www.jhsph.edu/")
[1] TRUE
```

There are also metacharacters for matching the beginning and the end of a string which are "^" and "\$" respectively. Let's take a look at a few examples:

```
grep1("^a", c("bab", "aab"))
[1] FALSE TRUE
```

```
grep1("b$", c("bab", "aab"))
[1] TRUE TRUE
```

```
grep1("^ab]+$", c("bab", "aab", "abc"))
[1] TRUE TRUE FALSE
```

The last metacharacter we'll discuss is the OR metacharacter ("|"). The OR metacharacter matches either the regex on the left or the regex on the right side of this character. A few examples:

```
grep1("a|b", c("abc", "bcd", "cde"))
[1] TRUE TRUE FALSE
```

```
grep1("North|South", c("South Dakota", "North Carolina", "West Virginia"))
[1] TRUE TRUE FALSE
```

Finally we've learned enough to create a regular expression that matches all state names that both begin and end with a vowel:

1. We match the beginning of a string.
2. We create a character set of just capitalized vowels.
3. We specify one instance of that set.
4. Then any number of characters until:
5. A character set of just lowercase vowels.
6. We specify one instance of that set.
7. We match the end of a string.

```

start_end_vowel <- "^[AEIOU]{1}.[aeiou]{1}$"
vowel_state_lgl <- grep1(start_end_vowel, state.name)
head(vowel_state_lgl)
[1] TRUE TRUE TRUE FALSE FALSE

state.name[vowel_state_lgl]
[1] "Alabama"   "Alaska"    "Arizona"   "Idaho"     "Indiana"   "Iowa"
[7] "Ohio"       "Oklahoma"

```

Below is a table of several important metacharacters:

Metacharacter	Meaning
.	Any Character
\w	A Word
\W	Not a Word
\d	A Digit
\D	Not a Digit
\s	Whitespace
\S	Not Whitespace
[xyz]	A Set of Characters
[^xyz]	Negation of Set
[a-z]	A Range of Characters
^	Beginning of String
\$	End of String
\n	Newline
+	One or More of Previous
*	Zero or More of Previous
?	Zero or One of Previous
	Either the Previous or the Following
{5}	Exactly 5 of Previous
{2, 5}	Between 2 and 5 or Previous
{2, }	More than 2 of Previous

RegEx Functions in R

So far we've been using `grep1()` to see if a regex matches a string. There are a few other built in reged functions you should be aware of. First we'll review our workhorse of this chapter, `grep1()` which stands for “grep logical.”

```

grep1("[Ii]", c("Hawaii", "Illinois", "Kentucky"))
[1] TRUE TRUE FALSE

```

Then there's old fashioned `grep()` which returns the indices of the vector that match the regex:

```
grep("[Ii]", c("Hawaii", "Illinois", "Kentucky"))
[1] 1 2
```

The `sub()` function takes as arguments a regex, a “replacement,” and a vector of strings. This function will replace the first instance of that regex found in each string.

```
sub("[Ii]", "1", c("Hawaii", "Illinois", "Kentucky"))
[1] "Hawa1i"   "illinois" "Kentucky"
```

The `gsub()` function is nearly the same as `sub()` except it will replace every instance of the regex that is matched in each string.

```
gsub("[Ii]", "1", c("Hawaii", "Illinois", "Kentucky"))
[1] "Hawa11"   "111nois" "Kentucky"
```

The `strsplit()` function will split up strings according to the provided regex. If `strsplit()` is provided with a vector of strings it will return a list of string vectors.

```
two_s <- state.name[grep("ss", state.name)]
two_s
[1] "Massachusetts" "Mississippi"    "Missouri"      "Tennessee"
strsplit(two_s, "ss")
[[1]]
[1] "Ma"          "achusetts"

[[2]]
[1] "Mi"          "i"            "ippi"

[[3]]
[1] "Mi"          "ouri"

[[4]]
[1] "Tenne"        "ee"
```

The `stringr` Package

The `stringr` package, written by Hadley Wickham, is part of the `Tidyverse` group of R packages. This package takes a “data first” approach to functions involving regex, so usually the string is the first argument and the regex is the second argument. The majority of the function names in `stringr` begin with `str_`.

The `str_extract()` function returns the sub-string of a string that matches the provided regular expression.

```
library(stringr)
state_tbl <- paste(state.name, state.area, state.abb)
head(state_tbl)
[1] "Alabama 51609 AL"      "Alaska 589757 AK"      "Arizona 113909 AZ"
[4] "Arkansas 53104 AR"     "California 158693 CA"    "Colorado 104247 CO"
str_extract(state_tbl, "[0-9]+")
[1] "51609"   "589757"  "113909"  "53104"   "158693"  "104247"  "5009"
[8] "2057"    "58560"   "58876"   "6450"    "83557"   "56400"   "36291"
[15] "56290"   "82264"   "40395"   "48523"   "33215"   "10577"   "8257"
[22] "58216"   "84068"   "47716"   "69686"   "147138"  "77227"   "110540"
[29] "9304"    "7836"    "121666"  "49576"   "52586"   "70665"   "41222"
[36] "69919"   "96981"   "45333"   "1214"    "31055"   "77047"   "42244"
[43] "267339"  "84916"   "9609"    "40815"   "68192"   "24181"   "56154"
[50] "97914"
```

The `str_order()` function returns a numeric vector that corresponds to the alphabetical order of the strings in the provided vector.

```
head(state.name)
[1] "Alabama"      "Alaska"       "Arizona"      "Arkansas"     "California"
[6] "Colorado"
str_order(state.name)
[1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
[24] 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46
[47] 47 48 49 50

head(state.abb)
[1] "AL" "AK" "AZ" "AR" "CA" "CO"
str_order(state.abb)
[1] 2 1 4 3 5 6 7 8 9 10 11 15 12 13 14 16 17 18 21 20 19 22 23
[24] 25 24 26 33 34 27 29 30 31 28 32 35 36 37 38 39 40 41 42 43 44 46 45
[47] 47 49 48 50
```

The `str_pad()` function pads strings with other characters which is often useful when the string is going to be eventually printed for a person to read.

```
str_pad("Thai", width = 8, side = "left", pad = "-")
[1] "----Thai"
str_pad("Thai", width = 8, side = "right", pad = "-")
[1] "Thai---"
str_pad("Thai", width = 8, side = "both", pad = "-")
[1] "--Thai--"
```

The `str_to_title()` function acts just like `tolower()` and `toupper()` except it puts strings into Title Case.

```
cases <- c("CAPS", "low", "Title")
str_to_title(cases)
[1] "Caps"   "Low"    "Title"
```

The `str_trim()` function deletes whitespace from both sides of a string.

```
to_trim <- c("  space", "the      ", "      final frontier  ")
str_trim(to_trim)
[1] "space"           "the"            "final frontier"
```

The `str_wrap()` function inserts newlines in strings so that when the string is printed each line's length is limited.

```
pasted_states <- paste(state.name[1:20], collapse = " ")

cat(str_wrap(pasted_states, width = 80))
Alabama Alaska Arizona Arkansas California Colorado Connecticut Delaware Florida
Georgia Hawaii Idaho Illinois Indiana Iowa Kansas Kentucky Louisiana Maine
Maryland
cat(str_wrap(pasted_states, width = 30))
Alabama Alaska Arizona
Arkansas California Colorado
Connecticut Delaware Florida
Georgia Hawaii Idaho Illinois
Indiana Iowa Kansas Kentucky
Louisiana Maine Maryland
```

The `word()` function allows you to index each word in a string as if it were a vector.

```
a_tale <- "It was the best of times it was the worst of times it was the age of wisdom it was the ag\
e of foolishness"

word(a_tale, 2)
[1] "was"

word(a_tale, end = 3)
[1] "It was the"

word(a_tale, start = 11, end = 15)
[1] "of times it was the"
```

Summary

String manipulation in R is useful for data cleaning, plus it can be fun! For prototyping your first regular expressions I highly recommend checking out <http://regexp.com/>. If you're interested in what some people call a more “humane” way of constructing regular expressions you should check out the `rex` package by [Kevin Ushey](#) and [Jim Hester](#). If you'd like to find out more about text analysis I highly recommend reading [Tidy Text Mining in R](#) by [Julia Silge](#) and [David Robinson](#).

1.8 The Role of Physical Memory

The learning objectives of this section are to:

- Describe how memory is used in R sessions to store R objects

Generally speaking, R stores and manipulates all objects in the physical memory of your computer (i.e. the RAM). Therefore, it's important to be aware of the limits of your computing environment with respect to available memory and how that may affect your ability to use R. In the event that your computer's physical memory is insufficient for some of your work, there have been some developments that allow R users to deal with objects out of physical memory and we will discuss them below.

The first thing that is worth keeping in mind as you use R is how much physical memory your computer actually has. Typically, you can figure this out by looking at your operating system's settings. For example, as of this writing, Roger has a 2015-era Macbook with 8 GB of RAM. Of course, the amount of RAM available to R will be quite a bit less than that, but it's a useful upper bound. If you plan to read into R an object that is 16 GB on this computer, you're going to have ask Roger for a new computer.

The `pryr` package provides a number of useful functions for interrogating the memory usage of your R session. Perhaps the most basic is the `mem_used()` function, which tells you how much memory your current R session is using.

```
library(pryr)
mem_used()
127 MB
```

The primary use of this function is to make sure your memory usage in R isn't getting too big. If the output from `mem_used()` is in the neighborhood of 75%-80% of your total physical RAM, you might need to consider a few things.

First, you might consider removing a few very large objects in your workspace. You can see the memory usage of objects in your workspace by calling the `object_size()` function.

```
ls() ## Show objects in workspace
[1] "a"                      "a_tale"                  "andrew_tracks"
[4] "b"                      "cases"                  "check_months"
[7] "check_tracks"           "check_weekdays"        "denver"
[10] "ext_tracks"             "ext_tracks_colnames"  "ext_tracks_file"
[13] "ext_tracks_widths"      "input"                  "join_funcs"
[16] "katrina"                "katrina_reduced"       "knots_to_mph"
[19] "logdates"                "logs"                   "m"
[22] "mc_tibl"                 "meso_url"               "miami"
[25] "msg"                     "old"                    "pasted_states"
[28] "readr_functions"         "regular_expression"   "shapes"
[31] "start_end_vowel"        "state_tbl"              "string_to_search"
```

```
[34] "team_standings"      "teams"                  "to_trim"
[37] "two_cities"          "two_s"                  "VADeaths"
[40] "vowel_state_lgl"    "wc_table"               "worldcup"
[43] "x"                   "y"                      "zika_brazil"
[46] "zika_file"
object_size(worldcup)
61.2 kB
```

The `object_size()` function will print the number of bytes (or kilobytes, or megabytes) that a given object is using in your R session. If you want see what the memory usage of the largest 5 objects in your workspace is, you can use the following code.

```
library(magrittr)
sapply(ls(), function(x) object.size(get(x))) %>% sort %>% tail(5)
worldcup      denver check_tracks   ext_tracks      miami
61424        222768     239848       1842472      13121608
```

Note: We have had to use the `object.size()` function here (see note below) because the current version of `object_size()` in `pryr` throws an error for certain types of objects.

Here we can see that the `miami` and `ext_tracks` objects (created in previous chapters of this book) are currently taking up the most memory in our R session. Since we no longer need those objects, we can remove them from the workspace and free up some memory.

```
mem_used()
127 MB
rm(ext_tracks, miami)
mem_used()
125 MB
```

Here you can see how much memory we save by deleting these two objects. But you may be wondering why there isn't a larger savings, given the number reported by `object_size()`. This has to do with the internal representation of the `miami` object, which is of the class `ggmap`. Occasionally, certain types of R objects can appear to take up more memory than they actually do, in which case functions like `object_size()` will get confused.

Viewing the change in memory usage by executing an R expression can actually be simplified using the `mem_change()` function. We can see what happens when we remove the next three largest objects.

```
mem_change(rm(check_tracks, denver, b))
-459 kB
```

Here the decrease is about 400 KB.

R has a built in function called `object.size()` that also calculates the size of an object, but it uses a slightly different calculation than `object_size()` in `pryr`. While the two functions will generally agree for most objects, for things like functions and formulas, which have enclosing environments attached to them, they will differ. Similarly, objects with shared elements (i.e. character vectors) may result in different computations of their size. The `compare_size()` function in `pryr` allows you to see how the two functions compare in their calculations. We will discuss these concepts more in the next chapter.

Back of the Envelope Calculations

When reading in large datasets or creating large R objects, it's often useful to do a back of the envelope calculation of how much memory the object will occupy in the R session (ideally *before* creating the object). To do this it's useful to know roughly how much memory different types of atomic data types in R use.

It's difficult to generalize how much memory is used by data types in R, but on most 64 bit systems today, integers are 32 bits (4 bytes) and double-precision floating point numbers (numerics in R) are 64 bits (8 bytes). Furthermore, character data are usually 1 byte per character. Because most data come in the form of numbers (integer or numeric) and letters, just knowing these three bits of information can be useful for doing many back of the envelope calculations.

For example, an integer vector is roughly 4 bytes times the number of elements in the vector. We can see that for a zero-length vector, that still requires some memory to represent the data structure.

```
object_size(integer(0))  
40 B
```

However, for longer vectors, the overhead stays roughly constant, and the size of the object is determined by the number of elements.

```
object_size(integer(1000)) ## 4 bytes per integer  
4.04 kB  
object_size(numeric(1000)) ## 8 bytes per numeric  
8.04 kB
```

If you are reading in tabular data of integers and floating point numbers, you can roughly estimate the memory requirements for that table by multiplying the number of rows by the memory required for each of the columns. This can be a useful exercise to do before reading in large datasets. If you accidentally read in a dataset that requires more memory than your computer has available, you may end up freezing your R session (or even your computer).

The `.Machine` object in R (found in the `base` package) can give you specific details about how your computer/operation system stores different types of data.

```
str(.Machine)
List of 18
 $ double.eps      : num 2.22e-16
 $ double.neg.eps  : num 1.11e-16
 $ double.xmin     : num 2.23e-308
 $ double.xmax     : num 1.8e+308
 $ double.base     : int 2
 $ double.digits   : int 53
 $ double.rounding : int 5
 $ double.guard    : int 0
 $ double.ulp.digits: int -52
 $ double.neg.ulp.digits: int -53
 $ double.exponent : int 11
 $ double.min.exp  : int -1022
 $ double.max.exp  : int 1024
 $ integer.max     : int 2147483647
 $ sizeof.long      : int 8
 $ sizeof.longlong  : int 8
 $ sizeof.longdouble: int 16
 $ sizeof.pointer   : int 8
```

The floating point representation of a decimal number contains a set of bits representing the *exponent* and another set of bits representing the *significand* or the *mantissa*. Here the number of bits used for the exponent is 11, from `double.exponent`, and the number of bits for the significand is 53, from the `double.digits` element. Together, each double precision floating point number requires 64 bits, or 8 bytes to store.

For integers, we can see that the maximum integer indicated by the `integer.max` is 2147483647, we can take the base 2 log of that number and see that it requires 31 bits to encode. Because we need another bit to encode the sign of the number, the total number of bits for an integer is 32, or 4 bytes.

Much of the point of this discussion of memory is to determine if your computer has sufficient memory to do the work you want to do. If you determine that the data you're working with cannot be completely stored in memory for a given R session, then you may need to resort to alternate tactics. We discuss one such alternative in the section below, "Working with large datasets".

Internal Memory Management in R

If you're familiar with other programming languages like C, you'll notice that you do not need to explicitly allocate and de-allocate memory for objects in R. This is because R has a garbage collection system that recycles unused memory and gives it back to R. This happens automatically without the need for user intervention.

Roughly, R will periodically cycle through all of the objects that have been created and see if there are still any references to the object somewhere in the session. If there are no references, the object is garbage-collected and the memory returned. Under normal usage, the garbage collection is not noticeable, but occasionally, when working with very large R

objects, you may notice a “hiccup” in your R session when R triggers a garbage collection to reclaim unused memory. There’s not really anything you can do about this except not panic when it happens.

The `gc()` function in the `base` package can be used to explicitly trigger a garbage collection in R. Calling `gc()` explicitly is never actually needed, but it does produce some output that is worth understanding.

```
gc()
  used (Mb) gc trigger (Mb) max used (Mb)
Ncells 1692783 90.5    2637877 140.9  2637877 140.9
Vcells 3766734 28.8    11515884  87.9  18887069 144.1
```

The `used` column gives you the amount of memory currently being used by R. The distinction between `Ncells` and `Vcells` is not important—the `mem_used()` function in `pryr` essentially gives you the sum of this column. The `gc trigger` column gives you the amount of memory that can be used before a garbage collection is triggered. Generally, you will see this number go up as you allocate more objects and use more memory. The `max used` column shows the maximum space used since the last call to `gc(reset = TRUE)` and is not particularly useful.

1.9 Working with Large Datasets

The learning objectives of this section are to:

- Read and manipulate large datasets

R now offers a variety of options for working with large datasets. We won’t try to cover all these options in detail here, but rather give an overview of strategies to consider if you need to work with a large dataset, as well as point you to additional resources to learn more about working with large datasets in R.

While there are a variety of definitions of how large a dataset must be to qualify as “large”, in this section we don’t formally define a limit. Instead, this section is meant to give you some strategies anytime you work with a dataset large enough that you notice it’s causing problems. For example, data large enough for R to be noticeably slow to read or manipulate the data, or large enough it’s difficult to store the data locally on your computer.

In-memory strategies

In this section, we introduce the basics of why and how to use `data.table` to work with large datasets in R. We have included a video demonstration online showing how functions from the `data.table` package can be used to load and explore a large dataset more efficiently.

The `data.table` package can help you read a large dataset into R and explore it more efficiently. The `fread` function in this package, for example, can read large flat files in much more quickly than comparable base R packages. Since all of the `data.table` functions will work with smaller datasets, as well, we'll illustrate using `data.table` with the Zika data accessed from GitHub in an earlier section of this chapter. We've saved that data locally to illustrate how to read it in and work with it using `data.table`.

First, to read this data in using `fread`, you can run:

```
library(data.table)
brazil_zika <- fread("data/COES_Microcephaly-2016-06-25.csv")
head(brazil_zika, 2)
#> #> report_date      location location_type          data_field
#> 1: 2016-06-25    Brazil-Acre           state microcephaly_confirmed
#> 2: 2016-06-25 Brazil-Alagoas           state microcephaly_confirmed
#> #> data_field_code time_period time_period_type value   unit
#> 1:          BR0002        NA             NA       2 cases
#> 2:          BR0002        NA             NA      75 cases
class(brazil_zika)
[1] "data.table" "data.frame"
```

If you are working with a very large dataset, `data.table` will provide a status bar showing your progress towards loading the code as you read it in using `fread`.

If you have a large dataset for which you only want to read in certain columns, you can save time when using `data.table` by only reading in the columns you want with the `select` argument in `fread`. This argument takes a vector of either the names or positions of the columns that you want to read in:

```
fread("data/COES_Microcephaly-2016-06-25.csv",
      select = c("location", "value", "unit")) %>%
  dplyr::slice(1:3)
#> #> location value   unit
#> 1  Brazil-Acre    2 cases
#> 2 Brazil-Alagoas  75 cases
#> 3  Brazil-Amapa     7 cases
```

Many of the `fread` arguments are counterparts to arguments in the `read.table` family of functions in base R (for example, `na.strings`, `sep`, `skip`, `colClasses`). One that is particular useful is `nrows`. If you're working with data that takes a while to read in, using `nrows = 20`

or some other small number will allow you to make sure you have set all of the arguments in `fread` appropriately for the dataset before you read in the full dataset.

If you already have a dataset loaded to your R session, you can use the `data.table` function to convert a data frame into a `data.table` object. (Note: if you use `fread`, the data is automatically read into a `data.table` object.) A `data.table` object also has the class `data.frame`; this means that you can use all of your usual methods for manipulating a data frame with a `data.table` object. However, for extra speed, use `data.table` functions to manipulate, clean, and explore the data in a `data.table` object. You can find out more about using `data.table` functions at the [data.table wiki](#).

Many of the functions in `data.table`, like many in `dplyr`, use non-standard evaluation. This means that, while they'll work fine in interactive programming, you'll need to take some extra steps when you use them to write functions for packages. We'll cover non-standard evaluation in the context of developing packages in a later section.

When you are working with datasets that are large, but can still fit in-memory, you'll want to optimize your code as much as possible. There are more details on profiling and optimizing code in a later chapter, but one strategy for speeding up R code is to write some of the code in C++ and connect it to R using the `Rcpp` package. Since C++ is a compiled rather than an interpreted language, it runs much faster than similar code written in R. If you are more comfortable coding in another compiled language (C or FORTRAN, for example), you can also use those, although the `Rcpp` package is very nicely written and well-maintained, which makes C++ an excellent first choice for creating compiled code to speed up R.

Further, a variety of R packages have been written that help you run R code in parallel, either locally or on a cluster. Parallel strategies may be work pursuing if you are working with very large datasets, and if the coding tasks can be split to run in parallel. To get more ideas and find relevant packages, visit CRAN's [High-Performance and Parallel Computing with R task view](#).

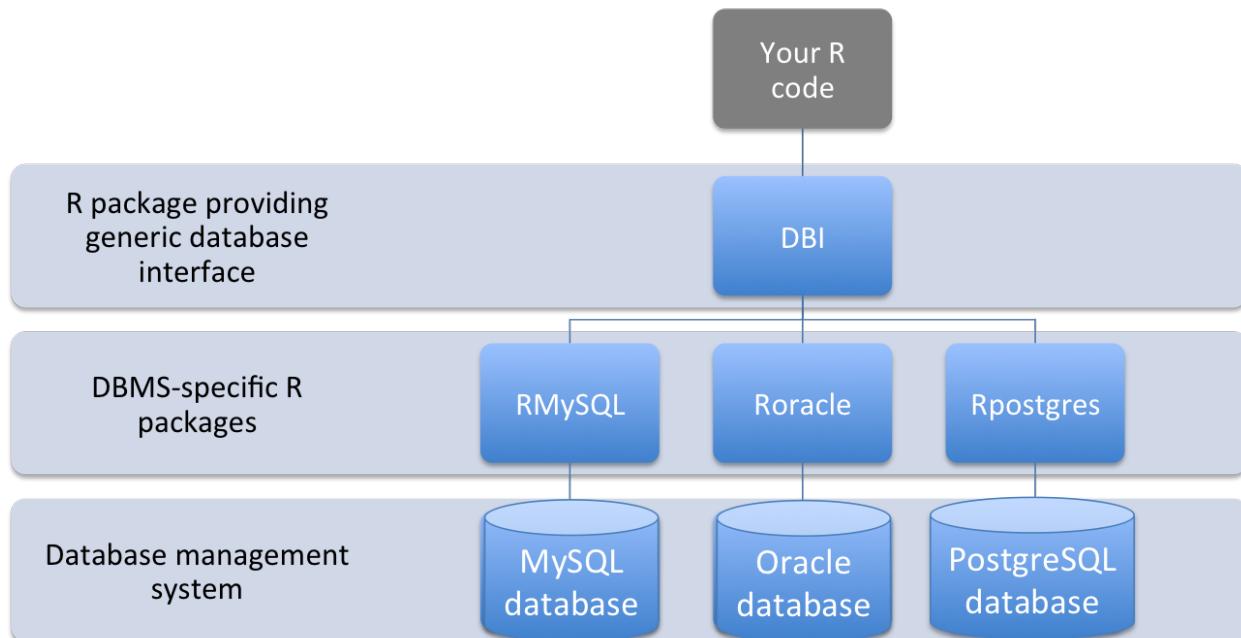
Out-of-memory strategies

If you need to work with a very large dataset, there are also some options to explore and model the dataset without ever loading it into R, while still using R commands and working from the R console or an R script. These options can make working with large datasets more efficient, because they let other software handle the heavy lifting of sifting through the data and / or avoid loading large datasets into RAM, instead using data stored on hard drive.

For example, database management systems are optimized to more efficiently store and better search through large sets of data; popular examples include Oracle, MySQL, and PostgreSQL. There are several R packages that allow you to connect your R session to a database. With these packages, you can use functions from the R console or an R script to search and subset data without loading the whole dataset into R, and so take advantage of

the improved efficiency of the database management system in handling data, as well as work with data too big to fit in memory.

The `DBI` package is particularly convenient for interfacing R code with a database management system, as it provides a top-level interface to a number of different database management systems, with system-specific code applied by a lower-level, more specific R package (Figure @ref(fig:rdbi)).



Structure of interface between code in an R script and data stored in a database management system using DBI-compliant packages

The `DBI` package therefore allows you to use the same commands for working with database-stored data in R, without worrying about details specific to the exact type of database management system you're connecting to. The following table outlines the `DBI` functions you can use to perform a variety of tasks when working with data stored in a database:

Task	DBI Function
Create a new driver object for an instance of a database	<code>dbDriver</code>
Connect to database instance	<code>dbConnect</code>
Find available tables in a connected database instance	<code>dbListTables</code>
Find available fields within a table	<code>dbListFields</code>
Query a connected database instance	<code>dbSendQuery</code>
Pull a data frame into R from a query result	<code>dbFetch</code>
Jointly query and pull data from a database instance	<code>dbGetQuery</code>
Close result set from a query	<code>dbClearResult</code>
Write a new table in a database instance	<code>dbWriteTable</code>
Remove a table from a database instance	<code>dbRemoveTable</code>
Disconnect from a database instance	<code>dbDisconnect</code>

The `DBI` package depends on lower-level R packages to translate its generic commands to work for specific database management systems. DBI-compliant R packages have not been written for every database management system, so there are some databases for which DBI commands will not work. DBI-compliant R packages that are available include:

Database Management System	R packages
Oracle	<code>ROracle</code>
MySQL	<code>RMySQL</code>
Microsoft SQL Server	<code>RSQLServer</code>
PostgreSQL	<code>RPostgres</code>
SQLite	<code>RSQLite</code>



For more on the `DBI` package, including its history, see [the package's GitHub README page](#).

The packages for working with database management systems require you to send commands to the database management system in that system’s command syntax (e.g., SQL). You can, however, do “SELECT” database queries directly using `dplyr` syntax for some database systems, rather than with SQL syntax. While this functionality is limited to “SELECT” calls, often this is all you’ll need within a data analysis script. For more details, see the [dplyr database vignette](#).

In addition to database management systems, there are other options for working with large data out-of-memory in R. For example, the `bigmemory` and associated packages can be used to access and work with large matrices stored on hard drive rather than in RAM, by storing the data in a C++ matrix structure and loading to R pointers to the data, rather than the full dataset. This family of packages includes packages that can be used to summarize and model the data (`biglm`, `bigglm`, `biganalytics`, `bigtabulate`, `bigalgebra`). One limitation is that these packages only work with matrices, not data frames; matrices require all elements share a class (e.g., all numeric).

Finally, there are some packages that allow you to write R code that uses other software to load and work with data through an R API provided by the other software. For example, the `h2o` package allows you to write R code to load and fit machine learning models in `H2O`, which is open-source software that facilitates distributed machine learning. `H2O` includes functions to fit and evaluate numerous machine learning models, including ensemble models, which would take quite a while to fit within R with a large training dataset. Since processing is done using compiled code, models can be fit on large datasets more quickly. However, while the `h2o` package allows you to use R-like code from within an R console to explore and model your data, it is not actually running R, but instead is using the R code, through the R API, to run Java-encoded functions. As a result, you only have access to a small subset of R’s total functionality, since you can only run the R-like functions written into `H2O`’s own software.

1.10 Diagnosing Problems

The learning objectives of this section are to:

- Describe how to diagnose programming problems and to look up answers from the web or forums

Inevitably, no matter what your level of expertise, you will get to a point in your R programming where you're stuck. It happens to us every single day.

The first question is always "How do you know you have a problem?" Two things must be satisfied in this situation:

1. You had a certain expectation for what was supposed to happen
2. Something *other* than that expectation actually happened

While it might seem overly didactic to separate out these two things, one common mistake is to only focus on the second part, i.e. what actually happened. Typically, we see an error message or a warning or some other bad sign and we intuitively know that there is a problem. While it's important to recognize these warning signs, it's equally important to be able to say specifically what your expectation was. What output were you expecting to see? What did you think the answer was going to be?

The more specific you can be with your expectation, the more likely you'll be able to figure out what went wrong. In particular, in many cases it might be that your expectations were incorrect. For example, you might think it's a bug that the `log()` function returns `NaN` when called on a negative number. If you were expecting there to be an error in this situation, then your expectation is incorrect because the `log()` function was specifically designed to return the `NaN` value (indicating an undefined operation) and give a warning when called with negative numbers.

There are two basic approaches to diagnosing and solving problems.

1. Googling
2. Asking a human

Before asking a human, it's usually best to see if you can Google your way out. This can be a real timesaver for all involved. We discuss both approaches below.

How to Google Your Way Out of a Jam

Like with any other programming language, it's essential that you know how to Google your way out of a jam. A related resource in this situation is the [Stack Overflow](#) web site, which is a popular Q&A web site for programming related questions. However, often results from Google will simply point you to Stack Overflow, so Google can serve as useful wrapper around a variety of web site like this.

While we don't exactly have an algorithm for getting unstuck from a jam, here are few tips.

- If you get an error message, **copy and paste the entire error message into Google**. Why? Because, almost surely, someone else has gotten this very same error and has asked a question about it on some forum that Google has indexed. Chances are, that person copy-and-pasted the error message into that forum posting and, presto! You have your answer. Or something close to it.
- For working with certain high-level functions, you can simply Google the name of the function, perhaps with the phrase “R function” following it in case it is a somewhat generic function name. This will usually bring up the help page for the function first, but it will also commonly bring up various tutorials that people have written that use this function. Often, seeing how other people use a certain function can be very helpful in understanding how a function works.
- If you’re trying to learn a new R package, Google “[package name] vignette” and “[package name] tutorial”. Often, someone will have written course slides, a blog post, or a document that walks you through how to use the package.
- If you are struggling with how to write the code for a plot, try using Google Images. Google “r [name or description of plot]” (e.g., “r pareto plot”) and then choose the “Images” tab in the results. Scroll through to find something that looks like the plot you want to create, and then check the image’s website. It will often include the R code used to create the image.

Asking for Help

In the event that Googling around does not find you an answer, you may need to wade into a forum like Stack Overflow, Reddit, or perhaps the R-help mailing list to get help with a problem. When asking questions on a forum, there are some general rules that are always worth following.

- Read the posting guide for the forum, if there is one. This may cover the rules of posting to the forum and will save you a bit of grief later on.
- If the forum has a FAQ, read it. The answer to your question may already be there.
- State the problem you’re trying to solve, along with the approach that you took that lead to your problem. In particular, **state what you were expecting to see from your code**. Sometimes the source of your problem lies higher up the chain than you might think. In particular, **your expectations may be incorrect**.
- Show that you’ve done your homework and have tried to diagnose the problem yourself, read the help page, Googled for answers, etc.
- **Provide a reproducible example of your problem.** This cannot be stressed enough. In order for others to help you, it’s critical that they can reproduce the problem on their own machines. Otherwise, they will have to diagnose your problem from afar, and much like with human beings, this is often very difficult to do. If your problem involves massive amounts of computation, try to come up with a simple example that reproduces the same problem. Other people will not download your 100 GB dataset just so they can reproduce your error message.

2. Advanced R Programming

This course covers advanced topics in R programming that are necessary for developing powerful, robust, and reusable data science tools. Topics covered include functional programming in R, robust error handling, object oriented programming, profiling and benchmarking, debugging, and proper design of functions. Upon completing this course you will be able to identify and abstract common data analysis tasks and to encapsulate them in user-facing functions. Because every data science environment encounters unique data challenges, there is always a need to develop custom software specific to your organization's mission. You will also be able to define new data types in R and to develop a universe of functionality specific to those data types to enable cleaner execution of data science tasks and stronger reusability within a team.

The learning objectives of the chapter are:

- Describe the control flow of an R program
- Write a function that abstracts a single concept/procedure
- Describe functional programming concepts
- Write functional programming code using the `purrr` package
- Manipulate R expressions to “compute on the language”
- Describe the semantics of R environments
- Implement exception handling routines in R functions
- Design and Implement a new S3, S4, or reference class with generics and methods
- Apply debugging tools to identify bugs in R programs
- Apply profiling and timing tools to optimize R code
- Describe the principles of tidyverse functions

2.1 Control Structures

Note: Some of the material in this section is adapted from [R Programming for Data Science](#).

The learning objectives of the section are:

- Describe the control flow of an R program

Control structures in R allow you to control the flow of execution of a series of R expressions. Basically, control structures allow you to put some “logic” into your R code, rather than just always executing the same R code every time. Control structures allow you to respond to inputs or to features of the data and execute different R expressions accordingly.

Commonly used control structures are

- `if` and `else`: testing a condition and acting on it
- `for`: execute a loop a fixed number of times
- `break`: break the execution of a loop
- `next`: skip an iteration of a loop

Most control structures are not used in interactive sessions, but rather when writing functions or longer expressions. However, these constructs do not have to be used in functions and it's a good idea to become familiar with them before we delve into functions.

`if-else`

The `if-else` combination is probably the most commonly used control structure in R (or perhaps any language). This structure allows you to test a condition and act on it depending on whether it's true or false.

For starters, you can just use the `if` statement.

```
if(<condition>) {  
    ## do something  
}  
## Continue with rest of code
```

The above code does nothing if the condition is false. If you have an action you want to execute when the condition is false, then you need an `else` clause.

```
if(<condition>) {  
    ## do something  
} else {  
    ## do something else  
}
```

You can have a series of tests by following the initial `if` with any number of `else ifs`.

```
if(<condition1>) {  
    ## do something  
} else if(<condition2>) {  
    ## do something different  
} else {  
    ## do something different  
}
```

Here is an example of a valid if/else structure.

```
## Generate a uniform random number
x <- runif(1, 0, 10)
if(x > 3) {
    y <- 10
} else {
    y <- 0
}
```

The value of `y` is set depending on whether `x > 3` or not.

Of course, the `else` clause is not necessary. You could have a series of if clauses that always get executed if their respective conditions are true.

```
if(<condition1>) {
}
if(<condition2>) {
}
```

for Loops

For loops are pretty much the only looping construct that you will need in R. While you may occasionally find a need for other types of loops, in most data analysis situations, there are very few cases where a for loop isn't sufficient.

In R, for loops take an iterator variable and assign it successive values from a sequence or vector. For loops are most commonly used for iterating over the elements of an object (list, vector, etc.)

```
numbers <- rnorm(10)
for(i in 1:10) {
    print(numbers[i])
}
[1] -0.9567815
[1] 1.347491
[1] -0.03158058
[1] 0.5960358
[1] 1.133312
[1] -0.7085361
[1] 1.525453
[1] 1.114152
[1] -0.1214943
[1] -0.2898258
```

This loop takes the `i` variable and in each iteration of the loop gives it values 1, 2, 3, ..., 10, executes the code within the curly braces, and then the loop exits.

The following three loops all have the same behavior.

```
x <- c("a", "b", "c", "d")

for(i in 1:4) {
  ## Print out each element of 'x'
  print(x[i])
}

[1] "a"
[1] "b"
[1] "c"
[1] "d"
```

The `seq_along()` function is commonly used in conjunction with for loops in order to generate an integer sequence based on the length of an object (in this case, the object `x`).

```
## Generate a sequence based on length of 'x'
for(i in seq_along(x)) {
  print(x[i])
}

[1] "a"
[1] "b"
[1] "c"
[1] "d"
```

It is not necessary to use an index-type variable.

```
for(letter in x) {
  print(letter)
}

[1] "a"
[1] "b"
[1] "c"
[1] "d"
```

For one line loops, the curly braces are not strictly necessary.

```
for(i in 1:4) print(x[i])
[1] "a"
[1] "b"
[1] "c"
[1] "d"
```

However, curly braces are sometimes useful even for one-line loops, because that way if you decide to expand the loop to multiple lines, you won't be burned because you forgot to add curly braces (and you *will* be burned by this).

Nested for loops

for loops can be nested inside of each other.

```
x <- matrix(1:6, 2, 3)

for(i in seq_len(nrow(x))) {
  for(j in seq_len(ncol(x))) {
    print(x[i, j])
  }
}
```

Nested loops are commonly needed for multidimensional or hierarchical data structures (e.g. matrices, lists). Be careful with nesting though. Nesting beyond 2 to 3 levels often makes it difficult to read or understand the code. If you find yourself in need of a large number of nested loops, you may want to break up the loops by using functions (discussed later).

next, break

`next` is used to skip an iteration of a loop.

```
for(i in 1:100) {
  if(i <= 20) {
    ## Skip the first 20 iterations
    next
  }
  ## Do something here
}
```

`break` is used to exit a loop immediately, regardless of what iteration the loop may be on.

```
for(i in 1:100) {
  print(i)

  if(i > 20) {
    ## Stop loop after 20 iterations
    break
  }
}
```

Summary

- Control structures like `if-else` and `for` allow you to control the flow of an R program.
- Control structures mentioned here are primarily useful for writing programs; for command-line interactive work, the “apply” functions are typically more useful.

2.2 Functions

The learning objectives of the section are:

- Write a function that abstracts a single concept/procedure

The development of functions in R represents the next level of R programming, beyond executing commands at the command line and writing scripts containing multiple R expressions. When writing R functions, one has to consider the following things:

1. Functions are used to **encapsulate** a sequence of expressions that are executed together to achieve a specific goal. A single function typically does “one thing well”—often taking some input and then generating output that can potentially be handed off to another function for further processing. Drawing the lines where functions begin and end is a key skill for writing functions. When writing a function, it’s important to ask yourself *what do I want to encapsulate?*
2. There is going to be a **user** who will desire the ability to modify certain aspects of your code to match their specific needs or application. Aspects of your code that can be modified often become *function arguments* that can be specified by the user. This user can range from yourself (at a later date) to people you have never met using your code for purposes you never dreamed of. When writing any function it’s important to ask *what will the user want to modify in this function?* Ultimately, the answer to this question will lead to the function’s **interface**.

Code

Often we start out analyzing data by writing straight R code at the console. This code is designed to accomplish a single task—whatever it is that we are trying to do *right now*. For example, consider the following code that operates on download logs published by RStudio from their mirror of the Comprehensive R Archive Network (CRAN). This code counts the number of times the `filehash` package was downloaded on July 20, 2016.

```
library(readr)
library(dplyr)

## Download data from RStudio (if we haven't already)
if(!file.exists("data/2016-07-20.csv.gz")) {
  download.file("http://cran-logs.rstudio.com/2016/2016-07-20.csv.gz",
                "data/2016-07-20.csv.gz")
}
cran <- read_csv("data/2016-07-20.csv.gz", col_types = "ccicccccci")
cran %>% filter(package == "filehash") %>% nrow
[1] 179
```

This computation is fairly straightforward and if one were only interested in knowing the number of downloads for this package on this day, there would be little more to say about the code. However, there are a few aspects of this code that one might want to modify or expand on:

- the **date**: this code only reads data for July 20, 2016. But what about data from other days? Note that we would first need to obtain that data if we were interested in knowing download statistics from other days.
- the **package**: this code only returns the number of downloads for the `filehash` package. However, there are many other packages on CRAN and we may want to know how many times these other packages were downloaded.

Once we've identified which aspects of a block of code we might want to modify or vary, we can take those things and abstract them to be arguments of a function.

Function interface

The following function has two arguments:

- `pkgname`, the name of the package as a character string
- `date`, a character string indicating the date for which you want download statistics, in year-month-day format

Given the date and package name, the function downloads the appropriate download logs from the RStudio server, reads the CSV file, and then returns the number of downloads for the package.

```
library(dplyr)
library(readr)

## pkgname: package name (character)
## date: YYYY-MM-DD format (character)
num_download <- function(pkgname, date) {
  ## Construct web URL
  year <- substr(date, 1, 4)
  src <- sprintf("http://cran-logs.rstudio.com/%s/%s.csv.gz",
                 year, date)

  ## Construct path for storing local file
  dest <- file.path("data", basename(src))

  ## Don't download if the file is already there!
  if(!file.exists(dest))
    download.file(src, dest, quiet = TRUE)

  cran <- read_csv(dest, col_types = "ccicccccci", progress = FALSE)
  cran %>% filter(package == pkgname) %>% nrow
}
```

Now we can call our function using whatever date or package name we choose.

```
num_download("filehash", "2016-07-20")
[1] 179
```

We can look up the downloads for a different package on a different day.

```
num_download("Rcpp", "2016-07-19")
[1] 13572
```

Note that for this date, the CRAN log file had to be downloaded separately because it had not yet been downloaded.

Default values

The way that the `num.download()` function is currently specified, the user must enter the date and package name each time the function is called. However, it may be that there is a logical “default date” for which we always want to know the number of downloads, for any package. We can set a **default value** for the date argument, for example, to be July 20, 2016. In that case, if the `date` argument is not explicitly set by the user, the function can use the default value. The revised function might look as follows:

```
num_download <- function(pkgname, date = "2016-07-20") {
  year <- substr(date, 1, 4)
  src <- sprintf("http://cran-logs.rstudio.com/%s/%s.csv.gz",
                year, date)
  dest <- file.path("data", basename(src))
  if(!file.exists(dest))
    download.file(src, dest, quiet = TRUE)
  cran <- read_csv(dest, col_types = "ccicccccci", progress = FALSE)
  cran %>% filter(package == pkgname) %>% nrow
}
```

Now we can call the function in the following manner. Notice that we do not specify the `date` argument.

```
num_download("Rcpp")
[1] 14761
```

Default values play a critical role in R functions because R functions are often called *interactively*. When using R in interactive mode, it can be a pain to have to specify the value of every argument in every instance of calling the function. Sometimes we want to call a function multiple times while varying a single argument (keeping the other arguments at a sensible default).

Also, function arguments have a tendency to proliferate. As functions mature and are continuously developed, one way to add more functionality is to increase the number of

arguments. But if these new arguments do not have sensible default values, then users will generally have a harder time using the function.

As a function author, you have tremendous influence over the user's behavior by specifying defaults, so take care in choosing them. However, just note that a judicious use of default values can greatly improve the user experience with respect to your function.

Re-factoring code

Now that we have a function written that handles the task at hand in a more general manner (i.e. it can handle any package and any date), it is worth taking a closer look at the function and asking whether it is written in the most useful possible manner. In particular, it could be argued that this function does too many things:

1. Construct the path to the remote and local log file
2. Download the log file (if it doesn't already exist locally)
3. Read the log file into R
4. Find the package and return the number of downloads

It might make sense to abstract the first two things on this list into a separate function. For example, we could create a function called `check_for_logfile()` to see if we need to download the log file and then `num_download()` could call this function.

```
check_for_logfile <- function(date) {
  year <- substr(date, 1, 4)
  src <- sprintf("http://cran-logs.rstudio.com/%s/%s.csv.gz",
                year, date)
  dest <- file.path("data", basename(src))
  if(!file.exists(dest)) {
    val <- download.file(src, dest, quiet = TRUE)
    if(!val)
      stop("unable to download file ", src)
  }
  dest
}
```

This file takes the original download code from `num_download()` and adds a bit of error checking to see if `download.file()` was successful (if not, an error is thrown with `stop()`).

Now the `num_download()` function is somewhat simpler.

```
num_download <- function(pkgname, date = "2016-07-20") {
  dest <- check_for_logfile(date)
  cran <- read_csv(dest, col_types = "ccicccccci", progress = FALSE)
  cran %>% filter(package == pkgname) %>% nrow
}
```

In addition to being simpler to read, another key difference is that the `num_download()` function does not need to know anything about downloading or URLs or files. All it knows is that there is a function `check_for_logfile()` that just deals with getting the data to your computer. From there, we can just read the data with `read_csv()` and get the information we need. This is the value of abstraction and writing functions.

Dependency Checking

The `num_downloads()` function depends on the `readr` and `dplyr` packages. Without them installed, the function won't run. Sometimes it is useful to check to see that the needed packages are installed so that a useful error message (or other behavior) can be provided for the user.

We can write a separate function to check that the packages are installed.

```
check_pkg_deps <- function() {
  if(!require(readr)) {
    message("installing the 'readr' package")
    install.packages("readr")
  }
  if(!require(dplyr))
    stop("the 'dplyr' package needs to be installed first")
}
```

There are a few things to note about this function. First, it uses the `require()` function to attempt to load the `readr` and `dplyr` packages. The `require()` function is similar to `library()`, however `library()` stops with an error if the package cannot be loaded whereas `require()` returns `TRUE` or `FALSE` depending on whether the package can be loaded or not. For both functions, if the package is available, it is loaded and attached to the `search()` path.

Typically, `library()` is good for interactive work because you usually can't go on without a specific package (that's why you're loading it in the first place!). On the other hand, `require()` is good for programming because you may want to engage in different behaviors depending on which packages are not available.

For example, in the above function, if the `readr` package is not available, we go ahead and install the package for the user (along with providing a message). However, if we cannot load the `dplyr` package we throw an error. This distinction in behaviors for `readr` and `dplyr` is a bit arbitrary in this case, but it illustrates the flexibility that is afforded by using `require()` versus `library()`.

Now, our updated function can check for package dependencies.

```
num_download <- function(pkgname, date = "2016-07-20") {
  check_pkg_deps()
  dest <- check_for_logfile(date)
  cran <- read_csv(dest, col_types = "ccicccccci", progress = FALSE)
  cran %>% filter(package == pkgname) %>% nrow
}
```

Vectorization

One final aspect of this function that is worth noting is that as currently written it is not *vectorized*. This means that each argument must be a single value—a single package name and a single date. However, in R, it is a common paradigm for functions to take vector arguments and for those functions to return vector or list results. Often, users are bitten by unexpected behavior because a function is assumed to be vectorized when it is not.

One way to vectorize this function is to allow the `pkgname` argument to be a character vector of package names. This way we can get download statistics for multiple packages with a single function call. Luckily, this is fairly straightforward to do. The two things we need to do are

1. Adjust our call to `filter()` to grab rows of the data frame that fall within a vector of package names
2. Use a `group_by() %>% summarize()` combination to count the downloads *for each* package.

```
## 'pkgname' can now be a character vector of names
num_download <- function(pkgname, date = "2016-07-20") {
  check_pkg_deps()
  dest <- check_for_logfile(date)
  cran <- read_csv(dest, col_types = "ccicccccci", progress = FALSE)
  cran %>% filter(package %in% pkgname) %>%
    group_by(package) %>%
    summarize(n = n())
}
```

Now we can call the following

```
num_download(c("filehash", "weathermetrics"))
# A tibble: 2 × 2
  package     n
  <chr>   <int>
1 filehash    179
2 weathermetrics    7
```

Note that the output of `num_download()` has changed. While it previously returned an integer vector, the vectorized function returns a data frame. If you are authoring a function that

is used by many people, it is usually wise to give them some warning before changing the nature of the output.

Vectorizing the `date` argument is similarly possible, but it has the added complication that for each date you need to download another log file. We leave this as an exercise for the reader.

Argument Checking

Checking that the arguments supplied by the reader are proper is a good way to prevent confusing results or error messages from occurring later on in the function. It is also a useful way to enforce documented requirements for a function.

In this case, the `num_download()` function is expecting both the `pkgname` and `date` arguments to be character vectors. In particular, the `date` argument should be a character vector of length 1. We can check the class of an argument using `is.character()` and the length using the `length()` function.

The revised function with argument checking is as follows.

```
num_download <- function(pkgname, date = "2016-07-20") {
  check_pkg_deps()

  ## Check arguments
  if(!is.character(pkgname))
    stop("'pkgname' should be character")
  if(!is.character(date))
    stop("'date' should be character")
  if(length(date) != 1)
    stop("'date' should be length 1")

  dest <- check_for_logfile(date)
  cran <- read_csv(dest, col_types = "ccicccccci",
    progress = FALSE)
  cran %>% filter(package %in% pkgname) %>%
    group_by(package) %>%
    summarize(n = n())
}
```

Note that here, we chose to `stop()` and throw an error if the argument was not of the appropriate type. However, an alternative would have been to simply coerce the argument to be of character type using the `as.character()` function.

```
num_download("filehash", c("2016-07-20", "2016-0-21"))
Error in num_download("filehash", c("2016-07-20", "2016-0-21")): 'date' should be length 1
```

R package

R packages are collections of functions that together allow one to conduct a series of related operations. We will not go into detail about R packages here, but we bring them up only to

indicate that they are the natural evolution of writing many functions. R packages similarly have an interface or API which specifies to the user what functions he/she can call in their own code. The development and maintenance of R packages is the major focus of the next chapter.

When Should I Write a Function?

Deciding when to write a function depends on the context in which you are programming in R. For a one-off type of activity, it's probably not worth considering the design of a function or set of functions. However, in our experience, there are relatively few one-off scenarios. In particular, such a scenario implies that whatever you did worked on the very first try.

In reality, we often have to repeat certain tasks or we have to share code with others. Sometimes those “other people” are simply ourselves 3 months later. As the great [Karl Broman](#) once famously said

Your closest collaborator is you six months ago, but you don't reply to emails.

This comment relates to the general question of whether some code will ever have any **users**, including yourself later on. If the code will likely have more than one user, they will benefit from the abstraction and simplification afforded by encapsulating the code in functions and providing a clean interface.

In Roger's book, *Executive Data Science*, he writes about when to write a function:

- If you're going to do something **once** (that does happen on occasion), just write some code and *document it very well*. The important thing is that you want to make sure that you understand what the code does, and so that requires both writing the code well and writing documentation. You want to be able to reproduce it later on if you ever come back to it, or if someone else comes back to it.
- If you're going to do something **twice**, write a function. This allows you to abstract a small piece of code, and it forces you to define an interface, so you have well defined inputs and outputs.
- If you're going to do something **three** times or more, you should think about writing a small package. It doesn't have to be commercial level software, but a small package which encapsulates the set of operations that you're going to be doing in a given analysis. It's also important to write some real documentation so that people can understand what's supposed to be going on, and can apply the software to a different situation if they have to.

Summary

Developing functions is a key aspect of programming in R and typically involves a bottom-up process.

- Code is written to accomplish a specific task or a specific instance of a task.

- The code is examined to identify key aspects that may be modified by other users; these aspects are abstracted out of the code and made into arguments of a function.
- Functions are written to accomplish more general versions of a task; specific instances of the task are indicated by setting values of function arguments.
- Function code can be re-factored to provide better modularity and to divide functions into specific sub-tasks.
- Functions can be assembled and organized into R packages.

2.3 Functional Programming

The learning objectives of the section are:

- Describe functional programming concepts
- Write functional programming code using the `purrr` package

What is Functional Programming?

Functional programming is a programming philosophy based on [lambda calculus](#). Lambda calculus was created by [Alonzo Church](#), the PhD adviser to [Alan Turing](#) who is known for his role in cracking the encryption of the Nazi's Enigma machine during World War Two. Functional programming has been a popular approach ever since it helped bring down the Third Reich.

Functional programming concentrates on four constructs:

1. Data (numbers, strings, etc)
2. Variables (function arguments)
3. Functions
4. Function Applications (evaluating functions given arguments and/or data)

By now you're used to treating variables inside of functions as data, whether they're values like numbers and strings, or they're data structures like lists and vectors. With functional programming you can also consider the possibility that you can provide a function as an argument to another function, and a function can return another function as its result.

If you've used functions like `sapply()` or `args()` then it's easy to imagine how functions as arguments to other functions can be used. In the case of `sapply()` the provided function is applied to data, and in the case of `args()` information about the function is returned. What's rarer to see is a function that returns a function when it's evaluated. Let's look at a small example of how this can work:

```

adder_maker <- function(n){
  function(x){
    n + x
  }
}

add2 <- adder_maker(2)
add3 <- adder_maker(3)

add2(5)
[1] 7
add3(5)
[1] 8

```

In the example above the function `adder_maker()` returns a function with no name. The function returned adds `n` to its only argument `x`.

Core Functional Programming Functions

There are groups of functions that are essential for functional programming. In most cases they take a function and a data structure as arguments, and that function is applied to that data structure in some way. The `purrr` library contains many of these functions and we'll be using it throughout this section. Function programming is concerned mostly with lists and vectors. I may refer to just lists or vectors, but you should know that what applies for lists generally applies for vectors and vice-versa.

Map

The map family of functions applies a function to the elements of a data structure, usually a list or a vector. The function is evaluated once for each element of the vector with the vector element as the first argument to the function. The return value is the same kind of data structure (a list or vector) but with every element replaced by the result of the function being evaluated with the corresponding element as the argument to the function. In the `purrr` package the `map()` function returns a list, while the `map_lgl()`, `map_chr()`, and `map_db1()` functions return vectors of logical values, strings, or numbers respectively. Let's take a look at a few examples:

```

library(purrr)

map_chr(c(5, 4, 3, 2, 1), function(x){
  c("one", "two", "three", "four", "five")[x]
})
[1] "five"   "four"   "three"  "two"    "one"

map_lgl(c(1, 2, 3, 4, 5), function(x){
  x > 3
})
[1] FALSE FALSE FALSE  TRUE  TRUE

```

Think about evaluating each function above with just one of the arguments in the specified numeric vector, and then combining all of those function results into one vector.

The `map_if()` function takes as its arguments a list or vector containing data, a predicate function, and then a function to be applied. A predicate function is a function that returns `TRUE` or `FALSE` for each element in the provided list or vector. In the case of `map_if()`: if the predicate functions evaluates to `TRUE`, then the function is applied to the corresponding vector element, however if the predicate function evaluates to `FALSE` then the function is not applied. The `map_if()` function always returns a list, so I'm piping the result of `map_if()` to `unlist()` so it look prettier:

```
map_if(1:5, function(x){
  x %% 2 == 0
},
function(y){
  y^2
}) %>% unlist()
[1] 1 4 3 16 5
```

Notice how only the even numbers are squared, while the odd numbers are left alone.

The `map_at()` function only applies the provided function to elements of a vector specified by their indexes. `map_at()` always returns a list so like before I'm piping the result to `unlist()`:

```
map_at(seq(100, 500, 100), c(1, 3, 5), function(x){
  x - 10
}) %>% unlist()
[1] 90 200 290 400 490
```

Like we expected to happen the provided function is only applied to the first, third, and fifth element of the vector provided.

In each of the examples above we have only been mapping a function over one data structure, however you can map a function over two data structures with the `map2()` family of functions. The first two arguments should be two vectors of the same length, followed by a function which will be evaluated with an element of the first vector as the first argument and an element of the second vector as the second argument. For example:

```
map2_chr(letters, 1:26, paste)
[1] "a 1"  "b 2"  "c 3"  "d 4"  "e 5"  "f 6"  "g 7"  "h 8"  "i 9"  "j 10"
[11] "k 11" "l 12" "m 13" "n 14" "o 15" "p 16" "q 17" "r 18" "s 19" "t 20"
[21] "u 21" "v 22" "w 23" "x 24" "y 25" "z 26"
```

The `pmap()` family of functions is similar to `map2()`, however instead of mapping across two vectors or lists, you can map across any number of lists. The list argument is a list of lists that the function will map over, followed by the function that will applied:

```
pmap_chr(list(
  list(1, 2, 3),
  list("one", "two", "three"),
  list("uno", "dos", "tres")
), paste)
[1] "1 one uno"     "2 two dos"    "3 three tres"
```

Mapping is a powerful technique for thinking about how to apply computational operations to your data.

Reduce

List or vector reduction iteratively combines the first element of a vector with the second element of a vector, then that combined result is combined with the third element of the vector, and so on until the end of the vector is reached. The function to be applied should take at least two arguments. Where mapping returns a vector or a list, reducing should return a single value. Some examples using `reduce()` are illustrated below:

```
reduce(c(1, 3, 5, 7), function(x, y){
  message("x is ", x)
  message("y is ", y)
  message("")
  x + y
})
x is 1
y is 3

x is 4
y is 5

x is 9
y is 7

[1] 16
```

On the first iteration `x` has the value 1 and `y` has the value 3, then the two values are combined (they're added together). On the second iteration `x` has the value of the result from the first iteration (4) and `y` has the value of the third element in the provided numeric vector (5). This process is repeated for each iteration. Here's a similar example using string data:

```
reduce(letters[1:4], function(x, y){  
  message("x is ", x)  
  message("y is ", y)  
  message("")  
  paste0(x, y)  
})  
x is a  
y is b  
  
x is ab  
y is c  
  
x is abc  
y is d  
  
[1] "abcd"
```

By default `reduce()` starts with the first element of a vector and then the second element and so on. In contrast the `reduce_right()` function starts with the last element of a vector and then proceeds to the second to last element of a vector and so on:

```
reduce_right(letters[1:4], function(x, y){  
  message("x is ", x)  
  message("y is ", y)  
  message("")  
  paste0(x, y)  
})  
x is d  
y is c  
  
x is dc  
y is b  
  
x is dcba  
y is a  
  
[1] "dcba"
```

Search

You can search for specific elements of a vector using the `contains()` and `detect()` functions. `contains()` will return `TRUE` if a specified element is present in a vector, otherwise it returns `FALSE`:

```
contains(letters, "a")
[1] TRUE
contains(letters, "A")
[1] FALSE
```

The `detect()` function takes a vector and a predicate function as arguments and it returns the first element of the vector for which the predicate function returns `TRUE`:

```
detect(20:40, function(x){
  x > 22 && x %% 2 == 0
})
[1] 24
```

The `detect_index()` function takes the same arguments, however it returns the index of the provided vector which contains the first element that satisfies the predicate function:

```
detect_index(20:40, function(x){
  x > 22 && x %% 2 == 0
})
[1] 5
```

Filter

The group of functions that includes `keep()`, `discard()`, `every()`, and `some()` are known as filter functions. Each of these functions takes a vector and a predicate function. For `keep()` only the elements of the vector that satisfy the predicate function are returned while all other elements are removed:

```
keep(1:20, function(x){
  x %% 2 == 0
})
[1]  2  4  6  8 10 12 14 16 18 20
```

The `discard()` function works similarly, it only returns elements that don't satisfy the predicate function:

```
discard(1:20, function(x){
  x %% 2 == 0
})
[1]  1  3  5  7  9 11 13 15 17 19
```

The `every()` function returns `TRUE` only if every element in the vector satisfies the predicate function, while the `some()` function returns `TRUE` if at least one element in the vector satisfies the predicate function:

```
every(1:20, function(x){  
  x %% 2 == 0  
})  
  
some(1:20, function(x){  
  x %% 2 == 0  
})
```

Compose

Finally, the `compose()` function combines any number of functions into one function:

```
n_unique <- compose(length, unique)  
# The composition above is the same as:  
# n_unique <- function(x){  
#   length(unique(x))  
# }  
  
rep(1:5, 1:5)  
[1] 1 2 2 3 3 4 4 4 4 5 5 5 5 5  
  
n_unique(rep(1:5, 1:5))  
[1] 5
```

Functional Programming Concepts

Partial Application

Partial application of functions can allow functions to behave a little like data structures. Using the `partial()` function from the `purrr` package you can specify some of the arguments of a function, and then `partial()` will return a function that only takes the unspecified arguments. Let's take a look at a simple example:

```
library(purrr)  
  
mult_three_n <- function(x, y, z){  
  x * y * z  
}  
  
mult_by_15 <- partial(mult_three_n, x = 3, y = 5)  
  
mult_by_15(z = 4)  
[1] 60
```

By using partial application you can bind some data to the arguments of a function before using that function elsewhere.

Side Effects

Side effects of functions occur whenever a function interacts with the “outside world” – reading or writing data, printing to the console, and displaying a graph are all side effects. The results of side effects are one of the main motivations for writing code in the first place! Side effects can be tricky to handle though, since the order in which functions with side effects are executed often matters and there are variables that are external to the program (the relative location of some data). If you want to evaluate a function across a data structure you should use the `walk()` function from `purrr`. Here’s a simple example:

```
library(purrr)

walk(c("Friends, Romans, countrymen",
      "lend me your ears;",
      "I come to bury Caesar",
      "not to praise him."), message)
Friends, Romans, countrymen,
lend me your ears;
I come to bury Caesar,
not to praise him.
```

Recursion

Recursion is very powerful tool, both mentally and in software development, for solving problems. Recursive functions have two main parts: a few easy to solve problems called “base cases,” and then a case for more complicated problems where **the function is called inside of itself**. The central philosophy of recursive programming is that problems can be broken down into simpler parts, and then combining those simple answers results in the answer to a complex problem.

Imagine you wanted to write a function that adds together all of the numbers in a vector. You could of course accomplish this with a loop:

```
vector_sum_loop <- function(v){
  result <- 0
  for(i in v){
    result <- result + i
  }
  result
}

vector_sum_loop(c(5, 40, 91))
[1] 136
```

You could also think about how to solve this problem recursively. First ask yourself: what’s the base case of finding the sum of a vector? If the vector only contains one element, then the sum is just the value of that element. In the more complex case the vector has more than one element. We can remove the first element of the vector, but then what should we

do with the rest of the vector? Thankfully we have a function for computing the sum of all of the elements of a vector because we're writing that function right now! So we'll add the value of the first element of the vector to whatever the cumulative sum is of the rest of the vector. The resulting function is illustrated below:

```
vector_sum_rec <- function(v){  
  if(length(v) == 1){  
    v  
  } else {  
    v[1] + vector_sum_rec(v[-1])  
  }  
}  
  
vector_sum_rec(c(5, 40, 91))  
[1] 136
```

Another useful exercise for thinking about applications for recursion is computing the Fibonacci sequence. The Fibonacci sequence is a sequence of integers that starts: 0, 1, 1, 2, 3, 5, 8 where each proceeding integer is the sum of the previous two integers. This fits into a recursive mental framework very nicely since each subsequent number depends on the previous two numbers.

Let's write a function to computes the nth digit of the Fibonacci sequence such that the first number in the sequence is 0, the second number is 1, and then all proceeding numbers are the sum of the n - 1 and the n - 2 Fibonacci number. It is immediately evident that there are three base cases:

1. n must be greater than 0.
2. When n is equal to 1, return 0.
3. When n is equal to 2, return 1.

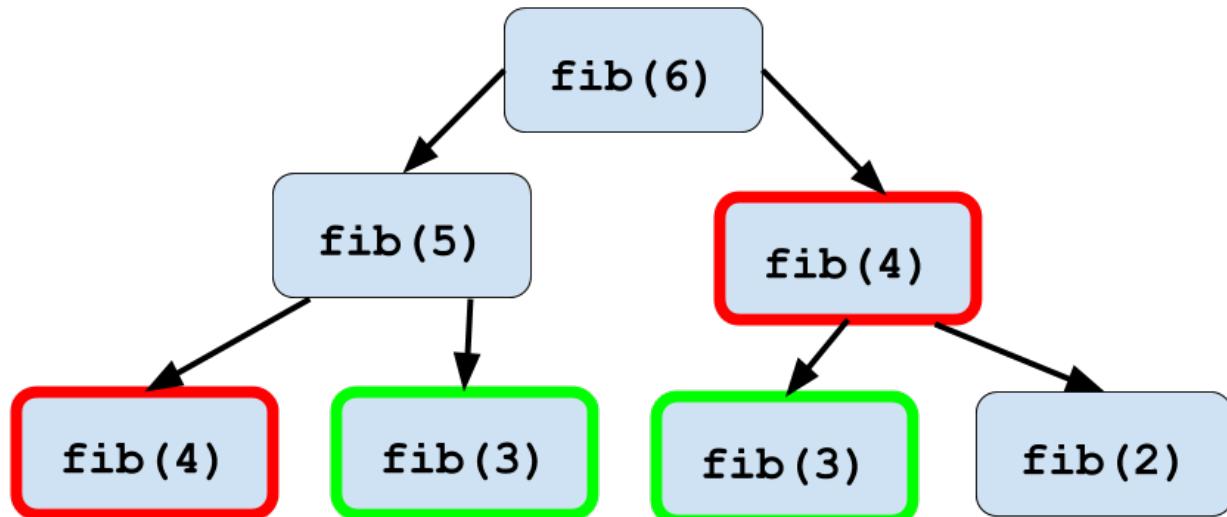
And then the recursive case:

- Otherwise return the sum of the n - 1 Fibonacci number and the n - 2 Fibonacci number.

Let's turn those words into code:

```
fib <- function(n){  
  stopifnot(n > 0)  
  if(n == 1){  
    0  
  } else if(n == 2){  
    1  
  } else {  
    fib(n - 1) + fib(n - 2)  
  }  
}  
  
fib(1)  
[1] 0  
fib(2)  
[1] 1  
fib(3)  
[1] 1  
fib(4)  
[1] 2  
fib(5)  
[1] 3  
fib(6)  
[1] 5  
fib(7)  
[1] 8  
  
map_dbl(1:12, fib)  
[1]  0  1  1  2  3  5  8 13 21 34 55 89
```

Looks like it's working well! There is one optimization that we could apply here which comes up in recursive programming often. When you execute the function `fib(6)`, within that function you'll execute `fib(5)` and `fib(4)`. Then within the execution of `fib(5)`, `fib(4)` will be executed again. An illustration of this phenomenon is below:



Memoization of fib() function

This duplication of computation slows down your program significantly as you calculate larger numbers in the Fibonacci sequence. Thankfully you can use a technique called memoization in order to speed this computation up. Memoization stores the value of each calculated Fibonacci number in table so that once a number is calculated you can look it up instead of needing to recalculate it!

Below is an example of a function that can calculate the first 25 Fibonacci numbers. First we'll create a very simple table which is just a vector containing 0, 1, and then 23 NAs. First the `fib_mem()` function will check if the number is in the table, and if it is then it is returned. Otherwise the Fibonacci number is recursively calculated and stored in the table. Notice that we're using the complex assignment operator `<-<` in order to modify the table outside the scope of the function. You'll learn more about the complex operator in the section titled *Expressions & Environments*.

```

fib_tbl <- c(0, 1, rep(NA, 23))

fib_mem <- function(n){
  stopifnot(n > 0)

  if(!is.na(fib_tbl[n])){
    fib_tbl[n]
  } else {
    fib_tbl[n - 1] <-< fib_mem(n - 1)
    fib_tbl[n - 2] <-< fib_mem(n - 2)
    fib_tbl[n - 1] + fib_tbl[n - 2]
  }
}

map_dbl(1:12, fib_mem)
[1]  0  1  1  2  3  5  8 13 21 34 55 89
  
```

It works! But is it any faster than the original `fib()`? Below I'm going to use the `microbenchmark` package in order assess whether `fib()` or `fib_mem()` is faster:

```
library(purrr)
library(microbenchmark)
library(tidyr)
library(magrittr)
library(dplyr)

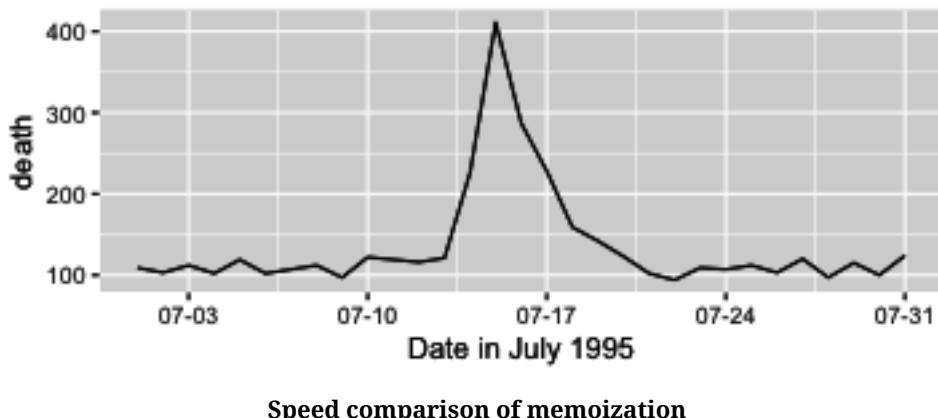
fib_data <- map(1:10, function(x){microbenchmark(fib(x), times = 100)$time})
names(fib_data) <- paste0(letters[1:10], 1:10)
fib_data <- as.data.frame(fib_data)

fib_data %>>%
  gather(num, time) %>%
  group_by(num) %>%
  summarise(med_time = median(time))

memo_data <- map(1:10, function(x){microbenchmark(fib_mem(x))$time})
names(memo_data) <- paste0(letters[1:10], 1:10)
memo_data <- as.data.frame(memo_data)

memo_data %>>%
  gather(num, time) %>%
  group_by(num) %>%
  summarise(med_time = median(time))

plot(1:10, fib_data$med_time, xlab = "Fibonacci Number", ylab = "Median Time (Nanoseconds)",
     pch = 18, bty = "n", xaxt = "n", yaxt = "n")
axis(1, at = 1:10)
axis(2, at = seq(0, 350000, by = 50000))
points(1:10 + .1, memo_data$med_time, col = "blue", pch = 18)
legend(1, 300000, c("Not Memoized", "Memoized"), pch = 18,
       col = c("black", "blue"), bty = "n", cex = 1, y.intersp = 1.5)
```



As you can see as higher Fibonacci numbers are calculated the time it takes to calculate a number with `fib()` grows exponentially, while the time it takes to do the same task with `fib_mem()` stays constant.

Summary

- Functional programming is based on lambda calculus.
- This approach concentrates on data, variables, functions, and function applications.
- It's possible for functions to be able to return other functions.
- The core functional programming concepts can be summarized in the following categories: map, reduce, search, filter, and compose.
- Partial application of functions allows functions to be used like data structures.
- Side effects are difficult to debug although they motivate a huge fraction of computer programming.
- The most important part of understanding recursion is understanding recursion.

2.4 Expressions & Environments

The learning objectives of this section are:

- Manipulate R expressions to “compute on the language”
- Describe the semantics of R environments

Expressions

Expressions are encapsulated operations that can be executed by R. This may sound complicated, but using expressions allows you manipulate code with code! You can create an expression using the `quote()` function. For that function’s argument, just type whatever you would normally type into the R console. For example:

```
two_plus_two <- quote(2 + 2)
two_plus_two
2 + 2
```

You can execute this expressions using the `eval()` function:

```
eval(two_plus_two)
[1] 4
```

You might encounter R code that is stored as a string that you want to evaluate with `eval()`. You can use `parse()` to transform a string into an expression:

```
tpt_string <- "2 + 2"

tpt_expression <- parse(text = tpt_string)

eval(tpt_expression)
[1] 4
```

You can reverse this process and transform an expression into a string using `deparse()`:

```
deparse(two_plus_two)
[1] "2 + 2"
```

One interesting feature about expressions is that you can access and modify their contents like you a `list()`. This means that you can change the values in an expression, or even the function being executed in the expression before it is evaluated:

```
sum_expr <- quote(sum(1, 5))
eval(sum_expr)
[1] 6
sum_expr[[1]]
sum
sum_expr[[2]]
[1] 1
sum_expr[[3]]
[1] 5
sum_expr[[1]] <- quote(paste0)
sum_expr[[2]] <- quote(4)
sum_expr[[3]] <- quote(6)
eval(sum_expr)
[1] "46"
```

You can compose expressions using the `call()` function. The first argument is a string containing the name of a function, followed by the arguments that will be provided to that function.

```
sum_40_50_expr <- call("sum", 40, 50)
sum_40_50_expr
sum(40, 50)
eval(sum_40_50_expr)
[1] 90
```

You can capture the the expression an R user typed into the R console when they executed a function by including `match.call()` in the function the user executed:

```
return_expression <- function(...){
  match.call()
}

return_expression(2, col = "blue", FALSE)
return_expression(2, col = "blue", FALSE)
```

You could of course then manipulate this expression inside of the function you're writing. The example below first uses `match.call()` to capture the expression that the user entered. The first argument of the function is then extracted and evaluated. If the first expression is a number, then a string is returned describing the first argument, otherwise the string "The first argument is not numeric." is returned.

```
first_arg <- function(...){
  expr <- match.call()
  first_arg_expr <- expr[[2]]
  first_arg <- eval(first_arg_expr)
  if(is.numeric(first_arg)){
    paste("The first argument is", first_arg)
  } else {
    "The first argument is not numeric."
  }
}

first_arg(2, 4, "seven", FALSE)
[1] "The first argument is 2"

first_arg("two", 4, "seven", FALSE)
[1] "The first argument is not numeric."
```

Expressions are a powerful tool for writing R programs that can manipulate other R programs.

Environments

Environments are data structures in R that have special properties with regard to their role in how R code is executed and how memory in R is organized. You may not realize it but you're probably already familiar with one environment called the global environment. Environments formalize relationships between variable names and values. When you enter `x <- 55` into the R console what you're saying is: assign the value of 55 to a variable called `x`, and store this assignment in the global environment. The global environment is therefore where most R users do most of their programming and analysis.

You can create a new environment using `new.env()`. You can assign variables in that environment in a similar way to assigning a named element of a list, or you can use `assign()`. You can retrieve the value of a variable just like you would retrieve the named element of a list, or you can use `get()`. Notice that `assign()` and `get()` are opposites:

```

my_new_env <- new.env()
my_new_env$x <- 4
my_new_env$x
[1] 4

assign("y", 9, envir = my_new_env)
get("y", envir = my_new_env)
[1] 9
my_new_env$y
[1] 9

```

You can get all of the variable names that have been assigned in an environment using `ls()`, you can remove an association between a variable name and a value using `rm()`, and you can check if a variable name has been assigned in an environment using `exists()`:

```

ls(my_new_env)
[1] "x" "y"
rm(y, envir = my_new_env)
exists("y", envir = my_new_env)
[1] TRUE
exists("x", envir = my_new_env)
[1] TRUE
my_new_env$x
[1] 4
my_new_env$y
NULL

```

Environments are organized in parent/child relationships such that every environment keeps track of its parent, but parents are unaware of which environments are their children. Usually the relationships between environments is not something you should try to directly control. You can see the parents of the global environment using the `search()` function:

```

search()
[1] ".GlobalEnv"           "package:magrittr"
[3] "package:tidyverse"     "package:microbenchmark"
[5] "package:purrr"         "package:dplyr"
[7] "package:readr"          "package:parallel"
[9] "package:knitr"          "package:stats"
[11] "package:graphics"       "package:grDevices"
[13] "package:utils"          "package:datasets"
[15] "Autoloads"              "package:base"

```

As you can see `package:magrittr` is the parent of `.GlobalEnv`, and `package:tidyverse` is parent of `package:magrittr`, and so on. In general the parent of `.GlobalEnv` is always the last package that was loaded using `library()`. Notice that after I load the `ggplot2` package, that package becomes the parent of `.GlobalEnv`:

```
library(ggplot2)
search()
[1] ".GlobalEnv"           "package:ggplot2"
[3] "package:magrittr"      "package:tidyverse"
[5] "package:microbenchmark" "package:purrr"
[7] "package:dplyr"          "package:readr"
[9] "package:parallel"        "package:knitr"
[11] "package:stats"          "package:graphics"
[13] "package:grDevices"       "package:utils"
[15] "package:datasets"        "Autoloads"
[17] "package:base"
```

Execution Environments

Although there may be several cases where you need to create a new environment using `new.env()`, you will more often create new environments whenever you execute functions. An execution environment is an environment that exists temporarily within the scope of a function that is being executed. For example if we have the following code:

```
x <- 10

my_func <- function(){
  x <- 5
  return(x)
}

my_func()
```

What do you think will be the result of `my_func()`? Make your guess and then take a look at the executed code below:

```
x <- 10

my_func <- function(){
  x <- 5
  return(x)
}

my_func()
[1] 5
```

So what exactly is happening above? First the name `x` is brought assigned the value 10 in the global environment. Then the name `my_func` is being assigned the value of the function `function(){x <- 5; return(x)}` in the global environment. When `my_func()` is executed, a new environment is created called the execution environment which only exists while `my_func()` is running. Inside of the execution environment the name `x` is assigned the value 5. When `return()` is executed it looks first in the execution environment for a value that is assigned to `x`. Then the value 5 is returned. In contrast to the situation above, take a look at this variation:

```
x <- 10

another_func <- function(){
  return(x)
}

another_func()
[1] 10
```

In this situation the execution environment inside of `another_func()` does not contain an assignment for the name `x`, so R looks for an assignment in the parent environment of the execution environment which is the global environment. Since `x` is assigned the value 10 in the global environment 10 is returned.

After seeing the cases above you may be curious if it's possible for an execution environment to manipulate the global environment. You're already familiar with the assignment operator `<-`, however you should also be aware that there's another assignment operator called the *complex assignment operator* which looks like `<<-`. You can use the complex assignment operator to re-assign or even create name-value bindings in the global environment from within an execution environment. In this first example, the function `assign1()` will change the value associated with the name `x`:

```
x <- 10
x
[1] 10

assign1 <- function(){
  x <<- "Wow!"
}

assign1()
x
[1] "Wow!"
```

You can see that the value associated with `x` has been changed from 10 to "Wow!" in the global environment. You can also use `<<-` to assign names to values that have not been yet been defined in the global environment *from inside a function*:

```
a_variable_name
Error in eval(expr, envir, enclos): object 'a_variable_name' not found
exists("a_variable_name")
[1] FALSE

assign2 <- function(){
  a_variable_name <<- "Magic!"
}

assign2()
```

```
exists("a_variable_name")
[1] TRUE
a_variable_name
[1] "Magic!"
```

If you want to see a case for using `<->` in action, see the section of this book about functional programming and the discussion there about memoization.

Summary

- Expressions are a powerful tool for manipulating and executing R code.
- Environments record associations between names and values.
- Execution environments create a scope for variable names inside of functions.

2.5 Error Handling and Generation

The learning objectives of this section are:

- Implement exception handling routines in R functions

What is an error?

Errors most often occur when code is used in a way that it is not intended to be used. For example adding two strings together produces the following error:

```
"hello" + "world"
Error in "hello" + "world": non-numeric argument to binary operator
```

The `+` operator is essentially a function that takes two numbers as arguments and finds their sum. Since neither `"hello"` nor `"world"` are numbers, the R interpreter produces an error. Errors will stop the execution of your program, and they will (hopefully) print an error message to the R console.

In R there are two other constructs which are related to errors: warnings and messages. Warnings are meant to indicate that something seems to have gone wrong in your program that should be inspected. Here's a simple example of a warning being generated:

```
as.numeric(c("5", "6", "seven"))
Warning: NAs introduced by coercion
[1] 5 6 NA
```

The `as.numeric()` function attempts to convert each string in `c("5", "6", "seven")` into a number, however it is impossible to convert `"seven"`, so a warning is generated. Execution of the code is not halted, and an `NA` is produced for `"seven"` instead of a number.

Messages simply print to the R console, though they are generated by an underlying mechanism that is similar to how errors and warning are generated. Here's a small function that will generate a message:

```
f <- function(){
  message("This is a message.")
}

f()
This is a message.
```

Generating Errors

There are a few essential functions for generating errors, warnings, and messages in R. The `stop()` function will generate an error. Let's generate an error:

```
stop("Something erroneous has occurred!")
```

```
Error: Something erroneous has occurred!
```

If an error occurs inside of a function then the name of that function will appear in the error message:

```
name_of_function <- function(){
  stop("Something bad happened.")
}

name_of_function()
Error in name_of_function(): Something bad happened.
```

The `stopifnot()` function takes a series of logical expressions as arguments and if any of them are false an error is generated specifying which expression is false. Let's take a look at an example:

```
error_if_n_is_greater_than_zero <- function(n){
  stopifnot(n <= 0)
  n
}

error_if_n_is_greater_than_zero(5)
Error: n <= 0 is not TRUE
```

The `warning()` function creates a warning, and the function itself is very similar to the `stop()` function. Remember that a warning does not stop the execution of a program (unlike an error.)

```
warning("Consider yourself warned!")
Warning: Consider yourself warned!
```

Just like errors, a warning generated inside of a function will include the name of the function in which it was generated:

```
make_NA <- function(x){
  warning("Generating an NA.")
  NA
}

make_NA("Sodium")
Warning in make_NA("Sodium"): Generating an NA.
[1] NA
```

Messages are simpler than errors or warnings; they just print strings to the R console. You can issue a message with the `message()` function:

```
message("In a bottle.")
In a bottle.
```

When to generate errors or warnings

Stopping the execution of your program with `stop()` should only happen in the event of a catastrophe - meaning only if it is impossible for your program to continue. If there are conditions that you can anticipate that would cause your program to create an error then you should document those conditions so whoever uses your software is aware. Common failure conditions like providing invalid arguments to a function should be checked at the beginning of your program so that the user can quickly realize something has gone wrong. Checking function inputs is a typical use of the `stopifnot()` function.

You can think of a function as kind of contract between you and the user: if the user provides specified arguments, your program will provide predictable results. Of course it's impossible for you to anticipate all of the potential uses of your program, so the results of executing a function can only be predictable with regard to the type of the result. It's appropriate to create a warning when this contract between you and the user is violated. A perfect example of this situation is the result of `as.numeric(c("5", "6", "seven"))`, which we saw before. The user expects a vector of numbers to be returned as the result of `as.numeric()` but "seven" is coerced into being NA, which is not completely intuitive.

R has largely been developed according to the [Unix Philosophy](#) (which is further discussed in Chapter 3), which generally discourages printing text to the console unless something unexpected has occurred. Languages that commonly run on Unix systems like C, C++, and Go are rarely used interactively, meaning that they usually underpin computer infrastructure (computers “talking” to other computers). Messages printed to the console are therefore not very useful since nobody will ever read them and it's not straightforward for other programs to capture and interpret them. In contrast R code is frequently executed

by human beings in the R console, which serves as an interactive environment between the computer and person at the keyboard. If you think your program should produce a message, make sure that the output of the message is primarily meant for a human to read. You should avoid signaling a condition or the result of your program to another program by creating a message.

How should errors be handled?

Imagine writing a program that will take a long time to complete because of a complex calculation or because you're handling a large amount of data. If an error occurs during this computation then you're liable to lose all of the results that were calculated before the error, or your program may not finish a critical task that a program further down your pipeline is depending on. If you anticipate the possibility of errors occurring during the execution of your program then you can design your program to handle them appropriately.

The `tryCatch()` function is the workhorse of handling errors and warnings in R. The first argument of this function is any R expression, followed by conditions which specify how to handle an error or a warning. The last argument, `finally`, specifies a function or expression that will be executed after the expression no matter what, even in the event of an error or a warning.

Let's construct a simple function I'm going to call `beera` that catches errors and warnings gracefully.

```
beera <- function(expr){  
  tryCatch(expr,  
    error = function(e){  
      message("An error occurred:\n", e)  
    },  
    warning = function(w){  
      message("A warning occurred:\n", w)  
    },  
    finally = {  
      message("Finally done!")  
    })  
}
```

This function takes an expression as an argument and tries to evaluate it. If the expression can be evaluated without any errors or warnings then the result of the expression is returned and the message `Finally done!` is printed to the R console. If an error or warning is generated then the functions that are provided to the `error` or `warning` arguments are printed. Let's try this function out with a few examples.

```
beera({  
  2 + 2  
})  
Finally done!  
[1] 4  
  
beera({  
  "two" + 2  
})  
An error occurred:  
Error in "two" + 2: non-numeric argument to binary operator  
  
Finally done!  
  
beera({  
  as.numeric(c(1, "two", 3))  
})  
A warning occurred:  
simpleWarning in doTryCatch(return(expr), name, parentenv, handler): NAs introduced by coercion  
  
Finally done!
```

Notice that we've effectively transformed errors and warnings into messages.

Now that you know the basics of generating and catching errors you'll need to decide when your program should generate an error. My advice to you is to limit the number of errors your program generates as much as possible. Even if you design your program so that it's able to catch and handle errors, the error handling process slows down your program by orders of magnitude. Imagine you wanted to write a simple function that checks if an argument is an even number. You might write the following:

```
is_even <- function(n){  
  n %% 2 == 0  
}  
  
is_even(768)  
[1] TRUE  
  
is_even("two")  
Error in n%%2: non-numeric argument to binary operator
```

You can see that providing a string causes this function to raise an error. You could imagine though that you want to use this function across a list of different data types, and you only want to know which elements of that list are even numbers. You might think to write the following:

```
is_even_error <- function(n){
  tryCatch(n %% 2 == 0,
    error = function(e){
      FALSE
    })
}

is_even_error(714)
[1] TRUE

is_even_error("eight")
[1] FALSE
```

This appears to be working the way you intended, however when applied to more data this function will be seriously slow compared to alternatives. For example I could check that `n` is numeric before treating `n` like a number:

```
is_even_check <- function(n){
  is.numeric(n) && n %% 2 == 0
}

is_even_check(1876)
[1] TRUE

is_even_check("twelve")
[1] FALSE
```



Notice that by using `is.numeric()` before the “AND” operator (`&&`) the expression `n %% 2 == 0` is never evaluated. This is a programming language design feature called “short circuiting.” The expression can never evaluate to `TRUE` if the left hand side of `&&` evaluates to `FALSE`, so the right hand side is ignored.

To demonstrate the difference in the speed of the code we’ll use the `microbenchmark` package to measure how long it takes for each function to be applied to the same data.

```
library(microbenchmark)
microbenchmark(sapply(letters, is_even_check))
```

```
Unit: microseconds
          expr     min      lq      mean    median      uq      max neval
sapply(letters, is_even_check) 46.224 47.7975 61.43616 48.6445 58.4755 167.091   100
```

```
microbenchmark(sapply(letters, is_even_error))

Unit: microseconds
      expr     min      lq      mean    median      uq      max neval
sapply(letters, is_even_error) 640.067 678.0285 906.3037 784.4315 1044.501 2308.931   100
```

The error catching approach is nearly 15 times slower!

Proper error handling is an essential tool for any software developer so that you can design programs that are error tolerant. Creating clear and informative error messages is essential for building quality software. One closing tip I recommend is to put documentation for your software online, including the meaning of the errors that your software can potentially throw. Often a user's first instinct when encountering an error is to search online for that error message, which should lead them to your documentation!

Summary

- Errors, warnings, and messages can be generated within R code using the functions `stop`, `stopifnot`, `warning`, and `message`.
- Catching errors, and providing useful error messaging, can improve user experience with functions but can also slow down code substantially.

2.6 Debugging

The learning objectives of this section are:

- Apply debugging tools to identify bugs in R programs

Debugging is the process of getting your expectations to converge with reality. When writing software in any language, we develop a certain set of expectations about how the software should behave and what it should do. But inevitably, when we run the software, it does something *different* from what we expected. In these situations, we need to engage in a process to determine if

1. Our expectations were incorrect, based on the documented behavior of the software; or
2. There is a problem with the code, such that the programming is not done in a way that will match expectations.

This is the process of debugging.

In the previous section, we discussed what to do when software generates conditions (errors, warnings, messages) in a manner that is completely *expected*. In those cases, we

know that certain functions will generate errors and we want to handle them in a manner that is not the usual way.

This section describes the tools for debugging your software in R. R comes with a set of built-in tools for interactive debugging that can be useful for tracking down the source of problems. These functions are

- `browser()`: an interactive debugging environment that allows you to step through code one expression at a time
- `debug()` / `debugonce()`: a function that initiates the browser within a function
- `trace()`: a function that allows you to temporarily insert pieces of code into other functions to modify their behavior
- `recover()`: a function for navigating the function call stack after a function has thrown an error
- `traceback()`: a function that prints out the function call stack after an error occurs but does nothing if there's no error

`traceback()`

If an error occurs, the easiest thing to do is to immediately call the `traceback()` function. This function returns the function call stack just before the error occurred so that you can see what level of function calls the error occurred. If you have many functions calling each other in succession, the `traceback()` output can be useful for identifying where to go digging first.

For example, the following code gives an error.

```
check_n_value <- function(n) {
  if(n > 0) {
    stop("n should be <= 0")
  }
}
error_if_n_is_greater_than_zero <- function(n){
  check_n_value(n)
  n
}
error_if_n_is_greater_than_zero(5)
Error in check_n_value(n): n should be <= 0
```

Running the `traceback()` function immediately after getting this error would give us

```
traceback()
3: stop("n should be <= 0") at #2
2: check_n_value(n) at #2
1: error_if_n_is_greater_than_zero(5)
```

From the traceback, we can see that the error occurred in the `check_n_value()` function. Put another way, the `stop()` function was called from within the `check_n_value()` function.

Browsing a Function Environment

From the traceback output, it is often possible to determine in which function and on which line of code an error occurs. If you are the author of the code in question, one easy thing to do is to insert a call to the `browser()` function in the vicinity of the error (ideally, *before* the error occurs). The `browser()` function takes no arguments and is just placed wherever you want in the function. Once it is called, you will be in the browser environment, which is much like the regular R workspace environment except that you are inside a function.

```
check_n_value <- function(n) {
  if(n > 0) {
    browser() ## Error occurs around here
    stop("n should be <= 0")
  }
}
```

Now, when we call `error_if_n_is_greater_than_zero(5)`, we will see the following.

```
error_if_n_is_greater_than_zero(5)
Called from: check_n_value(n)
Browse[1]>
```

Tracing Functions

If you have easy access to the source code of a function (and can modify the code), then it's usually easiest to insert `browser()` calls directly into the code as you track down various bugs. However, if you do not have easy access to a function's code, or perhaps a function is inside a package that would require rebuilding after each edit, it is sometimes easier to make use of the `trace()` function to make temporary code modifications.

The simplest use of `trace()` is to just call `trace()` on a function without any other arguments.

```
trace("check_n_value")
Error in trace("check_n_value"): could not find function "check_n_value"
```

Now, whenever `check_n_value()` is called by any other functions, you will see a message printed to the console indicating that the function was called.

```
error_if_n_is_greater_than_zero(5)
Error in check_n_value(n): n should be <= 0
```

Here we can see that `check_n_value()` was called once before the error occurred. But we can do more with `trace()`, such as inserting a call to `browser()` in a specific place, such as right before the call to `stop()`.

We can obtain the expression numbers of each part of a function by calling `as.list()` on the `body()` of a function.

```
as.list(body(check_n_value))
[[1]]
`{`

[[2]]
if (n > 0) {
  stop("n should be <= 0")
}
```

Here, the `if` statement is the second expression in the function (the first “expression” being the very beginning of the function). We can further break down the second expression as follows.

```
as.list(body(check_n_value)[[2]])
[[1]]
`if`


[[2]]
n > 0


[[3]]
{
  stop("n should be <= 0")
}
```

Now we can see the call to `stop()` is the third sub-expression within the second expression of the overall function. We can specify this to `trace()` by passing an integer vector wrapped in a list to the `at` argument.

```
trace("check_n_value", browser, at = list(c(2, 3)))
Error in getFunction(what, where = whereF): no function 'check_n_value' found
```

The `trace()` function has a side effect of modifying the function and converting into a new object of class “function”.

```
check_n_value
function(n) {
  if(n > 0) {
    stop("n should be <= 0")
  }
}
<environment: 0x7fba0097e788>
```

You can see the internally modified code by calling

```
body(check_n_value)
{
  if (n > 0) {
    stop("n should be <= 0")
  }
}
```

Here we can see that the code has been altered to add a call to `browser()` just before the call to `stop()`.

We can add more complex expressions to a function by wrapping them in a call to `quote()` within the the `trace()` function. For example, we may only want to invoke certain behaviors depending on the local conditions of the function.

```
trace("check_n_value", quote({
  if(n == 5) {
    message("invoking the browser")
    browser()
  }
}), at = 2)
Error in getFunction(what, where = whereF): no function 'check_n_value' found
```

Here, we only invoke the `browser()` if `n` is specifically 5.

```
body(check_n_value)
{
  if (n > 0) {
    stop("n should be <= 0")
  }
}
```

Debugging functions within a package is another key use case for `trace()`. For example, if we wanted to insert tracing code into the `glm()` function within the `stats` package, the only addition to the `trace()` call we would need is to provide the namespace information via the `where` argument.

```
trace("glm", browser, at = 4, where = asNamespace("stats"))
Tracing function "glm" in package "namespace:stats"
[1] "glm"
```

Here we show the first few expressions of the modified `glm()` function.

```
body(stats::glm)[1:5]
{
  call <- match.call()
  if (is.character(family))
    family <- get(family, mode = "function", envir = parent.frame())
  {
    .doTrace(browser(), "step 4")
    if (is.function(family))
      family <- family()
  }
  if (is.null(family$family)) {
    print(family)
    stop("'family' not recognized")
  }
}
```

Using `debug()` and `debugonce()`

The `debug()` and `debugonce()` functions can be called on other functions to turn on the “debugging state” of a function. Calling `debug()` on a function makes it such that when that function is called, you immediately enter a browser and can step through the code one expression at a time.

```
## Turn on debugging state for 'lm' function
debug(lm)
```

A call to `debug(f)` where `f` is a function is basically equivalent to `trace(f, browser)` which will call the `browser()` function upon entering the function.

The debugging state is persistent, so once a function is flagged for debugging, it will remain flagged. Because it is easy to forget about the debugging state of a function, the `debugonce()` function turns on the debugging state the next time the function is called, but then turns it off after the browser is exited.

`recover()`

The `recover()` function is not often used but can be an essential tool when debugging complex code. Typically, you do not call `recover()` directly, but rather set it as the function to invoke anytime an error occurs in code. This can be done via the `options()` function.

```
options(error = recover)
```

Usually, when an error occurs in code, the code stops execution and you are brought back to the usual R console prompt. However, when `recover()` is in use and an error occurs, you are given the function call stack and a menu.

```
error_if_n_is_greater_than_zero(5)
Error in check_n_value(n) : n should be <= 0

Enter a frame number, or 0 to exit

1: error_if_n_is_greater_than_zero(5)
2: #2: check_n_value(n)

Selection:
```

Selecting a number from this menu will bring you into that function on the call stack and you will be placed in a browser environment. You can exit the browser and then return to this menu to jump to another function in the call stack.

The `recover()` function is very useful if an error is deep inside a nested series of function calls and it is difficult to pinpoint exactly where an error is occurring (so that you might use `browser()` or `trace()`). In such cases, the `debug()` function is often of little practical use because you may need to step through many many expressions before the error actually occurs. Another scenario is when there is a stochastic element to your code so that errors occur in an unpredictable way. Using `recover()` will allow you to browse the function environment only when the error eventually does occur.

Final Thoughts on Debugging

The debugging tools in any programming language can be essential for tracking down problems in code, especially when the code becomes complex and spans many lines. However, one should not lean on them too heavily so that they become a regular part of the programming process. It is easy to get into a situation where you “throw some code out there” and then let the debugger catch it before something bad happens. If you find yourself coding up a function and then immediately calling `debug()` on it, you are in this situation.

A better approach is to think carefully about what a function should do and then consider how to code it up. A few minutes of careful forethought can often save the hapless programmer hours of debugging.

Summary

- Debugging in R is facilitated with the functions `browser`, `debug`, `trace`, `recover`, and `traceback`.
- These debugging tools should not be used as a crutch when developing functions.

2.7 Profiling and Benchmarking

The learning objectives of this section are:

- Apply profiling and timing tools to optimize R code

Some of the R code that you write will be slow. Slow code often isn't worth fixing in a script that you will only evaluate a few times, as the time it will take to optimize the code will probably exceed the time it takes the computer to run it. However, if you are writing functions that will be used repeatedly, it is often worthwhile to identify slow sections of the code so you can try to improve speed in those sections.

In this section, we will introduce the basics of profiling R code, using functions from two packages, `microbenchmark` and `profvis`. The `profvis` package is fairly new and requires recent versions of both R (version 3.0 or higher) and RStudio. If you are having problems running either package, you should try updating both R and RStudio (the Preview version of RStudio, which will provide full functionality for `profvis`, is available for download [here](#)).

`microbenchmark`

The `microbenchmark` package is useful for running small sections of code to assess performance, as well as for comparing the speed of several functions that do the same thing. The `microbenchmark` function from this package will run code multiple times (100 times is the default) and provide summary statistics describing how long the code took to run across those iterations. The process of timing a function takes a certain amount of time itself. The `microbenchmark` function adjusts for this overhead time by running a certain number of “warm-up” iterations before running the iterations used to time the code.

You can use the `times` argument in `microbenchmark` to customize how many iterations are used. For example, if you are working with a function that is a bit slow, you might want to run the code fewer times when benchmarking (although with slower or more complex code, it likely will make more sense to use a different tool for profiling, like `profvis`).

You can include multiple lines of code within a single call to `microbenchmark`. However, to get separate benchmarks of line of code, you must separate each line by a comma:

```
library(microbenchmark)
microbenchmark(a <- rnorm(1000),
               b <- mean(rnorm(1000)))
Unit: microseconds
      expr     min      1q    mean   median      uq     max
a <- rnorm(1000) 78.286  96.8885 101.5368 100.0850 103.7450 198.914
b <- mean(rnorm(1000)) 84.680 108.2950 114.0244 113.2595 119.3025 241.997
neval
100
100
```

The `microbenchmark` function is particularly useful for comparing functions that take the same inputs and return the same outputs. As an example, say we need a function that can identify days that meet two conditions: (1) the temperature equals or exceeds a threshold temperature (27 degrees Celsius in the examples) and (2) the temperature equals or exceeds the hottest temperature in the data before that day. We are aiming for a function that can input a data frame that includes a column named `temp` with daily mean temperature in Celsius, like this data frame:

date	temp
2015-07-01	26.5
2015-07-02	27.2
2015-07-03	28.0
2015-07-04	26.9
2015-07-05	27.5
2015-07-06	25.9
2015-07-07	28.0
2015-07-08	28.2

and outputs a data frame that has an additional binary `record_temp` column, specifying if that day meet the two conditions, like this:

date	temp	record_temp
2015-07-01	26.5	FALSE
2015-07-02	27.2	TRUE
2015-07-03	28.0	TRUE
2015-07-04	26.9	FALSE
2015-07-05	27.5	FALSE
2015-07-06	25.9	FALSE
2015-07-07	28.0	TRUE
2015-07-08	28.2	TRUE

Below are two example functions that can perform these actions. Since the `record_temp` column depends on temperatures up to that day, one option is to use a loop to create this value. The first function takes this approach. The second function instead uses tidyverse functions to perform the same tasks.

```
# Function that uses a loop
find_records_1 <- function(datafr, threshold){
  highest_temp <- c()
  record_temp <- c()
  for(i in 1:nrow(datafr)){
    highest_temp <- max(highest_temp, datafr$temp[i])
    record_temp[i] <- datafr$temp[i] >= threshold &
      datafr$temp[i] >= highest_temp
  }
  datafr <- cbind(datafr, record_temp)
  return(datafr)
}

# Function that uses tidyverse functions
find_records_2 <- function(datafr, threshold){
  datafr <- datafr %>%
    mutate_(over_threshold = ~ temp >= threshold,
           cummax_temp = ~ temp == cummax(temp),
           record_temp = ~ over_threshold & cummax_temp) %>%
    select_.dots = c("-over_threshold", "-cummax_temp"))
  return(as.data.frame(datafr))
}
```

If you apply the two functions to the small example data set, you can see that they both create the desired output:

```
example_data <- data_frame(date = c("2015-07-01", "2015-07-02",
                                    "2015-07-03", "2015-07-04",
                                    "2015-07-05", "2015-07-06",
                                    "2015-07-07", "2015-07-08"),
                           temp = c(26.5, 27.2, 28.0, 26.9,
                                   27.5, 25.9, 28.0, 28.2))

(test_1 <- find_records_1(example_data, 27))
  date temp record_temp
1 2015-07-01 26.5 FALSE
2 2015-07-02 27.2 TRUE
3 2015-07-03 28.0 TRUE
4 2015-07-04 26.9 FALSE
5 2015-07-05 27.5 FALSE
6 2015-07-06 25.9 FALSE
7 2015-07-07 28.0 TRUE
8 2015-07-08 28.2 TRUE

(test_2 <- find_records_2(example_data, 27))
  date temp record_temp
1 2015-07-01 26.5 FALSE
2 2015-07-02 27.2 TRUE
3 2015-07-03 28.0 TRUE
4 2015-07-04 26.9 FALSE
5 2015-07-05 27.5 FALSE
6 2015-07-06 25.9 FALSE
7 2015-07-07 28.0 TRUE
8 2015-07-08 28.2 TRUE

all.equal(test_1, test_2)
[1] TRUE
```

The performance of these two functions can be compared using `microbenchmark`:

```
record_temp_perf <- microbenchmark(find_records_1(example_data, 27),
                                    find_records_2(example_data, 27))
record_temp_perf
Unit: microseconds
      expr     min      1q     mean    median      max      uq     neval
find_records_1(example_data, 27) 505.33 811.4225 908.3908 891.236
find_records_2(example_data, 27) 838.85 1334.7440 1623.8295 1459.676
                                     uq      max
956.0905 2443.247   100
1637.8945 11156.739   100
```

This output gives summary statistics (`min`, `1q`, `mean`, `median`, `uq`, and `max`) describing the time it took to run the two function over the 100 iterations of each function call. By default, these

times are given in a reasonable unit, based on the observed profiling times (units are given in microseconds in this case).

It's useful to check next to see if the relative performance of the two functions is similar for a bigger data set. The `chicagoNMMAPS` data set from the `dlnm` package includes temperature data over 15 years in Chicago, IL. Here are the results when we benchmark the two functions with that data (note, this code takes a minute or two to run):

```
library(dlnm)
data("chicagoNMMAPS")

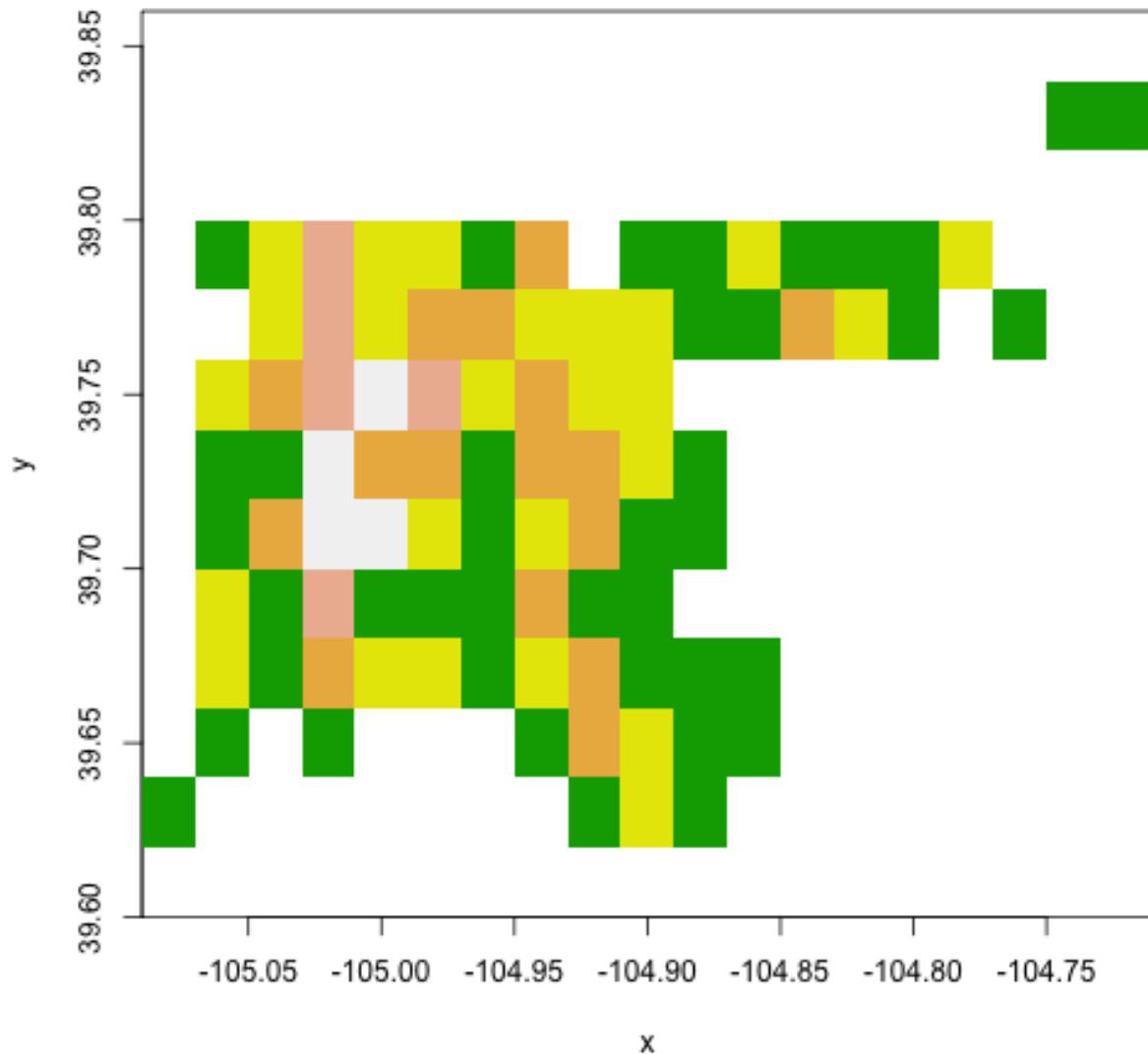
record_temp_perf_2 <- microbenchmark(find_records_1(chicagoNMMAPS, 27),
                                      find_records_2(chicagoNMMAPS, 27))
record_temp_perf_2
Unit: milliseconds

      expr       min        1q      mean
find_records_1(chicagoNMMAPS, 27) 127.363049 151.735575 182.203054
find_records_2(chicagoNMMAPS, 27)   1.438063   1.773797   2.708233
median      uq      max neval
176.103582 198.89603 298.57544   100
  2.235095   2.79743  11.84717   100
```

While the function with the loop (`find_records_1`) performed better with the very small sample data, the function that uses tidyverse functions (`find_records_2`) performs much, much better with a larger data set.

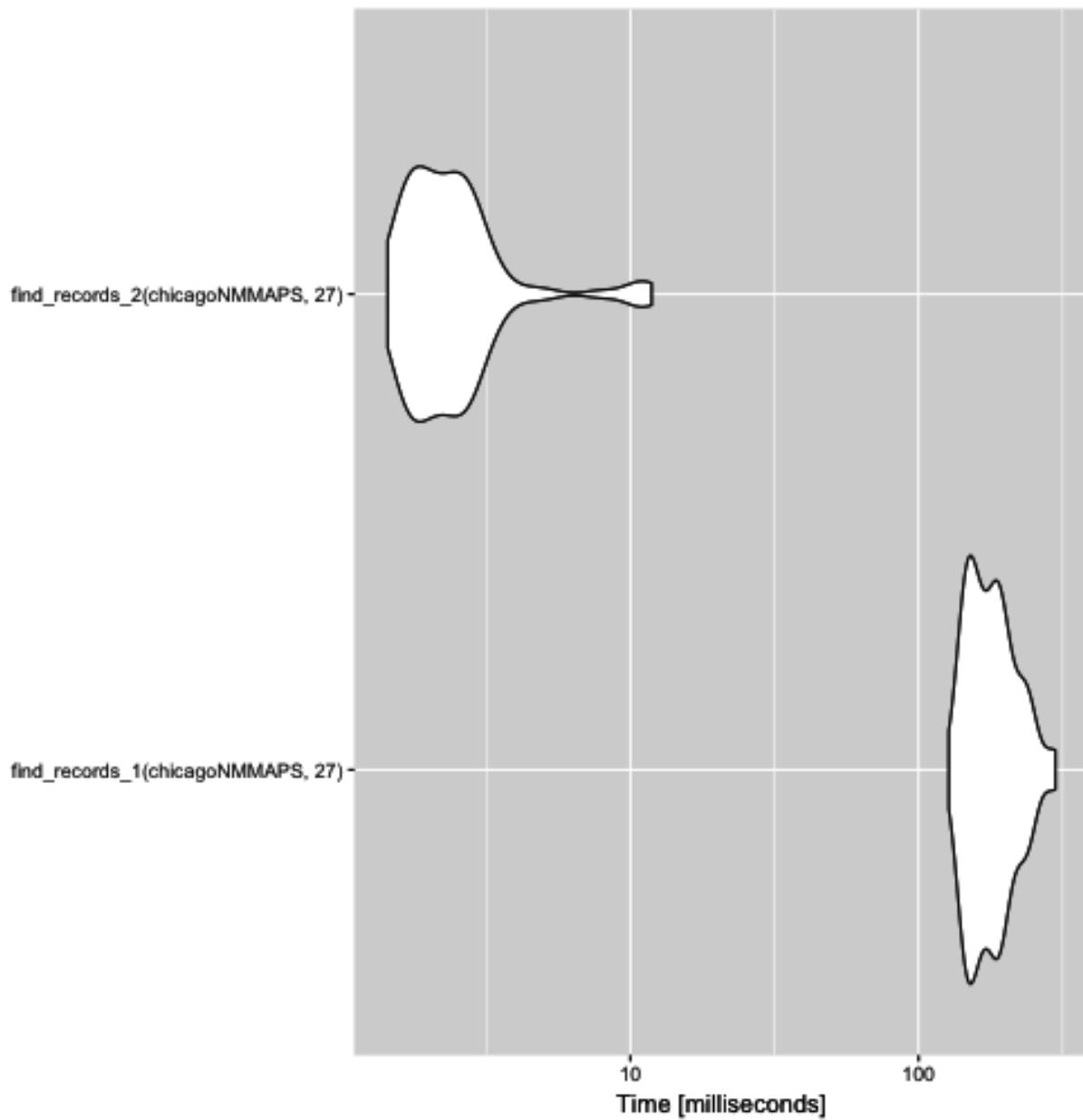
The `microbenchmark` function returns an object of the “`microbenchmark`” class. This class has two methods for plotting results, `autoplot.microbenchmark` and `boxplot.microbenchmark`. To use the `autoplot` method, you will need to have `ggplot2` loaded in your R session.

```
library(ggplot2)
# For small example data
autoplot(record_temp_perf)
```



Timing comparison of find records functions

```
# For larger data set  
autoplots(record_temp_perf_2)
```



Timing comparison of find records functions

By default, this plot gives the “Time” axis on a log scale. You can change this with the argument `log = FALSE`.

profvis

Once you’ve identified slower code, you’ll likely want to figure out which parts of the code are causing bottlenecks. The `profvis` function from the `profvis` package is very useful for this type of profiling. This function uses the `RProf` function from base R to profile code, and then

displays it in an interactive visualization in RStudio. This profiling is done by sampling, with the `RProf` function writing out the call stack every 10 milliseconds while running the code.

To profile code with `profvis`, just input the code (in braces if it is multi-line) into `profvis` within RStudio. For example, we found that the `find_records_1` function was slow when used with a large data set. To profile the code in that function, run:

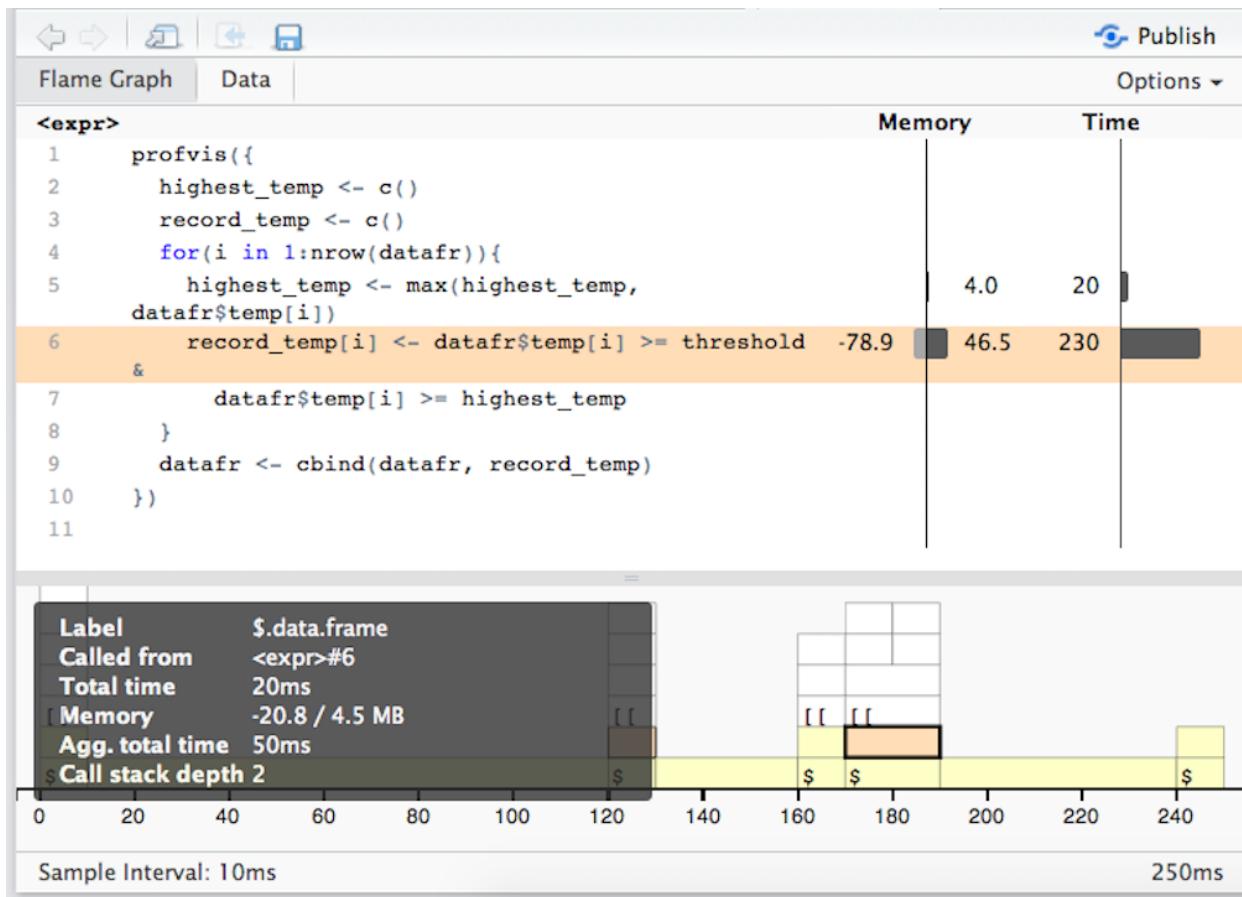
```
library(profvis)
datafr <- chicagoNMMAPS
threshold <- 27

profvis({
  highest_temp <- c()
  record_temp <- c()
  for(i in 1:nrow(datafr)){
    highest_temp <- max(highest_temp, datafr$temp[i])
    record_temp[i] <- datafr$temp[i] >= threshold &
      datafr$temp[i] >= highest_temp
  }
  datafr <- cbind(datafr, record_temp)
})
```

The `profvis` output gives you two options for visualization: “Flame Graph” or “Data” (a button to toggle between the two is given in the top left of the `profvis` visualization created when you profile code). The “Data” output defaults to show you the time usage of each first-level function call. Each of these calls can be expanded to show deeper and deeper functions calls within the call stack. This expandable interface allows you to dig down within a call stack to determine what calls are causing big bottlenecks. For functions that are part of a package you have loaded with `devtools::load_all`, this output includes a column with the file name where a given function is defined. This functionality makes this “Data” output pane particularly useful in profiling functions in a package you are creating.

The “Flame Graph” view in `profvis` output gives you two panels. The top panel shows the code called, with bars on the right to show memory use and time spent on the line. The bottom panel also visualizes the time used by each line of code, but in this case it shows time use horizontally and shows the full call stack at each time sample, with initial calls shown at the bottom of the graph, and calls deeper in the call stack higher in the graph. Clicking on a block in the bottom panel will show more information about a call, including which file it was called from, how much time it took, how much memory it took, and its depth in the call stack.

Figure @ref(fig:profvisexample) shows example output from profiling the code in the `find_records_1` function defined earlier in this section.



Example of profvis output for a function to find record temperatures.

Based on this visualization, most of the time is spent on line 6, filling in the `record_temp` vector. Now that we know this, we could try to improve the function, for example by doing a better job of initializing vectors before running the loop.

The `profvis` visualization can be used to profile code in functions you're writing as part of a package. If some of the functions in the code you are profiling are in a package currently loaded with `loaded` with `devtools::load_all`, the top panel in the Flame Graph output will include the code defining those functions, which allows you to explore speed and memory use within the code for each function. You can also profile code within functions from other packages— for more details on the proper set-up, see the “FAQ” section of [RStudio’s profvis documentation](#).

The `profvis` function will not be able to profile code that runs too quickly. Trying to profile functions that are too fast will give you the following error message:

```
Error in parse_rprof(prof_output, expr_source) :
  No parsing data available. Maybe your function was too fast?
```

You can use the argument `interval` in `profvis` to customize the sampling interval. The default is to sample every 10 milliseconds (`interval = 0.01`), but you can decrease this sampling

interval. In some cases, you may be able to use this option to profile faster-running code. However, you should avoid using an interval smaller than about 5 milliseconds, as below that you will get inaccurate estimates with `profvis`. If you are running very fast code, you're better off profiling with `microbenchmark`, which can give accurate estimates at finer time intervals.

Here are some tips for optimizing your use of `profvis`:

- You may find it convenient to use the “Show in new window” button on the RStudio pane with profiling results to expand this window while you are interpreting results.
- An “Options” button near the top right gives different options for how to display the profiling results, including whether to include memory profiling results and whether to include lines of code with zero time.
- You can click-and-drag results in the bottom visualization panel, as well as pan in and out.
- You may need to update your version of RStudio to be able to use the full functionality of `profvis`. You can download a Preview version of RStudio [here](#).
- If you'd like to share code profiling results from `profvis` publicly, you can do that by using the “Publish” button on the top right of the rendered profile visualization to publish the visualization to RPubs. The “FAQ” section of RStudio's `profvis` documentation includes more tips for sharing a code profile visualization online.
- If you get a lot of blocks labeled “<Anonymous>”, try updating your version of R. In newer versions of R, functions called using `package::function()` syntax or `list$function()` syntax are labeled in profiling blocks in a more meaningful way. This is likely to be a particular concern if you are profiling code in a package you are developing, as you will often be using `package::function()` syntax extensively to pass CRAN checks.

Find out more

If you'd like to learn more about profiling R code, or improving performance of R code once you've profiled, you might find these resources helpful:

- RStudio's `profvis` documentation
- Section on performant code in Hadley Wickham's *Advanced R* book
- “FasteR! HigheR! StrongeR! - A Guide to Speeding Up R Code for Busy People”, an article by Noam Ross

Summary

- Profiling can help you identify bottlenecks in R code.
- The `microbenchmark` package helps you profile short pieces of code and compare functions with each other. It runs the code many times and provides summary statistics across the iterations.
- The `profvis` package allows you to visualize performance across more extensive code. It can be used to profile code within functions being developed for a package, as long as the package source code has been loaded locally using `devtools::load_all`.

2.8 Non-standard Evaluation

Functions from packages like `dplyr`, `tidyr`, and `ggplot2` are excellent for creating efficient and easy-to-read code that cleans and displays data. However, they allow shortcuts in calling columns in data frames that allow some room for ambiguity when you move from evaluating code interactively to writing functions for others to use. The non-standard evaluation used within these functions mean that, if you use them as you would in an interactive session, you'll get a lot of "no visible bindings" warnings when you run CRAN checks on your package. These warnings will look something like this:

```
map_counties: no visible binding for global variable 'fips'
map_counties: no visible binding for global variable 'storm_dist'
map_counties: no visible binding for global variable 'tot_precip'
Undefined global functions or variables:
  fips storm_dist tot_precip
```

When you write a function for others to use, you need to avoid non-standard evaluation and so avoid all of these functions (culprits include many `dplyr` and `tidyr` functions— including `mutate`, `select`, `filter`, `group_by`, `summarize`, `gather`, `spread`— but also some functions in `ggplot2`, including `aes`). Fortunately, these functions all have standard evaluation alternatives, which typically have the same function name followed by an underscore (for example, the standard evaluation version of `mutate` is `mutate_`).

The input to the function call will need to be a bit different for standard evaluation versions of these functions. In many cases, this change is as easy as using formula notation (`~`) within the call, but in some cases it requires something more complex, including using the `.dots` argument.

Here is a table with examples of non-standard evaluation calls and their standard evaluation alternatives (these are all written assuming that the function is being used as a step in a piping flow, where the input data frame has already been defined earlier in the piping sequence):

Non-standard evaluation version	Standard evaluation version
<code>filter(fips %in% counties)</code>	<code>filter_(~ fips %in% counties)</code>
<code>mutate(max_rain = max(tot_precip))</code>	<code>mutate_(max_rain = ~ max(tot_precip))</code>
<code>summarize(tot_precip = sum(precip))</code>	<code>summarize_(tot_precip = ~ sum(precip))</code>
<code>group_by(storm_id, fips)</code>	<code>group_by_(~ storm_id, ~ fips)</code>
<code>aes(x = long, y = lat)</code>	<code>aes_(x = ~ long, y = ~ lat)</code>
<code>select(-start_date, -end_date)</code>	<code>select_(.dots = c('start_date', 'end_date'))</code>
<code>select(-start_date, -end_date)</code>	<code>select_(.dots = c('-start_date', '-end_date'))</code>
<code>spread(key, mean)</code>	<code>spread_(key_col = 'key', value_col = 'mean')</code>
<code>gather(key, mean)</code>	<code>gather_(key_col = 'key', value_col = 'mean')</code>

If you have any non-standard evaluation in your package code (which you'll notice because of the "no visible bindings" warnings you'll get when you check the package), go through and change any instances to use standard evaluation alternatives. This change prevents these warnings when you check your package and will also ensure that the functions behave like you expect them to when they are run by other users in their own R sessions.

In this section, we've explained only how to convert from functions that use non-standard evaluation to those that use standard evaluation, to help in passing CRAN checks as you go from coding scripts to writing functions for packages. If you would like to learn more about non-standard evaluation in R, you should check out the [chapter on non-standard evaluation](#) in Hadley Wickham's *Advanced R* book.

Summary

- Functions that use non-standard evaluation can cause problems within functions written for a package.
- The NSE functions in tidyverse packages all have standard evaluation analogues that should be used when writing functions that will be used by others.

2.9 Object Oriented Programming

The learning objectives of this section are:

- Design and Implement a new S3, S4, or reference class with generics and methods

Introduction

Object oriented programming is one of the most successful and widespread philosophies of programming and is a cornerstone of many programming languages including Java, Ruby, Python, and C++. R has three object oriented systems because the roots of R date back to 1976, when the idea of object orientated programming was barely [four years old](#). New object oriented paradigms were added to R as they were invented, so some of the ideas in R about object oriented programming have gone stale in the years since. It's still important to understand these older systems since a huge amount of R code is written with them, and they're still useful and interesting! Long time object oriented programmers reading this book may find these old ideas refreshing.

The two older object oriented systems in R are called S3 and S4, and the modern system is called RC which stands for "reference classes." Programmers who are already familiar with object oriented programming will feel at home using RC.

Object Oriented Principles

There are several key principles in object oriented programming which span across R's object systems and other programming languages. The first are the ideas of a **class** and an **object**.

The world is made up of physical objects - the chair you're sitting in, the clock next to your bed, the bus you ride every day, etc. Just like the world is full of physical objects, your programs can be made of objects as well. A class is a blueprint for an object: it describes the parts of an object, how to make an object, and what the object is able to do. If you were to think about a class for a bus (as in the public buses that roam the roads) this class would describe attributes for the bus like the number of seats on the bus, the number of windows, the top speed of the bus, and the maximum distance the bus can drive on one tank of gas.

Buses in general can perform the same actions, and these actions are also described in the class: a bus can open and close its doors, the bus can steer, and the accelerator or the brake can be used to slow down or speed up the bus. Each of these actions can be described as a **method** which is a **function** that is associated with a particular class. We'll be using this class in order to create individual bus objects, so we should provide a **constructor** which is a method where we can specify attributes of the bus as arguments. This constructor method will then return an individual bus object with the attributes that we specified.

You could also imagine that after making the bus class you might want to make a special kind of class for a **party bus**. Party buses have all of the same attributes and methods as our bus class, but they also have additional attributes and methods like the number of refrigerators, window blinds that can be opened and closed, and smoke machines that can be turned on and off. Instead of rewriting the entire bus class and then adding new attributes and methods, it is possible for the party bus class to **inherit** all of the attributes and methods from the bus class. In this framework of inheritance, we talk about the bus class as the super-class of the party bus, and the party bus is the sub-class of the bus. What this relationship means is that the party bus has all of the same attributes and methods as the bus class plus additional attributes and methods.

S3

Conveniently everything in R is an object. By “everything” I mean every single “thing” in R including numbers, functions, strings, data frames, lists, etc. If you want to know the class of an object in R you can simply use the `class()` function:

```
class(2)
[1] "numeric"
class("is in session.")
[1] "character"
class(class)
[1] "function"
```

Now it's time to wade into some of the quirks of R's object oriented systems. In the S3 system you can arbitrarily assign a class to any object, which goes against most of what we discussed in the *Object Oriented Principles* section. Class assignments can be made using the `structure()` function, or you can assign the class using `class()` and `<-:`:

```
special_num_1 <- structure(1, class = "special_number")
class(special_num_1)
[1] "special_number"

special_num_2 <- 2
class(special_num_2)
[1] "numeric"
class(special_num_2) <- "special_number"
class(special_num_2)
[1] "special_number"
```

This is completely legal R code, but if you want to have a better behaved S3 class you should create a constructor which returns an S3 object. The `shape_S3()` function below is a constructor that returns a `shape_S3` object:

```
shape_s3 <- function(side_lengths){
  structure(list(side_lengths = side_lengths), class = "shape_S3")
}

square_4 <- shape_s3(c(4, 4, 4, 4))
class(square_4)
[1] "shape_S3"

triangle_3 <- shape_s3(c(3, 3, 3))
class(triangle_3)
[1] "shape_S3"
```

We've now made two `shape_S3` objects: `square_4` and `triangle_3`, which are both instantiations of the `shape_S3` class. Imagine that you wanted to create a method that would return `TRUE` if a `shape_S3` object was a square, `FALSE` if a `shape_S3` object was not a square, and `NA` if the object provided as an argument to the method was not a `shape_s3` object. This can be achieved using R's **generic methods** system. A generic method can return different values based depending on the class of its input. For example `mean()` is a generic method that can find the average of a vector of numbers or it can find the “average day” from a vector of dates. The following snippet demonstrates this behavior:

```
mean(c(2, 3, 7))
[1] 4
mean(c(as.Date("2016-09-01"), as.Date("2016-09-03")))
[1] "2016-09-02"
```

Now let's create a generic method for identifying `shape_S3` objects that are squares. The creation of every generic method uses the `UseMethod()` function in the following way with only slight variations:

```
[name of method] <- function(x) UseMethod("[name of method]")
```

Let's call this method `is_square`:

```
is_square <- function(x) UseMethod("is_square")
```

Now we can add the actual function definition for detecting whether or not a shape is a square by specifying `is_square.shape_S3`. By putting a dot (.) and then the name of the class after `is_square`, we can create a method that associates `is_square` with the `shape_S3` class:

```
is_square.shape_S3 <- function(x){
  length(x$side_lengths) == 4 &&
  x$side_lengths[1] == x$side_lengths[2] &&
  x$side_lengths[2] == x$side_lengths[3] &&
  x$side_lengths[3] == x$side_lengths[4]
}

is_square(square_4)
[1] TRUE
is_square(triangle_3)
[1] FALSE
```

Seems to be working well! We also want `is_square()` to return `NA` when its argument is not a `shape_S3`. We can specify `is_square.default` as a last resort if there is no method associated with the object passed to `is_square()`.

```
is_square.default <- function(x){
  NA
}

is_square("square")
[1] NA
is_square(c(1, 1, 1, 1))
[1] NA
```

Let's try printing `square_4`:

```
print(square_4)
$side_lengths
[1] 4 4 4 4

attr(, "class")
[1] "shape_S3"
```

Doesn't that look ugly? Lucky for us `print()` is a generic method, so we can specify a print method for the `shape_S3` class:

```

print.shape_S3 <- function(x){
  if(length(x$side_lengths) == 3){
    paste("A triangle with side lengths of", x$side_lengths[1],
           x$side_lengths[2], "and", x$side_lengths[3])
  } else if(length(x$side_lengths) == 4) {
    if(is_square(x)){
      paste("A square with four sides of length", x$side_lengths[1])
    } else {
      paste("A quadrilateral with side lengths of", x$side_lengths[1],
             x$side_lengths[2], x$side_lengths[3], "and", x$side_lengths[4])
    }
  } else {
    paste("A shape with", length(x$side_lengths), "slides.")
  }
}

print(square_4)
[1] "A square with four sides of length 4"
print(triangle_3)
[1] "A triangle with side lengths of 3 3 and 3"
print(shape_s3(c(10, 10, 20, 20, 15)))
[1] "A shape with 5 slides."
print(shape_s3(c(2, 3, 4, 5)))
[1] "A quadrilateral with side lengths of 2 3 4 and 5"

```

Since printing an object to the console is one of the most common things to do in R, nearly every class has an associated print method! To see all of the methods associated with a generic like `print()` use the `methods()` function:

```

head(methods(print), 10)
[1] "print,ANY-method"          "print,diagonalMatrix-method"
[3] "print,sparseMatrix-method" "print.acf"
[5] "print.anova"              "print.anova.gam"
[7] "print.anova.lme"          "print.aov"
[9] "print.aovlist"            "print.ar"

```

One last note on S3 with regard to inheritance. In the previous section we discussed how a sub-class can inherit attributes and methods from a super-class. Since you can assign any class to an object in S3, you can specify a super class for an object the same way you would specify a class for an object:

```

class(square_4)
[1] "shape_S3"
class(square_4) <- c("shape_S3", "square")
class(square_4)
[1] "shape_S3" "square"

```

To check if an object is a sub-class of a specified class you can use the `inherits()` function:

```
inherits(square_4, "square")
[1] TRUE
```

Example: S3 Class/Methods for Polygons

The S3 system doesn't have a formal way to define a class but typically, we use a list to define the class and elements of the list serve as data elements.

Here is our definition of a polygon represented using Cartesian coordinates. The class contains an element called `xcoord` and `ycoord` for the x- and y-coordinates, respectively. The `make_poly()` function is the “constructor” function for polygon objects. It takes as arguments a numeric vector of x-coordinates and a corresponding numeric vector of y-coordinates.

```
## Constructor function for polygon objects
## x a numeric vector of x coordinates
## y a numeric vector of y coordinates
make_poly <- function(x, y) {
  if(length(x) != length(y))
    stop("'x' and 'y' should be the same length")

  ## Create the "polygon" object
  object <- list(xcoord = x, ycoord = y)

  ## Set the class name
  class(object) <- "polygon"
  object
}
```

Now that we have a class definition, we can develop some methods for operating on objects from that class.

The first method we'll define is the `print()` method. The `print()` method should just show some simple information about the object and should not be too verbose—just enough information that the user knows what the object is.

Here the `print()` method just shows the user how many vertices the polygon has. It is a convention for `print()` methods to return the object `x` invisibly.

```
## Print method for polygon objects
## x an object of class "polygon"
print.polygon <- function(x, ...) {
  cat("a polygon with", length(x$xcoord),
      "vertices\n")
  invisible(x)
}
```

Next is the `summary()` method. The `summary()` method typically shows a bit more information and may even do some calculations. This `summary()` method computes the ranges of the x- and y-coordinates.

The typical approach for `summary()` methods is to allow the summary method to compute something, but to *not* print something. The strategy is

1. The `summary()` method returns an object of class “summary_‘class name’”
2. There is a separate `print()` method for “summary_‘class name’” objects.

For example, here is the `summary()` method.

```
## Summary method for polygon objects
## object an object of class "polygon"

summary.polygon <- function(object, ...) {
  object <- list(rng.x = range(object$xcoord),
                 rng.y = range(object$ycoord))
  class(object) <- "summary_polygon"
  object
}
```

Note that it simply returns an object of class `summary_polygon`. Now the corresponding `print()` method.

```
## Print method for summary.polygon objects
## x an object of class "summary_polygon"
print.summary_polygon <- function(x, ...) {
  cat("x:", x$rng.x[1], "-->", x$rng.x[2], "\n")
  cat("y:", x$rng.y[1], "-->", x$rng.y[2], "\n")
  invisible(x)
}
```

Now we can make use of our new class and methods.

```
## Construct a new "polygon" object
x <- make_poly(1:4, c(1, 5, 2, 1))
```

We can use the `print()` to see what the object is.

```
print(x)
a polygon with 4 vertices
```

And we can use the `summary()` method to get a bit more information about the object.

```
out <- summary(x)
class(out)
[1] "summary_polygon"
print(out)
x: 1 --> 4
y: 1 --> 5
```

Because of auto-printing we can just call the `summary()` method and let the results auto-print.

```
summary(x)
$rng.x
[1] 1 4

$rng.y
[1] 1 5

attr(,"class")
[1] "summary_polygon"
```

From here, we could build other methods for interacting with our `polygon` object. For example, it may make sense to define a `plot()` method or maybe methods for intersecting two polygons together.

S4

The S4 system is slightly more restrictive than S3, but it's similar in many ways. To create a new class in S4 you need to use the `setClass()` function. You need to specify two or three arguments for this function: `Class` which is the name of the class as a string, `slots`, which is a named list of attributes for the class with the class of those attributes specified, and optionally `contains` which includes the super-class of they class you're specifying (if there is a super-class). Take look at the class definition for a `bus_S4` and a `party_bus_S4` below:

```
setClass("bus_S4",
  slots = list(n_seats = "numeric",
              top_speed = "numeric",
              current_speed = "numeric",
              brand = "character"))
setClass("party_bus_S4",
  slots = list(n_subwoofers = "numeric",
               smoke_machine_on = "logical"),
  contains = "bus_S4")
```

Now that we've created the `bus_S4` and the `party_bus_S4` classes we can create bus objects using the `new()` function. The `new()` function's arguments are the name of the class and values for each “slot” in our S4 object.

```
my_bus <- new("bus_S4", n_seats = 20, top_speed = 80,
               current_speed = 0, brand = "Volvo")
my_bus
An object of class "bus_S4"
Slot "n_seats":
[1] 20

Slot "top_speed":
[1] 80

Slot "current_speed":
[1] 0

Slot "brand":
[1] "Volvo"
my_party_bus <- new("party_bus_S4", n_seats = 10, top_speed = 100,
                      current_speed = 0, brand = "Mercedes-Benz",
                      n_subwoofers = 2, smoke_machine_on = FALSE)
my_party_bus
An object of class "party_bus_S4"
Slot "n_subwoofers":
[1] 2

Slot "smoke_machine_on":
[1] FALSE

Slot "n_seats":
[1] 10

Slot "top_speed":
[1] 100

Slot "current_speed":
[1] 0

Slot "brand":
[1] "Mercedes-Benz"
```

You can use the `@` operator to access the slots of an S4 object:

```
my_bus@n_seats
[1] 20
my_party_bus@top_speed
[1] 100
```

This is essentially the same as using the `$` operator with a list or an environment.

S4 classes use a generic method system that is similar to S3 classes. In order to implement a new generic method you need to use the `setGeneric()` function and the `standardGeneric()` function in the following way:

```
setGeneric("new_generic", function(x){
  standardGeneric("new_generic")
})
```

Let's create a generic function called `is_bus_moving()` to see if a `bus_S4` object is in motion:

```
setGeneric("is_bus_moving", function(x){
  standardGeneric("is_bus_moving")
})
[1] "is_bus_moving"
```

Now we need to actually define the function which we can do with `setMethod()`. The `setMethod()` function takes as arguments the name of the method as a string, the method signature which specifies the class of each argument for the method, and then the function definition of the method:

```
setMethod("is_bus_moving",
  c(x = "bus_S4"),
  function(x){
    x@current_speed > 0
  })
[1] "is_bus_moving"

is_bus_moving(my_bus)
[1] FALSE
my_bus@current_speed <- 1
is_bus_moving(my_bus)
[1] TRUE
```

In addition to creating your own generic methods, you can also create a method for your new class from an existing generic. First use the `setGeneric()` function with the name of the existing method you want to use with your class, and then use the `setMethod()` function like in the previous example. Let's make a `print()` method for the `bus_S4` class:

```
setGeneric("print")
[1] "print"

setMethod("print",
  c(x = "bus_S4"),
  function(x){
    paste("This", x@brand, "bus is traveling at a speed of", x@current_speed)
  })
[1] "print"

print(my_bus)
[1] "This Volvo bus is traveling at a speed of 1"
print(my_party_bus)
[1] "This Mercedes-Benz bus is traveling at a speed of 0"
```

Reference Classes

With reference classes we leave the world of R's old object oriented systems and enter the philosophies of other prominent object oriented programming languages. We can use the `setRefClass()` function to define a class' fields, methods, and super-classes. Let's make a reference class that represents a student:

```
Student <- setRefClass("Student",
  fields = list(name = "character",
    grad_year = "numeric",
    credits = "numeric",
    id = "character",
    courses = "list"),
  methods = list(
    hello = function(){
      paste("Hi! My name is", name)
    },
    add_credits = function(n){
      credits <- credits + n
    },
    get_email = function(){
      paste0(id, "@jhu.edu")
    }
  )
)
```

To recap: we've created a class definition called `Student` which defines the student class. This class has five fields and three methods. To create a `Student` object use the `new()` method:

```
brooke <- Student$new(name = "Brooke", grad_year = 2019, credits = 40,
  id = "ba123", courses = list("Ecology", "Calculus III"))
roger <- Student$new(name = "Roger", grad_year = 2020, credits = 10,
  id = "rp456", courses = list("Puppetry", "Elementary Algebra"))
```

You can access the fields and methods of each object using the `$` operator:

```
brooke$credits
[1] 40
roger$hello()
[1] "Hi! My name is Roger"
roger$get_email()
[1] "rp456@jhu.edu"
```

Methods can change the state of an object, for instance in the case of the `add_credits()` function:

```
brooke$credits
[1] 40
brooke$add_credits(4)
brooke$credits
[1] 44
```

Notice that the `add_credits()` method uses the complex assignment operator (`<-<-`). You need to use this operator if you want to modify one of the fields of an object with a method. You'll learn more about this operator in the Expressions & Environments section.

Reference classes can inherit from other classes by specifying the `contains` argument when they're defined. Let's create a sub-class of `Student` called `Grad_Student` which includes a few extra features:

```
Grad_Student <- setRefClass("Grad_Student",
                           contains = "Student",
                           fields = list(thesis_topic = "character"),
                           methods = list(
                             defend = function(){
                               paste0(thesis_topic, ". QED.")
                             }
                           ))
jeff <- Grad_Student$new(name = "Jeff", grad_year = 2021, credits = 8,
                         id = "j155", courses = list("Fitbit Repair",
                                         "Advanced Base Graphics"),
                         thesis_topic = "Batch Effects")
jeff$defend()
[1] "Batch Effects. QED."
```

Summary

- R has three object oriented systems: S3, S4, and Reference Classes.
- Reference Classes are the most similar to classes and objects in other programming languages.
- Classes are blueprints for an object.
- Objects are individual instances of a class.
- Methods are functions that are associated with a particular class.
- Constructors are methods that create objects.
- Everything in R is an object.
- S3 is a liberal object oriented system that allows you to assign a class to any object.
- S4 is a more strict object oriented system that builds upon ideas in S3.
- Reference Classes are a modern object oriented system that is similar to Java, C++, Python, or Ruby.

2.10 Gaining Your ‘tidyverse’ Citizenship

The learning objectives of this section are:

- Describe the principles of tidyverse functions

Many of the tools that we discuss in this book revolve around the so-called “tidyverse” set of tools. These tools, largely developed by Hadley Wickham but also including a diverse community of developers, have a set of principles that are adhered to when they are being developed. Hadley Wickham laid out these principles in his [Tidy Tools Manifesto](#), a vignette within the `tidyverse` package.

The four basic principles of the tidyverse are:

Reuse existing data structures

R has a number of data structures (data frames, vectors, etc.) that people have grown accustomed to over the many years of R’s existence. While it is often tempting to develop custom data structures, for example, by using S3 or S4 classes, it is often worthwhile to consider reusing a commonly used structure. You’ll notice that many tidyverse functions make heavy use of the data frame (typically as their first argument), because the data frame is a well-known, well-understood structure used by many analysts. Data frames have a well-known and reasonably standardized corresponding file format in the CSV file.

While common data structures like the data frame may not be perfectly suited to your needs as you develop your own software, it is worth considering using them anyway because the enormous value to the community that is already familiar with them. If the user community feels familiar with the data structures required by your code, they are likely to adopt them quicker.

Compose simple functions with the pipe

One of the [original principles of the Unix operating system](#) was that every program should do “one thing well”. The limitation of only doing one thing (but well!) was removed by being able to easily pipe the output of one function to be the input of another function (the pipe operator on Unix was the `|` symbol). Typical Unix commands would contain long strings of commands piped together to (eventually) produce some useful output. On Unix systems, the unifying concept that allowed programs to pipe to each other was the use of [textual formats]. All data was rendered in textual formats so that if you wrote a new program, you would not need to worry about decoding some obscure proprietary format.

Much like the original Unix systems, the tidyverse eschews building monolithic functions that have many bells and whistles. Rather, once you are finished writing a simple function, it is better to start afresh and work off the input of another function to produce new output (using the `%>%` operator, for example). The key to this type of development is having *clean interfaces* between functions and an expectation that the output of every function may

serve as the input to another function. This is why the first principle (reuse existing data structures) is important, because the reuse of data structures that are well-understood and characterized lessens the burden on other developers who are developing new code and would prefer not to worry about new-fangled data structures at every turn.

Embrace functional programming

This can be a tough principle for people coming from other non-functional programming languages. But the reality is, R is a functional programming language (with its roots in Scheme) and it's best not to go against the grain. In our section on Functional Programming, we outlined many of the principles that are fundamental to functional-style programming. In particular, the `purrr` package implements many of those ideas.

One benefit to functional programming is that it can at times be easier to reason about when simply looking at the code. The inability to modify arguments enables us to predict what the output of a function will be given a certain input, allowing for things like memoization. Functional programming also allows for simple parallelization, so that we can quickly parallelize any code that uses `lapply()` or `map()`.

Design for humans

Making your code *readable* and *usable* by people is goal that is overlooked surprisingly often. The result is things like function names that are obscure and do not actually communicate what they do. When writing code, using things like good, explicit, function names, with descriptive arguments, can allow for users to quickly learn your API. If you have a set of functions with a similar purpose, they might share a prefix (see e.g. `geom_point()`, `geom_line()`, etc.). If you have an argument like `color` that could either take arguments 1, 2, and 3, or `black`, `red`, and `green`, think about which set of arguments might be easier for humans to handle.

3. Building R Packages

This section covers building R packages. Writing good code for data science is only part of the job. In order to maximize the usefulness and reusability of data science software, code must be organized and distributed in a manner that adheres to community-based standards and provides a good user experience. This section covers the primary means by which R software is organized and distributed to others. We cover R package development, writing good documentation and vignettes, writing robust software, cross-platform development, continuous integration tools, and distributing packages via CRAN and GitHub. Learners will produce R packages that satisfy the criteria for submission to CRAN.

The Learning objectives for this section are:

- Recognize the basic structure and purpose of an R package
- Create a simple R package skeleton using the devtools package
- Recognize the key directives in a NAMESPACE file
- Create R function documentation using roxygen2
- Create vignettes using knitr and R Markdown
- Create an R package that contains data (and associated documentation)
- Create unit tests for an R package using the testthat package
- Categorize errors in the R CMD check process
- Recall the principles of open source software
- Recall two open source licenses
- Create a GitHub repository for an R package
- Create an R package that is tested and deployed on Travis
- Create an R package that is tested and deployed on Appveyor
- Recognize characteristics of R packages that are not cross-platform

3.1 Before You Start

Building R packages requires a toolchain that must be in place before you begin developing. If you are developing packages that contain only R code, then the tools you need come with R and RStudio. However, if you want to build packages with compiled C, C++, or Fortran code (or which to build other people's packages with such code), then you will need to install additional tools. Which tools you install depends on what platform you are running.

Mac OS

For developing in Mac OS, you will first need to download the Xcode development environment. While you do not need the IDE that comes with Xcode to develop R packages you need many of the tools that come with it, including the C compiler (`clang`). Xcode can be obtained from either the [Mac App Store](#) or from Apple's [Xcode developer's page](#). Once this is installed you will have the C compiler as well as a number of additional Unix shell tools. You will also have necessary header files for compiling C code.

While it's unlikely that you will be building your own packages with Fortran code, many older packages (including R itself) contain Fortran code. Therefore, in order to build these packages, you need a Fortran compiler. Mac OS does not come with one by default and so you can download the [GNU Fortran Compiler](#) from the R for Mac tools page.

There are more details provided on the [R for Mac tools package](#) maintained by Simon Urbanek, particularly for older versions of Mac OS.

Windows

On Windows, the R Core has put together a package of tools that you can download all at once and install via a simple installer tool. The [Rtools](#) package comes in different versions, depending on the version of R that you are using. Make sure to get the version of Rtools that matches your version of R. Once you have installed this, you will have most of the tools needed to build R packages. You can optionally install a few other tools, documented [here](#).

Unix/Linux

If you are using R on a Unix-like system then you may have already have the tools for building R packages. In particular, if you built R from the sources, then you already have a C compiler and Fortran compiler. If, however, you installed R from a package management system, then you may need to install the compilers, as well as the header files. These usually come in packages with the suffix `-devel`. For example, the header files for the `readline` package may come in the package `readline-devel`. The catch is that these `-devel` packages are not needed to run R, only to build R packages from the sources.

3.2 R Packages

The objectives of this section are:

- Recognize the basic structure and purpose of an R package
- Recognize the key directives in a NAMESPACE file

An R package is a mechanism for extending the basic functionality of R. It is the natural extension of writing functions that each do a specific thing well. In the previous chapter, we discussed how writing functions abstracts the behavior of a set of R expressions by

providing a defined interface, with inputs (i.e. function arguments) and outputs (i.e. return values). The use of functions simplifies things for the user because the user no longer needs to be knowledgeable of the details of the underlying code. They only need to understand the inputs and outputs.

Once one has developed many functions, it becomes natural to group them in to collections of functions that are aimed at achieving an overall goal. This collection of functions can be assembled into an R package. R packages represent another level of abstraction, where the interface presented to the user is a set of **user-facing functions**. These functions provide access to the underlying functionality of the package and simplify the user experience because the one does not need to be concerned with the many other helper functions that are required.

R packages are a *much* better way to distribute code to others because they provide a clean and uniform user experience for people who want to interact with your code. R packages require documentation in a standardized format, and the various tools that come with R (and RStudio) help to check your packages so that they do not contain inconsistencies or errors. R users are already familiar with how to use R packages, and so they will be able to quickly adopt your code if is presented in this format.

This chapter highlights the key elements of building R packages. The fine details of building a package can be found in the [Writing R Extensions](#) manual.

Basic Structure of an R Package

An R package begins life as a directory on your computer. This directory has a specific layout with specific files and sub-directories. The two required sub-directories are

- `R`, which contains all of your R code files
- `man`, which contains your documentation files.

At the top level of your package directory you will have a `DESCRIPTION` file and a `NAMESPACE` file. This represents the minimal requirements for an R package. Other files and sub-directories can be added and will discuss how and why in the sections below.



While RStudio is not required to build R packages, it contains a number of convenient features that make the development process easier and faster. That said, in order to use RStudio for package development, you must setup the environment properly. Details of how to do this can be found in Roger's [RStudio package development pre-flight check list](#).

DESCRIPTION File

The `DESCRIPTION` file is an essential part of an R package because it contains key metadata for the package that is used by repositories like CRAN and by R itself. In particular, this

file contains the package name, the version number, the author and maintainer contact information, the license information, as well as any dependencies on other packages.

As an example, here is the DESCRIPTION file for the `mvtsplot` package on CRAN. This package provides a function for plotting multivariate time series data.

```
Package: mvtsplot
Version: 1.0-3
Date: 2016-05-13
Depends: R (>= 3.0.0)
Imports: splines, graphics, grDevices, stats, RColorBrewer
Title: Multivariate Time Series Plot
Author: Roger D. Peng <rpeng@jhsph.edu>
Maintainer: Roger D. Peng <rpeng@jhsph.edu>
Description: A function for plotting multivariate time series data.
License: GPL (>= 2)
URL: https://github.com/rdpeng/mvtsplot
```

NAMESPACE File

The NAMESPACE file specifies the interface to the package that is presented to the user. This is done via a series of `export()` statements, which indicate which functions in the package are exported to the user. Functions that are not exported cannot be called directly by the user (although see below). In addition to exports, the NAMESPACE file also specifies what functions or packages are *imported* by the package. If your package depends on functions from another package, you must import them via the NAMESPACE file.

Below is the NAMESPACE file for the `mvtsplot` package described above.

```
export("mvtsplot")

import(splines)
import(RColorBrewer)
importFrom("grDevices", "colorRampPalette", "gray")
importFrom("graphics", "abline", "axis", "box", "image", "layout",
          "lines", "par", "plot", "points", "segments", "strwidth",
          "text", "Axis")
importFrom("stats", "complete.cases", "lm", "na.exclude", "predict",
          "quantile")
```

Here we can see that only a single function is exported from the package (the `mvtsplot()` function). There are two types of import statements:

- `import()`, simply takes a package name as an argument, and the interpretation is that all exported functions from that external package will be accessible to your package
- `importFrom()`, takes a package and a series of function names as arguments. This directive allows you to specify exactly which function you need from an external package. For example, this package imports the `colorRampPalette()` and `gray()` functions from the `grDevices` package.

Generally speaking, it is better to use `importFrom()` and to be specific about which function you need from an external package. However, in some cases when you truly need almost every function in a package, it may be more efficient to simply `import()` the entire package.

With respect to exporting functions, it is important to think through carefully which functions you want to export. First and foremost, exported functions must be documented and supported. Users will generally expect exported functions to be there in subsequent iterations of the package. It's usually best to limit the number of functions that you export (if possible). It's always possible to export something later if it is needed, but removing an exported function once people have gotten used to having it available can result in upset users. Finally, exporting a long list of functions has the effect of cluttering a user's namespace with function names that may conflict with functions from other packages. Minimizing the number of exports reduces the chances of a conflict with other packages (using more package-specific function names is another way).

Namespace Function Notation

As you start to use many packages in R, the likelihood of two functions having the same name increases. For example, the commonly used `dplyr` package has a function named `filter()`, which is also the name of a function in the `stats` package. If one has both packages loaded (a more than likely scenario) how can one specify exactly which `filter()` function they want to call?

In R, every function has a full name, which includes the package namespace as part of the name. This format is along the lines of

```
<package name>::<exported function name>
```

For example, the `filter()` function from the `dplyr` package can be referenced as `dplyr::filter()`. This way, there is no confusion over which `filter()` function we are calling. While in principle every function can be referenced in this way, it can be tiresome for interactive work. However, for programming, it is often safer to reference a function using the full name if there is even a chance that there might be confusion.

It is possible to call functions that are *not* exported by package by using the namespace notation. The `:::` operator can be used for this purpose, as in `<package name>:::<unexported function name>`. This can be useful for examining the code of an unexported function (e.g. for debugging purposes) or for temporarily accessing some unexported feature of a package. However, it's not a good idea to make this a habit as such unexported functions may change or even be eliminated in future versions of the package. Furthermore, use of the `:::` operator is not allowed for packages that reside on CRAN.

Loading and Attaching a Package Namespace

When dealing with R packages, it's useful to understand the distinction between *loading* a package namespace and *attaching* it. When package A imports the namespace of package B, package A loads the namespace of package B in order to gain access to the exported functions

of package B. However, when the namespace of package B is loaded, it is only available to package A; it is not placed on the search list and is not visible to the user or to other packages.

Attaching a package namespace places that namespace on the search list, making it visible to the user and to other packages. Sometimes this is needed because certain functions need to be made visible to the user and not just to a given package.

The R Sub-directory

The R sub-directory contains all of your R code, either in a single file, or in multiple files. For larger packages it's usually best to split code up into multiple files that logically group functions together. The names of the R code files do not matter, but generally it's not a good idea to have spaces in the file names.

The man Sub-directory

The man sub-directory contains the documentation files for all of the exported objects of a package. With older versions of R one had to write the documentation of R objects directly into the man directory using a LaTeX-style notation. However, with the development of the roxygen2 package, we no longer need to do that and can write the documentation directly into the R code files. Therefore, you will likely have little interaction with the man directory as all of the files in there will be auto-generated by the roxygen2 package.

Summary

R packages provide a convenient and standardized mechanism for distributing R code to a wide audience. As part of building an R package you design an interface to a collection of functions that users can access to make use of the functionality you provide. R packages are directories containing R code, documentation files, package metadata, and export/import information. Exported functions are functions that are accessible by the user; imported functions are functions in other packages that are used by your package.

3.3 The devtools Package

The objective of this section is

- Create a simple R package skeleton using the devtools package

R package development has become substantially easier in recent years with the introduction of a package by Hadley Wickham called devtools. As the package name suggests, this includes a variety of functions that facilitate software development in R.



Hands down, the best resource for mastering the devtools package is the book *R Packages* by Hadley Wickham. The full book is available online for free at <http://r-pkgs.had.co.nz>. It is also available as a hard copy book published by O'Reilly. If you plan to develop a lot of R packages, it is well worth your time to read this book closely.

Key `devtools` functions

Here are some of the key functions included in `devtools` and what they do, roughly in the order you are likely to use them as you develop an R package:

Function	Use
<code>create</code>	Create the file structure for a new package
<code>load_all</code>	Load the code for all functions in the package
<code>document</code>	Create <code>\man</code> documentation files and the “NAMESPACE” file from <code>roxygen2</code> code
<code>use_data</code>	Save an object in your R session as a dataset in the package
<code>use_package</code>	Add a package you’re using to the <code>DESCRIPTION</code> file
<code>use_vignette</code>	Set up the package to include a vignette
<code>use_readme_rmd</code>	Set up the package to include a README file in R Markdown format
<code>use_build_ignore</code>	Specify files that should be ignored when building the R package (for example, if you have a folder where you’re drafting a journal article about the package, you can include all related files in a folder that you set to be ignored during the package build)
<code>check</code>	Check the full R package for any ERRORS, WARNINGS, or NOTES
<code>build_win</code>	Build a version of the package for Windows and send it to be checked on a Windows machine. You’ll receive an email with a link to the results.
<code>use_travis</code>	Set the package up to facilitate using Travis CI with the package
<code>use_cran_comments</code>	Create a file where you can add comments to include with your CRAN submission.
<code>submit_cran</code>	Submit the package to CRAN
<code>use_news_md</code>	Add a file to the package to give news on changes in new versions

Some of these functions you’ll only need to use once for a package. The one-time (per package) functions are mostly those that set up a certain type of infrastructure for the package. For example, if you want to use R Markdown to create a README file for a package you are posting to GitHub, you can create the proper infrastructure with the `use_readme_rmd` function. This function adds a starter README file in the main directory of the package with the name “`README.Rmd`”. You can edit this file and render it to Markdown to provide GitHub users more information about your package. However, you will have problems with your CRAN checks if there is a README file in this top-level directory of the package, so the `use_readme_rmd` function also adds the files names for the R Markdown README file, and the Markdown file it creates, in the “`.Rbuildignore`” file, so it is not included when the package is built.

Creating a package

The earliest infrastructure function you will use from the `devtools` package is `create`. This function inputs the filepath for the directory where you would like to create the package

and creates the initial package structure (as a note, this directory should not yet exist). You will then add the elements (code, data, etc.) for the package within this structure. As an alternative to `create`, you can also initialize an R package in RStudio by selecting “File” -> “New Project” -> “New Direction” -> “R Package”.



In addition to starting a package using `create` or by creating a new project in RStudio, you could also create the package by hand, creating and then filling a directory. However, it's hard to think of any circumstances where there would be a good reason to do that rather than using some of the more convenient tools offered by `devtools` and RStudio.

Figure @ref(fig:initialpackagestructure) gives an example of what the new package directory will look like after you create an initial package structure with `create` or via the RStudio “New Project” interface. This initial package directory includes an `R` subdirectory, where you will save R scripts with all code defining R functions for the package. It also includes two files that will store metadata and interface information about your package (`DESCRIPTION` and `NAMESPACE`), as well as an R project file (`.Rproj` extension) that saves some project options for the directory. Finally, the initial package structure includes two files that can be used to exclude some files in the directory from either being followed by git (`.gitignore`) or included when the package is built (`.Rbuildignore`). These two files have names that start with a dot, so they may not be listed if you look at the package directory structure in a file manager like “Finder” on Macs. These “dot-files” will, however, be listed in the “Files” tab that shows up in one of the RStudio panes when you open an R project like a package directory, as shown in this figure.

	..	
	<code>.gitignore</code>	29 B
	<code>.Rbuildignore</code>	28 B
	<code>DESCRIPTION</code>	329 B
	<code>examplepackage.Rproj</code>	312 B
	<code>NAMESPACE</code>	96 B
	<code>R</code>	

Example of the directory contents of the initial package structure created with `devtools`.

Other functions

In contrast to the `devtools` infrastructure functions that you will only use once per package, there are other `devtools` functions you'll use many times as you develop a package. Two of the workhorses of `devtools` are `load_all` and `document`. The `load_all` function loads the entire package (by default, based on the current working directory, although you can also give the filepath to load a package directory elsewhere). In addition to loading all R functions, it also loads all package data and compiles and connects C, C++, and FORTRAN code in the package. As you add to a package, you can use `load_all` to ensure you're using the latest version of all package functions and data. The `document` function rewrites the help files and `NAMESPACE` file based on the latest version of the `roxygen2` comments for each function (writing `roxygen2` is covered in more detail in the next section).



RStudio has created a very helpful Package Development Cheatsheet that covers many of the `devtools` functions. A pdf of this cheatsheet is available [here](#).

Summary

The `devtools` package contains functions that help with R package development. These functions include `create`, which creates the initial structure for a new package, as well as a number of functions for adding useful infrastructure to the package directory and functions to load and document the package.

3.4 Documentation

The objectives of this section are:

- Create R function documentation using `roxygen2`
- Create vignettes using `knitr` and R Markdown

There are two main types of documentation you may want to include with packages:

- Longer documents that give tutorials or overviews for the whole package
- Shorter, function-specific help files for each function or group of related functions

You can create the first type of document using package vignettes, README files, or both. For the function-specific help files, the easiest way to create these is with the `roxygen2` package.

In this section, we'll cover why and how to create this documentation. In addition, vignette / README documentation can be done using `knitr` to create R Markdown documents that mix R code and text, so we'll include more details on that process.

Vignettes and README files

You will likely want to create a document that walks users through the basics of how to use your package. You can do this through two formats:

- Vignette: This document is bundled with your R package, so it becomes locally available to a user once they install your package from CRAN. They will also have it available if they install the package from GitHub, as long as they use the `build_vignettes = TRUE` option when running `install_github`.
- README file: If you have your package on GitHub, this document will show up on the main page of the repository.

A package likely only needs a README file if you are posting the package to GitHub. For any GitHub repository, if there is a README.md file in the top directory of the repository, it will be rendered on the main GitHub repository page below the listed repository content. For an example, visit <https://github.com/geanders/countytimezones> and scroll down. You'll see a list of all the files and subdirectories included in the package repository and below that is the content in the package's README.md file, which gives a tutorial on using the package.

If the README file does not need to include R code, you can write it directly as an .md file, using Markdown syntax, which is explained in more detail in the next section. If you want to include R code, you should start with a README.Rmd file, which you can then render to Markdown using knitr. You can use the devtools package to add either a README.md or README.Rmd file to a package directory using `use_readme_md` or `use_readme_rmd`, respectively. These functions will add the appropriate file to the top level of the package directory and will also add the file name to ".Rbuildignore", since having one of these files in the top level of the package directory could otherwise cause some problems when building the package.

The README file is a useful way to give GitHub users information about your package, but it will not be included in builds of the package or be available through CRAN for packages that are posted there. Instead, if you want to create tutorials or overview documents that are included in a package build, you should do that by adding one or more package vignettes. Vignettes are stored in a vignettes subdirectory within the package directory.

To add a vignette file, saved within this subdirectory (which will be created if you do not already have it), use the `use_vignette` function from `devtools`. This function takes as arguments the file name of the vignette you'd like to create and the package for which you'd like to create it (the default is the package in the current working directory). For example, if you are currently working in your package's top-level directory and you would like to add a vignette called "model_details", you can do that with the code:

```
use_vignette("model_details")
```

You can have more than one vignette per package, which can be useful if you want to include one vignette that gives a more general overview of the package as well as a few vignettes that go into greater detail about particular aspects or applications.



Once you create a vignette with `use_vignette`, be sure to update the Vignette Index Entry in the vignette's YAML (the code at the top of an R Markdown document). Replace “Vignette Title” there with the actual title you use for the vignette.

Knitr / Markdown

Both vignettes and README files can be written as R Markdown files, which will allow you to include R code examples and results from your package. One of the most exciting tools in R is the `knitr` system for combining code and text to create a reproducible document. In terms of the power you get for time invested in learning a tool, `knitr` probably can't be beat. Everything you need to know to create and “knit” a reproducible document can be learned in about 20 minutes, and while there is a lot more you can do to customize this process if you want to, probably 80% of what you'll ever want to do with `knitr` you'll learn in those first 20 minutes.

R Markdown files are mostly written using Markdown. To write R Markdown files, you need to understand what markup languages like Markdown are and how they work. In Word and other word processing programs you have used, you can add formatting using buttons and keyboard shortcuts (e.g., “Ctrl-B” for bold). The file saves the words you type. It also saves the formatting, but you see the final output, rather than the formatting markup, when you edit the file (WYSIWYG – what you see is what you get). In markup languages, on the other hand, you markup the document directly to show what formatting the final version should have (e.g., you type `**bold**` in the file to end up with a document with **bold**). Examples of markup languages include:

- HTML (HyperText Markup Language)
- LaTex
- Markdown (a “lightweight” markup language)

Common Markdown formatting elements

To write a file in Markdown, you'll need to learn the conventions for creating formatting. This table shows what you would need to write in a flat file for some common formatting choices:

Code	Rendering	Explanation
<code>**text**</code>	<code>text</code>	boldface
<code>*text*</code>	<code>text</code>	italicized
<code>[text](www.google.com)</code>	<code>text</code>	hyperlink
<code># text</code>		first-level header
<code>## text</code>		second-level header

Some other simple things you can do in Markdown include:

- Lists (ordered or bulleted)

- Equations
- Tables
- Figures from files
- Block quotes
- Superscripts

The start of a Markdown file gives some metadata for the file (authors, title, format) in a language called YAML. For example, the YAML section of a package vignette might look like this:

```
---
```

```
title: "Model Details for example_package"
author: "Jane Doe"
date: "2016-12-08"
output: rmarkdown::html_vignette
vignette: >
  \%VignetteIndexEntry{Model Details for example_package}
  \%VignetteEngine{knitr::rmarkdown}
  \%VignetteEncoding{UTF-8}
```

```
---
```

When creating R Markdown documents using the RStudio toolbar, much of this YAML will be automatically generated based on your specifications when opening the initial file. However, this is not the case with package vignettes, for which you'll need to go into the YAML and add the authors and title yourself. Leave the vignette engine, vignette encoding, output, and date as their default values.

For more Markdown conventions, see [RStudio's R Markdown Reference Guide](#) (link also available through “Help” in RStudio).

R Markdown files work a lot like Markdown files, but add the ability to include R code that will be run before rendering the final document. This functionality is based on *literate programming*, an idea developed by Donald Knuth, to mix executable code with regular text. The files you create can then be “knitted”, to run any embedded code. The final output will have results from your code and the regular text.

The basic steps of opening and rendering an R Markdown file in RStudio are:

- To open a new R Markdown file, go to “File” -> “New File” -> “RMarkdown...”. To start, choose a “Document” in “HTML” format.
- This will open a new R Markdown file in RStudio. The file extension for R Markdown files is “.Rmd”.
- The new file comes with some example code and text. You can run the file as-is to try out the example. You will ultimately delete this example code and text and replace it with your own.
- Once you “knit” the R Markdown file, R will render an HTML file with the output. This is automatically saved in the same directory where you saved your .Rmd file.

- Write everything besides R code using Markdown syntax.

The `knit` function from the `knitr` package works by taking a document in R Markdown format (among a few possible formats), reading through it for any markers of the start of R code, running any of the code between that “start” marker and a marker showing a return to regular Markdown, writing any of the relevant results from R code into the Markdown file in Markdown format, and then passing the entire document to software that can render from Markdown to the desired output format (for example, compile a pdf, Word, or HTML document).

This means that all a user needs to do to include R code within a document is to properly separate it from other parts of the document through the appropriate markers. To indicate R code in an RMarkdown document, you need to separate off the code chunk using the following syntax:

```
```{r}
my_vec <- 1:10
````
```

This syntax tells R how to find the start and end of pieces of R code (*code chunks*) when the file is rendered. R will walk through, find each piece of R code, run it and create output (printed output or figures, for example), and then pass the file along to another program to complete rendering (e.g., Tex for pdf files).

You can specify a name for each chunk, if you’d like, by including it after “r” when you begin your chunk. For example, to give the name `load_mtcars` to a code chunk that loads the `mtcars` dataset, specify that name in the start of the code chunk:

```
```{r load_mtcars}
data(mtcars)
````
```



Here are a couple of tips for naming code chunks:

- Chunk names must be unique across a document.
- Any chunks you don’t name are given ordered numbers by `knitr`.

You do not have to name each chunk. However, there are some advantages:

- It will be easier to find any errors.
- You can use the chunk labels in referencing for figure labels.
- You can reference chunks later by name.

Common knitr chunk options

You can also add options when you start a chunk. Many of these options can be set as TRUE / FALSE and include:

| Option | Action |
|----------|---|
| echo | Print out the R code? |
| eval | Run the R code? |
| messages | Print out messages? |
| warnings | Print out warnings? |
| include | If FALSE, run code, but don't print code or results |

Other chunk options take values other than TRUE / FALSE. Some you might want to include are:

| Option | Action |
|------------|---|
| results | How to print results (e.g., <code>hide</code> runs the code, but doesn't print the results) |
| fig.width | Width to print your figure, in inches (e.g., <code>fig.width = 4</code>) |
| fig.height | Height to print your figure |

To include any of these options, add the option and value in the opening brackets and separate multiple options with commas:

```
```{r messages = FALSE, echo = FALSE}
mtcars[1, 1:3]
```

```

You can set “global” options at the beginning of the document. This will create new defaults for all of the chunks in the document. For example, if you want `echo`, `warning`, and `message` to be FALSE by default in all code chunks, you can run:

```
```{r global_options}
knitr::opts_chunk$set(echo = FALSE, message = FALSE,
warning = FALSE)
```

```

If you set both global and local chunk options that you set specifically for a chunk will take precedence over global options. For example, running a document with:

```
```{r global_options}
knitr::opts_chunk$set(echo = FALSE, message = FALSE,
warning = FALSE)
```

```

```
```{r check_mtcars, echo = TRUE}
head(mtcars, 1)
```

```

would print the code for the `check_mtcars` chunk, because the option specified for that specific chunk (`echo = TRUE`) would override the global option (`echo = FALSE`).

You can also include R output directly in your text (“inline”) using backticks: \bigskip
“There are `r nrow(mtcars)` observations in the `mtcars` data set. The average miles per gallon is `r mean(mtcars\$mpg, na.rm = TRUE)`.”

\bigskip

Once the file is rendered, this gives: \bigskip

“There are 32 observations in the `mtcars` data set. The average miles per gallon is 20.090625.”
\bigskip



Here are some tips that will help you diagnose some problems rendering R Markdown files:

- Be sure to save your R Markdown file before you run it.
- All the code in the file will run “from scratch”—as if you just opened a new R session.
- The code will run using, as a working directory, the directory where you saved the R Markdown file.
- To use the latest version of functions in a package you are developing in an R Markdown document, rebuild the package before knitting the document. You can build a package using the “Build” tab in one of the RStudio panes.

You’ll want to try out pieces of your code as you write an R Markdown document. There are a few ways you can do that:

- You can run code in chunks just like you can run code from a script (Ctrl-Return or the “Run” button).
- You can run all the code in a chunk (or all the code in all chunks) using the different options under the “Run” button in RStudio.
- All the “Run” options have keyboard shortcuts, so you can use those.



Two excellent books for learning more about creating reproducible documents with R are *Dynamic Documents with R and knitr* by Yihui Xie (the creator of `knitr`) and *Reproducible Research with R and RStudio* by Christopher Gandrud. The first goes into the technical details of how `knitr` and related code works, which gives you the tools to extensively customize a document. The second provides an extensive view of how to use tools from R and other open source software to conduct, write up, and present research in a reproducible and efficient way. RStudio’s [R Markdown Cheatsheet](#) is another very useful reference.

Help files and roxygen2

In addition to writing tutorials that give an overview of your whole package, you should also write specific documentation showing users how to use and interpret any functions you expect users to directly call.

These help files will ultimately go in a folder called `/man` of your package, in an R documentation format (`.Rd` file extensions) that is fairly similar to LaTeX. You used to have to write all of these files as separate files. However, the `roxygen2` package lets you put all of the help information directly in the code where you define each function. Further, `roxygen2` documentation allows you to include tags (`@export`, `@importFrom`) that will automate writing the package NAMESPACE file, so you don't need to edit that file by hand.

With `roxygen2`, you add the help file information directly above the code where you define each functions, in the R scripts saved in the `R` subdirectory of the package directory. You start each line of the `roxygen2` documentation with `#'` (the second character is an apostrophe, not a backtick). The first line of the documentation should give a short title for the function, and the next block of documentation should be a longer description. After that, you will use tags that start with `@` to define each element you're including. You should leave an empty line between each section of documentation, and you can use indentation for second and later lines of elements to make the code easier to read.

Here is a basic example of how this `roxygen2` documentation would look for a simple “Hello world” function:

```
#' Print "Hello world"
#
#' This is a simple function that, by default, prints "Hello world". You can
#' customize the text to print (using the \code{to_print} argument) and add
#' an exclamation point (\code{excited = TRUE}).
#
#' @param to_print A character string giving the text the function will print
#' @param excited Logical value specifying whether to include an exclamation
#'   point after the text
#
#' @return This function returns a phrase to print, with or without an
#'   exclamation point added. As a side effect, this function also prints out
#'   the phrase.
#
#' @examples
#' hello_world()
#' hello_world(excited = TRUE)
#' hello_world(to_print = "Hi world")
#
#' @export
hello_world <- function(to_print = "Hello world", excited = FALSE){
  if(excited) to_print <- paste0(to_print, "!")
  print(to_print)
}
```

You can run the `document` function from the `devtools` package at any time to render the latest version of these `roxygen2` comments for each of your functions. This will create function-specific help files in the package’s “man” subdirectory as well as update the package’s NAMESPACE file.

Common roxygen2 tags

Here are some of the common `roxygen2` tags to use in creating this documentation:

| Tag | Meaning |
|-----------------------------|---|
| <code>@return</code> | A description of the object returned by the function |
| <code>@parameter</code> | Explanation of a function parameter |
| <code>@inheritParams</code> | Name of a function from which to get parameter definitions |
| <code>@examples</code> | Example code showing how to use the function |
| <code>@details</code> | Add more details on how the function works (for example, specifics of the algorithm being used) |
| <code>@note</code> | Add notes on the function or its use |
| <code>@source</code> | Add any details on the source of the code or ideas for the function |
| <code>@references</code> | Add any references relevant to the function |
| <code>@importFrom</code> | Import a function from another package to use in this function (this is especially useful for inline functions like <code>%%</code> and <code>%within%</code>) |
| <code>@export</code> | Export the function, so users will have direct access to it when they load the package |

Here are a few things to keep in mind when writing help files using `roxygen2`:

- The tags `@example` and `@examples` do different things. You should *always* use the `@examples` (plural) tag for example code, or you will get errors when you build the documentation.
- The `@inheritParams` function can save you a lot of time, because if you are using the same parameters in multiple functions in your package, you can write and edit those parameter descriptions just in one place. However, keep in mind that you must point `@inheritParams` to the function where you originally define the parameters using `@param`, not another function where you use the parameters but define them using an `@inheritParams` pointer.
- If you want users to be able to directly use the function, you must include `@export` in your `roxygen2` documentation. If you have written a function but then find it isn’t being found when you try to compile a README file or vignette, a common culprit is that you have forgotten to export the function.

Common roxygen2 formatting tags

You can include formatting (lists, etc.) and equations in the `roxygen2` documentation. Here are some of the common formatting tags you might want to use:

| Tag | Meaning |
|-------------------------|--|
| <code>\code{}</code> | Format in a typeface to look like code |
| <code>\dontrun{}</code> | Use with examples, to avoid running the example code during package builds and testing |
| <code>\link{}</code> | Link to another R function |
| <code>\eqn{}}</code> | Include an inline equation |
| <code>\deqn{}}</code> | Include a display equation (i.e., shown on its own line) |
| <code>\itemize{}</code> | Create an itemized list |
| <code>\url{}</code> | Include a web link |
| <code>\href{}}</code> | Include a web link |

Some tips on using the R documentation format:

- Usually, you'll want you use the `\link` tag only in combination with the `\code` tag, since you're linking to another R function. Make sure you use these with `\code` wrapping `\link`, not the other way around (`\code{\link{other_function}}`), or you'll get an error.
- Some of the equation formatting, including superscripts and subscripts, won't parse in Markdown-based documentation (but will for pdf-based documentation). With the `\eqn` and `\deqn` tags, you can include two versions of an equation, one with full formatting, which will be fully compiled by pdf-based documentation, and one with a reduced form that looks better in Markdown-based documentation (for example, `\deqn{ \frac{X^2}{Y} }{ X2 / Y }`).
- For any examples in help files that take a while to run, you'll want to wrap the example code in the `\dontrun` tag.
- The tags `\url` and `\href` both include a web link. The difference between the two is that `\url` will print out the web address in the help documentation, `\href` allows you to use text other than the web address for the anchor text of the link. For example: "For more information, see `\url{www.google.com}`."; "For more information, `\href{www.google.com}{Google it}`".

In addition to document functions, you should also document any data that comes with your package. To do that, create a file in the `/R` folder of the package called “`data.R`” to use to documentation all of the package’s datasets. You can use `roxygen2` to document each dataset, and end each with the name of the dataset in quotation marks. There are more details on documenting package data using `roxygen2` in the next section.

Summary

You should include documentation to help others use your package, both longer-form documentation through vignettes or README files and function-specific help files. Longer-form documentation can be written using R Markdown files, which can include executable R code examples, while function-specific help files can be written using `roxygen2` comments within the R scripts where each function is defined.

3.5 Data Within a Package

The objective of this section is:

- Create an R package that contains data (and associated documentation)

Many R packages are designed to manipulate, visualize, and model data so it may be a good idea for you to include some data in your package. The primary reason most developers include data in their package is to demonstrate how to use the functions included in the package with the included data. Creating a package as a means to distribute data is also a method that is gaining popularity. Additionally you may want to include data that your package uses internally, but is not available to somebody who is using your package. When including data in your package consider the fact that your compressed package file should be smaller than 5MB, which is the largest package size that CRAN allows. If your package is larger than 5MB make sure to inform users in the instructions for downloading and installing your package.

Data for Demos

Data Objects

Including data in your package is easy thanks to the `devtools` package. To include datasets in a package, first create the objects that you would like to include in your package inside of the global environment. You can include any R object in a package, not just data frames. Then make sure you're in your package directory and use the `use_data()` function, listing each object that you want to include in your package. The names of the objects that you pass as arguments to `use_data()` will be the names of the objects when a user loads the package, so make sure you like the variable names that you're using.

You should then document each data object that you're including in the package. This way package users can use common R help syntax like `?dataset` to find out more information about the included data set. You should create one R file called `data.R` in the `R/` directory of your package. You can write the data documentation in the `data.R` file. Let's take a look at some documentation examples from the `minimap` package. First we'll look at the documentation for a data frame called `maple`:

```
#' Production and farm value of maple products in Canada
#'
#' @source Statistics Canada. Table 001-0008 - Production and farm value of
#'   maple products, annual. \url{http://www5.statcan.gc.ca/cansim/}
#' @format A data frame with columns:
#'   \describe{
#'     \item{Year}{A value between 1924 and 2015.}
#'     \item{Syrup}{Maple products expressed as syrup, total in thousands of gallons.}
#'     \item{CAD}{Gross value of maple products in thousands of Canadian dollars.}
#'     \item{Region}{Postal code abbreviation for territory or province.}
```

```
#' }
#' @examples
#' \dontrun{
#'   maple
#' }
"maple"
```

Data frames that you include in your package should follow the general schema above where the documentation page has the following attributes:

- An informative title describing the object.
- A `@source` tag describing where the data was found.
- A `@format` tag which describes the data in each column of the data frame.
- And then finally a string with the name of the object.

The `minimap` package also includes a few vectors. Let's look at the documentation for `mexico_abb`:

```
#' Postal Abbreviations for Mexico
#'
#' @examples
#' \dontrun{
#'   mexico_abb
#' }
"mexico_abb"
```

You should always include a title for a description of a vector or any other object. If you need to elaborate on the details of a vector you can include a description in the documentation or a `@source` tag. Just like with data frames the documentation for a vector should end with a string containing the name of the object.

Raw Data

A common task for R packages is to take raw data from files and to import them into R objects so that they can be analyzed. You might want to include some sample raw data files so you can show different methods and options for importing the data. To include raw data files in your package you should create a directory under `inst/extdata` in your R package. If you stored a data file in this directory called `response.json` in `inst/extdata` and your package is named `mypackage` then a user could access the path to this file with `system.file("extdata", "response.json", package = "mypackage")`. Include that line of code in the documentation to your package so that your users know how to access the raw data file.

Internal Data

Functions in your package may need to have access to data that you don't want your users to be able to access. For example the `swirl` package contains translations for menu items into languages other than English, however that data has nothing to do with the purpose of the `swirl` package and so it's hidden from the user. To add internal data to your package you can use the `use_data()` function from `devtools`, however you must specify the `internal = TRUE` argument. All of the objects you pass to `use_data(..., internal = TRUE)` can be referenced by the same name within your R package. All of these objects will be saved to one file called `R/sysdata.rda`.

Data Packages

There are several packages which were created for the sole purpose of distributing data including `janeaustenr`, `gapminder`, `babynames`, and `lego`. Using an R package as a means of distributing data has advantages and disadvantages. On one hand the data is extremely easy to load into R, as a user only needs to install and load the package. This can be useful for teaching folks who are new to R and may not be familiar with importing and cleaning data. Data packages also allow you document datasets using `roxygen2`, which provides a much cleaner and more programmer-friendly kind of code book compared to including a file that describes the data. On the other hand data in a data package is not accessible to people who are not using R, though there's nothing stopping you from distributing the data in multiple ways.

If you decide to create a data package you should document the process that you used to obtain, clean, and save the data. One approach to doing this is to use the `use_data_raw()` function from `devtools`. This will create a directory inside of your package called `data_raw`. Inside of this directory you should include any raw files that the data objects in your package are derived from. You should also include one or more R scripts which import, clean, and save those data objects in your R package. Theoretically if you needed to update the data package with new data files you should be able to just run these scripts again in order to rebuild your package.

Summary

Including data in a package is useful for showing new users how to use your package, using data internally, and sharing and documenting datasets. The `devtools` package includes several useful functions to help you add data to your package including `use_data()` and `use_data_raw()`. You can document data within your package just like you would document a function.

3.6 Software Testing Framework for R Packages

The objective of this section is:

- Create unit tests for an R package using the `testthat` package

Once you've written code for an R package and have gotten that code to a point where you believe it's working, it may be a good time to step back and consider a few things about your code.

- **How do you know it's working?** Given that you wrote the functions, you have a certain set of *expectations* about how the functions should behave. Specifically, for a given set of inputs you expect a certain output. Having these expectations clearly in mind is an important aspect of knowing whether code is "working".
- **Have you already tested your code?** Chances are, throughout the development of your code, you ran little tests to see if your functions were working. Assuming these tests were valid for the code you were testing, it's worth keeping these tests on hand and making them part of your package.

Setting up a battery of tests for the code in your package can play a big role in maintaining the ongoing smooth operation of the package in hunting down bugs in the code, should they arise. Over time, many aspects of a package can change. Specifically:

- As you actively develop your code, you may change/break older code without knowing it. For example, modifying a helper function that lots of other functions rely on may be better for some functions but may break behavior for other functions. Without a comprehensive testing framework, you might not know that some behavior is broken until a user reports it to you.
- The environment in which your package runs can change. The version of R, libraries, web sites and any other external resources, and packages can all change without warning. In such cases, your code may be unchanged, but because of an external change, your code may not produce the expected output given a set of inputs. Having tests in place that are run regularly can help to catch these changes even if your package isn't under active development.
- As you fix bugs in your code, it's often a good idea to include a specific test that addresses each bug so that you can be sure that the bug does not "return" in a future version of the package (this is also known as a regression).

Testing your code effectively has some implications for code design. In particular, it may be more useful to divide your code into smaller functions so that you can test individual pieces more effectively. For example, if you have one large function that returns `TRUE` or `FALSE`, it is easy to test this function, but ultimately it may not be possible to identify problems deep in the code by simply checking if the function returns the correct logical value. It may be better to divide up large function into smaller functions so that core elements of the function can be tested separately to ensure that they are behaving appropriately.

The `testthat` Package

The `testthat` package is designed to make it easy to setup a battery of tests for your R package. A nice introduction to the package can be found in Hadley Wickham's [article](#) in the *R Journal*. Essentially, the package contains a suite of functions for testing function/expression output with the expected output. The simplest use of the package is for testing a simple expression:

```
library(testthat)
expect_that(sqrt(3) * sqrt(3), equals(3))
```

Note that the `equals()` function allows for some numerical fuzz, which is why this expression actually passes the test. When a test fails, `expect_that()` throws an error and does not return something.

```
## Use a strict test of equality (this test fails)
expect_that(sqrt(3) * sqrt(3), is_identical_to(3))

Error: sqrt(3) * sqrt(3) not identical to 3.
Objects equal but not identical
```

The `expect_that()` function can be used to wrap many different kinds of test, beyond just numerical output. The table below provides a brief summary of the types of comparisons that can be made.

| Expectation | Description |
|---------------------------------|--|
| <code>equals()</code> | check for equality with numerical fuzz |
| <code>is_identical_to()</code> | strict equality via <code>identical()</code> |
| <code>is_equivalent_to()</code> | like <code>equals()</code> but ignores object attributes |
| <code>is_a()</code> | checks the class of an object (using <code>inherits()</code>) |
| <code>matches()</code> | checks that a string matches a regular expression |
| <code>prints_text()</code> | checks that an expression prints to the console |
| <code>shows_message()</code> | checks for a message being generated |
| <code>gives_warning()</code> | checks that an expression gives a warning |
| <code>throws_error()</code> | checks that an expression (properly) throws an error |
| <code>is_true()</code> | checks that an expression is <code>TRUE</code> |

A collection of calls to `expect_that()` can be put together with the `test_that()` function, as in

```
test_that("model fitting", {
  data(airquality)
  fit <- lm(Ozone ~ Wind, data = airquality)
  expect_that(fit, is_a("lm"))
  expect_that(1 + 1, equals(2))
})
```

Typically, you would put your tests in an R file. If you have multiple sets of tests that test

different domains of a package, you might put those tests in different files. Individual files can have their tests run with the `test_file()` function. A collection of tests files can be placed in a directory and tested all together with the `test_dir()` function.

In the context of an R package, it makes sense to put the test files in the `tests` directory. This way, when running `R CMD check` (see the next section) all of the tests will be run as part of the process of checking the entire package. If any of your tests fail, then the entire package checking process will fail and will prevent you from distributing buggy code. If you want users to be able to easily see the tests from an installed package, you can place the tests in the `inst/tests` directory and have a separate file in the `tests` directory to run all of the tests.

3.7 Passing CRAN checks

The objective of this section is:

- Categorize errors in the R CMD check process

Before submitting a package to CRAN, you must pass a battery of tests that are run by the R itself via the `R CMD check` program. In RStudio, if you are in an R Package “Project” you can run `R CMD check` by clicking the `Check` button in the build tab. This will run a series of tests that check the metadata in your package, the `NAMESPACE` file, the code, the documentation, run any tests, build any vignettes, and many others.

Here is an example of the output form `R CMD check` for the `filehash` package which currently passes all tests.

```
* using R version 3.3.2 (2016-10-31)
* using platform: x86_64-apple-darwin13.4.0 (64-bit)
* using session charset: UTF-8
* checking for file 'filehash/DESCRIPTION' ... OK
* this is package 'filehash' version '2.3'
* checking package namespace information ... OK
* checking package dependencies ... OK
* checking if this is a source package ... OK
* checking if there is a namespace ... OK
* checking for executable files ... OK
* checking for hidden files and directories ... OK
* checking for portable file names ... OK
* checking for sufficient/correct file permissions ... OK
* checking whether package 'filehash' can be installed ... OK
* checking installed package size ... OK
* checking package directory ... OK
* checking 'build' directory ... OK
* checking DESCRIPTION meta-information ... OK
* checking top-level files ... OK
* checking for left-over files ... OK
* checking index information ... OK
```

```
* checking package subdirectories ... OK
* checking R files for non-ASCII characters ... OK
* checking R files for syntax errors ... OK
* checking whether the package can be loaded ... OK
* checking whether the package can be loaded with stated dependencies ... OK
* checking whether the package can be unloaded cleanly ... OK
* checking whether the namespace can be loaded with stated dependencies ... OK
* checking whether the namespace can be unloaded cleanly ... OK
* checking loading without being on the library search path ... OK
* checking dependencies in R code ... OK
* checking S3 generic/method consistency ... OK
* checking replacement functions ... OK
* checking foreign function calls ... OK
* checking R code for possible problems ... OK
* checking Rd files ... OK
* checking Rd metadata ... OK
* checking Rd cross-references ... OK
* checking for missing documentation entries ... OK
* checking for code/documentation mismatches ... OK
* checking Rd \usage sections ... OK
* checking Rd contents ... OK
* checking for unstated dependencies in examples ... OK
* checking line endings in C/C++/Fortran sources/headers ... OK
* checking compiled code ... OK
* checking sizes of PDF files under 'inst/doc' ... OK
* checking installed files from 'inst/doc' ... OK
* checking files in 'vignettes' ... OK
* checking examples ... OK
* checking for unstated dependencies in 'tests' ... OK
* checking tests ...
OK
* checking for unstated dependencies in vignettes ... OK
* checking package vignettes in 'inst/doc' ... OK
* checking running R code from vignettes ...
  'filehash.Rnw' ... OK
OK
* checking re-building of vignette outputs ... OK
* checking PDF version of manual ... OK
* DONE
Status: OK
```

Here is an example from the `mvtspolt` package where we've deliberately introduced some problems to the package in order to show the check output. Checks that have passed are not shown below.

```
* checking foreign function calls ... OK
* checking R code for possible problems ... NOTE
drawImage: no visible global function definition for 'Axis'
drawImageMargin: no visible global function definition for 'lm'
drawImageMargin: no visible global function definition for 'Axis'
splineFillIn: no visible global function definition for 'lm'
Undefined global functions or variables:
  Axis lm
Consider adding
  importFrom("graphics", "Axis")
  importFrom("stats", "lm")
to your NAMESPACE file.
```

Here, it appears that the functions `Axis()` and `lm()` are needed by the package but are not available because they are not imported from their respective packages. In this case, R CMD check provides a suggestion of how you can modify the NAMESPACE package, but you are probably better off modifying the `roxygen2` documentation in the code file instead.

Moving on the rest of the checks, we see:

```
* checking for missing documentation entries ... OK
* checking for code/documentation mismatches ... WARNING
Codoc mismatches from documentation object 'mvtspplot':
mvtspplot
Code: function(x, group = NULL, xtime = NULL, norm = c("internal",
  "global"), levels = 3, smooth.df = NULL, margin =
  TRUE, sort = NULL, main = "", palette = "PRGn",
  rowstat = "median", xlim, bottom.ylim = NULL,
  right.xlim = NULL, gcol = 1)
Docs: function(y, group = NULL, xtime = NULL, norm = c("internal",
  "global"), levels = 3, smooth.df = NULL, margin =
  TRUE, sort = NULL, main = "", palette = "PRGn",
  rowstat = "median", xlim, bottom.ylim = NULL,
  right.xlim = NULL, gcol = 1)
Argument names in code not in docs:
  x
Argument names in docs not in code:
  y
Mismatches in argument names:
  Position: 1 Code: x Docs: y
```

Here the problem is that the code has the first argument named `x` while the documentation has the first argument named `y`.

```
* checking Rd \usage sections ... WARNING
Undocumented arguments in documentation object 'mvtspplot'
  'y'
Documented arguments not in \usage in documentation object 'mvtspplot':
  'x'

Functions with \usage entries need to have the appropriate \alias
entries, and all their arguments documented.
The \usage entries must correspond to syntactically valid R code.
See chapter 'Writing R documentation files' in the 'Writing R
Extensions' manual.
```

Because of the mismatch in code and documentation for the first argument, we have an argument that is not properly documented (*y*) and an argument that is documented but not used (*x*).

In case the checks fly by too quickly, you will receive a summary message the end saying what errors and warnings you got.

```
* DONE
Status: 2 WARNINGS, 1 NOTE
```

A package cannot be submitted to CRAN if there are any errors or warnings. If there is a NOTE, a package may be submitted if there is a Really Good Reason for that note.

3.8 Open Source Licensing

The objectives of this section are:

- Recall the principles of open source software
- Recall two open source licenses

You can specify how your R package is licensed in the package DESCRIPTION file under the `License:` section. How you license your R package is important because it provides a set of constraints for how other R developers use your code. If you're writing an R package to be used internally in your company then your company may choose to not share the package. In this case licensing your R package is less important since the package belongs to your company. In your package DESCRIPTION you can specify `License: file LICENSE`, and then create a text file called `LICENSE` which explains that your company reserves all rights to the package.

However if you (or your company) would like to publicly share your R package you should consider open source licensing. The philosophy of open source revolves around three principles:

1. The source code of the software can be inspected.

2. The source code of the software can be modified.
3. Modified versions of the software can be redistributed.

Nearly all open source licenses provide the protections above. Let's discuss three of the most popular open source licenses among R packages.

The General Public License

Known as the GPL, the GNU GPL, and GPL-3, the General Public License was originally written by [Richard Stallman](#). The GPL is known as a *copyleft license*, meaning that any software that is bundled with or originates from software licensed under the GPL must also be released under the GPL. The exact meaning of “bundle” will depend a bit on the circumstances. For example, software distributed with an operating system can be licensed under different licenses even if the operating system itself is licensed under the GPL. You can use the GPL-3 as the license for your R package by specifying `License: GPL-3` in the DESCRIPTION file.

It is worth noting that R itself is licensed under [version 2 of the GPL](#), or GPL-2, which is an earlier version of this license.

The MIT License

The MIT license is a more permissive license compared to the GPL. MIT licensed software can be modified or incorporated into software that is not open source. The MIT license protects the copyright holder from legal liability that might be incurred from using the software. When using the MIT license in a R package you should specify `License: MIT + file LICENSE` in the DESCRIPTION file. You should then add a file called LICENSE to your package which uses the following template exactly:

```
YEAR: [The current year]
COPYRIGHT HOLDER: [Your name or your organization's name]
```

The CC0 License

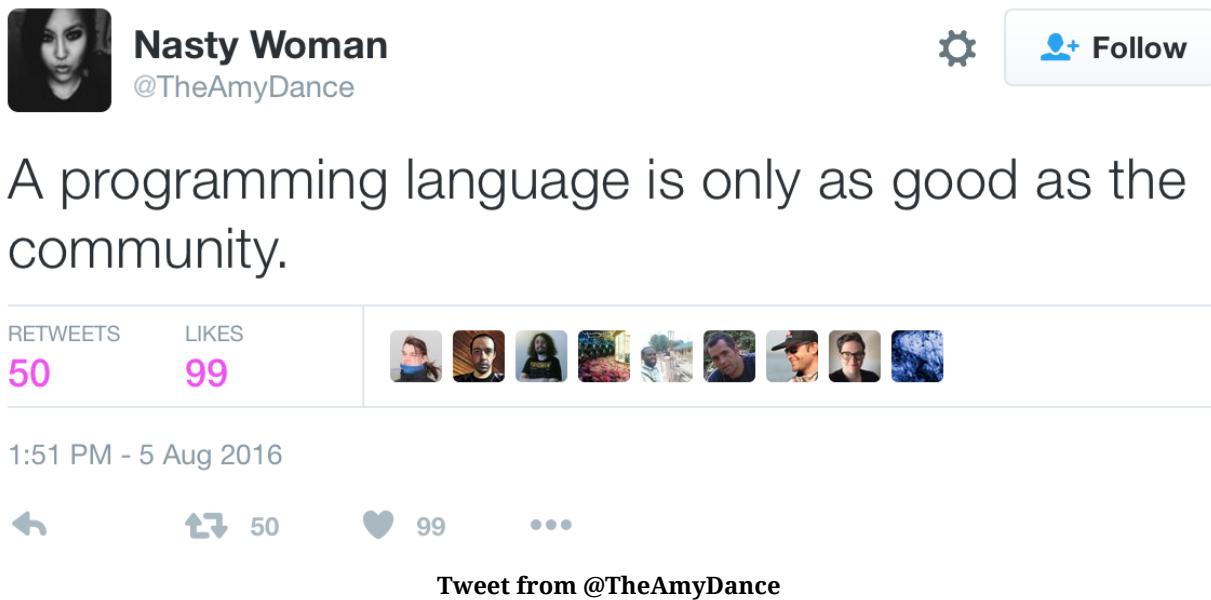
The [Creative Commons](#) licenses are usually used for artistic and creative works, however the CC0 license is also appropriate for software. The CC0 license dedicates your R package to the public domain, which means that you give up all copyright claims to your R package. The CC0 license allows your software to join other great works like *Pride and Prejudice*, *The Adventures of Huckleberry Finn*, and *The Scarlet Letter* in the public domain. You can use the CC0 license for your R package by specifying `License: cc0` in the DESCRIPTION file.

Why Open Source?

You've put weeks of sweat and mental anguish into writing a new R package, so why should you provide an open source license for software that you or your company owns by default? Let's discuss a few arguments for why open sourcing your software is a good idea.

Paying it Forward

Software development began in academic settings and the first computer programs with code that could be shared and run on multiple computers was shared between academics in the same way that academics share other kinds of scientific discoveries. The R programming language is open source, and there are hundreds of high-quality R packages that are also open source. A programming language can have lots of exciting features but the continued growth and improvement of a language is made possible by the people contributing to software written in that language. My colleague Amy [said it succinctly](#):



A screenshot of a Twitter post from user @TheAmyDance. The post contains the following text:

A programming language is only as good as the community.

The post has 50 retweets and 99 likes. Below the tweet, there are icons for retweeting, favoriting, and more options, along with the timestamp 1:51 PM - 5 Aug 2016.

Tweet from @TheAmyDance

So with that in mind, if you feel that the R language or the R community has contributed to your success or the success of your company consider open sourcing your software so that the greater R community can benefit from its availability.

Linus's Law

Now let's turn off the NPR pledge campaign and move our line of thinking from the Berkeley *Kumbaya* circle to the Stanford MBA classroom: as a business person why should you open source your software? One great reason is a concept called Linus's Law which refers to Linus Torvalds, the creator of Linux. The Linux operating system is a huge open source software project involving thousands of people. Linux has a reputation for security and for its lack of bugs which is in part a result of so many people looking at and being able to modify the source code. If the users of your software are able to view and modify the source code of your R package your package will likely be improved because of Linus's Law.

Hiring

Open source software's relationship with hiring is a two-way street: if you open source your software and other people send you improvements and contributions you can potentially

identify job candidates who you know are already familiar with your source code. On the other hand if you’re looking for a job your contributions to open source software can be a part of a compelling portfolio which showcases your software skills.

However there are pitfalls you should be aware of when weighing a candidate’s open source contributions. Many open source contributions are essentially “free work” - work that a candidate was able to do in their spare time. The best candidates often cannot afford to make open source contributions. The most meaningful ways that an individual contributes to their community usually has nothing to do with writing software.

Summary

Licensing and copyright laws vary between countries and jurisdictions. You shouldn’t consider any part of this chapter as legal advice. If you have questions about open source licensing software you’re building at work you should consult with your legal department. In most situations software that you write on your own time belongs to you, and software that you write while being paid by somebody else belongs to whoever is paying you. Open source licensing allows you to put restrictions on how your software can be used by others. The open source philosophy does not oppose the commercial sale of software. Many companies offer an open source version of their software that comes with limitations, while also offering a paid license for more expansive commercial use. This business model is used by companies like [RStudio](#) and [Highcharts](#).

3.9 Version Control and GitHub

The objective of this section is:

- Create a GitHub repository for an R package

GitHub allows you to post and interact with online code repositories, where all repositories are under git version control. You can post R packages on GitHub and, with the `install_github` function from the `devtools` package, install R packages directly from GitHub. GitHub can be particularly useful for collaborating with others on R packages, as it allows all collaborators to push and pull code between their personal computers and a GitHub repository. While git historically required you to leave R and run git functions at a command line, RStudio now has a number of features that make it easier to interface directly with GitHub.

When using git and GitHub, there are three levels of tasks you’ll need to do:

1. Initial set-up—these are things you will only need to do once (at least per computer).
 - Download git
 - Configure git with your user name and email
 - Set up a GitHub account

- Set up a SSH key to link RStudio on your personal computer with your GitHub account
2. Set-up of a specific repository— these are things you will need to do every time you create a new repository, but will only need to do once per repository.
- Initialize the directory on your personal computer as a git repository
 - Make an initial commit of files in the repository
 - Create an empty GitHub repository
 - Add the GitHub repository as a remote branch of the local repository
 - Push the local repository to the GitHub remote branch
 - (If you are starting from a GitHub repository rather than a local repository, either clone the repository or fork and clone the repository instead.)
3. Day-to-day workflow for a repository— these are things you will do regularly as you develop the code in a repository.
- Commit changes in files in the repository to save git history locally
 - Push committed changes to the GitHub remote branch
 - Pull the latest version of the GitHub remote branch to incorporate changes from collaborators into the repository code saved on your personal computer
 - Write and resolve “Issues” with the code in the repository
 - Fix any merge conflicts that come up between different collaborators’ code edits
 - If the repository is a fork, keep up-to-date with changes in the upstream branch

Each of these elements are described in detail in this section. More generally, this section describes how to use git and GitHub for version control and collaboration when building R packages.

git

Git is a *version control system*. When a repository is under git version control, information about all changes made, saved, and committed on any non-ignored file in a repository is saved. This allows you to revert back to previous versions of the repository and search through the history for all commits made to any tracked files in the repository. If you are working with others, using git version control allows you to see every change made to the code, who made it, and why (through the commit messages).

You will need git on your computer to create local git repositories that you can sync with GitHub repositories. Like R, git is open source. You can [download it](#) for different operating systems.

After downloading git but before you use it, you should configure it. For example, you should make sure it has your name and email address. You can configure git from a bash shell (for Macs, you can use “Terminal”, while for PCs you can use GitBash, which comes with the git installation).

You can use `git config` functions to configure your version of git. Two changes you should make are to include your name and email address as the `user.name` and `user.email`. For example, the following code, if run in a bash shell, would configure a git account for a user named “Jane Doe” who has a generic email address:

```
git config --global user.name "Jane Doe"
git config --global user.email "jane.doe@university.edu"
```

Once you've installed git, you should restart RStudio so RStudio can identify that git is now available. Often, just restarting RStudio will be enough. However, in some cases, you may need to take some more steps to activate git in RStudio. To do this, go to "RStudio" -> "Preferences" -> "Git/SVN". Choose "Enable version control". If RStudio doesn't automatically find your version of git in the "Git executable" box (you'll know if that box is blank), browse for your git executable file using the "Browse" button beside that box. If you aren't sure where your git executable is saved, try opening a bash shell and running `which git`, which should give you the filepath if you have git installed.

Initializing a git repository

You can initialize a git repository either using commands from a bash shell or directly from RStudio. First, to initialize a git repository from a bash shell, take the following steps:

1. Use a shell ("Terminal" on Macs) to navigate to that directory. You can use `cd` to do that (similar to `setwd` in R).
2. Once you are in the directory, first check that it is not already a git repository. To do that, run `git status`. If you get the message `fatal: Not a git repository` (or any of the parent directories): `.git`, it is not yet a git repository. If you do not get an error from `git status`, the directory is already a repository, so you do not need to initialize it.
3. If the directory is not already a git repository, run `git init` to initialize it as a repository.

For example, if I wanted to make a directory called "example_analysis", which is a direct subdirectory of my home directory, a git repository, I could open a shell and run:

```
cd ~/example_analysis
git init
```

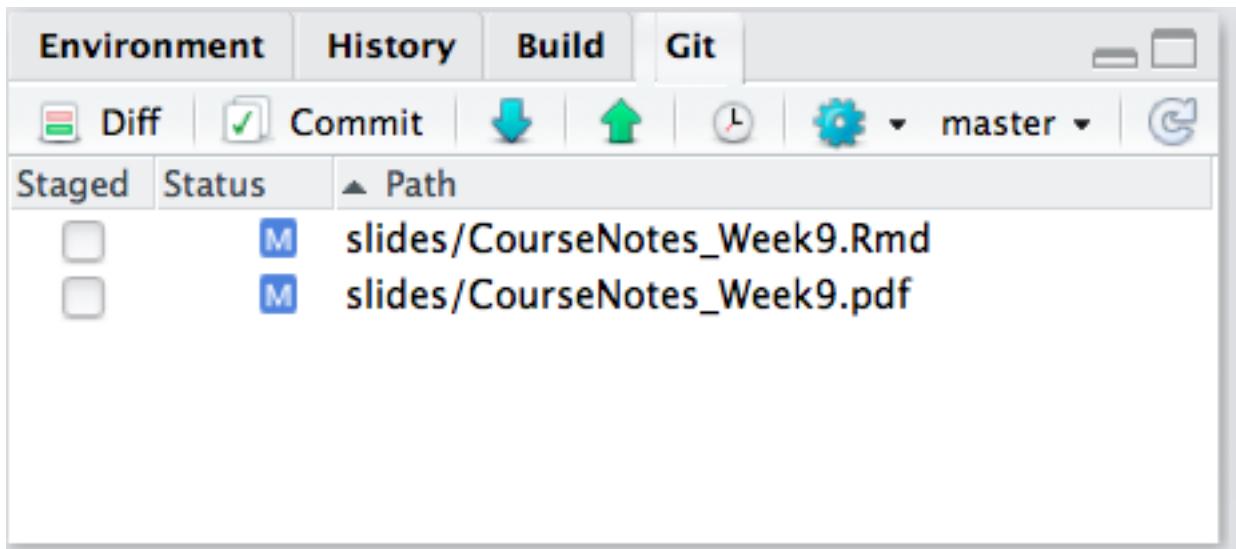
You can also initialize a directory as a git repository through R Studio. To do that, take the following steps:

1. Make the directory an R Project. If the directory is an R package, it likely already has an `.Rproj` file and so is an R Project. If the directory is not an R Project, you can make it one from RStudio by going to "File" -> "New Project" -> "Existing Directory", and then navigate to the directory you'd like to make an R project.
2. Open the R project.
3. Go to "Tools" -> "Version Control" -> "Project Setup".
4. In the box for "Version control system", choose "Git".



If you do not see “Git” in the box for “Version control system”, it means either that you do not have git installed on your computer or that RStudio was unable to find it. If so, see the earlier instructions for making sure that RStudio has identified the git executable.

Once you initialize the project as a git repository, you should have a “Git” window in one of your RStudio panes (top right pane by default). As you make and save changes to files, they will show up in this window for you to commit. For example, Figure @ref(fig:examplegitwindow) is what the Git window in RStudio looks like when there are changes to two files that have not yet been committed.

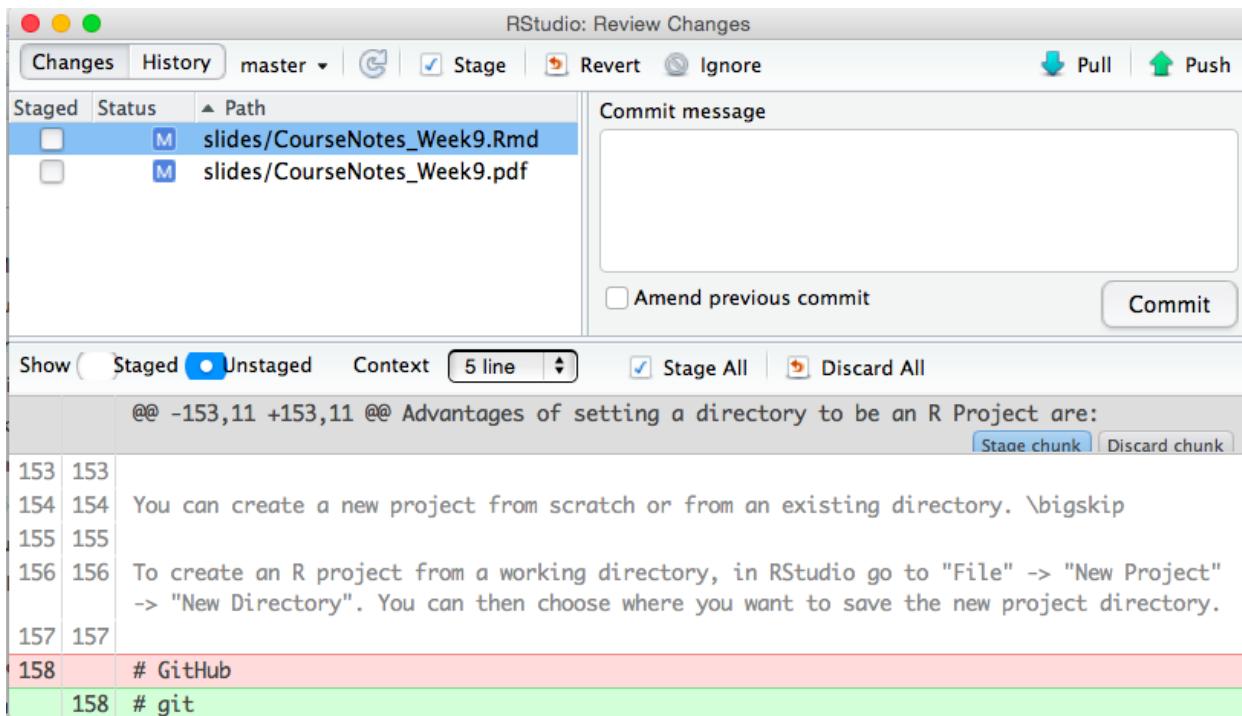


Example of a git window in RStudio when files in the repository have been changed and saved, but the changes haven’t yet been committed to git.

Committing

When you want git to record changes, you *commit* the files with the changes. Each time you commit, you have to include a short commit message with some information about the changes. You can make commits from a shell. However, the easiest workflow for an R project, including an R package directory, is to make git commits directly from the RStudio environment.

To make a commit from RStudio, click on the “Commit” button in the Git window. That will open a separate commit window that looks like Figure @ref(fig:examplecommitwindow).



Example of an RStudio commit window.

In this window, to commit changes:

1. Click on the boxes by the filenames in the top left panel to select the files to commit.
2. If you'd like, you can use the bottom part of the window to look through the changes you are committing in each file.
3. Write a message in the “Commit message” box in the top right panel. Keep the message to one line in this box if you can. If you need to explain more, write a short one-line message, skip a line, and then write a longer explanation.
4. Click on the “Commit” button on the right.

Once you commit changes to files, they will disappear from the Git window until you make and save more changes.

Browsing history

On the top left of the Commit window, you can toggle to “History”. This window allows you to explore the history of commits for the repository. Figure @ref(fig:examplehistorywindow) shows an example of this window. The top part of this window lists commits to the repository, from most recent to least. The commit message and author are shown for each commit. If you click on a commit, you can use the bottom panel to look through the changes made to that file with a specific commit.

| Subject | Author | Date | SHA |
|--|----------------------------------|------------|----------|
| HEAD origin/master origin/HEAD master Merge branch 'master' from 'origin/master' | Brooke Anderson <brooke.ander> | 2016-10-18 | a4acd950 |
| Add homework #4-5 | Brooke Anderson <brooke.ander> | 2016-10-18 | f5177dc7 |
| small ch. 9 edits | Rachel Severson <rachel.severso> | 2016-10-18 | e93ee9d2 |
| chapter eight edits | Rachel Severson <rachel.severso> | 2016-10-18 | f1d194ce |
| edits to chapter 7 | Rachel Severson <rachel.severso> | 2016-10-18 | 37151ee1 |

Commits 1-100 of 306

```

homework.Rmd
homework.Rmd
@@ -126,18 +126,21 @@ If you need them, here are some further tips:
126 126
127 127 ## Homework #4
128 128
129 **Due date: Oct. 26**
129 **Optional due date: Oct. 28**
130 130 Advanced R Markdown homework assignment
131

```

Example of the History window for exploring git commit history in RStudio.

Linking local repo to GitHub repo

GitHub allows you to host git repositories online. This allows you to:

- Work collaboratively on a shared repository
- Fork someone else's repository to create your own copy that you can use and change as you want
- Suggest changes to other people's repositories through pull requests

To do any of this, you will need a GitHub account. You can sign up at <https://github.com>. A free account is fine as long as you don't mind all of your repositories being "Public" (viewable by anyone).

The basic unit for working in GitHub is the repository. A repository is a directory of files with some supplemental files saving some additional information about the directory. While R Projects have this additional information saved as an ".RProj" file, git repositories have this information in a directory called ".git".



Because this pathname of the .git directory starts with a dot, it won't show up in many of the ways you list files in a directory. From a bash shell, you can see files that start with . by running `ls -a` from within that directory.

If you have a local directory that you would like to push to GitHub, these are the steps to do it. First, you need to make sure that the directory is under git version control. See the previous notes on initializing a repository. Next, you need to create an empty repository on GitHub to sync with your local repository. To do that:

1. In GitHub, click on the “+” in the upper right corner (“Create new”).
2. Choose “Create new repository”.
3. Give your repository the same name as the local directory you’d like to connect it to.
For example, if you want to connect it to a directory called “example_analysis” on your computer, name the repository “example_analysis”. (It is not required for your GitHub repository name to be identical to your local repository name, but it will make things easier.)
4. Leave everything else as-is (unless you’d like to add a short description in the “Description” box). Click on “Create repository” at the bottom of the page.

Now you are ready to connect the two repositories. First, you should change some settings in RStudio so GitHub will recognize that your computer can be trusted, rather than asking for you password every time. Do this by adding an SSH key from RStudio to your GitHub account with the following steps:

- In RStudio, go to “RStudio” -> “Preferences” -> “Git / svn”. Choose to “Create RSA key”.
- Click on “View public key”. Copy everything that shows up.
- Go to your GitHub account and navigate to “Settings”. Click on “SSH and GPG keys”.
- Click on “New SSH key”. Name the key something like “mylaptop”. Paste in your public key in the “Key box”.

Syncing RStudio and GitHub

Now you’re ready to push your local repository to the empty GitHub repository you created.

1. Open a shell and navigate to the directory you want to push. (You can open a shell from RStudio using the gear button in the Git window.)
2. Add the GitHub repository as a remote branch with the following command (this gives an example for adding a GitHub repository named “ex_repo” in my GitHub account, “geanders”):
`git remote add origin git@github.com:geanders/ex_repo.git` As a note, when you create a repository in GitHub, GitHub will provide suggested git code for adding the GitHub repository as the “origin” remote branch to a repository. That code is similar to the code shown above, but it uses “<https://github.com>” rather than “<git@github.com>”; the latter tends to work better with RStudio.
3. Push the contents of the local repository to the GitHub repository.
`git push -u origin master`

To pull a repository that already exists on GitHub and to which you have access (or that you’ve forked and so have access to the forked branch), first use `cd` from a bash shell on your personal computer to move into the directory where you want to put the repository. Then, use the `git clone` function to clone the repository locally. For example, to clone a GitHub repository called “ex_repo” posted in a GitHub account with the user name `janedoe`, you could run:

```
git clone git@github.com:janedoe/ex_repo.git
```

Once you have linked a local R project with a GitHub repository, you can push and pull commits using the blue down arrow (pull from GitHub) and green up arrow (push to GitHub) in the Git window in RStudio (see Figure @ref(fig:examplegitwindow) to see examples of these arrows).

GitHub helps you work with others on code. There are two main ways you can do this:

- **Collaborating:** Different people have the ability to push and pull directly to and from the same repository. When one person pushes a change to the repository, other collaborators can immediately get the changes by pulling the latest GitHub commits to their local repository.
- **Forking:** Different people have their own GitHub repositories, with each linked to their own local repository. When a person pushes changes to GitHub, it only makes changes to his own repository. The person must issue a pull request to another person's fork of the repository to share the changes.

Issues

Each original GitHub repository (i.e., not a fork of another repository) has a tab for “Issues”. This page works like a Discussion Forum. You can create new “Issue” threads to describe and discuss things that you want to change about the repository.

Issues can be closed once the problem has been resolved. You can close issues on the “Issue” page with the “Close issue” button. If a commit you make in RStudio closes an issue, you can automatically close the issue on GitHub by including “Close #[issue number]” in your commit message and then pushing to GitHub. For example, if issue #5 is “Fix typo in section 3”, and you make a change to fix that typo, you could make and save the change locally, commit that change with the commit message “Close #5”, and then push to GitHub, and issue #5 in “Issues” for that GitHub repository will automatically be closed, with a link to the commit that fixed the issue.

Pull request

You can use a *pull request* to suggest changes to a repository that you do not own or otherwise have the permission to directly change. Take the following steps to suggest changes to someone else’s repository:

1. Fork the repository
2. Make changes (locally or on GitHub)
3. Save your changes and commit them
4. Submit a pull request to the original repository
5. If there are not any conflicts and the owner of the original repository likes your changes, he or she can merge them directly into the original repository. If there are conflicts, these need to be resolved before the pull request can be merged.

You can also use pull requests within your own repositories. Some people will create a pull request every time they have a big issue they want to fix in one of their repositories.

In GitHub, each repository has a “Pull requests” tab where you can manage pull requests (submit a pull request to another fork or merge in someone else’s pull request for your fork).

Merge conflicts

At some point, if you are using GitHub to collaborate on code, you will get *merge conflicts*. These happen when two people have changed the same piece of code in two different ways at the same time.

For example, say two people are both working on local versions of the same repository, and the first person changes a line to `mtcars[1,]` while the second person changes the same line to `head(mtcars, 1)`. The second person pushes his commits to the GitHub version of the repository before the first person does. Now, when the first person pulls the latest commits to the GitHub repository, he will have a merge conflict for this line. To be able to commit a final version, the first person will need to decide which version of the code to use and commit a version of the file with that code.

If there are merge conflicts, they’ll show up like this in the file:

```
<<<<< HEAD
mtcars[1, ]
=====
head(mtcars, 1)
>>>> remote-branch
```

To fix them, search for all these spots in files with conflicts (Ctrl-F can be useful for this), pick the code you want to use, and delete everything else. For the example conflict, it could be resolved by changing the file from this:

```
<<<<< HEAD
mtcars[1, ]
=====
head(mtcars, 1)
>>>> remote-branch
```

To this:

```
head(mtcars, 1)
```

That merge conflict is now resolved. Once you resolve all merge conflicts in all files in the repository, you can save and commit the files.

These merge conflicts can come up in a few situations:

- You pull in commits from the GitHub branch of a repository you've been working on locally.
- Someone sends a pull request for one of your repositories, and you have updated some of the code between when the person forked the repository and submitted the pull request.

Summary

R code can be kept under version control using git, and RStudio offers convenient functionality for working with a directory under git version control. A directory under git version control can also be pushed to GitHub, which provides a useful platform for sharing and collaborating on code.

3.10 Software Design and Philosophy

Writing and designing software is a creative endeavor and like in other creative arts there are styles and guidelines that you can follow, however revolutions in the field can occur when those dogmas are broken properly. We're going to cover a few of the prominent ideas in software design in the last century. Above all of these suggestions I suggest one cardinal rule: Have empathy for your fellow human beings. Software is inherently complex, so set up your users to [fall into a pit of success](#).

The Unix Philosophy

The R programming language is open source software and many open source software packages draw some inspiration from the design of the Unix operating system which macOS and Linux are based on. [Ken Thompson](#) - one of the designers of Unix - first laid out this philosophy, and many Unix philosophy principles can be applied to R programs. The overarching philosophical theme of Unix programs is to **do one thing well**. Sticking to this rule accomplishes several objectives:

1. Since your program only does one thing the chance that your program contains many lines of code is reduced. This means that others can more easily read the code for your program so they can understand exactly how it works (if they need to know).
2. Simplicity in your program reduces the chance there will be major bugs in your program since fewer lines of code means fewer opportunities to make a mistake.
3. Your program will be easier for users to understand since the number of inputs and outputs are reduced for a program that only does one thing.
4. Programs built with other small programs have a higher chance of also being small. This ability to string several small programs together to make a more complex (but also small) program is called **composability**.

Unix command line programs are notable for their use of the pipe operator (`|`) and so the Unix philosophy also encourages programs to produce outputs that can be piped into

program inputs. Recently pipes in R have surged in popularity thanks to projects like the [magrittr](#) package. When it makes sense for your function to take data (usually a vector or a data frame) as an argument and then return data, you should consider making the data argument the first argument in your function so that your function can be part of a data pipeline.

One case where many R programs differ from the greater Unix philosophy is in terms of user interaction. Unix programs will usually only print a message to the user if a program produces an error or warning. Although this is a good guideline for your programs, many R programs print messages to the console even if the program works correctly. Many R users only use the language interactively, so showing messages to your users might make sense for your package. One issue with messages is that they produce output which is separate from the results of your program, and therefore messages are harder to capture.

Default Values

Every function argument is an opportunity for your function to fail the user by producing an error because of bad or unexpected inputs. Therefore you should provide as many default values for your functions as is reasonable. If there's an argument in your function that should only be one of a handful of values you should use the `match.arg()` function to check that one of the permitted values is provided:

```
multiply_by <- function(n, multiplier = c("two", "three", "four")){
  multiplier <- match.arg(multiplier)
  if(multiplier == "two"){
    n * 2
  } else if(multiplier == "three"){
    n * 3
  } else {
    n * 4
  }
}

multiply_by(5, "two")
[1] 10
multiply_by(5, "six")
Error in match.arg(multiplier): 'arg' should be one of "two", "three", "four"
```

Using `match.arg()` ensures that an error is thrown immediately if an erroneous argument value is provided.

Naming Things

Naming functions and variables is a challenge that programmers have always struggled with. Here are a few strategies you should use when naming things in R:

1. Use snake case and lowercase. Modern R packages use function and variable names like `geom_line()`, `bind_rows()`, and `unnest_token()` where words are separated by underscores (`_`) and all characters are lowercase. Once upon a time words were commonly separated by periods (`.`) but that scheme can cause confusion with regard to generic functions (see the object oriented programming chapter for more information).
2. Names should be short. A short name is faster to type and is more memorable than a long and complicated name. The length of a variable name has to be balanced with the fact that:
3. Names should be meaningful and descriptive. Function names should generally describe the actions they perform. Other object names should describe the data or attributes they encompass. In general you should avoid numbering variable names like `apple1`, `apple2`, and `apple3`. Instead you should create a data structure called `apples` so you can access each apple with `apple[[1]]`, `apple[[2]]`, and `apple[[3]]`.
4. Be sure that you're not assigning names that already exist and are common in R. For example `mean`, `summary`, and `rt` are already names of functions in R, so try to avoid overwriting them. You can check if a name is taken using the `apropos()` function:

```
apropos("mean")
[1] ".colMeans"      ".rowMeans"      "colMeans"       "kmeans"
[5] "mean"           "mean.Date"     "mean.default"   "mean.difftime"
[9] "mean.POSIXct"   "mean.POSIXlt"  "rowMeans"       "weighted.mean"
apropos("my_new_function")
character(0)
```

1. You might want to consider grouping similar functions together in families which all start with the same short prefix. For example in the `ggplot2` package the `aes_` family of functions set graphing aesthetics, the `gs_` family of functions interact with the Google Sheets API in the `googlesheets` package, and the `wq_` family of functions all write questions in the `swirlify` package.

Playing Well With Others

If you write a package with useful functions that are well designed then you may be lucky enough that your package becomes popular! Others may build upon your functions to extend or adapt thier features for other purposes. This means that when you establish a set of arguments for a function you're implicitly promising some amount of stability for the inputs and outputs of that function. Changing the order or the nature of function arguments or return values can break other people's code, creating work and causing pain for those who have chosen to use your software. For this reason you should think very carefully about function arguments and outputs to ensure that both can grow and change sustainably. You should seek to strike a balance between frustrating your users by making breaking changes and ensuring that your package follows up to date programming patterns and ideas. If you believe that the functions in a package you're developing are not yet stable you should make users aware of that fact so that they're warned if they choose to build on your work.

Summary

Most of software design is ensuring that your users stumble into their desired outcome. You may think you're writing the most intuitive package, but sitting down with a colleague and watching them use your package can teach you volumes about what users want and expect out of your package. There are libraries full of books written about software design and this chapter is only meant to serve as a jumping off point. If you happen to be looking for inspiration I highly recommend this talk Bret Victor called: *The Future of Programming*.

3.11 Continuous Integration

The objectives of this section are:

- Create an R package that is tested and deployed on Travis
- Create an R package that is tested and deployed on Appveyor

In modern software companies hundreds of people are simultaneously working on the source code of the same product while they develop different features for that product. At the same time those programmers are depending upon software that might be built by other teams within the company, or they may be using software built by other companies or individuals, which in turn is being actively developed and updated. The software development technique of continuous integration was developed to ensure that all of the components in this web of software are working together harmoniously.

R packages are usually not as big in terms of lines of code compared to software like Google's search engine, however it's plausible that your package may depend on several other packages which you want to make sure are still working the way you expected them to when you first included them in your code. When it comes to R packages continuous integration means ensuring that your package builds without any errors or warnings, and making sure that all of the tests that you've written for your package are passing. Building your R package will protect you against some big errors, but the best way that you can ensure continuous integration will be useful to you is if you build robust and complete tests for every function in your package.

Web Services for Continuous Integration

We'll discuss two services for continuous integration: the first is [Travis](#) which will test your package on Linux, and then there's [AppVeyor](#) which will test your package on Windows. Both of these services are free for R packages that are built in public GitHub repositories. These continuous integration services will run every time you push a new set of commits for your package repository. Both services integrate nicely with GitHub so you can see in GitHub's pull request pages whether or not your package is building correctly.

Using Travis

To start using Travis go to <https://travis-ci.org> and sign in with your GitHub account. Clicking on your name in the upper right hand corner of the site will bring up a list of your public GitHub repositories with a switch next to each repo. If you turn the switch on then the next time you push to that repository Travis will look for a `.travis.yml` file in the root of the repository, and it will run tests on your package accordingly.

Open up your R console and navigate to your R package repository. Now load the `devtools` package with `library(devtools)` and enter `use_travis()` into your R console. This command will set up a basic `.travis.yml` for your R package. You can now add, commit, and push your changes to GitHub, which will trigger the first build of your package on Travis. Go back to <https://travis-ci.org> to watch your package be built and tested at the same time! You may want to make some changes to your `.travis.yml` file, and you can see all of the options available in [this guide](#).

Once your package has been built for the first time you'll be able to obtain a badge, which is just a small image generated by Travis which indicates whether your package is building properly and passing all of your tests. You should display this badge in the `README.md` file of your package's GitHub repository so that you and others can monitor the build status of your package.

Using AppVeyor

You can start using AppVeyor by going to <https://www.appveyor.com/> and signing in with your GitHub account. After signing in click on “Projects” in the top navigation bar. If you have any GitHub repositories that use AppVeyor you'll be able to see them here. To add a new project click “New Project” and find the GitHub repo that corresponds to the R package you'd like to test on Windows. Click “Add” for AppVeyor to start tracking this repo.

Open up your R console and navigate to your R package repository. Now load the `devtools` package with `library(devtools)` and enter `use_appveyor()` into your R console. This command will set up a default `appveyor.yml` for your R package. You can now add, commit, and push your changes to GitHub, which will trigger the first build of your package on AppVeyor. Go back to <https://www.appveyor.com/> to see the result of the build. You may want to make some changes to your `appveyor.yml` file, and you can see all of the options available in the [r-appveyor guide](#) which is maintained by Kirill Müller. Like Travis, AppVeyor also generates badges that you should add to the `README.md` file of your package's GitHub repository.

Summary

Continuous integration is a strategy for testing new features and changes to your package as often as possible. Web services like Travis and AppVeyor make it possible to re-test your code on different platforms after every `git push`. Using continuous integration makes it easy for you and for others to simultaneously work on building an R package without breaking package features by mistake.

3.12 Cross Platform Development

The objective of this section is:

- Recognize characteristics of R packages that are not cross-platform

One of the great features about R is that you can run R code on multiple kinds of computers and operating systems and it will behave the same way on each one. Most of time you don't need to worry about what platform your R code is running on. The following sections discuss strategies and functions that you should use to ensure that your R code runs uniformly on every kind of system.

Handling Paths

Paths to files and folders can have big differences between operating systems. In general you should avoid constructing a path "by hand." For example if I wanted to access a file called `data.txt` that I know will be located on the user's desktop using the string "`~/Desktop/data.txt`" would not work if that code was run on a Windows machine. In general you should always use functions to construct and find paths to files and folders. The correct programmatic way to construct the path above is to use the `file.path()` function. So to get the file above I would do the following:

```
file.path("~/Desktop", "data.txt")
[1] "~/Desktop/data.txt"
```

Note that this book is probably being built on a Mac:

```
Sys.info()["sysname"]
sysname
"Darwin"
```

If the resulting line above says "Darwin" it's referring to the [core of macOS](#). If you don't have a Mac try running both lines of code above to see the resulting path and the type of system that you're running.

In general it's not guaranteed on any system that a particular file or folder you've looking for will exist — however if the user of your package has installed your package you can be sure that any files within your package exist on their machine. You can find the path to files included in your package using the `system.file()` function. Any files or folders in the `inst/` directory of your package will be copied one level up once your package is installed. If your package is called `ggplyr2` and there's file in your package under `inst/data/first.txt` you can get the path to that file with `system.file("data", "first.txt", package = "ggplyr2")`. Packaging files with your package is the best way to ensure that users have access to them when they're using your package.

In terms of constructing paths there are a few other functions you should be aware of. Remember that the results for many of these functions are contingent on this book being built on a Mac, so if you're using Windows I encourage you to run these functions yourself to see their result. The `path.expand()` function is usually used to find the absolute path name of a user's home directory when the tilde (~) is included in the path. The tilde is a shortcut for the path to the current user's home directory. Let's take a look at `path.expand()` in action:

```
path.expand("~")
[1] "/Users/rdpeng"
path.expand(file.path("~", "Desktop"))
[1] "/Users/rdpeng/Desktop"
```

The `normalizePath()` function is built on top of `path.expand()`, so it includes `path.expand()`'s features but it also creates full paths for other shortcuts like `..` which signifies the current working directory and `...` which signifies the directory above the current working directory. Let's take a look at some examples:

```
normalizePath(file.path("~", "R"))
```

```
[1] "/Users/sean/R"
```

```
normalizePath("..")
```

```
[1] "/Users/sean/books/msdr"
```

```
normalizePath("...")
```

```
[1] "/Users/sean/books"
```

To extract parts of a path you can use the `basename()` function to get the name of the file or the deepest directory in the path and you can use `dirname()` to get the part of the path that does not include either the file or the deepest directory. Let's take a look at some examples:

```

data_file <- normalizePath(file.path("~/", "data.txt"))
data_file
[1] "/Users/rdpeng/data.txt"
dirname(data_file)
[1] "/Users/rdpeng"
dirname(dirname(data_file))
[1] "/Users"
basename(data_file)
[1] "data.txt"

```

Saving Files & rappidirs

[CRAN's policy for R packages](#) contains the following statement:

Packages should not write in the users' home filesystem, nor anywhere else on the file system apart from the R session's temporary directory (or during installation in the location pointed to by TMPDIR: and such usage should be cleaned up). Installing into the system's R installation (e.g., scripts to its bin directory) is not allowed. Limited exceptions may be allowed in interactive sessions if the package obtains confirmation from the user.

In general you should strive to get the user's consent before you create or save files on their computer. With some functions consent is implicit, for example it's clear somebody using `write.csv()` consents to producing a csv file at a specified path. When it's not absolutely clear that the user will be creating a file or folder when they use your functions you should ask them specifically. Take a look at the code below for a skeleton of a function that asks for a user's consent:

```

#' A function for doing something
#'
#' This function takes some action. It also attempts to create a file on your
#' desktop called \code{data.txt}. If \code{data.txt} cannot be created a
#' warning is raised.
#'
#' @param force If set to \code{TRUE}, \code{data.txt} will be created on the
#' user's Desktop if their Desktop exists. If this function is used in an
#' interactive session the user will be asked whether or not \code{data.txt}
#' should be created. The default value is \code{FALSE}.
#'
#' @export
some_function <- function(force = FALSE){

  #
  # ... some code that does something useful ...
  #

  if(!dir.exists(file.path("~/", "Desktop"))){
    warning("No Desktop found.")
  }
}

```

```
    } else {
      if(!force && interactive()){
        result <- select.list(c("Yes", "No"),
                              title = "May this program create data.txt on your desktop?")
        if(result == "Yes"){
          file.create(file.path("~/Desktop", "data.txt"))
        }
      } else if(force){
        file.create(file.path("~/Desktop", "data.txt"))
      } else {
        warning("data.txt was not created on the Desktop.")
      }
    }
  }
}
```

The `some_function()` function above is a contrived example of how to ask for permission from the user to create a file on their hard drive. Notice that the description of the function clearly states that the function attempts to create the `data.txt` file. This function has a `force` argument which will create the `data.txt` file without asking the user first. By setting `force = FALSE` as the default, the user must set `force = TRUE`, which is one method to get consent from the user. The function above uses the `interactive()` function in order to determine whether the user is using this function in an R console or if this function is being run in a non-interactive session. If the user is in an interactive R session then using `select.list()` is a decent method to ask the user a question. You should strive to use `select.list()` and `interactive()` together in order to prevent an R session from waiting for input from a user that doesn't exist.

rappdirs

Even the contrived example above implicitly raises a good question: where should your package save files? The most obvious answer is to allow the user to provide an argument for the path where a file should be saved. This is a good idea as long as your package won't need to depend on the location of that file in the future, for example if your package is creating an output data file. But what if you need persistent and consistent access to a file? You might be tempted to use `path.package()` in order to find the directory that your package is installed in so you can store files there. This isn't a good idea because file access permissions often do not allow users to modify files where R packages are stored.

In order to find a location where you can read and write files that will persist on a user's computer you should use the `rappdirs` package. This package contains functions that will return paths to directories where your package can store files for future use. The `user_data_dir()` function will provide a user-specific path for your package, while the `site_data_dir()` function will return a directory path that is shared by all users. Let's take a look at `rappdirs` in action:

```
library(rappdirs)
Loading required package: methods
site_data_dir(appname = "ggplyr2")
[1] "/Library/Application Support/ggplyr2"
user_data_dir(appname = "ggplyr2")
[1] "/Users/rdpeng/Library/Application Support/ggplyr2"
```

Both of the examples above are probably the Mac-specific paths. We can get the Windows specific paths by specifying the `os` argument:

```
user_data_dir(appname = "ggplyr2", os = "win")
[1] "C:/Users/<username>/Local/ggplyr2/ggplyr2"
```

If you don't supply the `os` argument then the function will determine the operating system automatically. One feature about `user_data_dir()` you should note is the `roaming = TRUE` argument. Many Windows networks are configured so that any authorized user can log in to any computer on the network and have access to their desktop, settings, and files. Setting `roaming = TRUE` returns a special path so that R will have access to your packages files everywhere, but this requires the directory to be synced often. Make sure to only use `roaming = TRUE` if the files your package will storing with `rappdirs` are going to be small. For more information about `rappdirs` see <https://github.com/hadley/rappdirs>.

Options and Starting R

Several R Packages allow users to set global options that affect the behavior of the package using the `options()` function. The `options()` function returns a list, and named values in this list can be set using the following syntax: `options(key = value)`. It's a common feature for packages to allow a user to set options which may specify package defaults, or change the behavior of the package in some way. You should thoroughly document how your package is effected by which options are set.

When an R session begins a series of files are searched for and run if found as detailed in `help("Startup")`. One of those files is `.Rprofile`. The `.Rprofile` file is just a regular R file which is usually located in a user's home directory (which you can find with `normalizePath("~/")`). A user's `.Rprofile` is run every time they start an R session, so it's a good file for setting options that a user wants to be set when using R. If you want a user to be able to set an option that is related to your package that is unlikely to change (like a username or a key), then you should consider instructing them to create or make changes to their `.Rprofile`.

Package Installation

Your package documentation should prominently feature installation instructions. Many R packages that are distributed through GitHub recommend installing the `devtools` package, and then using `devtools::install_github()` to install the package. The `devtools` package is wonderful for developing R packages, but it has many dependencies which can make it

difficult for users to install. I recommend instructing folks to use the `ghit` package by [Thomas Leeper](#) and the `ghit::install_github()` function as a reliable alternative to `devtools`.

In cases where users might have a weak internet connection it's often easier for a user to download the source of your package as a zip file and then to install it using `install.packages()`. Instead of asking users to discern the path of zip file they've downloaded you should ask them to enter `install.packages(file.choose(), repos = NULL, type = "source")` into the R console and then they can interactively select the file they just downloaded. If a user is denied permission to modify their local package directory, they still may be able to use a package if they specify a directory they have access to with the `lib` argument for `install.packages()`.

Environmental Attributes

Occasionally you may need to know specific information about the hardware and software limitations of the computer that is running your R code. The environmental variables `.Platform` and `.Machine` are lists which contain named elements that can tell your program about the underlying machine. For example `.Platform$OS.type` is a good method for checking whether your program is in a Windows environment since the only values it can return are "windows" and "unix":

```
.Platform$OS.type  
[1] "unix"
```

For more information about information contained in `.Platform` see the help file: `help(".Platform")`.

The `.Machine` variable contains information specific to the computer architecture that your program is being run on. For example `.Machine$double.xmax` and `.Machine$double.xmin` are respectively the largest and smallest positive numbers that can be represented in R on your platform:

```
.Machine$double.xmax  
[1] 1.797693e+308  
.Machine$double.xmax + 100 == .Machine$double.xmax  
[1] TRUE  
.Machine$double.xmin  
[1] 2.225074e-308
```

You might also find `.Machine$double.eps` useful, which is the smallest number on a machine such that `1 + .Machine$double.eps != 1` evaluates to TRUE:

```
1 + .Machine$double.eps != 1  
[1] TRUE  
1 + .Machine$double.xmin != 1  
[1] FALSE
```

Summary

File and folder paths differ across platforms so R provides several functions to ensure that your program can construct paths correctly. The `rappdirs` package helps further by identifying locations where you can safely store files that your package can access. However before creating files anywhere on a user's disk you should always ask the user's permission. You should provide clear and easy instructions so people can easily install your package. The `.Platform` and `.Machine` variables can inform your program about hardware and software details.

4. Building Data Visualization Tools

The data science revolution has produced reams of new data from a wide variety of new sources. These new datasets are being used to answer new questions in ways never before conceived. Visualization remains one of the most powerful ways draw conclusions from data, but the influx of new data types requires the development of new visualization techniques and building blocks. This section provides you with the skills for creating those new visualization building blocks. We focus on the `ggplot2` framework and describe how to use and extend the system to suit the specific needs of your organization or team.

The objectives for this section are:

- Create data graphics using the `ggplot2` package
- Build graphics by combining multiple geoms
- Recognize the differences between different `geom_*` functions
- Recall the difference between mapping an aesthetic to a constant and mapping an aesthetic to a variable
- Create a scatterplot or a histogram using `ggplot2`
- Build plots using the six guidelines for good graphics
- Create simple maps using `ggplot2`
- Create dynamic maps using the `ggmap` package
- Build maps that have external data overlaid
- Create chloropleth maps of US counties
- Recognize the different graphical objects presented by the `grid` package
- Build build simple graphics using the `grid` package
- Create a `ggplot2` theme by modifying an existing theme
- Build a new geom function to implement a new feature or simplify a workflow

4.1 Basic Plotting With `ggplot2`

The `ggplot2` package allows you to quickly plot attractive graphics, to visualize and explore data. Objects created with `ggplot2` can also be extensively customized (more on that in the next subsection). While the structure of `ggplot2` code differs substantially from that of base R graphics, it offers a lot of power for the required effort. This subsection focuses on **useful**, rather than **attractive** graphs, since this subsection focuses on exploring rather than presenting data. Later sections will give more information about making more attractive or customized plots, as you'd want to do for final reports, papers, etc.

To show how to use basic `ggplot2`, we'll use a dataset of Titanic passengers, their characteristics, and whether or not they survived the sinking. This dataset has become fairly famous in

data science, because it's used, among other things, for one of Kaggle's long-term "learning" competitions, as well as in many tutorials and texts on building classification models.



Kaggle is a company that runs predictive modeling competitions, often sponsored by companies, with top competitors sometimes winning cash prizes or interviews at top companies. At any time, Kaggle is typically hosting several competitions, including some with no cash reward that are offered to help users get started with predictive modeling.

To get this dataset, you'll need to install and load the `titanic` package, and then you can load and rename the training datasets, which includes data on about two-thirds of the Titanic passengers:

```
# install.packages("titanic") # If you don't have the package installed
library(titanic)
data("titanic_train", package = "titanic")
titanic <- titanic_train
```

The other data example we'll use in this subsection is some data on players in the 2010 World Cup. This is available from the `faraway` package:

```
# install.packages("faraway") # If you don't have the package installed
library(faraway)
data("worldcup")
```



Unlike most data objects you'll work with, the data that comes with an R package will often have its own help file. You can access this using the `?operator`. For example, try running: `?worldcup`.

All of the plots we'll make today will use the `ggplot2` package (another member of the `tidyverse!`). If you don't already have that installed, you'll need to install it. You then need to load the package in your current session of R:

```
# install.packages("ggplot2") ## Uncomment and run if you don't have `ggplot2` installed
library(ggplot2)
```

The process of creating a plot using `ggplot2` follows conventions that are a bit different than most of the code you've seen so far in R (although it is somewhat similar to the idea of piping we introduced in an earlier course). The basic steps behind creating a plot with `ggplot2` are:

1. Create an object of the `ggplot` class, typically specifying the `data` and some or all of the `aesthetics`;
2. Add on `geoms` and other elements to create and customize the plot, using `+`.

You can add on one or many geoms and other elements to create plots that range from very simple to very customized. We'll focus on simple geoms and added elements, and then explore more detailed customization later.



If R gets to the end of a line and there is not some indication that the call is not over (e.g., `%>%` for piping or `+` for `ggplot2` plots), R interprets that as a message to run the call without reading in further code. A common error when writing `ggplot2` code is to put the `+` to add a geom or element at the beginning of a line rather than the end of a previous line— in this case, R will try to execute the call too soon. To avoid errors, be sure to end lines with `+`, don't start lines with it.

Initializing a `ggplot` object

The first step in creating a plot using `ggplot2` is to create a `ggplot` object. This object will not, by itself, create a plot with anything in it. Instead, it typically specifies the data frame you want to use and which aesthetics will be mapped to certain columns of that data frame (aesthetics are explained more in the next subsection).

Use the following conventions to initialize a `ggplot` object:

```
## Generic code
object <- ggplot(dataframe, aes(x = column_1, y = column_2))
## or, if you don't need to save the object
ggplot(dataframe, aes(x = column_1, y = column_2))
```

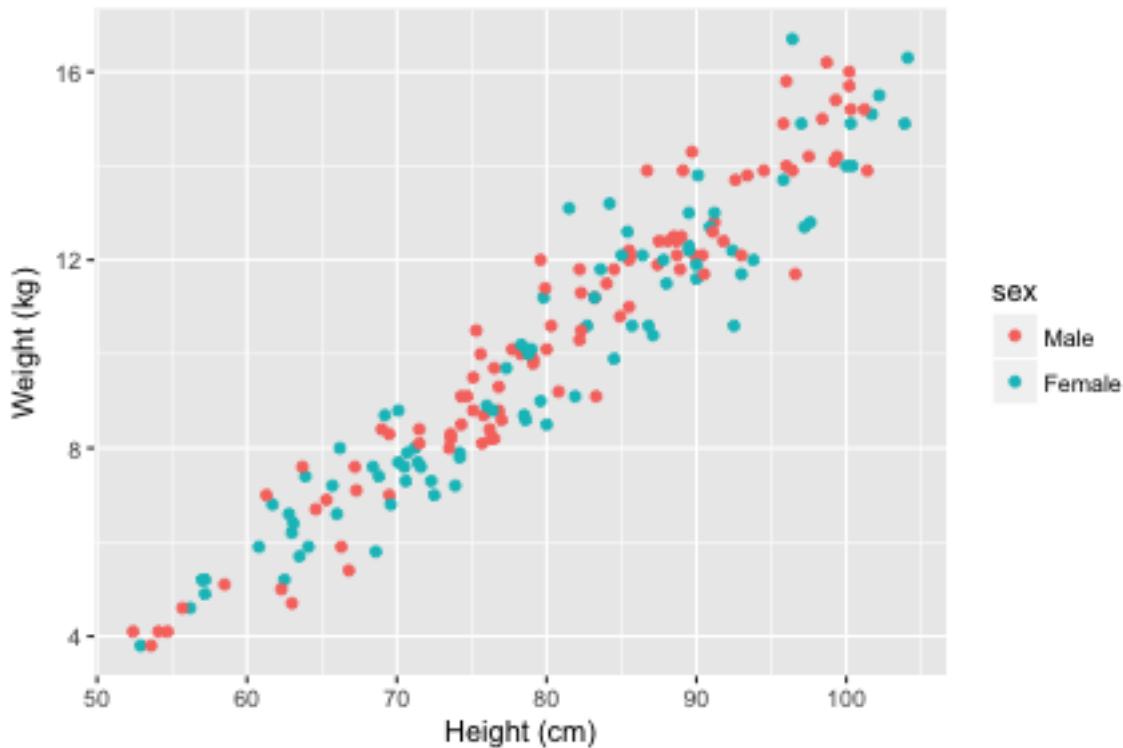
The `dataframe` is the first parameter in a `ggplot` call and, if you like, you can use the parameter definition with that call (e.g., `data = dataframe`). Aesthetics are defined within an `aes` function call that typically is used within the `ggplot` call.



In `ggplot2`, life is much easier if everything you want to plot is included in a `dataframe` as a column, and the first argument to `ggplot` must be a `dataframe`. This format has been a bit hard for some base R graphics users to adjust to, since base R graphics tends to plot based on vector, rather than `dataframe`, inputs. Trying to pass in a vector rather than a `dataframe` can be a common reason for `ggplot2` errors for all R users.

Plot aesthetics

Aesthetics are properties of the plot that can show certain elements of the data. For example, in Figure @ref(fig:aesmapex), color shows (i.e., is mapped to) gender, x-position shows height, and y-position shows weight in a sample data set of measurements of children in Nepal.



Example of how different properties of a plot can show different elements to the data. Here, color indicates gender, position along the x-axis shows height, and position along the y-axis shows weight. This example is a subset of data from the `nepali` dataset in the `faraway` package.



Any of these aesthetics could also be given a constant value, instead of being mapped to an element of the data. For example, all the points could be red, instead of showing gender. Later in this section, we will describe how to use these constant values for aesthetics. We'll discuss how to code this later in this section.

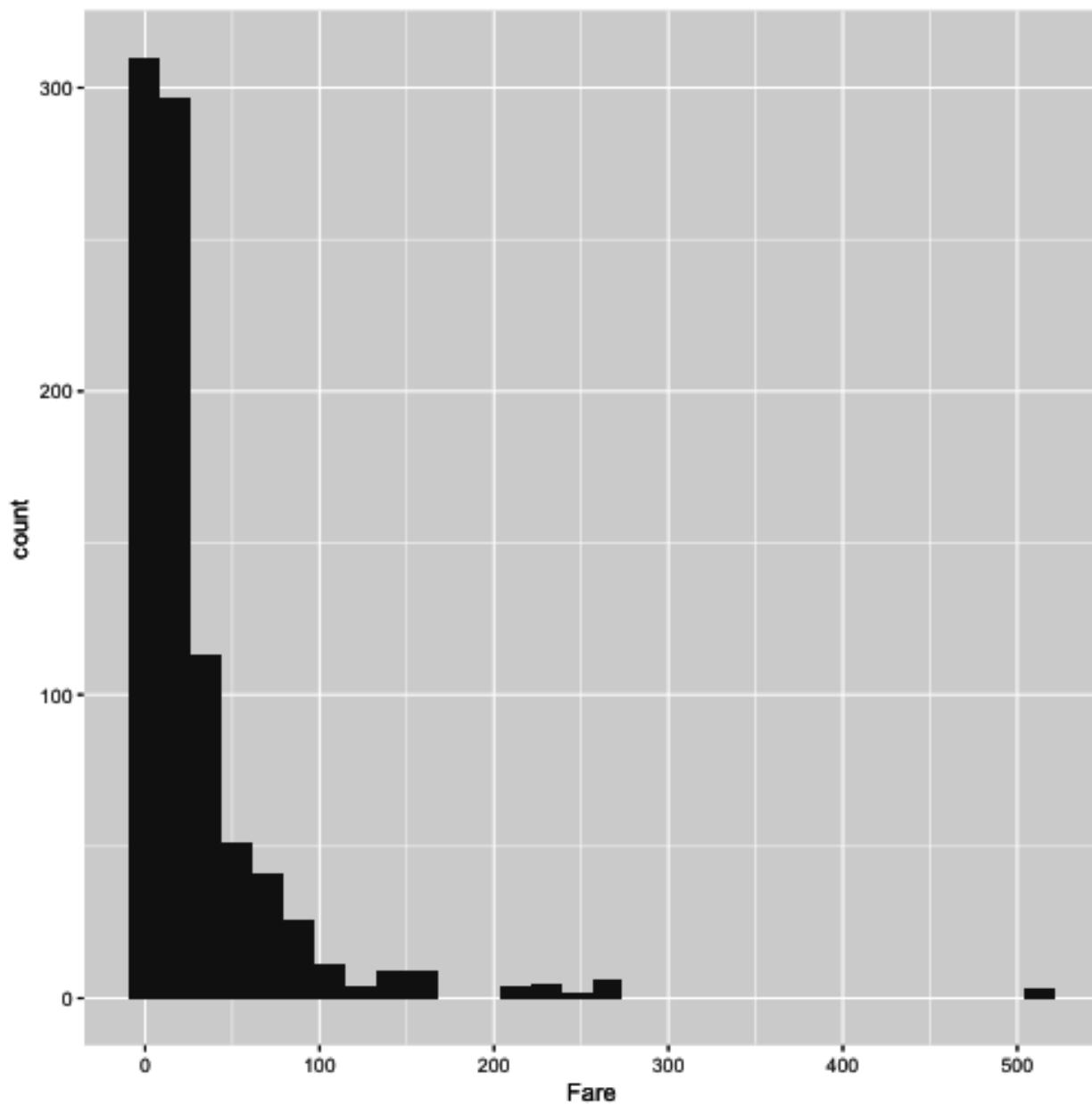
Which aesthetics are required for a plot depend on which geoms (more on those in a second) you're adding to the plot. You can find out the aesthetics you can use for a geom in the “Aesthetics” section of the geom’s help file (e.g., `?geom_point`). Required aesthetics are in bold in this section of the help file and optional ones are not. Common plot aesthetics you might want to specify include:

| Code | Description |
|-----------------------|--|
| <code>x</code> | Position on x-axis |
| <code>y</code> | Position on y-axis |
| <code>shape</code> | Shape |
| <code>color</code> | Color of border of elements |
| <code>fill</code> | Color of inside of elements |
| <code>size</code> | Size |
| <code>alpha</code> | Transparency (1: opaque; 0: transparent) |
| <code>linetype</code> | Type of line (e.g., solid, dashed) |

Creating a basic ggplot plot

The system of creating a `ggplot` object, mapping aesthetics to columns of the data, and adding geoms makes more sense once you try a few plots. For example, say you'd like to create a histogram showing the fares shown by passengers in the example Titanic data set. To plot the histogram, you'll first need to create a `ggplot` object using the dataframe with the column you want to print. In creating this `ggplot` object, you only need one aesthetic (`x`, the variable for which you want to create the histogram), and then you'll need to add a histogram geom. In code, this will be:

```
ggplot(data = titanic, aes(x = Fare)) +  
  geom_histogram()
```



plot of chunk unnamed-chunk-7

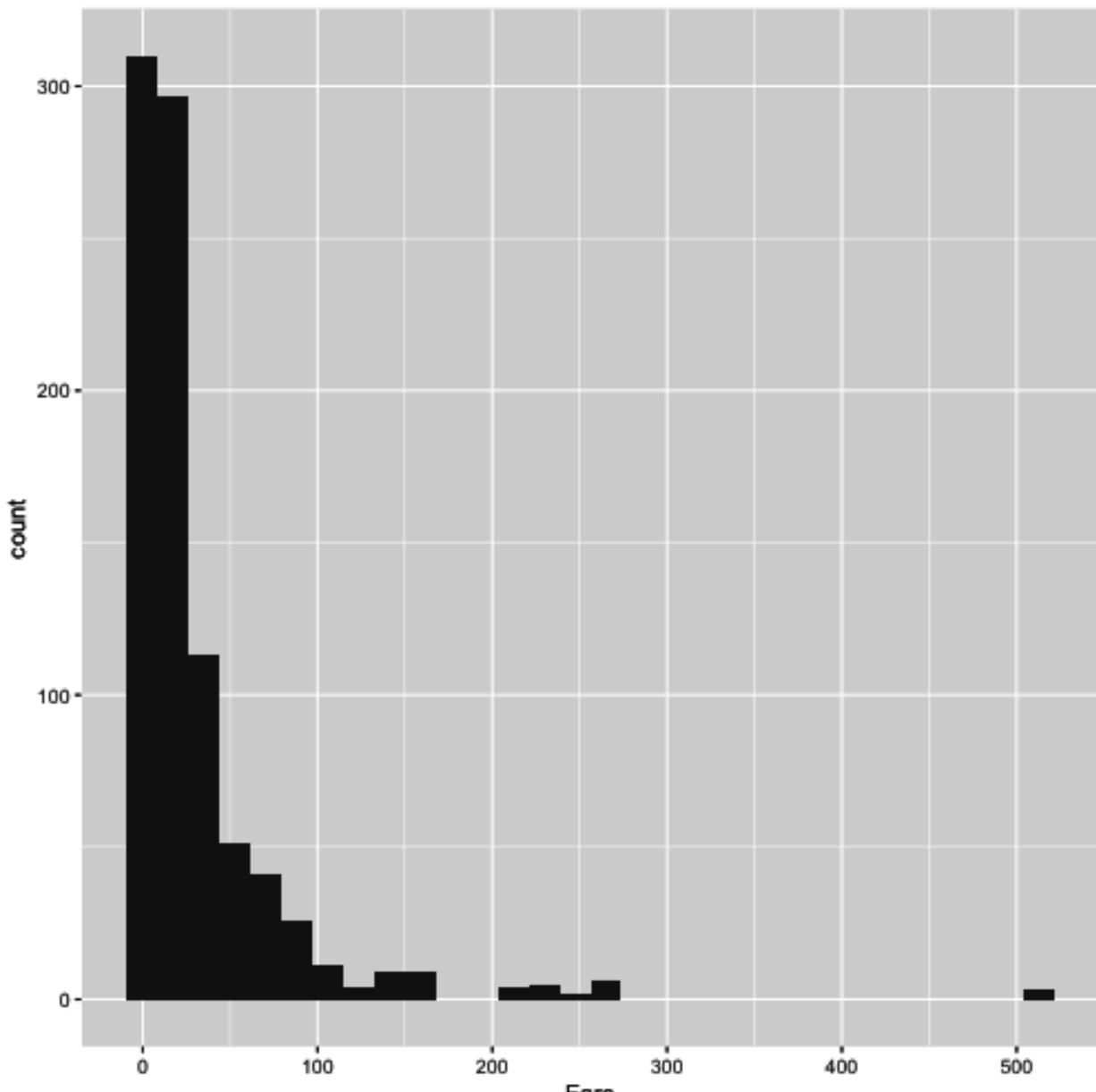
This code sets the dataframe as the `titanic` object in the user's working session, maps the values in the `Fare` column to the `x` aesthetic, and adds a histogram geom to generate a histogram.

There is some flexibility in writing the code to create this plot. For example, since the aesthetic mapping (showing `Fare` by position on the x-axis) only applies to one geom (`geom_histogram`), you could specify that aesthetic in an `aes` statement when adding the geom:

```
ggplot(data = titanic) +  
  geom_histogram(aes(x = Fare))
```

Similarly, you could specify the dataframe when adding the geom:

```
ggplot() +  
  geom_histogram(data = titanic, aes(x = Fare))  
`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



A basic ggplot plot

Finally, you can pipe your data into a `ggplot` call, since the `ggplot` function takes a dataframe as its first argument:

```
titanic %>%
  ggplot() +
  geom_histogram(aes(x = Fare))
# or
titanic %>%
  ggplot(aes(x = Fare)) +
  geom_histogram()
```

While all of these work, for simplicity we will use the syntax of specifying the data and aesthetics in the `ggplot` call for most examples in this subsection. Later, we'll show how this flexibility can be used to do things like use data from a different dataframe for some geoms or change aesthetics mappings between geoms.

A key thing to remember, however, is that `ggplot` is **not** flexible about whether you specify aesthetics within an `aes` call or not. We will discuss what happens if you do not later in the book, but it is very important that if you want to show values from a column of the data using aesthetics like color, size, shape, or position, you remember to make that specification within `aes`. Also, be sure that you specify the dataframe when or before you specify aesthetics (i.e., you can't specify aesthetics in the `ggplot` statement if you haven't specified the dataframe yet), and if you specify a dataframe within a geom, be sure to use `data =` syntax rather than relying on parameter position, as `data` is not the first parameter expected for geom functions.

Geoms

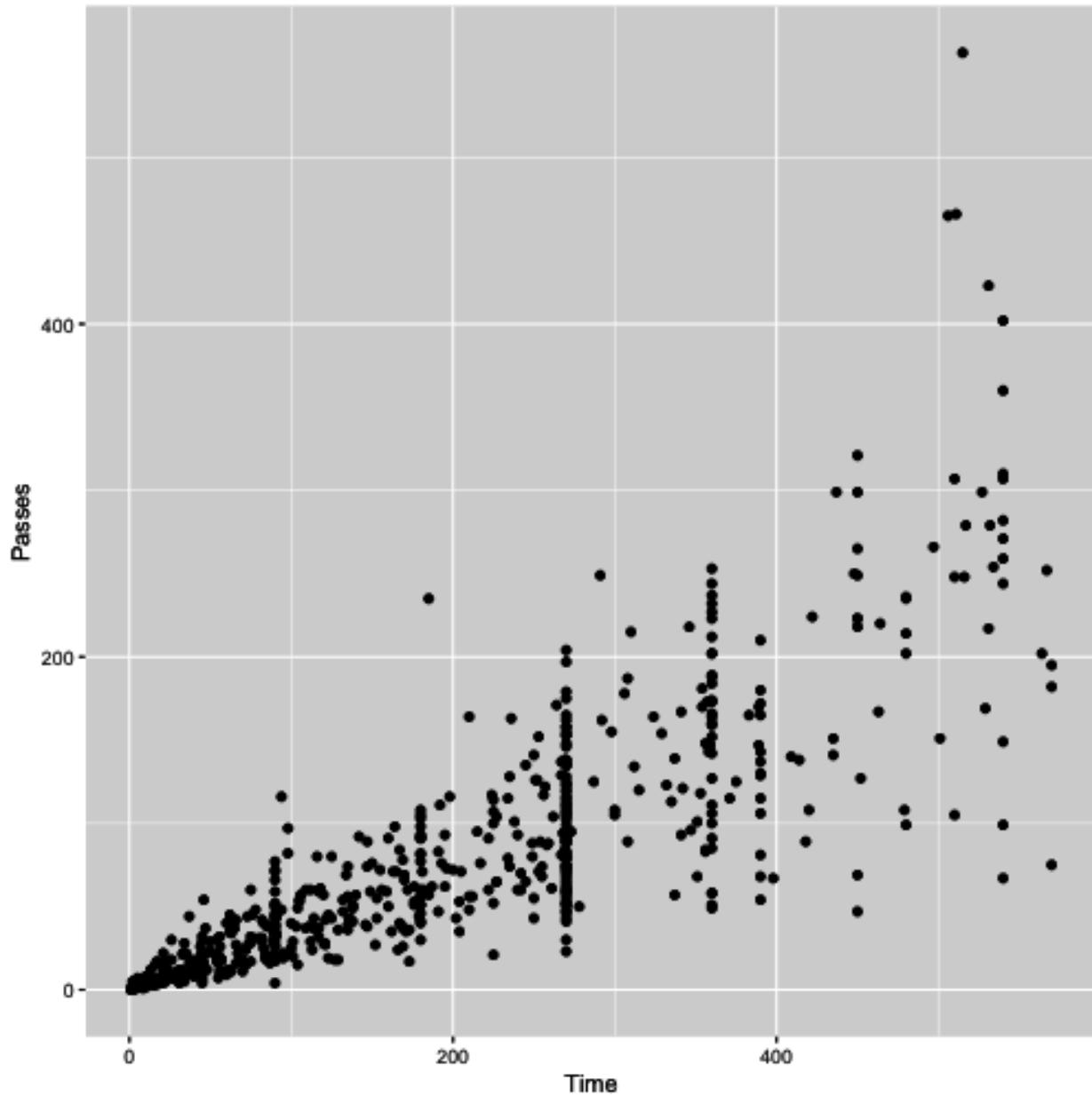
Geom functions actually plot elements of the plot; if you do not include at least one geom, you'll get a blank plot space. Geom functions have their own arguments to adjust how the graph is created. For example, you can use the `bins` argument to change the number of bins used to create a histogram—try:

```
ggplot(titanic, aes(x = Fare)) +
  geom_histogram(bins = 15)
```

As with any R functions, you can find out more about which arguments are available for geom functions by pulling up the function's help file (e.g., `?geom_histogram`).

Different geoms require different aesthetic inputs. For example, which `geom_histogram` only requires a single aesthetic (`x`). If you want to create a scatterplot, you'll need two aesthetics, `x` and `y`. In the `worldcup` dataset, the `Time` column gives the amount of time each player player and the `Passes` column gives the number of passes they made. To see the relationship between these two variables, you can run:

```
ggplot(worldcup, aes(x = Time, y = Passes)) +  
  geom_point()
```



Scatterplot of Time and Passes from `worldcup` data

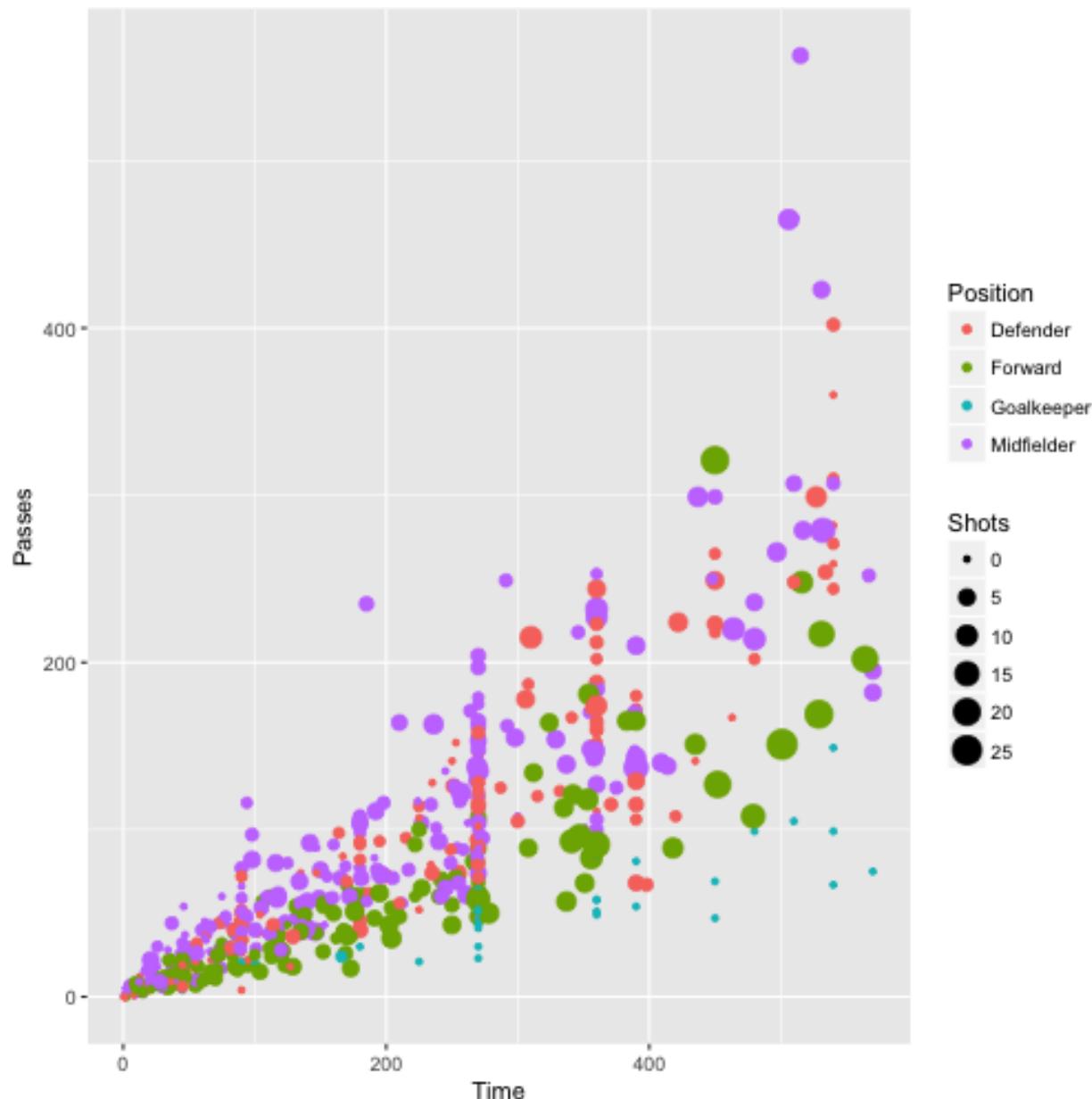
For any geom, there are both *required* and *accepted* aesthetics. For example, `geom_point` requires `x` and `y`, but will also accept `alpha` (transparency), `color`, `fill`, `group`, `size`, `shape`, and `stroke`. If you try to create a geom without one its required aesthetics, you will get an error:

```
ggplot(worldcup, aes(x = Time)) +  
  geom_point()  
Error: geom_point requires the following missing aesthetics: y
```

Plot with missing y aesthetic

You can, however, add accepted aesthetics to the geom if you'd like; for example, to use color to show player position and size to show shots on goal for the World Cup data, you could call:

```
ggplot(worldcup, aes(x = Time, y = Passes,
                      color = Position, size = Shots)) +
  geom_point()
```



Using color and size to show Position and Shots

The following table gives some of the more common geoms in `ggplot2`, along with the aesthetics commonly required for each and some of the most useful specific arguments for each (there are other useful arguments that can be applied to many different geom functions, which will be covered later). The functions of most are clear from the geom names (e.g., `geom_point` plots points; `geom_segment` plots segments).

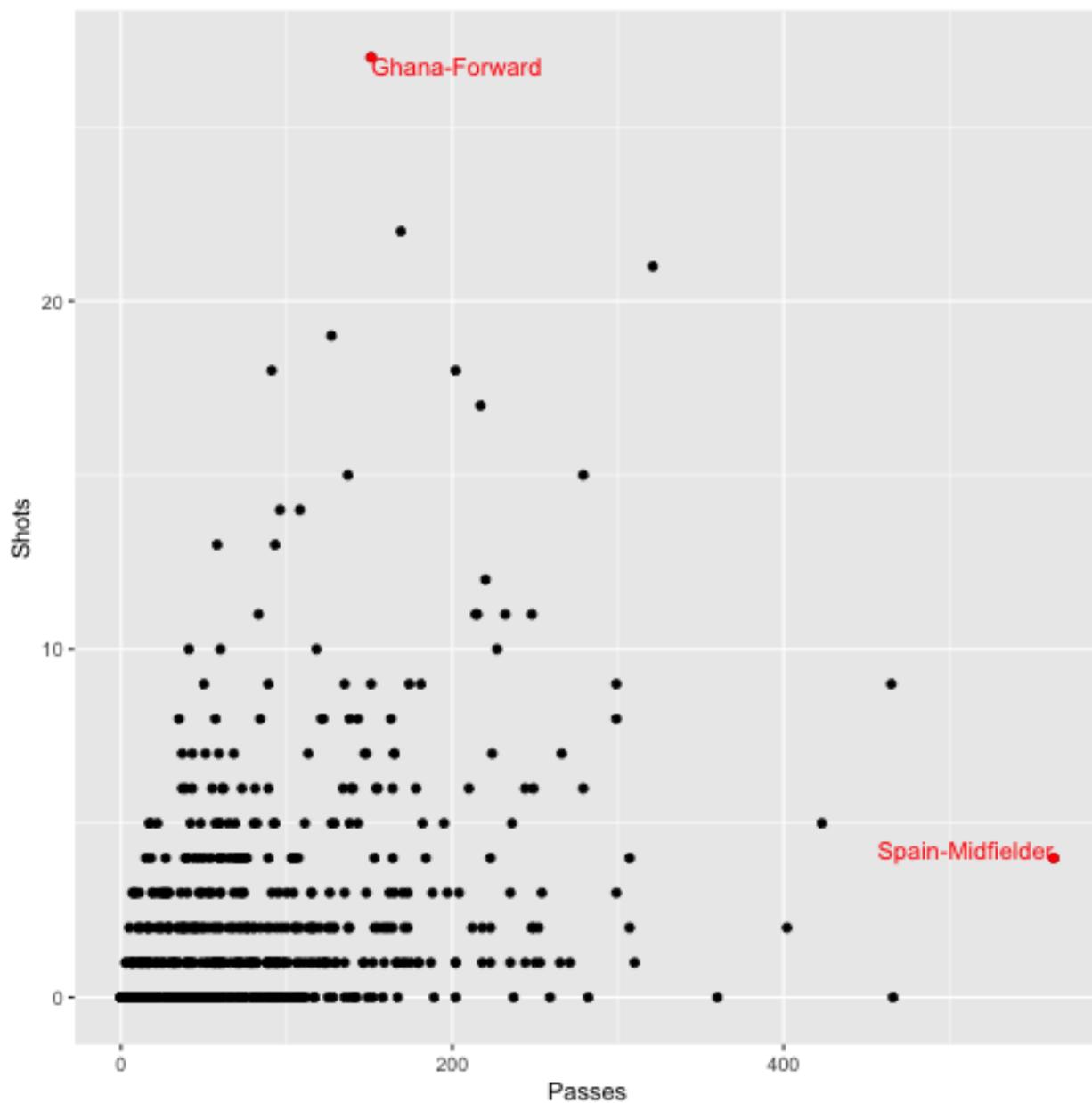
| Function | Common aesthetics | Common arguments |
|------------------|-------------------|----------------------------|
| geom_point() | x, y | |
| geom_line() | x, y | arrow, na.rm |
| geom_segment() | x, y, xend, yend | arrow, na.rm |
| geom_path() | x, y | na.rm |
| geom_polygon() | x, y | |
| geom_histogram() | x | bins, binwidth |
| geom_abline() | intercept, slope | |
| geom_hline() | yintercept | |
| geom_vline() | xintercept | |
| geom_boxplot() | x, y | outlier.color, notch, coef |
| geom_smooth() | x, y | method, se, span |
| geom_text() | x, y, label | parse, nudge_x, nudge_y |

Using multiple geoms

Several geoms can be added to the same `ggplot` object, which allows you to build up layers to create interesting graphs. For example, you could make the World Cup scatterplot of time versus shots more interesting by adding label points for noteworthy players with the player's team name and position:

```
library(dplyr)
noteworthy_players <- worldcup %>% filter(Shots == max(Shots) | 
                                             Passes == max(Passes)) %>%
  mutate(point_label = paste(Team, Position, sep = "-"))

ggplot(worldcup, aes(x = Passes, y = Shots)) +
  geom_point() +
  geom_text(data = noteworthy_players, aes(label = point_label),
            vjust = "inward", hjust = "inward", color = "red") +
  geom_point(data = noteworthy_players, color = "red")
```



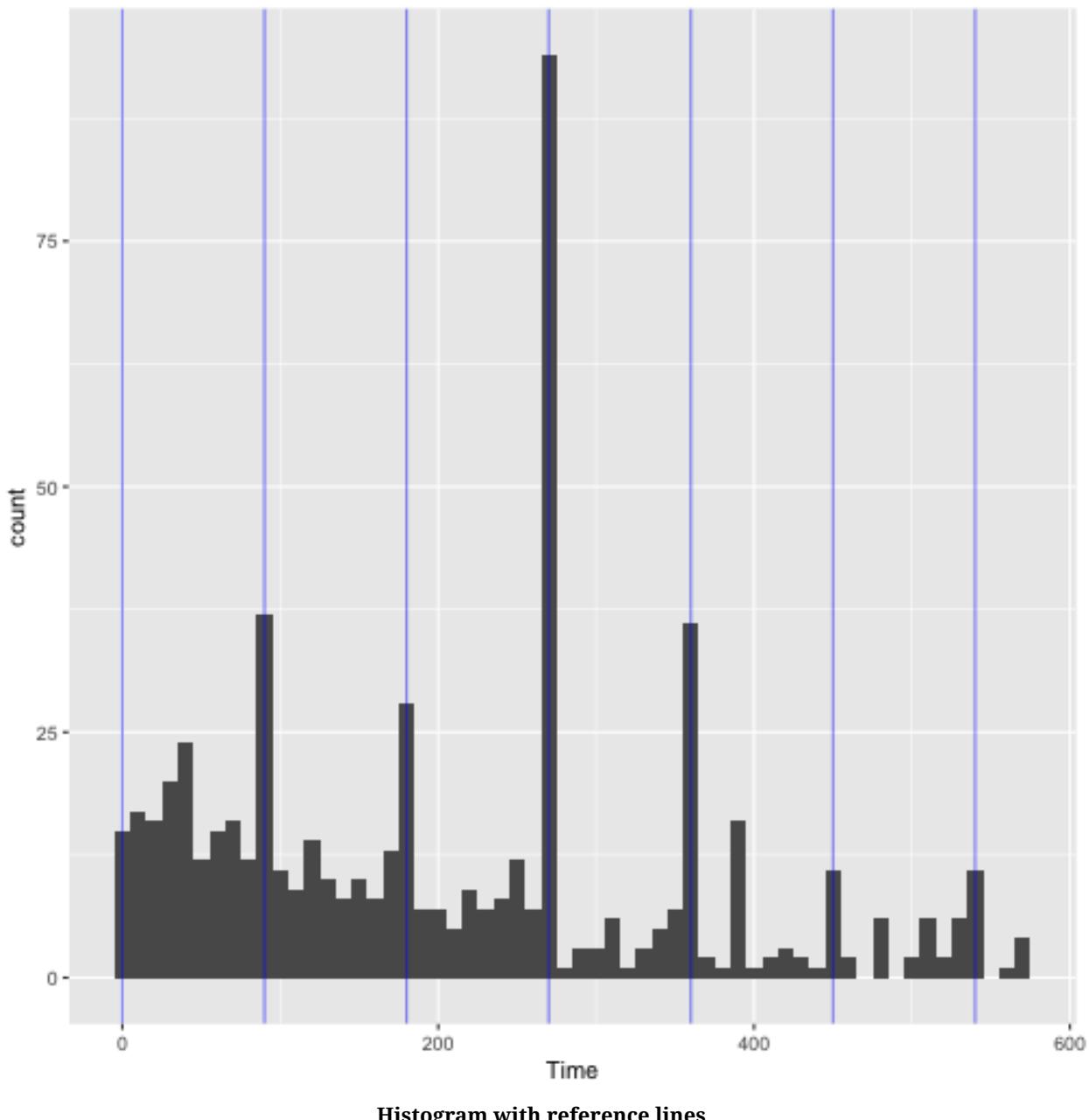
Adding label points to a scatterplot



In this example, we're using data from different dataframes for different geoms. We'll discuss how that works more later in this section.

As another example, there seemed to be some horizontal clustering in the scatterplot we made of player time versus passes made for the `worldcup` data. Soccer games last 90 minutes each, and different teams play a different number of games at the World Cup, based on how well they do. To check if horizontal clustering is at 90-minute intervals, you can plot a histogram of player time (`Time`), with reference lines every 90 minutes:

```
ggplot(worldcup, aes(x = Time)) +  
  geom_histogram(binwidth = 10) +  
  geom_vline(xintercept = 90 * 0:6,  
             color = "blue", alpha = 0.5)
```

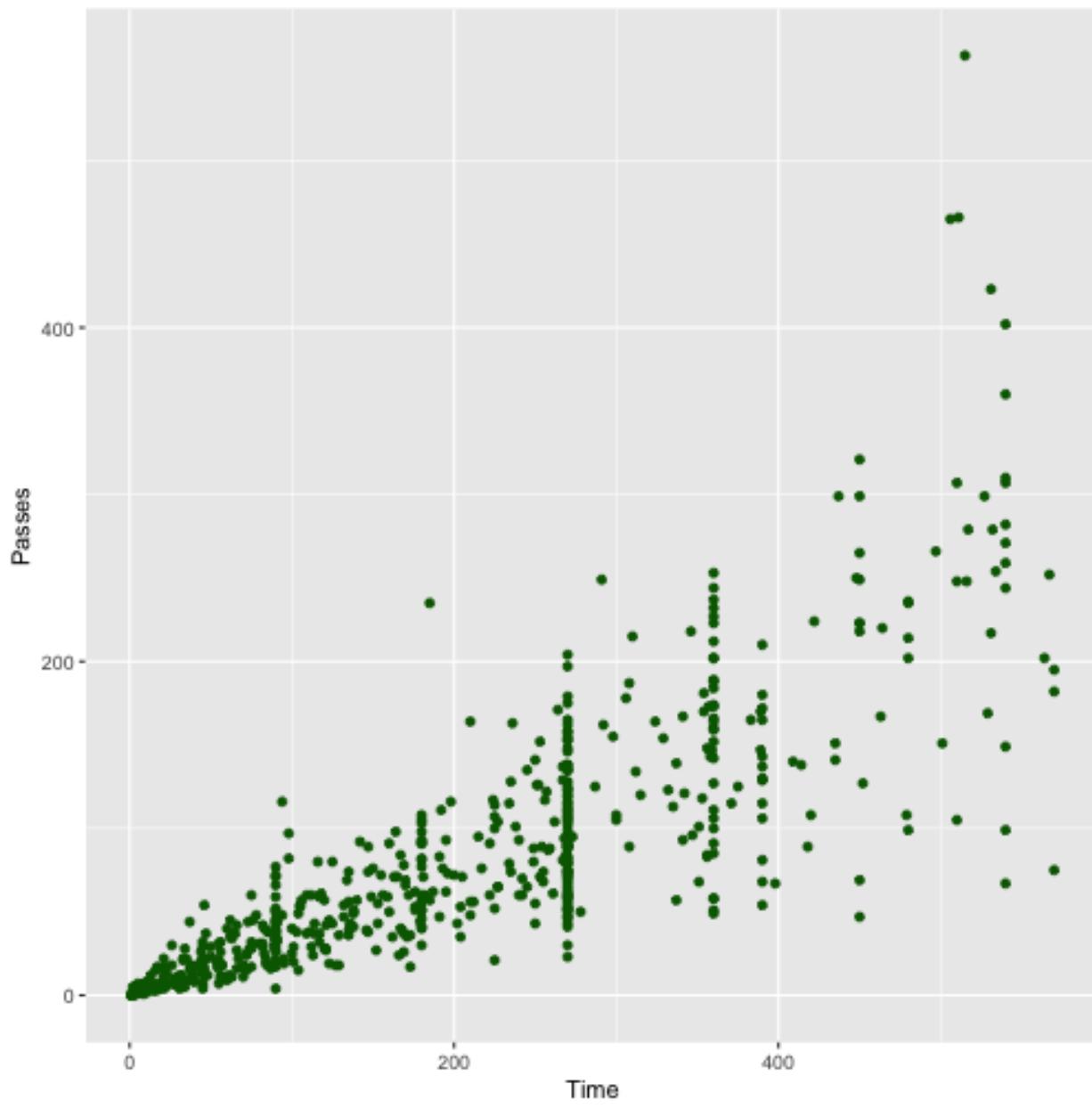


Based on this graph, player's times do cluster at 90-minute marks, especially at 270 minutes, which would be approximately after three games, the number played by all teams that fail to make it out of the group stage.

Constant aesthetics

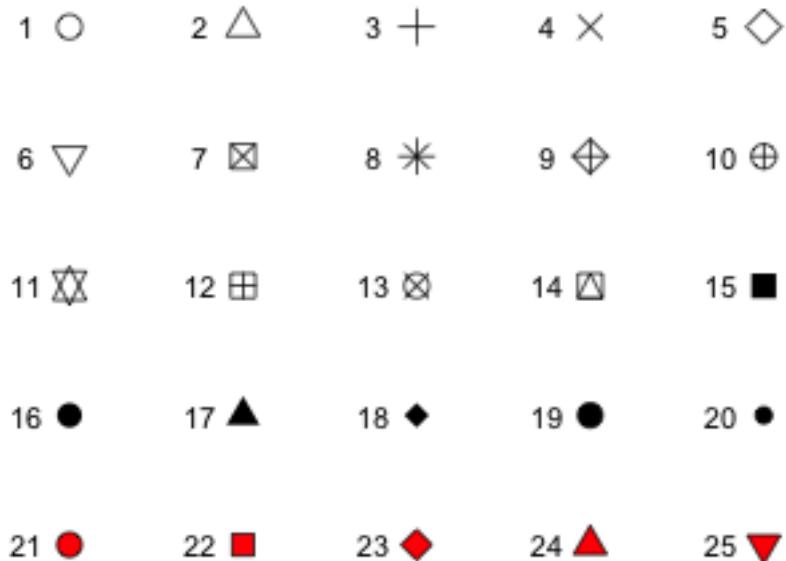
Instead of mapping an aesthetic to an element of your data, you can use a constant value for it. For example, you may want to make all the points green in the World Cup scatterplot:

```
ggplot(worldcup, aes(x = Time, y = Passes)) +  
  geom_point(color = "darkgreen")
```



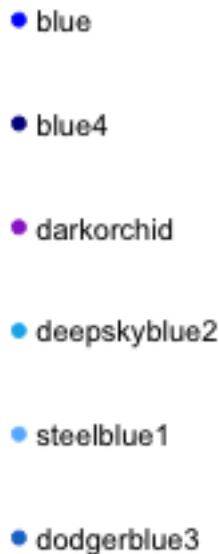
In this case, you'll define the color aesthetic when you add the geom, outside of an `aes` statement. You can do this with any of the aesthetics, including color, fill, shape, and size.

If you want to change the shape of points, in R, you specify the shape you want to use with a number. Figure @ref(fig:shapeexamples) shows the shapes that correspond to the numbers 1 to 25 in the `shape` aesthetic. This figure also provides an example of the difference between color (black for all these example points) and fill (red for these examples). You can see that some point shapes include a fill (21 for example), while some are either empty (1) or solid (19).



Examples of the shapes corresponding to different numeric choices for the `shape` aesthetic. For all examples, `color` is set to black and `fill` to red.

If you want to set color to be a constant value, you can do that in R using character strings for different colors. Figure @ref(fig:colorexamples) gives an example of a few of the different blues available in R. To find images that show all the named choices for colors that can be specified this way in R, google “R colors” and search by “Images” (for example, there is a pdf here: <http://www.stat.columbia.edu/~tzheng/files/Rcolor.pdf>).



Example of available shades of blue in R.

Other useful plot additions

There are also a number of elements besides geoms that you can add onto a `ggplot` object using `+`. A few that are used very frequently are:

| Element | Description |
|---------------------------------------|-------------------------|
| <code>ggtitle</code> | Plot title |
| <code>xlab</code> , <code>ylab</code> | x- and y-axis labels |
| <code>xlim</code> , <code>ylim</code> | Limits of x- and y-axis |

You can also use this syntax to customize plot scales and themes, which we will discuss later in this section.

Example plots

In this subsection, I'll show some examples of basic plots created with `ggplot2`. For the example plots in this subsection, I'll use a dataset in the `faraway` package called `nepali`. This gives data from a study of the health of a group of Nepalese children.

```
library(faraway)
data(nepali)
```

Each observation is a single measurement for a child; there can be multiple observations per child. I used the following code to select only the columns for child id, sex, weight, height, and age. I also used `distinct` to limit the dataset to only include one measurement for each child, the child's first measurement in the dataset.

```
nepali <- nepali %>%
  select(id, sex, wt, ht, age) %>%
  mutate(id = factor(id),
         sex = factor(sex, levels = c(1, 2),
                      labels = c("Male", "Female"))) %>%
  distinct(id, .keep_all = TRUE)
```

After this cleaning, the data looks like this:

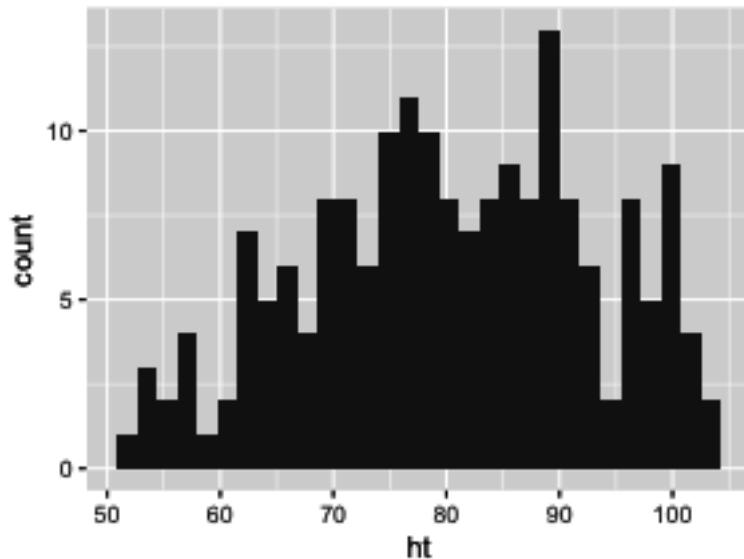
```
head(nepali)
#> #> #> #> #> #>
#> #> #> #> #> #>
#> #> #> #> #> #>
#> #> #> #> #> #>
#> #> #> #> #> #>
```

| | id | sex | wt | ht | age |
|---|--------|--------|------|-------|-----|
| 1 | 120011 | Male | 12.8 | 91.2 | 41 |
| 2 | 120012 | Female | 14.9 | 103.9 | 57 |
| 3 | 120021 | Female | 7.7 | 70.1 | 8 |
| 4 | 120022 | Female | 12.1 | 86.4 | 35 |
| 5 | 120023 | Male | 14.2 | 99.4 | 49 |
| 6 | 120031 | Male | 13.9 | 96.4 | 46 |

Histograms

Histograms show the distribution of a single variable. Therefore, `geom_histogram()` requires only one main aesthetic, `x`, the (numeric) vector for which you want to create a histogram. For example, to create a histogram of children's heights for the Nepali dataset (Figure @ref(fig:nepalihist1)), run:

```
ggplot(nepali, aes(x = ht)) +
  geom_histogram()
```



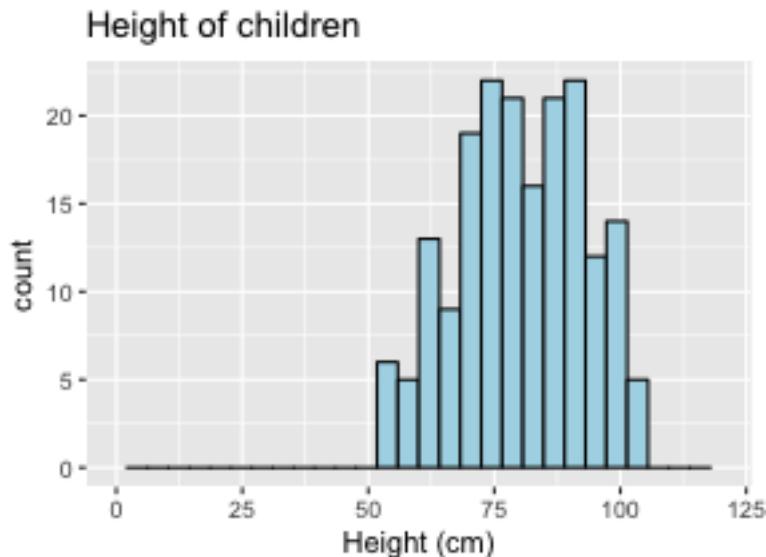
Basic example of plotting a histogram with `ggplot2`. This histogram shows the distribution of heights for the first recorded measurements of each child in the `nepali` dataset.



If you run the code with no arguments for `binwidth` or `bins` in `geom_histogram`, you will get a message saying “`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.”. This message is just saying that a default number of bins was used to create the histogram. You can use arguments to change the number of bins used, but often this default is fine. You may also get a message that observations with missing values were removed.

You can add some elements to the histogram to customize it a bit. For example (Figure @ref(fig:nepalihist2)), you can add a figure title (`ggtitle`) and clearer labels for the x-axis (`xlab`). You can also change the range of values shown by the x-axis (`xlim`).

```
ggplot(nepali, aes(x = ht)) +
  geom_histogram(fill = "lightblue", color = "black") +
  ggtitle("Height of children") +
  xlab("Height (cm)") + xlim(c(0, 120))
```

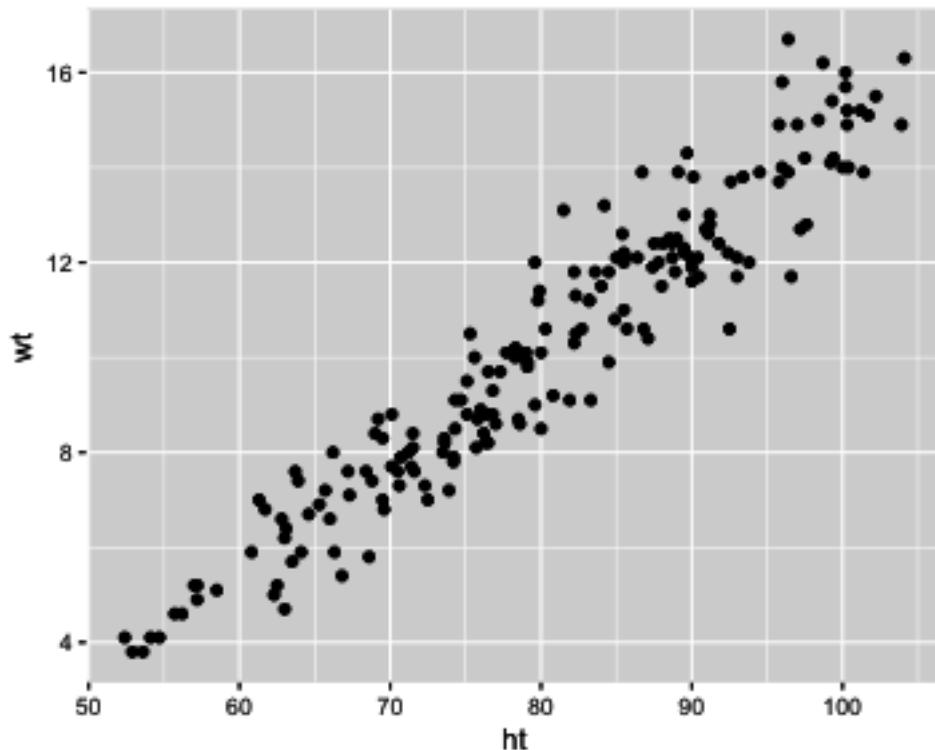


Example of adding ggplot elements to customize a histogram.

Scatterplots

A scatterplot shows how one variable changes as another changes. You can use the `geom_point` geom to create a scatterplot. For example, to create a scatterplot of height versus age for the Nepali data (Figure @ref(fig:nepaliscatter1)), you can run the following code:

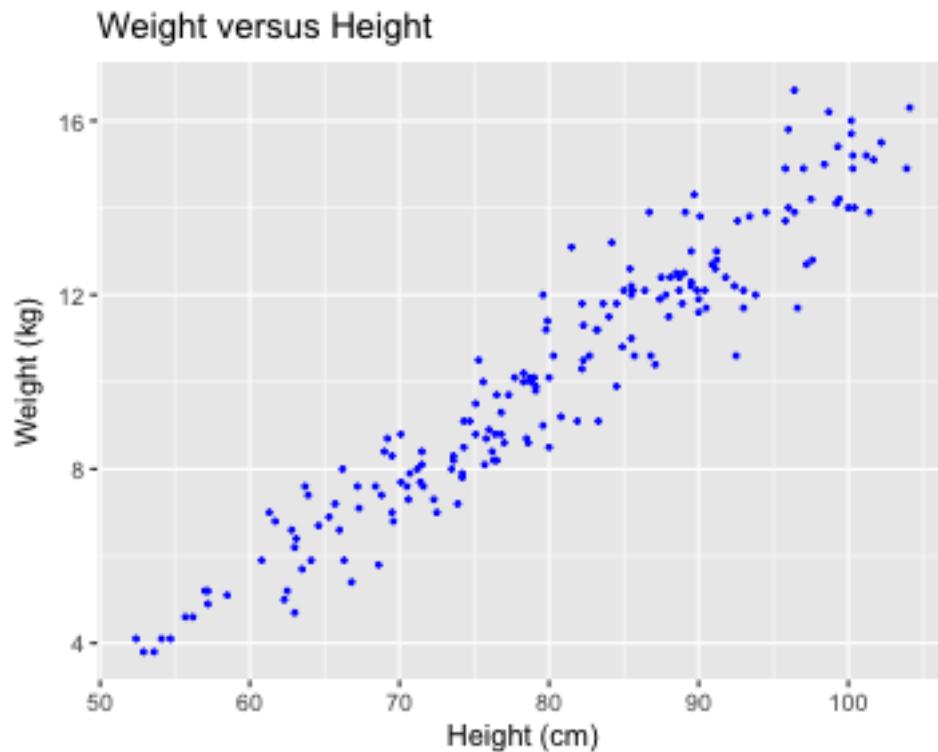
```
ggplot(nepali, aes(x = ht, y = wt)) +
  geom_point()
```



Example of creating a scatterplot. This scatterplot shows the relationship between children's heights and weights within the nepali dataset.

Again, you can use some of the options and additions to change the plot appearance. For example, to add a title, change the x- and y-axis labels, and change the color and size of the points on the scatterplot (Figure @ref(fig:nepaliscatter2)), you can run:

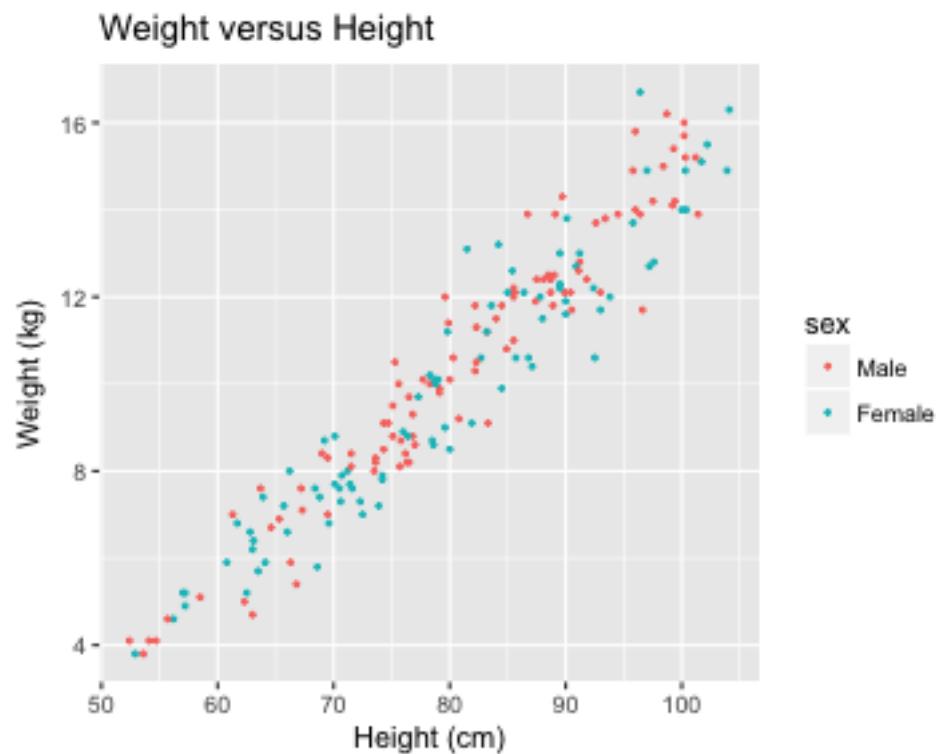
```
ggplot(nepali, aes(x = ht, y = wt)) +  
  geom_point(color = "blue", size = 0.5) +  
  ggtitle("Weight versus Height") +  
  xlab("Height (cm)") + ylab("Weight (kg)")
```



Example of adding ggplot elements to customize a scatterplot.

You can also try mapping another variable in the dataset to the `color` aesthetic. For example, to use color to show the sex of each child in the scatterplot (Figure @ref(fig:nepaliscatter3)), you can run:

```
ggplot(nepali, aes(x = ht, y = wt, color = sex)) +  
  geom_point(size = 0.5) +  
  ggtitle("Weight versus Height") +  
  xlab("Height (cm)") + ylab("Weight (kg)")
```

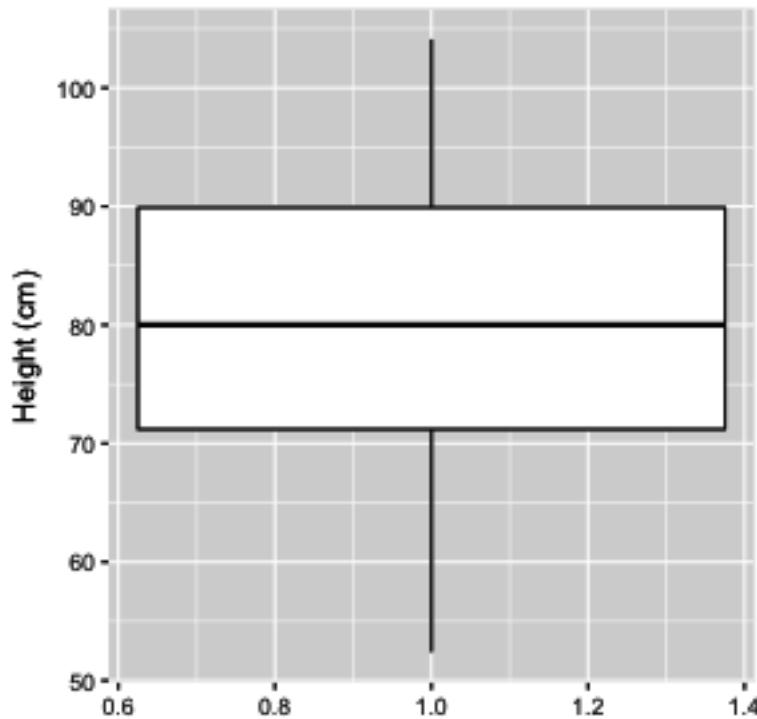


Example of mapping color to an element of the data in a scatterplot.

Boxplots

Boxplots can be used to show the distribution of a continuous variable. To create a boxplot, you can use the `geom_boxplot` geom. To plot a boxplot for a single, continuous variable, you can map that variable to `y` in the `aes` call, and map `x` to the constant `1`. For example, to create a boxplot of the heights of children in the Nepali dataset (Figure @ref(fig:nepaliboxplot1)), you can run:

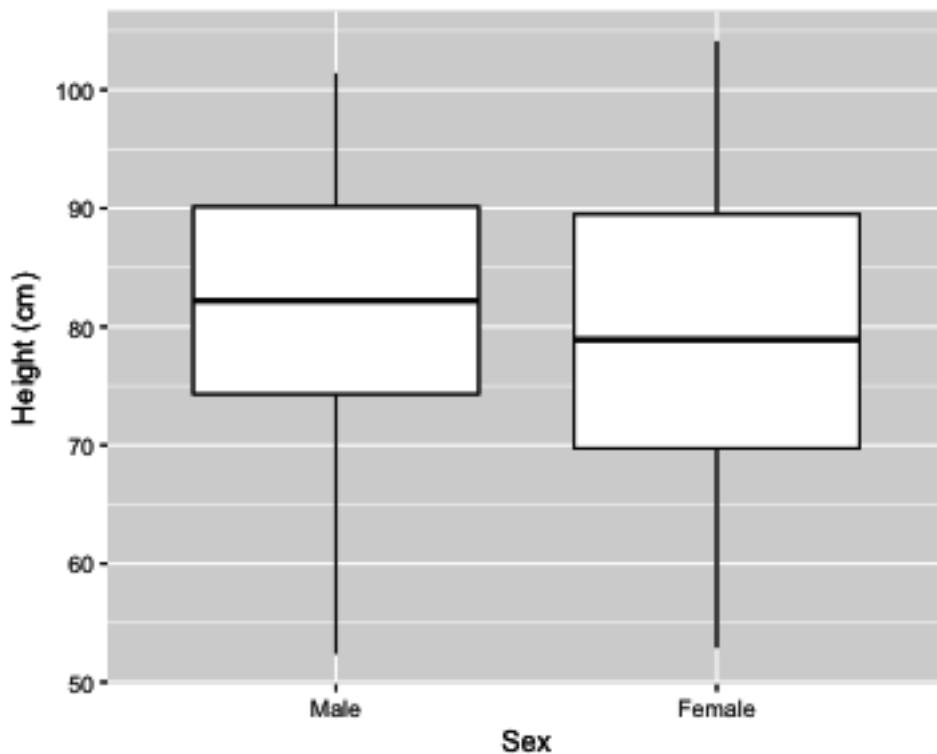
```
ggplot(nepali, aes(x = 1, y = ht)) +  
  geom_boxplot() +  
  xlab("") + ylab("Height (cm)")
```



Example of creating a boxplot. The example shows the distribution of height data for children in the nepali dataset.

You can also create separate boxplots, one for each level of a factor (Figure @ref(fig:nepaliboxplot2)). In this case, you'll need to include two aesthetics (`x` and `y`) when you initialize the `ggplot` object. The `y` variable is the variable for which the distribution will be shown, and the `x` variable should be a discrete (categorical or TRUE/FALSE) variable, and will be used to group the variable. This `x` variable should also be specified as the grouping variable, using `group` within the aesthetic call.

```
ggplot(nepali, aes(x = sex, y = ht, group = sex)) +  
  geom_boxplot() +  
  xlab("Sex") + ylab("Height (cm)")
```

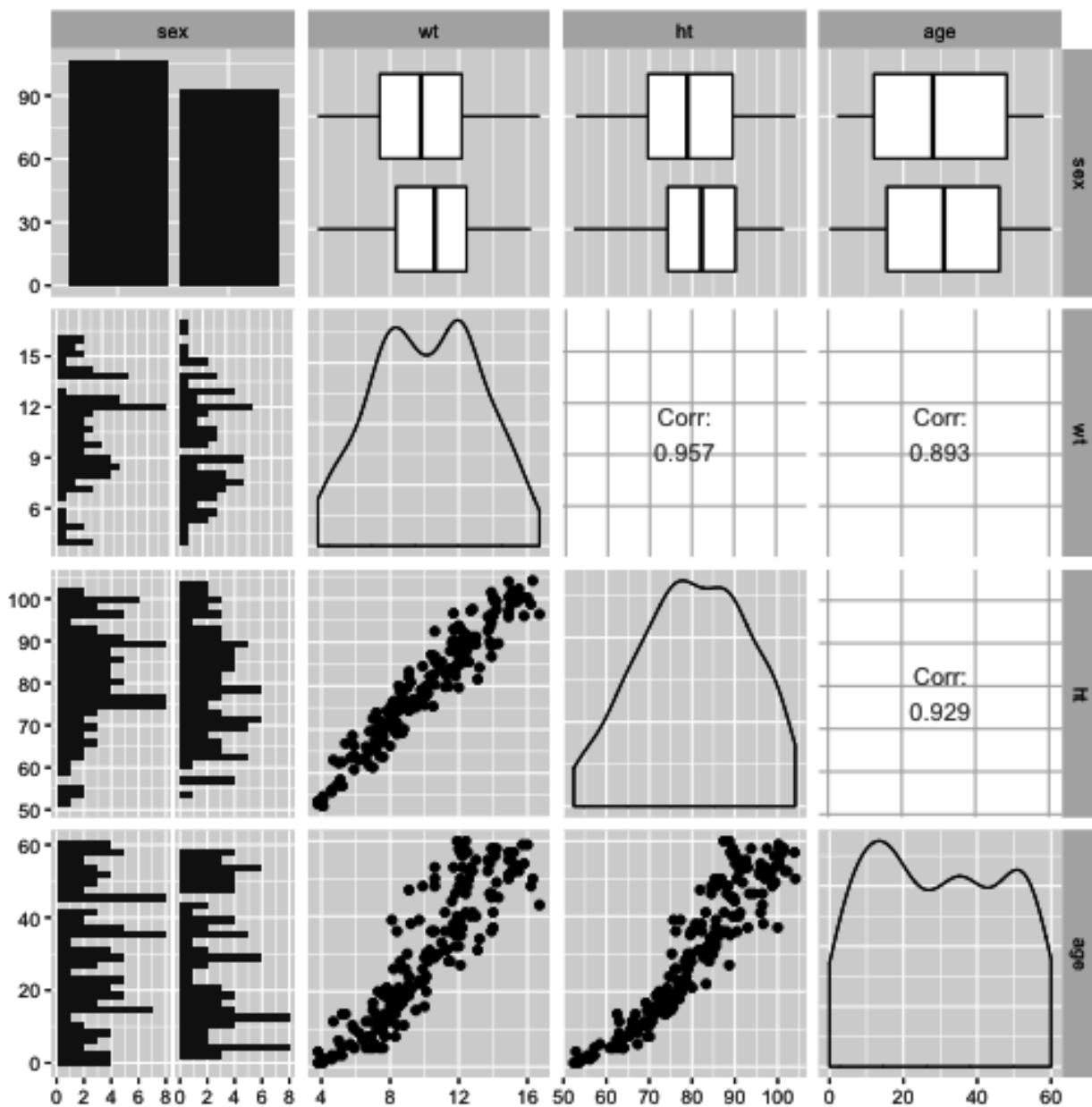


Example of creating separate boxplots, divided by a categorical grouping variable in the data.

Extensions of `ggplot2`

There are lots of R extensions for creating other interesting plots. For example, you can use the `ggpairs` function from the `GGally` package to plot all pairs of scatterplots for several variables (Figure @ref(fig:ggallyexample)).

```
library(GGally)
ggpairs(nepali %>% select(sex, wt, ht, age))
```



Example of using `ggpairs` from the `GGally` package for exploratory data analysis.

Notice how this output shows continuous and binary variables differently. For example, the center diagonal shows density plots for continuous variables, but a bar chart for the categorical variable.

See <https://www.ggplot2-exts.org> to find more `ggplot2` extensions. Later in this course, we will give an overview of how to make your own extensions.

4.2 Customizing ggplot2 Plots

With slightly more complex code, you can create very interesting and customized plots using `ggplot2`. This functionality is useful for creating plots that follow some of the guidelines for effective visualization, based for example on ideas from Edward Tufte and others. In this section, we'll provide an overview of some of the guidelines for creating good plots and show how you can customize `ggplot` objects to adhere to some of these guidelines. This will also give us the chance to go over how to customize `ggplot` objects, and we'll end by going over scales, themes, faceting, and color specifically.

Guidelines for good plots

There are a number of very thoughtful books and articles about creating graphics that effectively communicate information. Some of the authors I highly recommend (and from whose work I've pulled and aggregated the guidelines for good graphics we'll go over) are:

- Edward Tufte (**The Visual Display of Quantitative Information** is a classic)
- Howard Wainer
- Stephen Few
- Nathan Yau



While we overview some guidelines for effective plots here, this is mostly to provide a framework for showing how to customize `ggplot` objects. If you are interested in learning more about creating effective visualizations, you should read some of the thorough and thoughtful books written by the authors listed above.

In this section, we'll overview six guidelines for good graphics, based on the writings of these and other specialists in data display. The guidelines are:

1. Aim for high data density.
2. Use clear, meaningful labels.
3. Provide useful references.
4. Highlight interesting aspects of the data.
5. Consider using small multiples.
6. Make order meaningful.

For the examples in this subsection, we'll use `dplyr` for data cleaning and, for plotting, the packages `ggplot2`, `gridExtra`, and `ggthemes`.

```
library(tidyverse) ## Loads `dplyr` and `ggplot2`
library(gridExtra)
library(ggthemes)
```

You can load the data for the examples in this subsection with the following code:

```
library(faraway)
data(nepali)
data(worldcup)

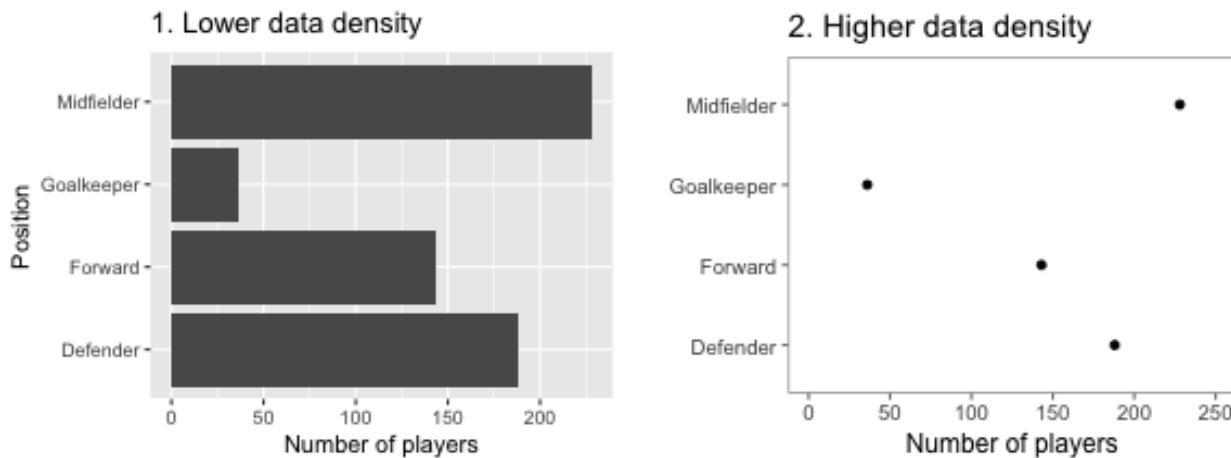
library(dlnm)
data(chicagoNMMAPS)
chic <- chicagoNMMAPS
chic_july <- chic %>%
  filter(month == 7 & year == 1995)
```

High data density

Guideline 1: Aim for high data density.

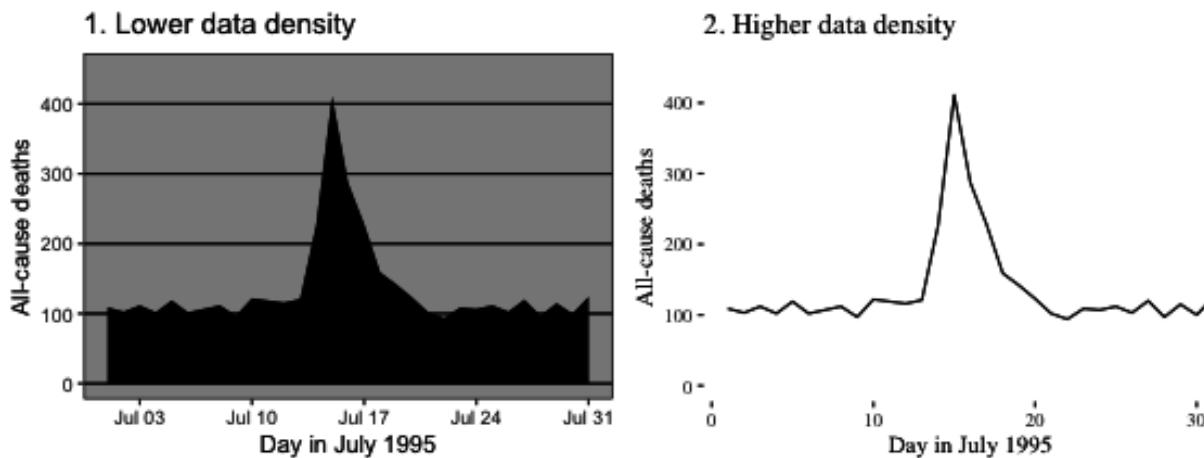
You should try to increase, as much as possible, the **data to ink ratio** in your graphs. This is the ratio of “ink” providing information to all ink used in the figure. In other words, if an element of the plot is redundant, take it out. If a graph uses up a lot of your printer’s ink, it should be packed with information.

The two graphs in Figure @ref(fig:datainkratio1) show the same information (“data”), but use very different amounts of ink. Each shows the number of players in each of four positions in the `worldcup` dataset. Notice how, in the plot on the right, a single dot for each category shows the same information that a whole filled bar is showing on the left. Further, the plot on the right has removed the gridded background, removing even more “ink” from the plot.



Example of plots with lower (left) and higher (right) data-to-ink ratios. Each plot shows the number of players in each position in the `worldcup` dataset from the `faraway` package.

Figure @ref(fig:datainkratio2) gives another example of two plots that show the same information but with very different data densities. This figure uses the `chicagoNMMAPS` data from the `dlnm` package, which includes daily mortality, weather, and air pollution data for Chicago, IL. Both plots show daily mortality counts during July 1995, when a very severe heat wave hit Chicago. Notice how many of the elements in the plot on the left, including the shading under the mortality time series and the colored background and grid lines, are unnecessary for interpreting the message from the data.

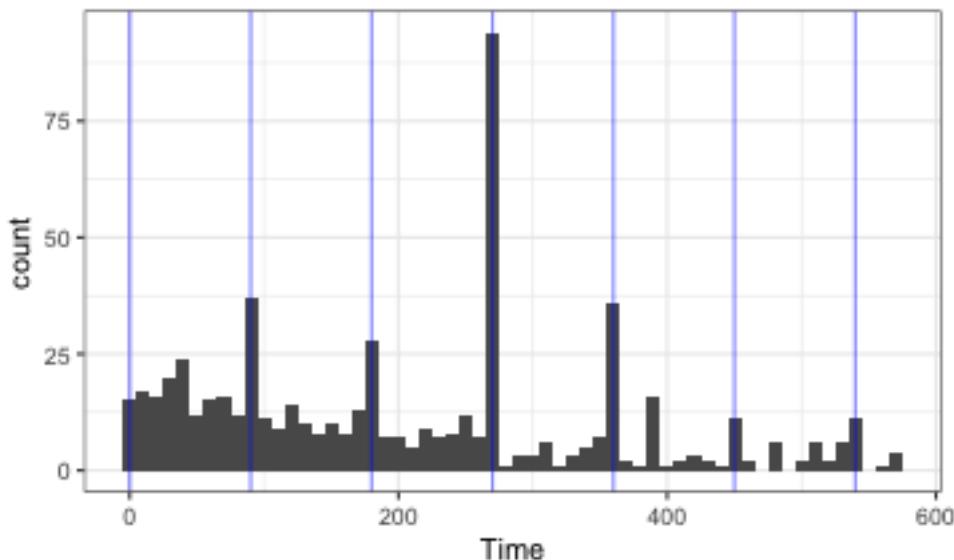


Example of plots with lower (left) and higher (right) data-to-ink ratios. Each plot shows daily mortality in Chicago, IL, in July 1995 using the `chicagoNMMAPS` data from the `dlnm` package.

By increasing the data-to-ink ratio in a plot, you can help viewers see the message of the data more quickly. A cluttered plot is harder to interpret. Further, you leave room to add some of the other elements we'll talk about, including elements to highlight interesting data and useful references. Notice how the plots on the left in Figures @ref(fig:datainkratio1) and @ref(fig:datainkratio2) are already cluttered and leave little room for adding extra elements, while the plots on the right of those figures have much more room for additions.

One quick way to increase data density in `ggplot2` is to change the *theme* for the plot. You can use themes in `ggplot` to quickly change several elements of the plot's appearance. There are several themes that come with `ggplot2`, including a black-and-white theme and a minimal theme. To use a theme, you can add a function with the theme to the `ggplot` object. For example, to use a minimal theme for the histogram of player time created in a previous subsection, you can run:

```
ggplot(worldcup, aes(x = Time)) +
  geom_histogram(binwidth = 10) +
  geom_vline(xintercept = 90 * 0:6,
             color = "blue", alpha = 0.5) +
  theme_bw()
```



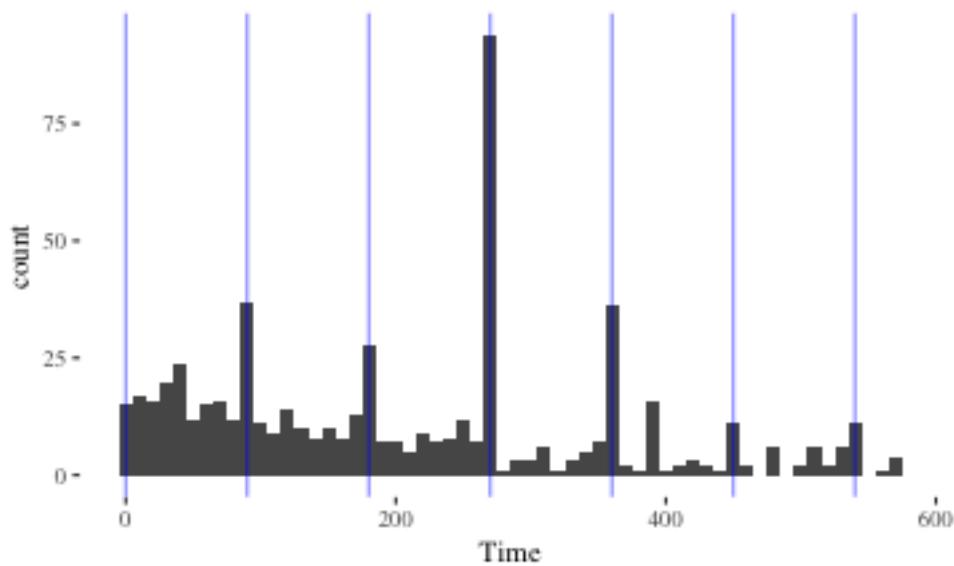
plot of chunk unnamed-chunk-24

Themes that come with `ggplot2` include:

- `theme_bw`
- `theme_minimal`
- `theme_void`

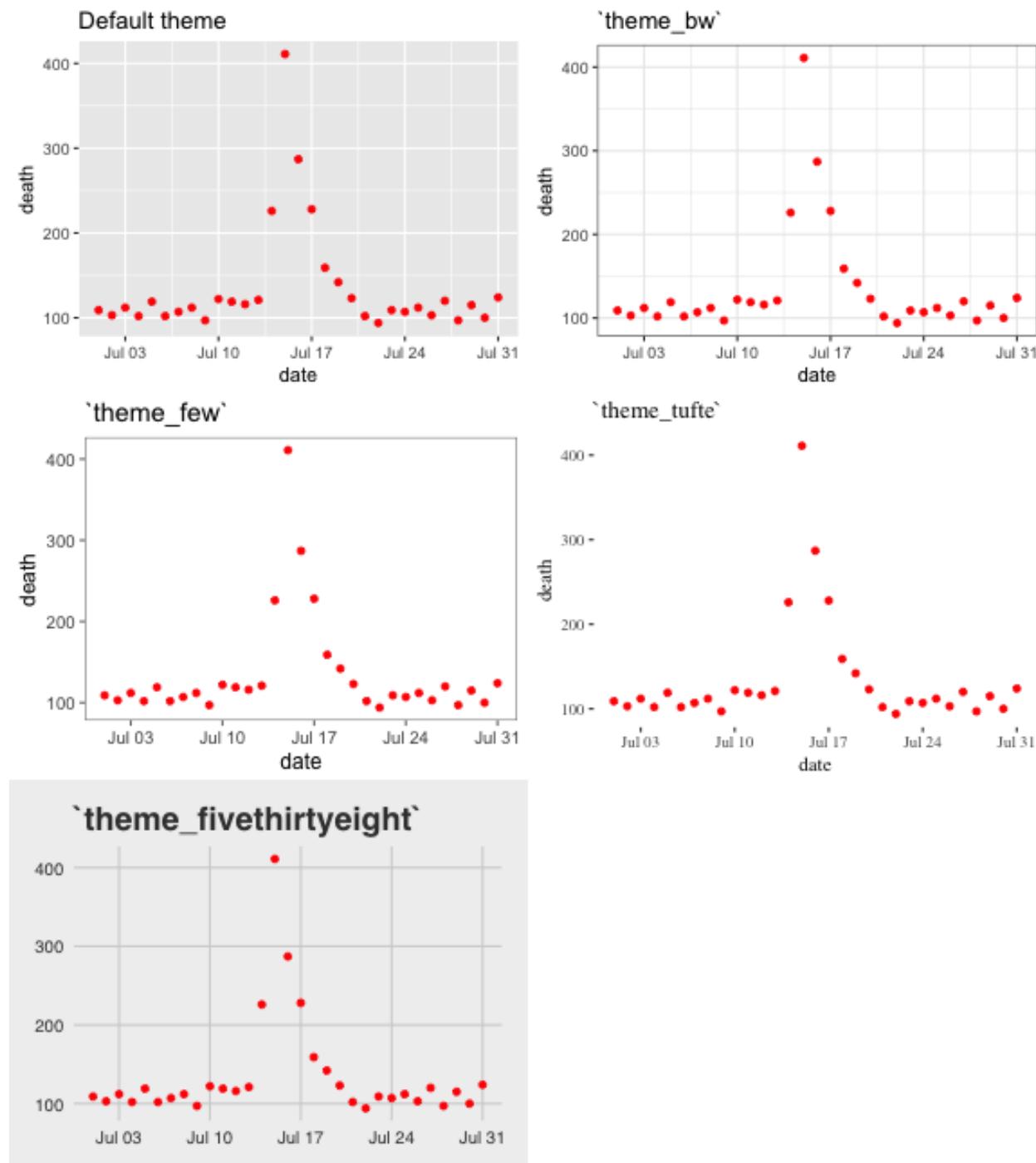
You can find more themes in packages that extend `ggplot2`. The `ggthemes` package, in particular, has some excellent additional themes. These include themes based on the graphing principles of Stephen Few (`theme_few`) and Edward Tufte (`theme_tufte`). Again, you can use one of these themes by adding it to a `ggplot` object:

```
library(ggthemes)
ggplot(worldcup, aes(x = Time)) +
  geom_histogram(binwidth = 10) +
  geom_vline(xintercept = 90 * 0:6,
             color = "blue", alpha = 0.5) +
  theme_tufte()
```



plot of chunk unnamed-chunk-25

The plots in Figure @ref(fig:themeexamples) shows some examples of the effects of using different themes. All show the same information— a plot of daily deaths in Chicago in July 1995. The top left graph shows the graph with the default theme. The other plots show the effects of adding different themes, including the black-and-white theme that comes with `ggplot2` (top right) and various themes from the `ggthemes` package.



Daily mortality in Chicago, IL, in July 1995. This figure gives an example of the plot using different themes.

We will teach you how to make your own theme later in the course.

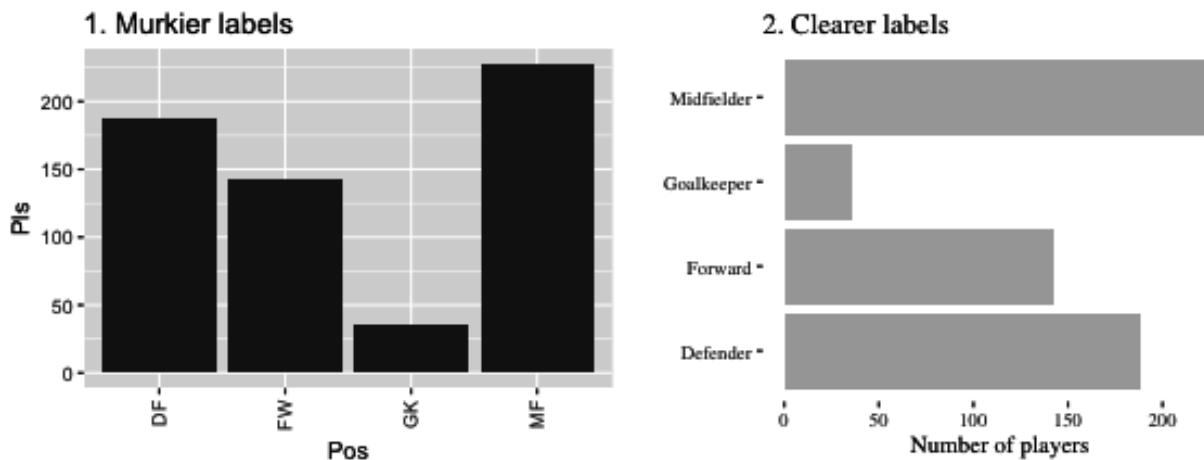
Meaningful labels

Guideline 2: Use clear, meaningful labels.

Graphs often default to use abbreviations for axis labels and other labeling. For example, the default is for `ggplot2` plots to use column names for the x- and y-axes of a scatterplot. While this is convenient for exploratory plots, it's often not adequate for plots for presentations and papers. You'll want to use short and easy-to-type column names in your `dataframe` to make coding easier, but you should use longer and more meaningful labeling in plots and tables that others need to interpret.

Furthermore, text labels can sometimes be aligned in a way that makes them hard to read. For example, when plotting a categorical variable along the x-axis, it can be difficult to fit labels for each category that are long enough to be meaningful.

Figure @ref(fig:labelsexample) gives an example of the same information (number of players in the World Cup data set by position) shown with labels that are harder to interpret (left) versus with clear, meaningful labels (right). Notice how the graph on the left is using abbreviations for the categorical variable (“DF” for “Defense”), abbreviations for axis labels (“Pos” for “Position” and “Pls” for “Number of players”), and has the player position labels in a vertical alignment. On the right graph, I have made the graph easier to quickly read and interpret by spelling out all labels and switching the x- and y-axes, so that there's room to fully spell out each position while still keeping the alignment horizontal, so the reader doesn't have to turn the page (or his head) to read the values.



The number of players in each position in the `worldcup` data from the `faraway` package. Both graphs show the same information, but the left graph has murkier labels, while the right graph has labels that are easier to read and interpret.

There are a few strategies you can use to make labels clearer when plotting with `ggplot2`:

- Add `xlab` and `ylab` elements to the plot, rather than relying on the column names in the original data. You can also relabel x- and y-axes with `scale` elements (e.g., `scale_x_continuous`), and the `scale` functions give you more power to also make other changes to the x- and y-axes (e.g., changing break points for the axis ticks). However, if you only need to change axis labels, `xlab` and `ylab` are often quicker.
- Include units of measurement in axis titles when relevant. If units are dollars or percent, check out the `scales` package, which allows you to add labels directly to

axis elements by including arguments like `labels = percent` in `scale` elements. See the helpfile for `scale_x_continuous` for some examples.

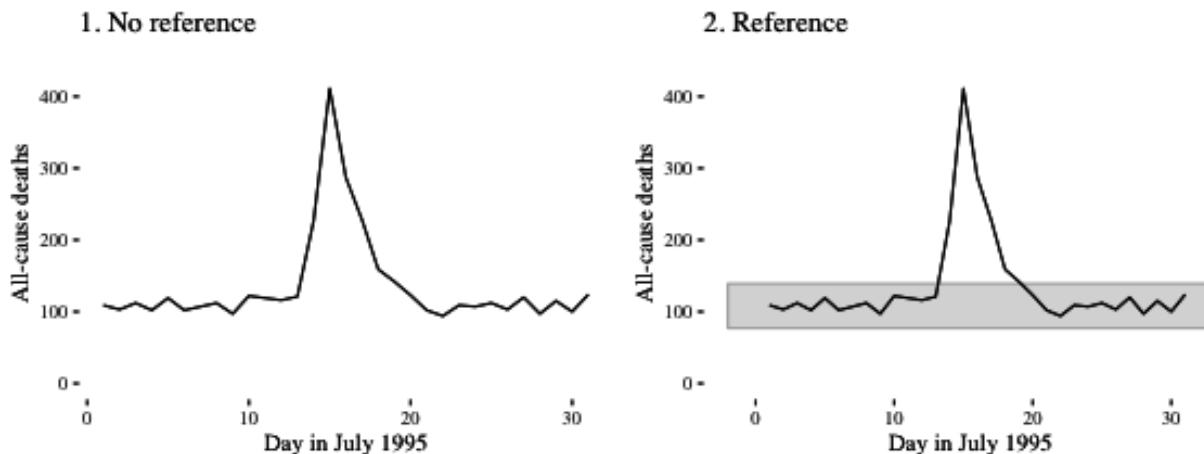
- If the x-variable requires longer labels, as is often the case with categorical data (for example, player positions Figure @ref(fig:labelsexample)), consider flipping the coordinates, rather than abbreviating or rotating the labels. You can use `coord_flip` to do this.

References

Guideline 3: Provide useful references.

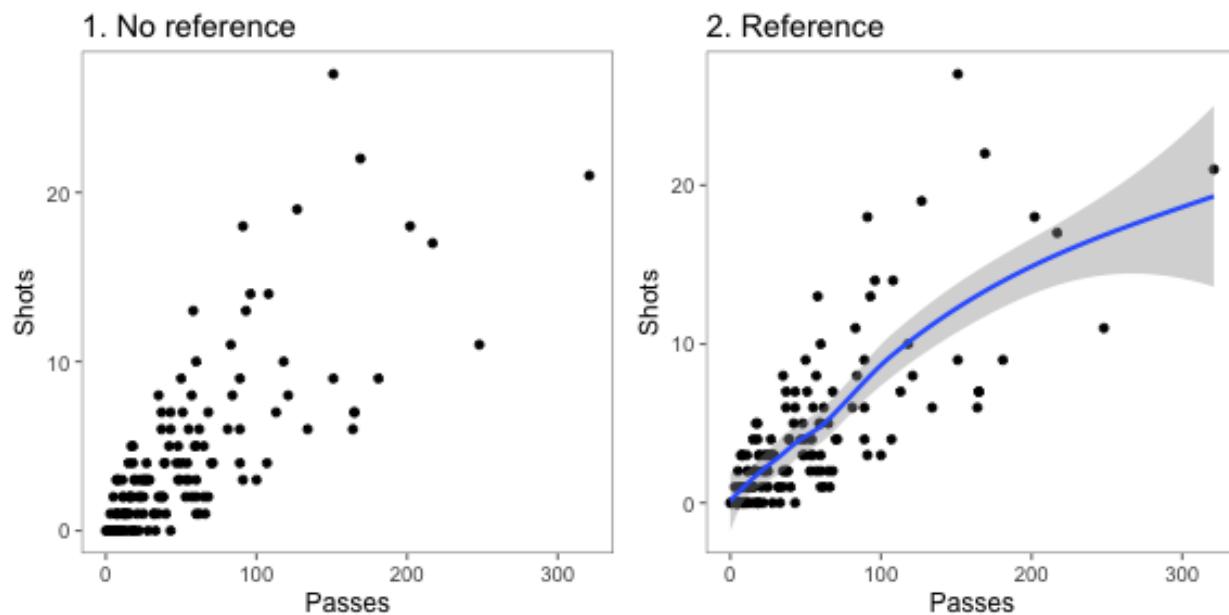
Data is easier to interpret when you add references. For example, if you show what is typical, it helps viewers interpret how unusual outliers are.

Figure @ref(fig:referenceexample1) shows daily mortality during July 1995 in Chicago, IL. The graph on the right has added shading showing the range of daily death counts in July in Chicago for neighboring years (1990–1994 and 1996–2000). This added reference helps clarify for viewers how unusual the number of deaths during the July 1995 heat wave was.



Daily mortality during July 1995 in Chicago, IL. In the graph on the right, I have added a shaded region showing the range of daily mortality counts for neighboring years, to show how unusual this event was.

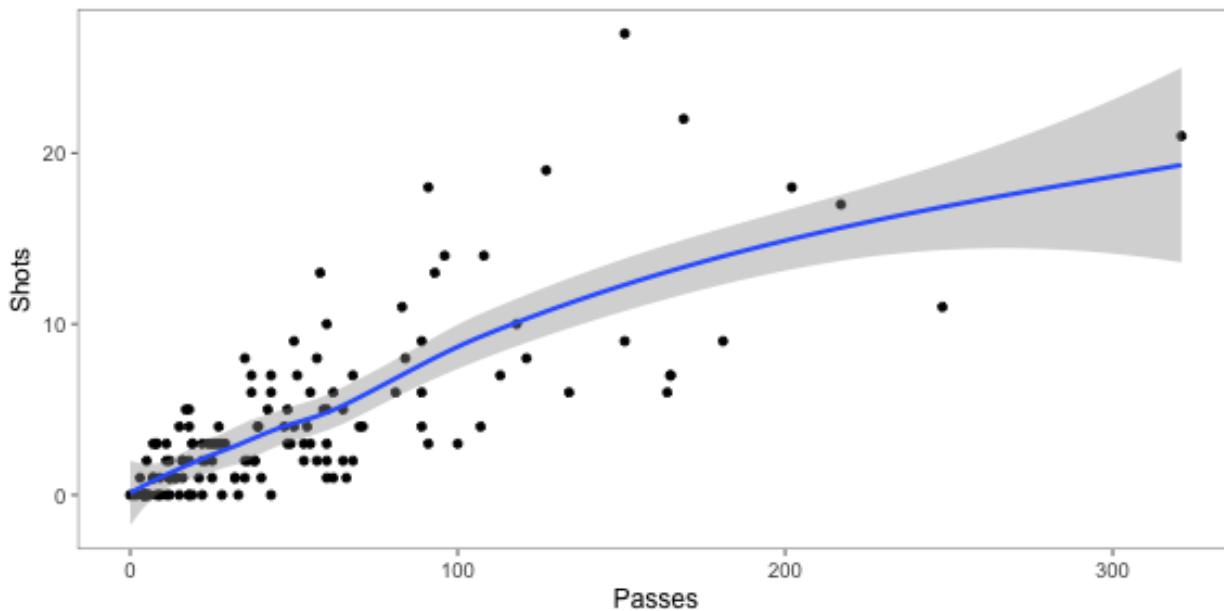
Another useful way to add references is to add a linear or smooth fit to the data, to help clarify trends in the data. Figure @ref(fig:referenceexample2) shows the relationship between passes and shots for Forwards in the `worldcup` dataset. The plot on the right has added a smooth function of the relationship between these two variables.



Relationship between passes and shots taken among Forwards in the worldcup dataset from the faraway package. The plot on the right has a smooth function added to help show the relationship between these two variables.

For scatterplots created with `ggplot2`, you can use the function `geom_smooth` to add a smooth or linear reference line. Here is the code that produces Figure @ref(fig:referenceexample3):

```
ggplot(filter(worldcup, Position == "Forward"), aes(x = Passes, y = Shots)) +
  geom_point(size = 1.5) +
  theme_few() +
  geom_smooth()
```



Relationship between passes and shots taken among Forwards in the worldcup dataset from the faraway package. The plot has a smooth function added to help show the relationship between these two variables.

The most useful `geom_smooth` parameters to know are:

- `method`: The default is to add a loess curve if the data includes less than 1000 points and a generalized additive model for 1000 points or more. However, you can change to show the fitted line from a linear model using `method = "lm"` or from a generalized linear model using `method = "glm"`.
- `span`: How wiggly or smooth the smooth line should be (smaller value: more flexible; larger value: more smooth)
- `se`: TRUE or FALSE, indicating whether to include shading for 95% confidence intervals.
- `level`: Confidence level for confidence interval (e.g., `0.90` for 90% confidence intervals)

Lines and polygons can also be useful for adding references, as in Figure @ref(fig:referenceexample1). Useful geoms for such shapes include:

- `geom_hline`, `geom_vline`: Add a horizontal or vertical line
- `geom_abline`: Add a line with an intercept and slope
- `geom_polygon`: Add a filled polygon
- `geom_path`: Add an unfilled polygon

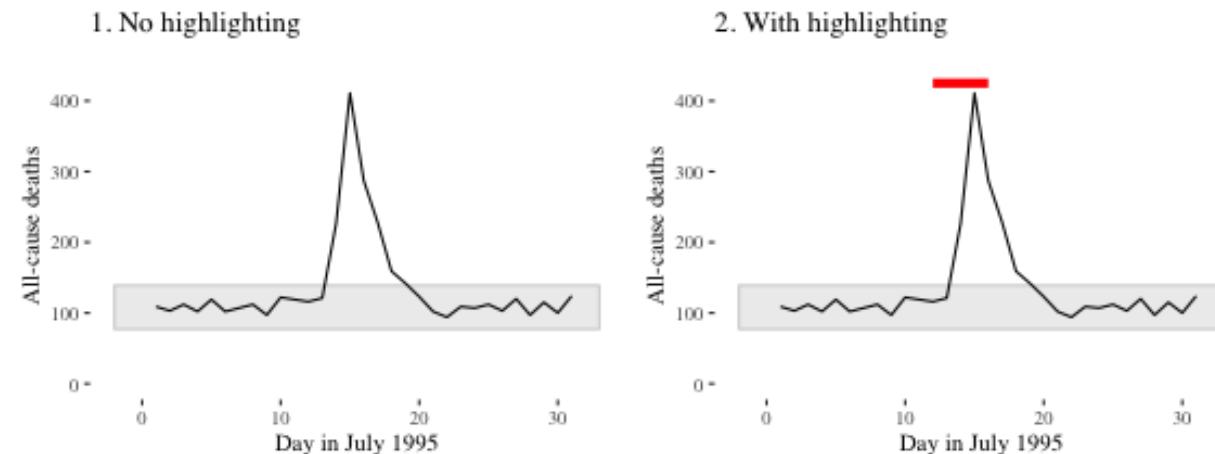
You want these references to support the main data shown in the plot, but not overwhelm it. When adding these references:

- Add reference elements first, so they will be plotted under the data, instead of on top of it.
- Use `alpha` to add transparency to these elements.
- Use colors that are unobtrusive (e.g., grays).
- For lines, consider using non-solid line types (e.g., `linetype = 3`).

Highlighting

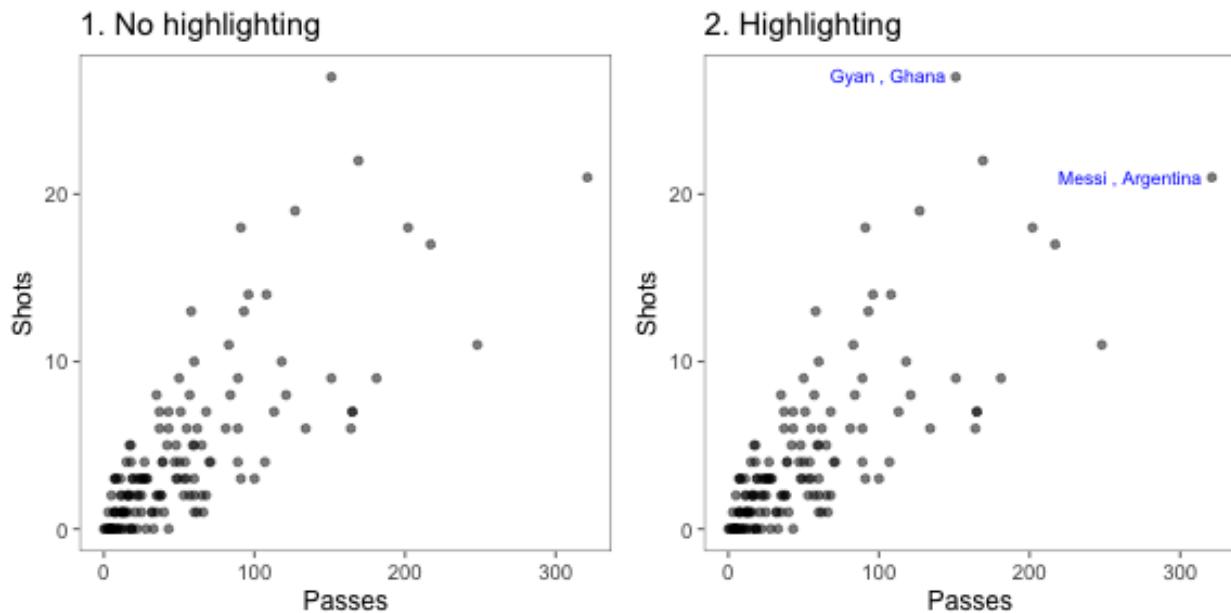
Guideline 4: Highlight interesting aspects.

Consider adding elements to highlight noteworthy elements of the data. For example, in the graph on the right of Figure @ref(fig:highlightexample1), the days of the heat wave (based on temperature measurements) have been highlighted over the mortality time series by using a thick red line.



Mortality in Chicago, July 1995. In the plot on the right, a thick red line has been added to show the dates of a heat wave.

In Figure @ref(fig:highlightpoints), the names of the players with the most shots and passes have been added to highlight these unusual points.



Passes versus shots for World Cup 2010 players. In the plot on the right, notable players have been highlighted.

You can add highlighting elements using geoms like `geom_text` and `geom_line`. Often, you will need to use a different dataframe for the geom. For example, you may want to create a subset of the original dataframe with points you want to highlight. You can pass that second dataframe to geoms using the `data` parameter in the geom. For example, to create the right plot in Figure @ref(fig:highlightpoints), we first created a subset with the players with the most shots and passes (including some pasting to create the label we want to use in the plot):

```

noteworthy_players <- worldcup %>%
  filter(Shots == max(Shots) | Passes == max(Passes)) %>%
  mutate(point_label = paste(Team, Position, sep = "-"))
noteworthy_players

```

| | Team | Position | Time | Shots | Passes | Tackles | Saves | point_label |
|---|-------|------------|------|-------|--------|---------|-------|------------------|
| 1 | Ghana | Forward | 501 | 27 | 151 | 1 | 0 | Ghana-Forward |
| 2 | Spain | Midfielder | 515 | 4 | 563 | 6 | 0 | Spain-Midfielder |

Now we can create a ggplot object based on the `worldcup` data, but use this subset for some of the geoms that highlight these players:

```
ggplot(worldcup, aes(x = Passes, y = Shots)) +  
  geom_point() +  
  geom_text(data = noteworthy_players, aes(label = point_label),  
            vjust = "inward", hjust = "inward", color = "red") +  
  geom_point(data = noteworthy_players, color = "red")
```

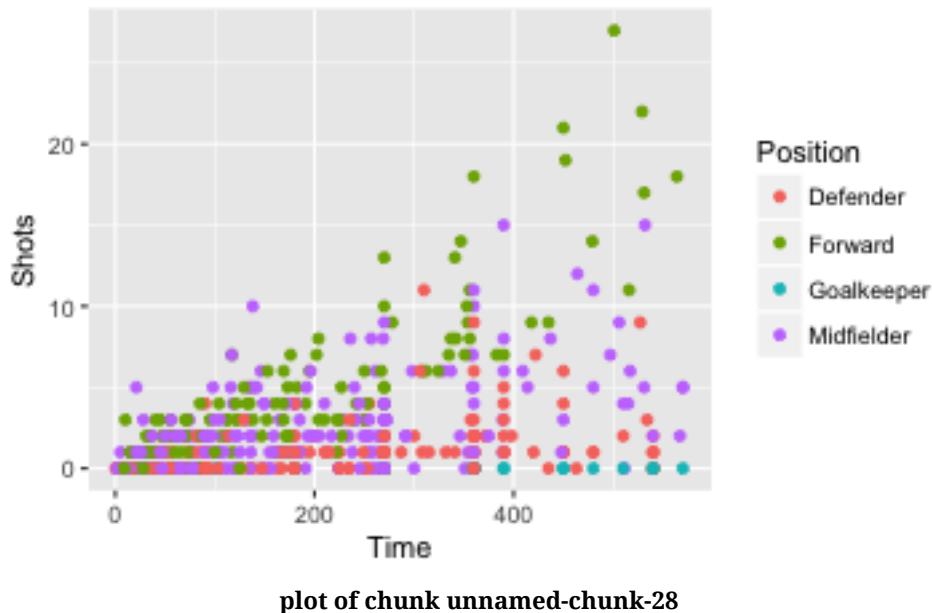
Small multiples

Guideline 5: When possible, use small multiples. \bigskip

Small multiples are graphs that use many small plots showing the mapping for different subsets of the data. Typically in small multiples, all plots use the same ranges for the x- and y-axes. This makes it easier to compare across plots, and it also allows you to save room by limiting axis annotation. Faceting creates multiple small plots, with each plot showing only data for one subset with the total data.

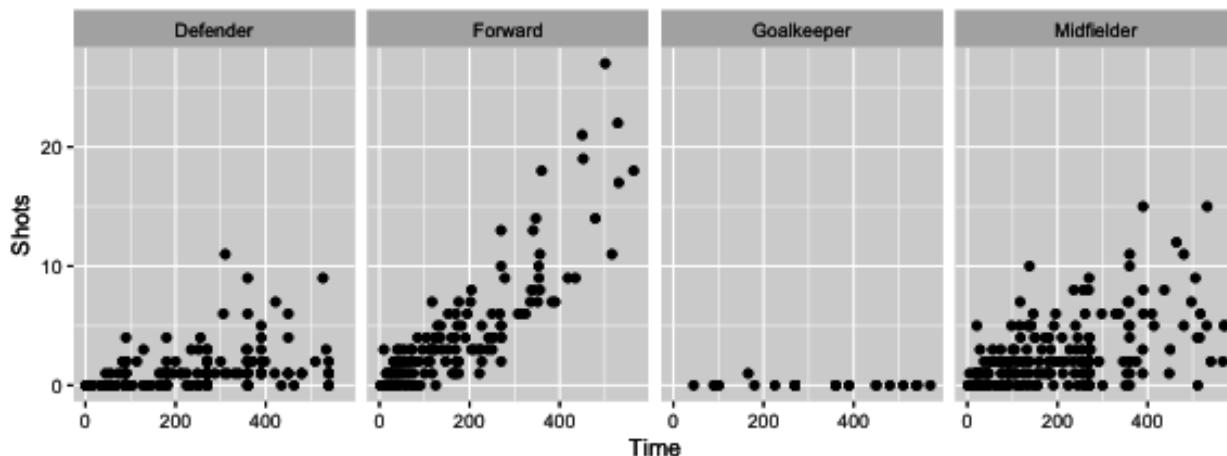
For example, the `worldcup` dataset used in earlier examples includes each player's position. If you want to explore a relationship (e.g., time played vs. shots on goal), you could try using color:

```
data(worldcup)
worldcup %>%
  ggplot(aes(x = Time, y = Shots, color = Position)) +
  geom_point()
```



However, often it's clearer to see relationships if you use faceting instead to create a small separate plot for each position. You can do this with either the `facet_grid` function or the `facet_wrap` function:

```
worldcup %>%
  ggplot(aes(x = Time, y = Shots)) +
  geom_point() +
  facet_grid(. ~ Position)
```



plot of chunk unnamed-chunk-29

You can create faceted plots using either `facet_grid` or `facet_wrap`. The two functions differ in whether it only allows facetting by one column, but allows the plots to wrap across several rows (`facet_wrap`) or whether it allows facetting by two columns but no wrapping (`facet_grid`)

The `facet_grid` function can facet by one or two variables. One will be shown by rows, and one by columns:

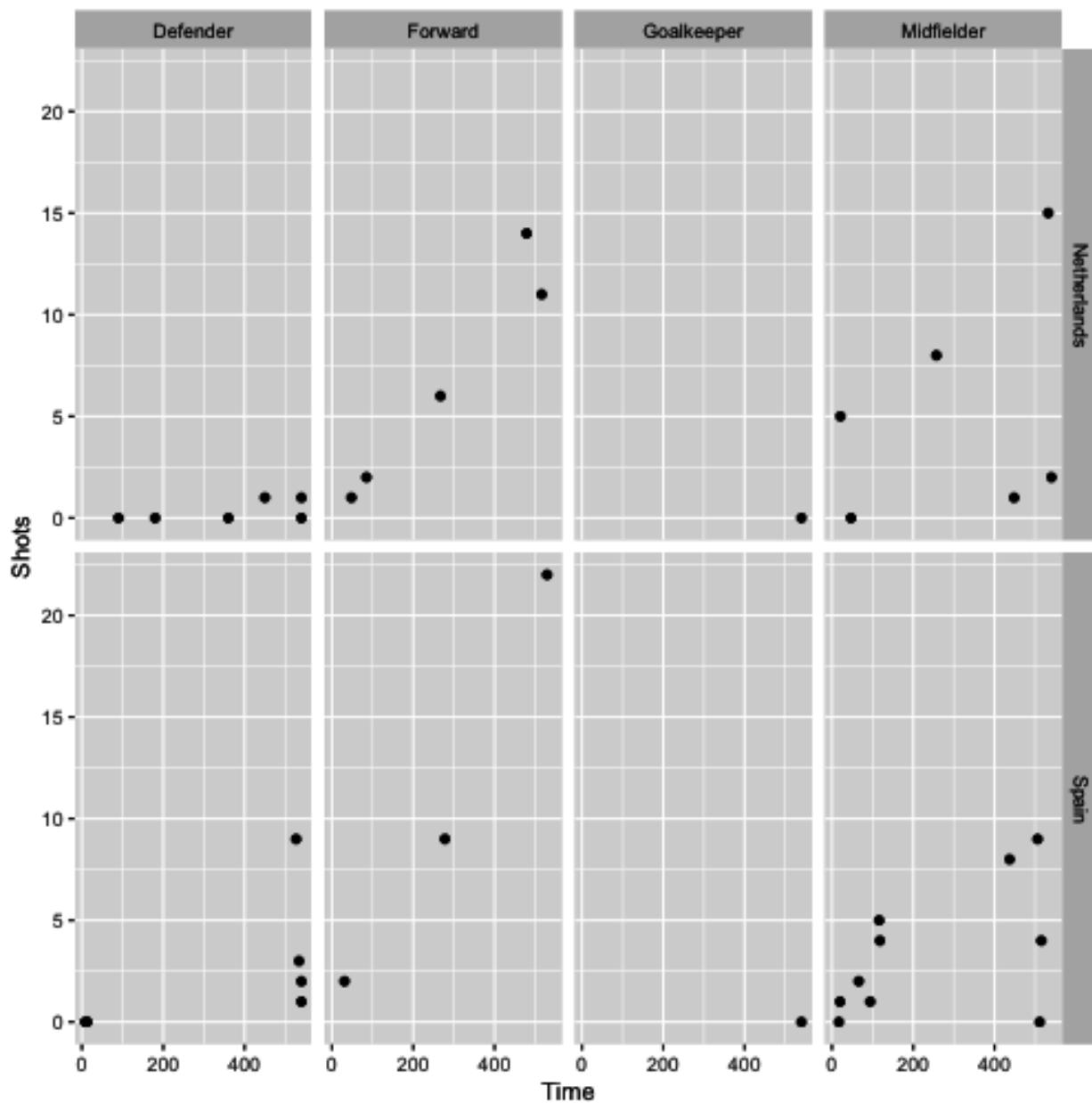
```
## Generic code
facet_grid([factor for rows] ~ [factor for columns])
```

The `facet_wrap()` function can only facet by one variable, but it can “wrap” the small graphs for that variable, so the don’t all have to be in one row or column:

```
## Generic code
facet_wrap(~ [factor for facetting], ncol = [number of columns])
```

For example, if you wanted to show relationships for the final two teams in World Cup 2010 (Spain and Holland), and facet by both position and team, you could run:

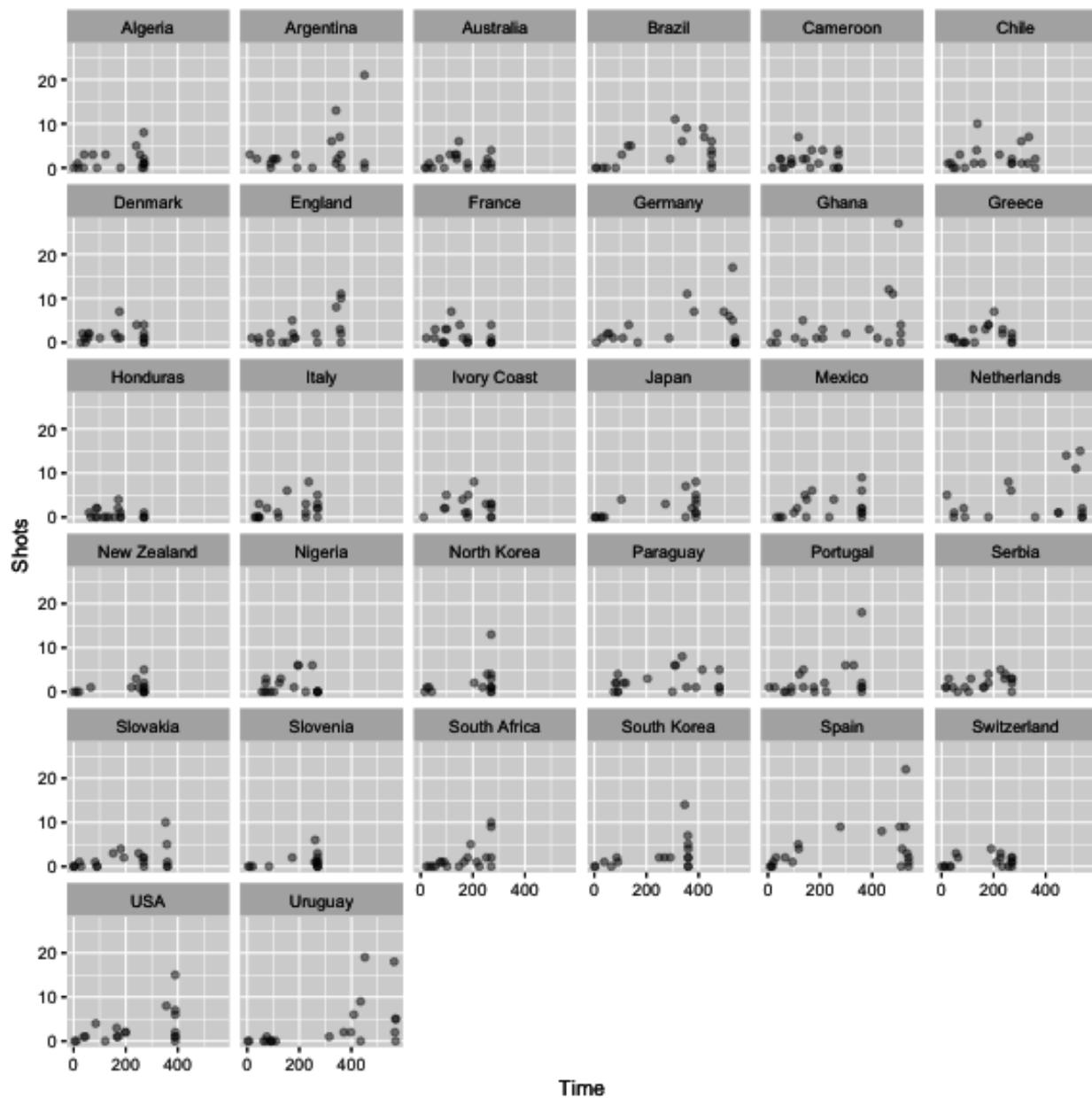
```
worldcup %>%
  filter(Team %in% c("Spain", "Netherlands")) %>%
  ggplot(aes(x = Time, y = Shots)) +
  geom_point() +
  facet_grid(Team ~ Position)
```



plot of chunk unnamed-chunk-32

With `facet_wrap`, you can only use a single variable for facetting. However, `facet_wrap` allows you to specify how many columns you want to use, which makes it useful if you want to facet across a variable with a lot of variables. For example, there are 32 teams in the World Cup. You can create a faceted graph of time played versus shots taken by team by running:

```
worldcup %>%
  ggplot(aes(x = Time, y = Shots)) +
  geom_point(alpha = 0.25) +
  facet_wrap(~ Team, ncol = 6)
```



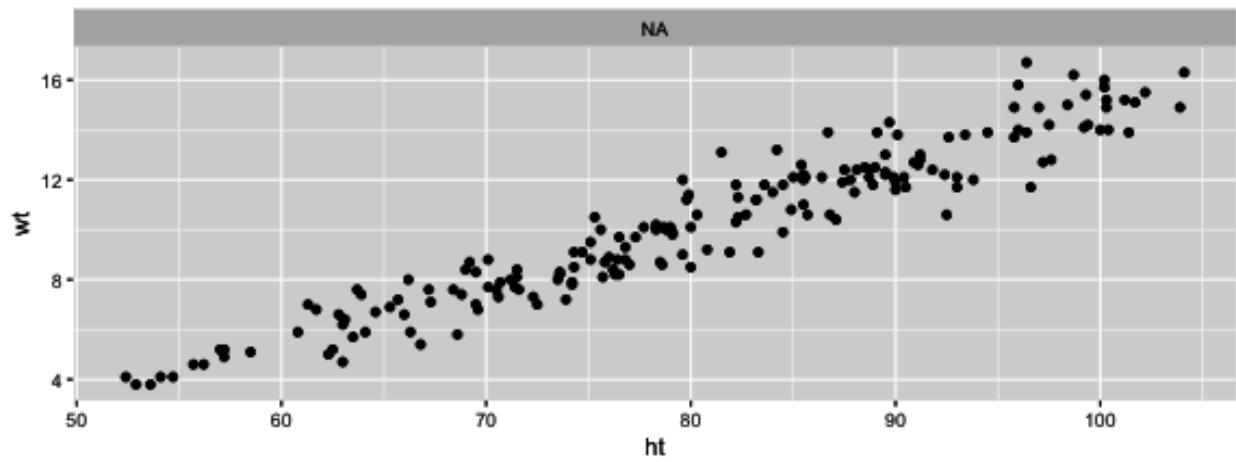
plot of chunk unnamed-chunk-33

Often, when you do faceting, you'll want to re-name your factors levels or re-order them. For this, you'll need to use the `factor()` function on the original vector, or use some of the tools from the `forcats` package. For example, to rename the `sex` factor levels from “1” and “2” to “Male” and “Female”, you can run:

```
nepali <- nepali %>%
  mutate(sex = factor(sex, levels = c(1, 2),
                      labels = c("Male", "Female")))
```

Notice that the labels for the two graphs have now changed:

```
ggplot(nepali, aes(ht, wt)) +
  geom_point() +
  facet_grid(. ~ sex)
```



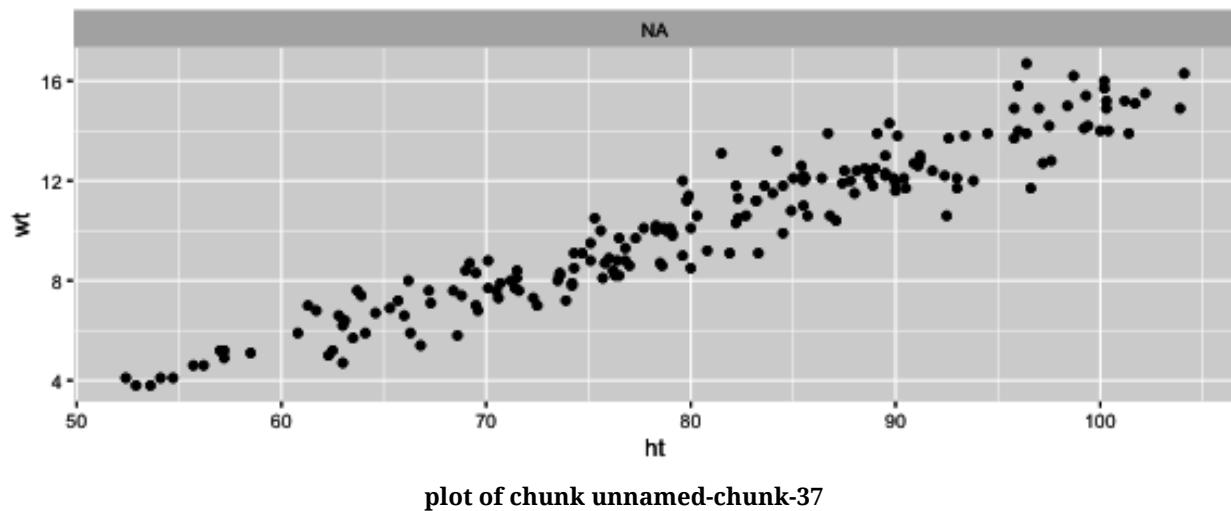
plot of chunk unnamed-chunk-35

To re-order the factor, and show the plot for “Female” first, you can use `factor` to change the order of the levels:

```
nepali <- nepali %>%
  mutate(sex = factor(sex, levels = c("Female", "Male")))
```

Now notice that the order of the plots has changed:

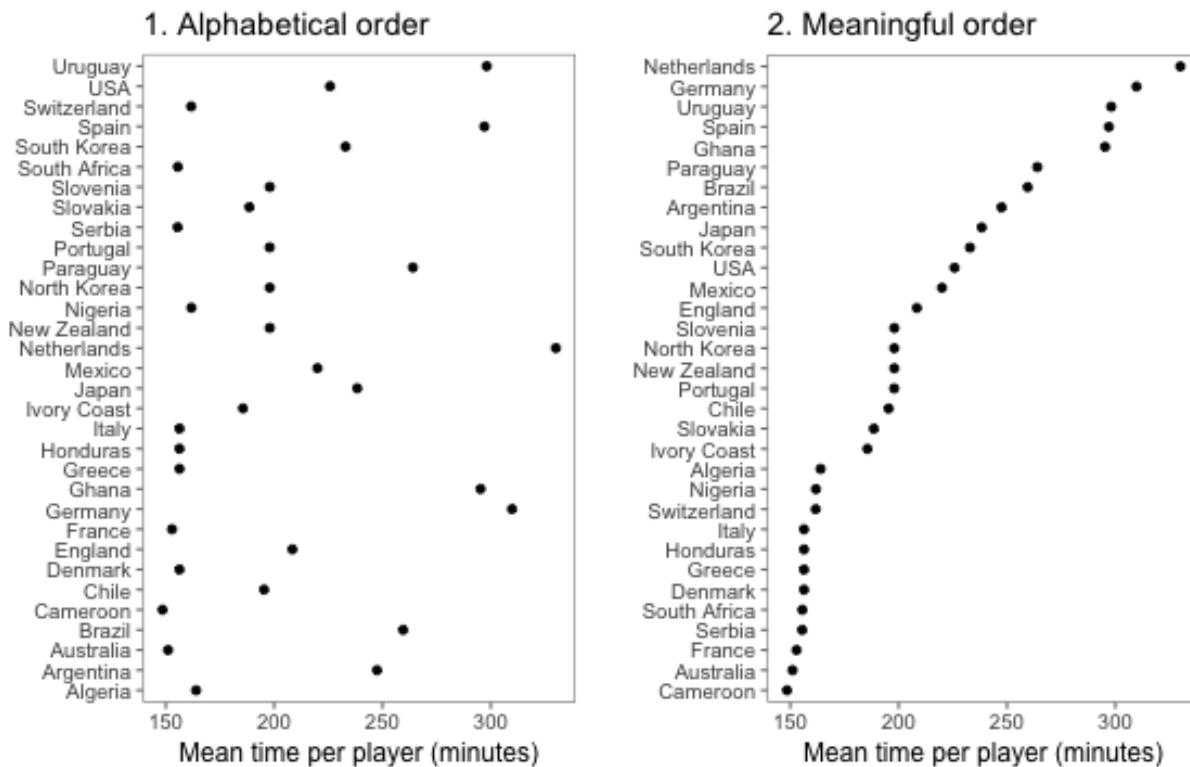
```
ggplot(nepali, aes(ht, wt)) +
  geom_point() +
  facet_grid(. ~ sex)
```



Order

Guideline 6: **Make order meaningful.**

Adding order to plots can help highlight interesting findings. Often, factor or categorical variables are ordered by something that is not interesting, like alphabetical order (Figure @ref(fig:plotorder), left plot).



Mean time per player in World Cup 2010 by team. The plot on the right has reordered teams to show patterns more clearly.

You can make the ranking of data clearer from a graph by using order to show rank (Figure @ref(fig:plotorder), right). You can re-order factor variables in a graph by resetting the factor using the `factor` function and changing the order that levels are included in the `levels` parameter. For example, here is the code for the two plots in Figure @ref(fig:plotorder):

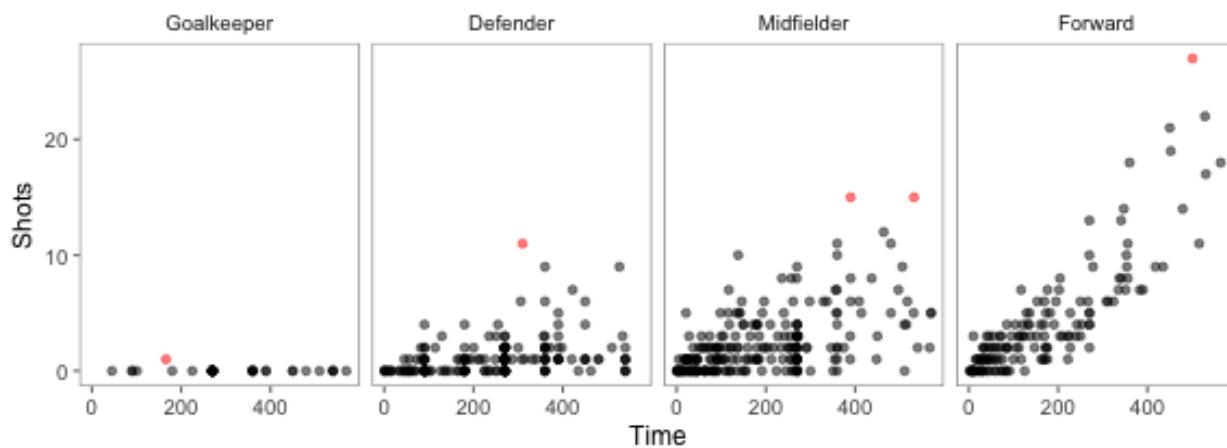
```
## Left plot
worldcup %>%
  group_by(Team) %>%
  summarize(mean_time = mean(Time)) %>%
  ggplot(aes(x = mean_time, y = Team)) +
  geom_point() +
  theme_few() +
  xlab("Mean time per player (minutes)") + ylab("")

## Right plot
worldcup %>%
  group_by(Team) %>%
  summarize(mean_time = mean(Time)) %>%
  arrange(mean_time) %>% # re-order and re-set
  mutate(Team = factor(Team, levels = Team)) %>% # factor levels before plotting
  ggplot(aes(x = mean_time, y = Team)) +
  geom_point() +
  theme_few() +
```

```
xlab("Mean time per player (minutes)") + ylab("")
```

As another example, you can customize the faceted plot created in the previous subsection to order these plots so from least to most average shots for a position using the following code. This example also has some added code to highlight the top players in each position in terms of shots on goal, as well as customizing colors and the theme.

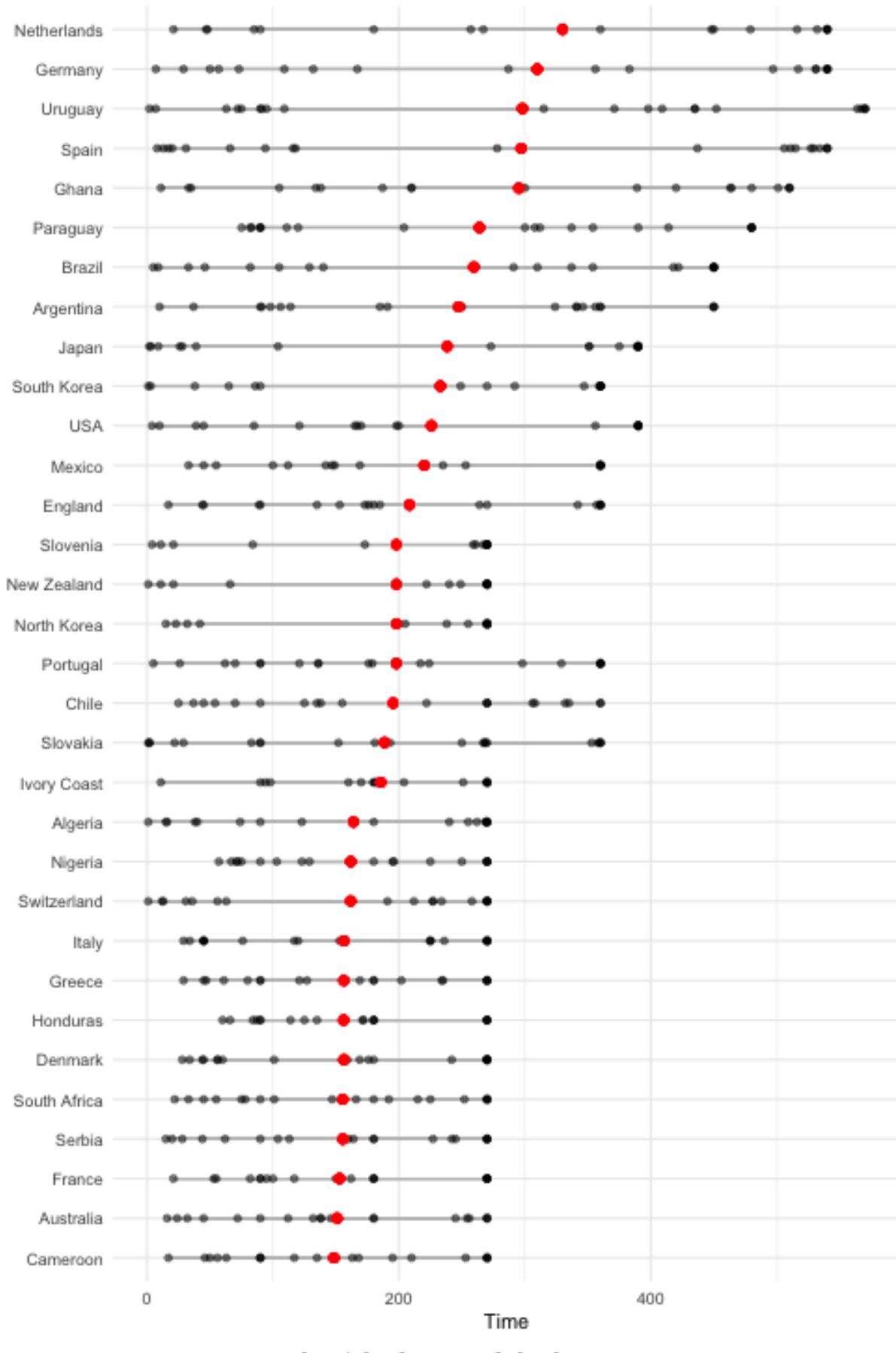
```
worldcup %>%
  dplyr::select(Position, Time, Shots) %>%
  dplyr::group_by(Position) %>%
  dplyr::mutate(ave_shots = mean(Shots),
               most_shots = Shots == max(Shots)) %>%
  dplyr::ungroup() %>%
  dplyr::arrange(ave_shots) %>%
  dplyr::mutate(Position = factor(Position, levels = unique(Position))) %>%
  ggplot(aes(x = Time, y = Shots, color = most_shots)) +
  geom_point(alpha = 0.5) +
  scale_color_manual(values = c("TRUE" = "red", "FALSE" = "black"),
                     guide = FALSE) +
  facet_grid(. ~ Position) +
  ggthemes::theme_few()
```



plot of chunk unnamed-chunk-39

As another example of ordering, suppose you wanted to show how playing times were distributed among players from each team for the World Cup data, with teams ordered by the average time for all their players. You can link up `dplyr` tools with `ggplot` to do this by using `group_by` to group the data by team, `mutate` to average player time within each team, `arrange` to order teams by that average player time, and `mutate` to reset the factor levels of the `Team` variable, using this new order, before plotting with `ggplot`:

```
worldcup %>%
  dplyr::select(Team, Time) %>%
  dplyr::group_by(Team) %>%
  dplyr::mutate(ave_time = mean(Time),
               min_time = min(Time),
               max_time = max(Time)) %>%
  dplyr::arrange(ave_time) %>%
  dplyr::ungroup() %>%
  dplyr::mutate(Team = factor(Team, levels = unique(Team))) %>%
  ggplot(aes(x = Time, y = Team)) +
  geom_segment(aes(x = min_time, xend = max_time, yend = Team),
               alpha = 0.5, color = "gray") +
  geom_point(alpha = 0.5) +
  geom_point(aes(x = ave_time), size = 2, color = "red", alpha = 0.5) +
  theme_minimal() +
  ylab("")
```



Scales and color

We'll finish this section by going into a bit more details about how to customize the scales and colors for ggplot objects, including more on scales and themes.

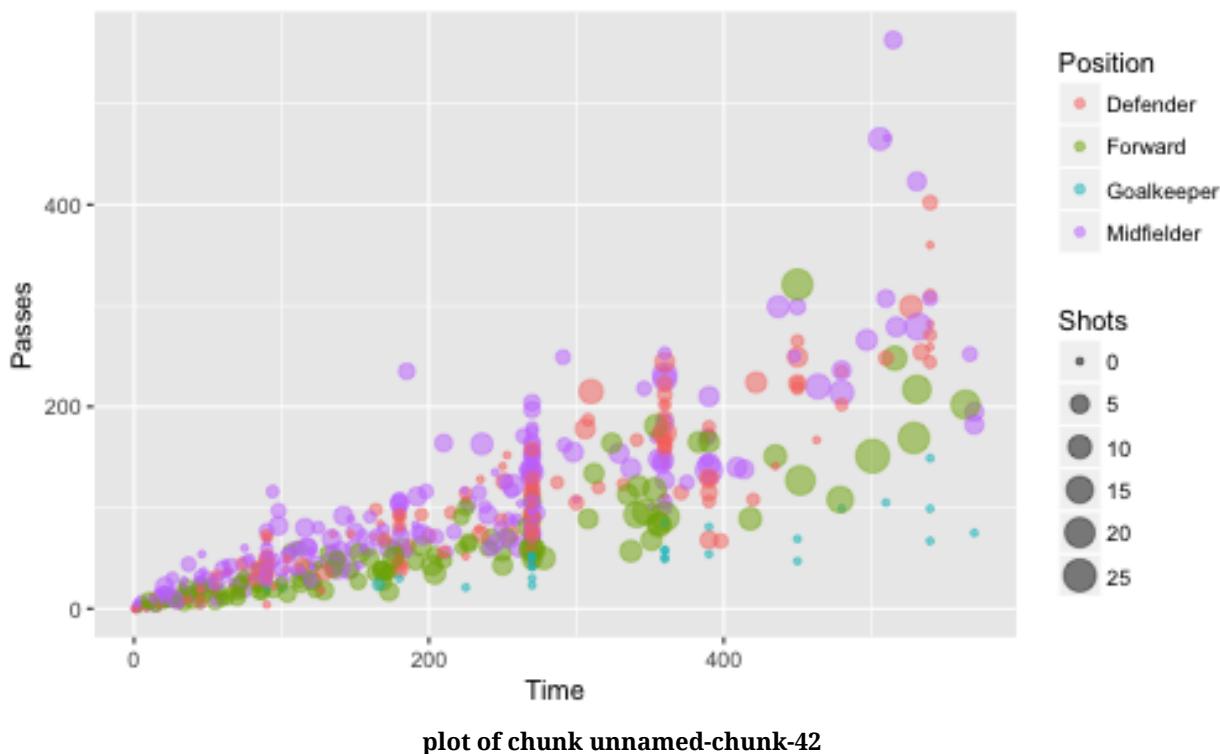
There are a number of different scale functions that allow you to customize the scales of ggplot objects. Because color is often mapped to an aesthetic, you can adjust colors in many ggplot objects using scales, as well (the exception is if you are using a constant color for an element). These functions follow the following convention:

```
## Generic code
scale_[aesthetic]_[vector type]
```

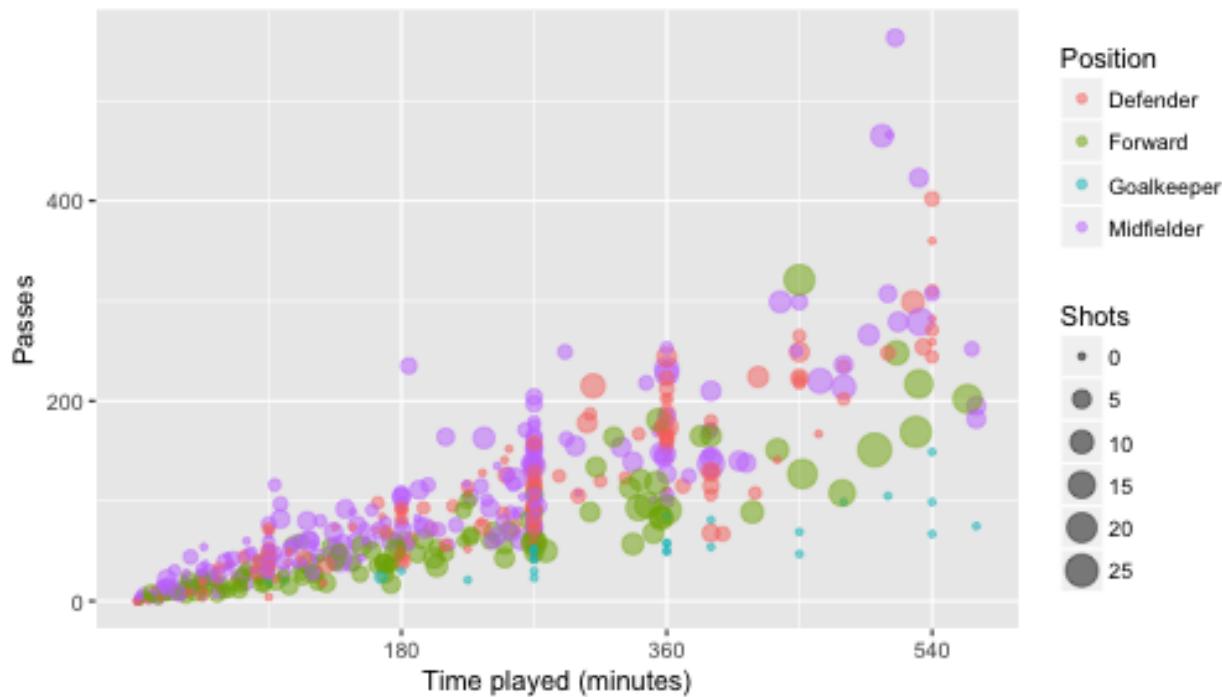
For example, to adjust the x-axis scale for a continuous variable, you'd use `scale_x_continuous`. You can use a `scale` function for an axis to change things like the axis label (which you could also change with `xlab` or `ylab`) as well as position and labeling of breaks.

For example, here is the default for plotting time versus passes for the `worldcup` dataset, with the number of shots taken shown by size and position shown by color:

```
ggplot(worldcup, aes(x = Time, y = Passes,
                      color = Position, size = Shots)) +
  geom_point(alpha = 0.5)
```



```
ggplot(worldcup, aes(x = Time, y = Passes,
                      color = Position, size = Shots)) +
  geom_point(alpha = 0.5) +
  scale_x_continuous(name = "Time played (minutes)",
                     breaks = 90 * c(2, 4, 6),
                     minor_breaks = 90 * c(1, 3, 5))
```



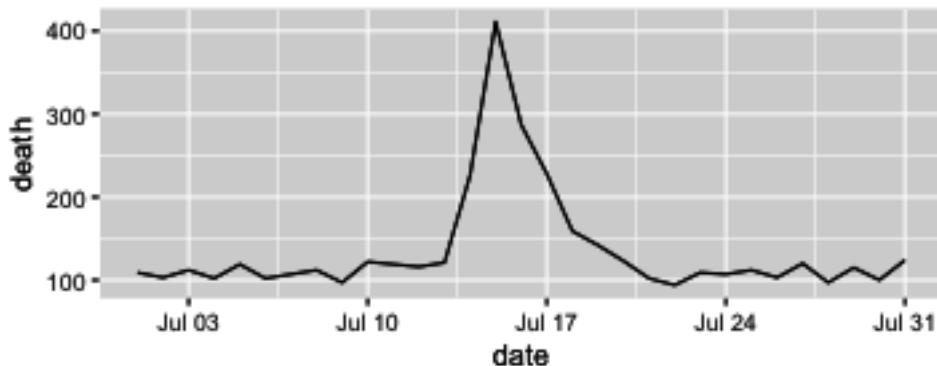
plot of chunk unnamed-chunk-43

Parameters you might find useful in `scale` functions include:

| Parameter | Description |
|--------------|---------------------------------|
| name | Label or legend name |
| breaks | Vector of break points |
| minor_breaks | Vector of minor break points |
| labels | Labels to use for each break |
| limits | Limits to the range of the axis |

For dates, you can use `scale` functions like `scale_x_date` and `scale_x_datetime`. For example, here's a plot of deaths in Chicago in July 1995 using default values for the x-axis:

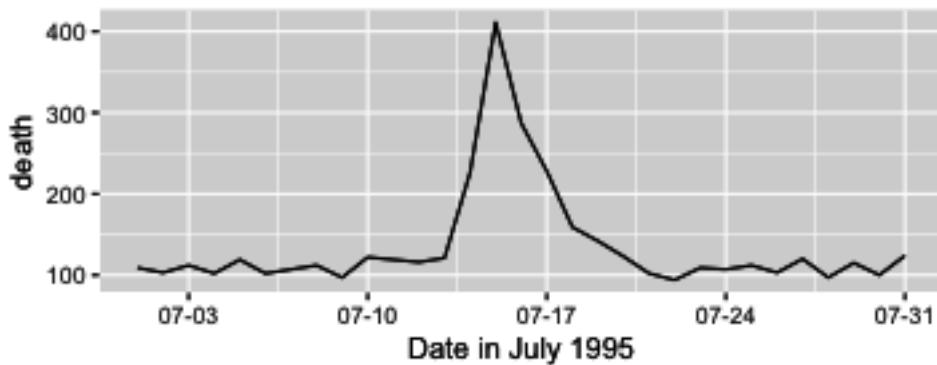
```
ggplot(chic_july, aes(x = date, y = death)) +
  geom_line()
```



plot of chunk unnamed-chunk-45

And here's an example of changing the formating and name of the x-axis:

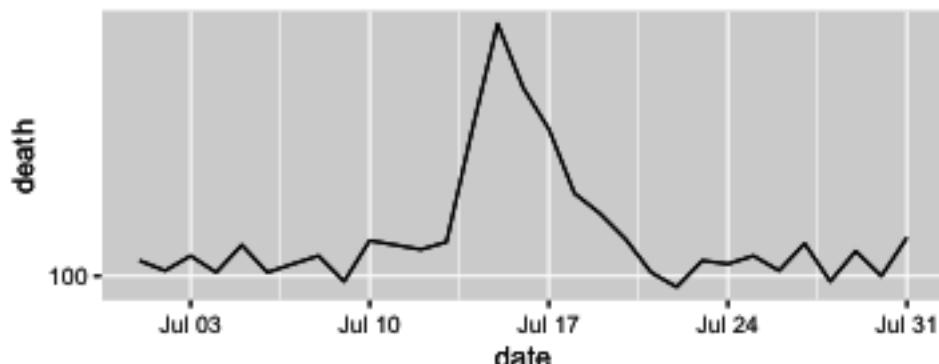
```
ggplot(chic_july, aes(x = date, y = death)) +  
  geom_line() +  
  scale_x_date(name = "Date in July 1995",  
               date_labels = "%m-%d")
```



plot of chunk unnamed-chunk-46

You can also use the `scale` functions to transform an axis. For example, to show the Chicago plot with "deaths" on a log scale, you can run:

```
ggplot(chic_july, aes(x = date, y = death)) +  
  geom_line() +  
  scale_y_log10()
```



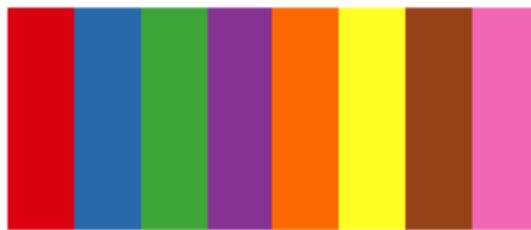
plot of chunk unnamed-chunk-47

For colors and fills, the conventions for the names of the `scale` functions can vary. For example, to adjust the color scale when you're mapping a discrete variable (i.e., categorical, like gender or animal breed) to color, you'd use `scale_color_hue`. To adjust the color scale for a continuous variable, like age, you'll use `scale_color_gradient`.

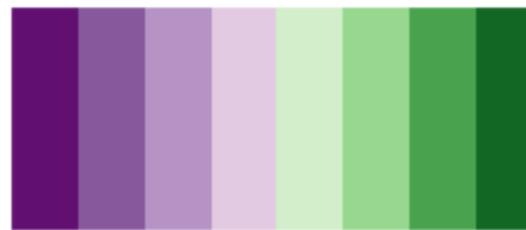
For any color scales, consider starting with `brewer` first (e.g., `scale_color_brewer`, `scale_color_distiller`). Scale functions from `brewer` allow you to set colors using different palettes. You can explore these palettes at <http://colorbrewer2.org/>.

The Brewer palettes fall into three categories: sequential, divergent, and qualitative. You should use sequential or divergent for continuous data and qualitative for categorical data. Use `display.brewer.pal` to show the palette for a given number of colors.

```
library(RColorBrewer)
display.brewer.pal(name = "Set1", n = 8)
display.brewer.pal(name = "PRGn", n = 8)
display.brewer.pal(name = "PuBuGn", n = 8)
```



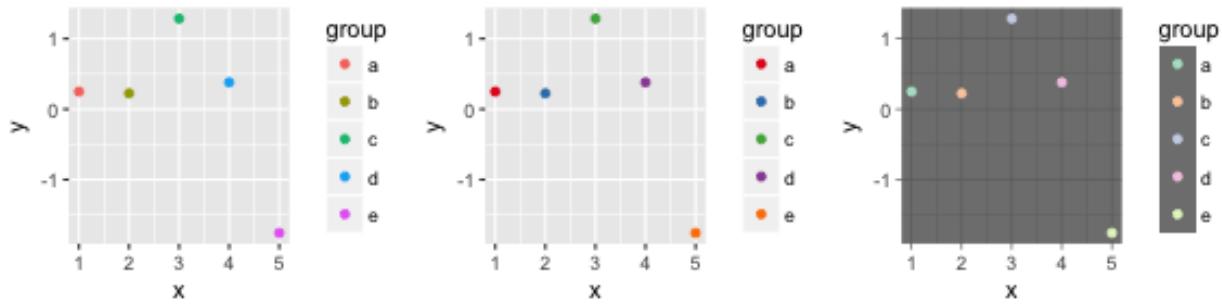
Set1 (qualitative)



PRGn (divergent)

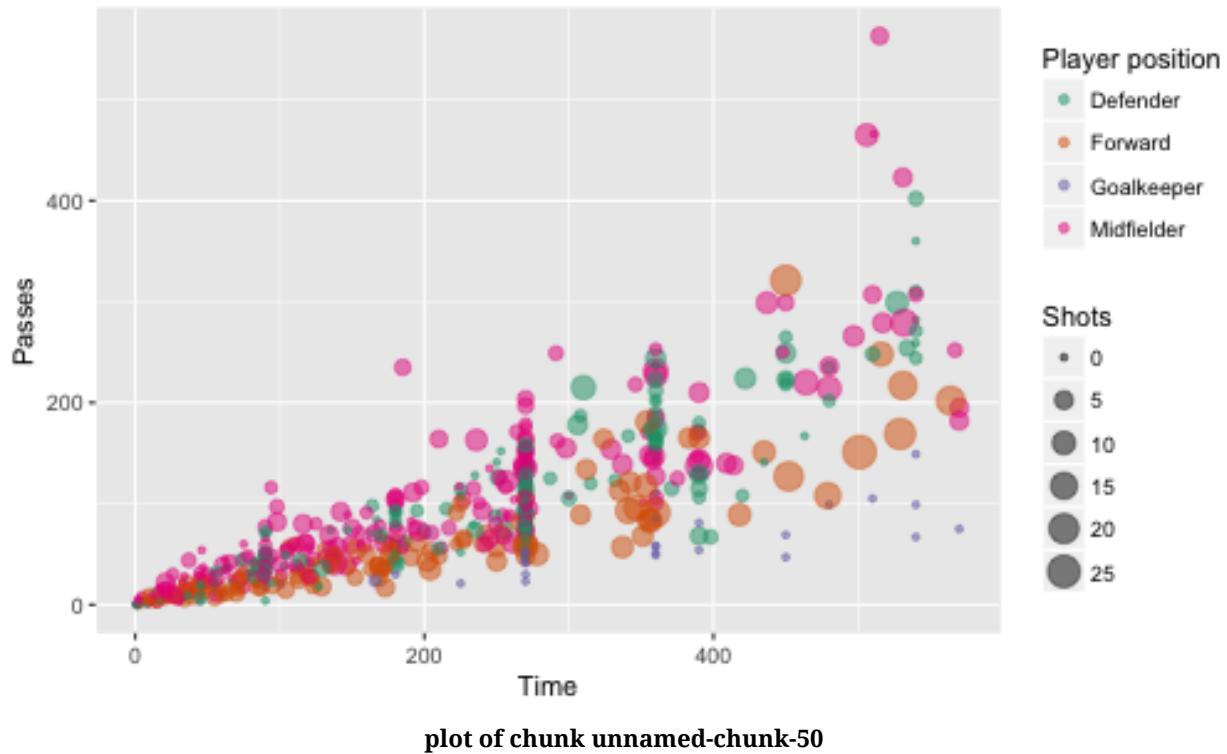
Use the `palette` argument within a `scales` function to customize the palette:

```
a <- ggplot(data.frame(x = 1:5, y = rnorm(5),
                        group = letters[1:5]),
             aes(x = x, y = y, color = group)) +
  geom_point()
b <- a + scale_color_brewer(palette = "Set1")
c <- a + scale_color_brewer(palette = "Pastel2") +
  theme_dark()
grid.arrange(a, b, c, ncol = 3)
```

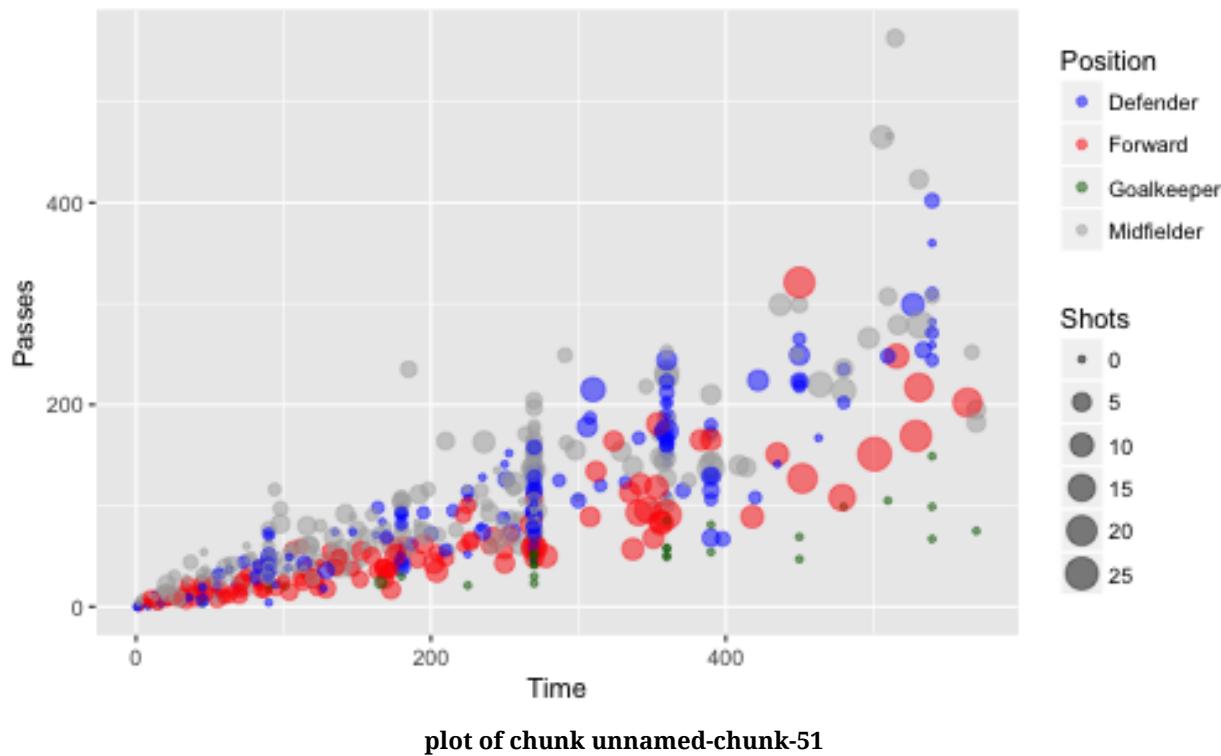


plot of chunk unnamed-chunk-49

```
ggplot(worldcup, aes(x = Time, y = Passes,
                      color = Position, size = Shots)) +
  geom_point(alpha = 0.5) +
  scale_color_brewer(palette = "Dark2",
                     name = "Player position")
```



You can also set colors manually:



Viridis color map

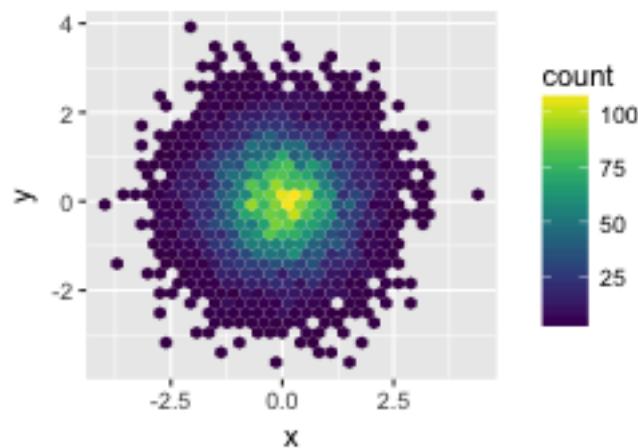
There is a package called `viridisLite` with some good color maps that are gaining popularity in visualization. From the package's GitHub repository:

“These four color maps are designed in such a way that they will analytically be perfectly perceptually-uniform, both in regular form and also when converted to black-and-white. They are also designed to be perceived by readers with the most common form of color blindness.”

- Viridis is now the default color map for `Matplotlib`, a key Python plotting library.
- The `viridisLite` package is a simpler version of the `viridis` package (and requires fewer dependencies).
- Several of the packages we'll look at today use Viridis as the default color map.

Here is an example of a hexbin graph of random values that uses the Viridis color map:

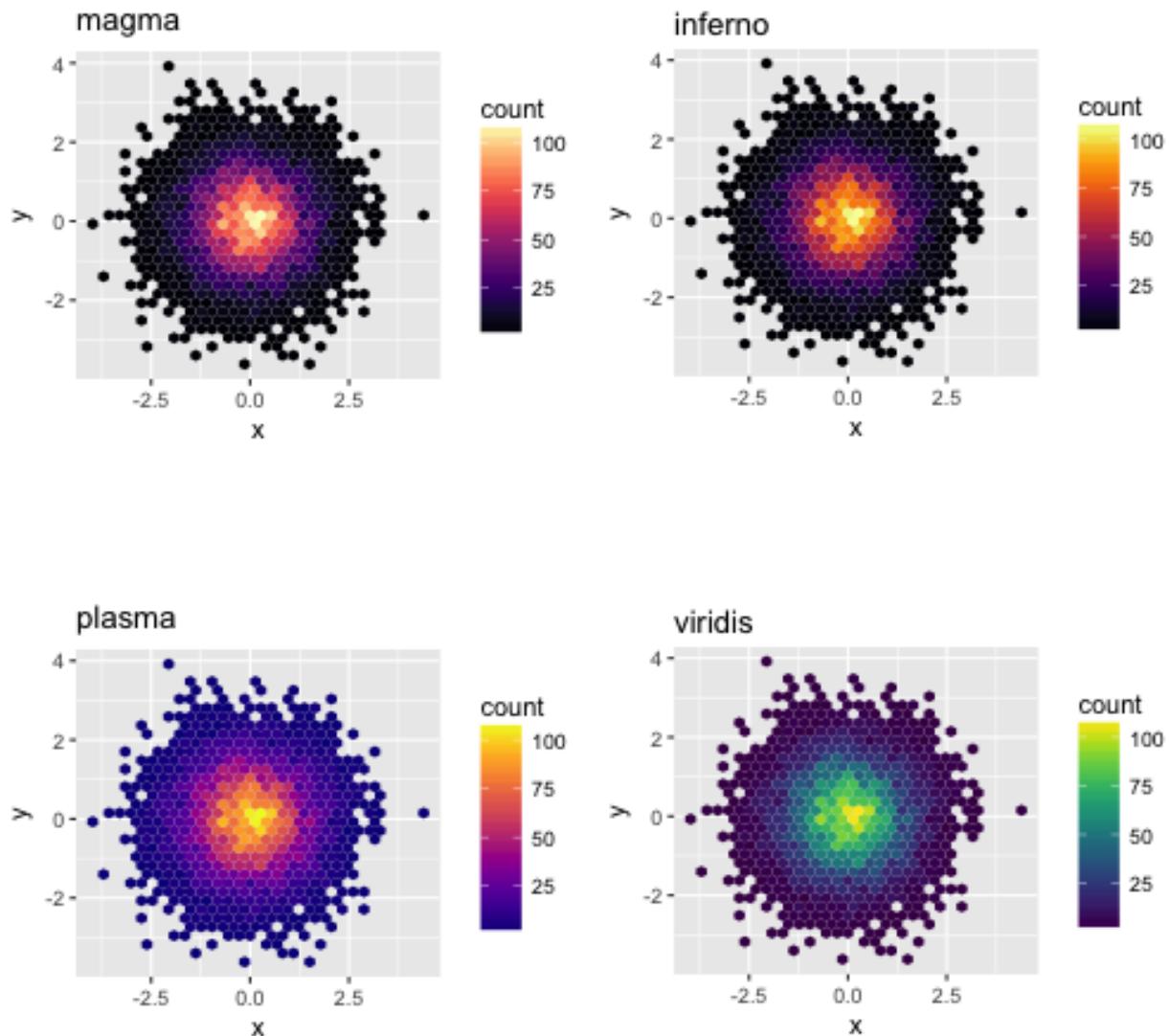
```
library(viridisLite)
library(hexbin)
Loading required package: methods
dat <- data.frame(x = rnorm(10000), y = rnorm(10000))
ggplot(dat, aes(x = x, y = y)) + geom_hex() + coord_fixed() +
  scale_fill_gradientn(colours = viridis(256, option = "D"))
```



plot of chunk unnamed-chunk-52

The `option` argument allows you to pick between four color maps: Magma, Inferno, Plasma, and Viridis. Here is the code to visualize each of those:

```
library(gridExtra)
ex_plot <- ggplot(dat, aes(x = x, y = y)) + geom_hex() + coord_fixed()
magma_plot <- ex_plot +
  scale_fill_gradientn(colours = viridis(256, option = "A")) +
  ggttitle("magma")
inferno_plot <- ex_plot +
  scale_fill_gradientn(colours = viridis(256, option = "B")) +
  ggttitle("inferno")
plasma_plot <- ex_plot +
  scale_fill_gradientn(colours = viridis(256, option = "C")) +
  ggttitle("plasma")
viridis_plot <- ex_plot +
  scale_fill_gradientn(colours = viridis(256, option = "D")) +
  ggttitle("viridis")
grid.arrange(magma_plot, inferno_plot, plasma_plot, viridis_plot, ncol = 2)
```



plot of chunk unnamed-chunk-53

To find out more

There are some excellent resources available for finding out more about creating plots using the `ggplot2` package.

If you want to get more practical tips on how to plot with `ggplot2`, these are good resources

- *R Graphics Cookbook* by Winston Chang: This “cookbook” style book is a useful reference to have to flip through when you have a specific task you want to figure

out how to do with ggplot2 (e.g., flip the coordinate axes, remove the figure legend).

- <http://www.cookbook-r.com/Graphs/>: Also created by Winston Chang, this website goes with the *R Graphics Cookbook* and is an excellent reference for quickly finding out how to do something specific in ggplot2.
- Google images: If you want to find example code for how to create a specific type of plot in R, try googling the name of the plot and “R”, and then search through the “Images” results. For example, if you wanted to plot a wind rose in R, google “wind rose r” and click on the “Images” tab. Often, the images that are returned will link back to a page that includes the example code to create the image (a blog post, for example).

For more technical details about plotting in R, check out the following resources:

- *ggplot2: Elegant Graphics for Data Analysis* by Hadley Wickham: Now in its second edition, this book was written by the creator of grid graphics and goes deeply into the details of why ggplot2 was created and how to use it.
- *R Graphics* by Paul Murrell: Also in its second edition, this book explains grid graphics, the graphics system that ggplot2 is built on. This course covers the basics of grid graphics in a later section to give you the tools to create your own ggplot2 extensions. However, if you want the full details on grid graphics, this book is where to find them.

4.3 Mapping

Often, data will include a spatial component, and you will want to map the data either for exploratory data analysis or to present interesting aspects of the data to others. R has a range of capabilities for mapping data. The simplest techniques involve using data that includes latitude and longitude values, and use these location values as the `x` and `y` aesthetics in a regular plot. R also has the ability to work with more complex spatial data objects and import shapefiles through extensions like the `sp` package.

In this section, we will cover the basics of mapping in R and touch on some of the more advanced possibilities. We will also present some useful packages for making quick but attractive maps in R. R also now has the capability to make interactive maps using the `plotly` and `leaflet` packages; in the end of this section, we'll present these packages and explain a bit more about `htmlWidgets` in general.

Basics of mapping

Point maps

It is very easy now to create point maps in R based on longitude and latitude values of specific locations. You can use the `map_data` function from the `ggplot2` package to pull data for maps at different levels (“usa”, “state”, “world”, “county”).

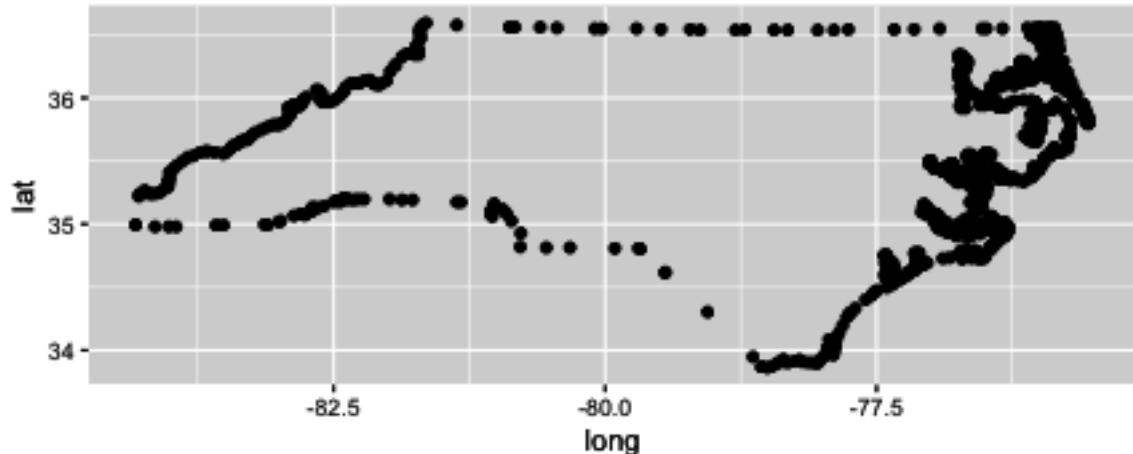
The maps you pull using `map_data` are just data to use to plot polygon shapes for areas like states and counties.

```
library(ggplot2)
us_map <- map_data("state")
head(us_map, 3)
  long      lat group order  region subregion
1 -87.46201 30.38968     1     1 alabama      <NA>
2 -87.48493 30.37249     1     2 alabama      <NA>
3 -87.52503 30.37249     1     3 alabama      <NA>
```

You can add points to these based on latitude and longitude.

Mapping uses the `long` and `lat` columns from this data for location:

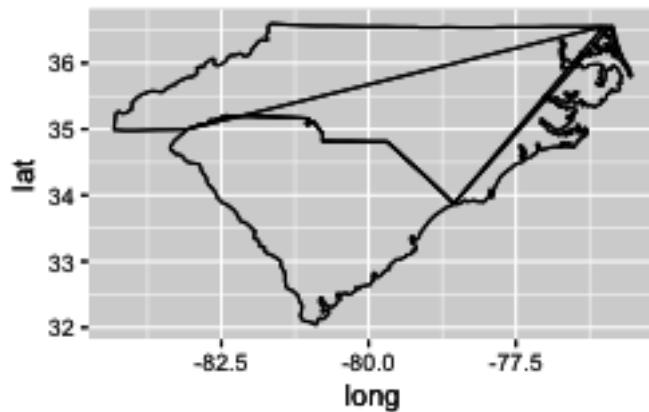
```
north_carolina <- us_map %>%
  filter(region == "north carolina")
ggplot(north_carolina, aes(x = long, y = lat)) +
  geom_point()
```



plot of chunk cav

If you try to plot lines, however, you'll have a problem:

```
library(stringr)
carolinias <- us_map %>%
  filter(str_detect(region, "carolina"))
ggplot(carolinias, aes(x = long, y = lat)) +
  geom_path()
```



plot of chunk caw

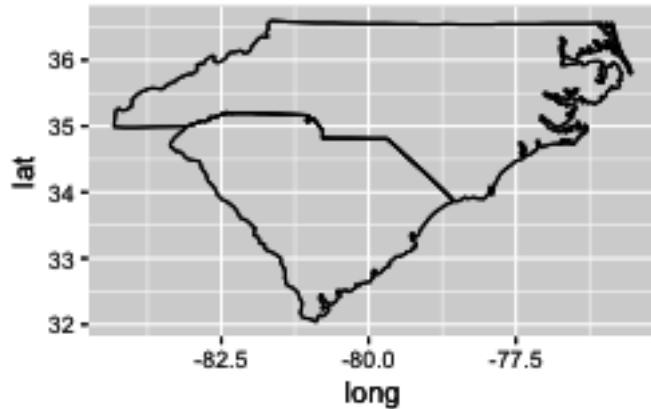
The `group` column fixes this problem. It will plot a separate path or polygon for each separate group. For mapping, this gives separate groupings for mainland versus islands and for different states:

```
carolinas %>%
  group_by(group) %>%
  slice(1)
Source: local data frame [4 x 6]
Groups: group [4]

  long      lat group order      region subregion
  <dbl>    <dbl> <dbl> <int>    <chr>   <chr>
1 -75.89399 36.55471     38  9549 north carolina knotts
2 -78.55824 33.86753     39  9587 north carolina main
3 -76.00285 36.55471     40 10321 north carolina spit
4 -83.10753 34.99053     47 11441 south carolina <NA>
```

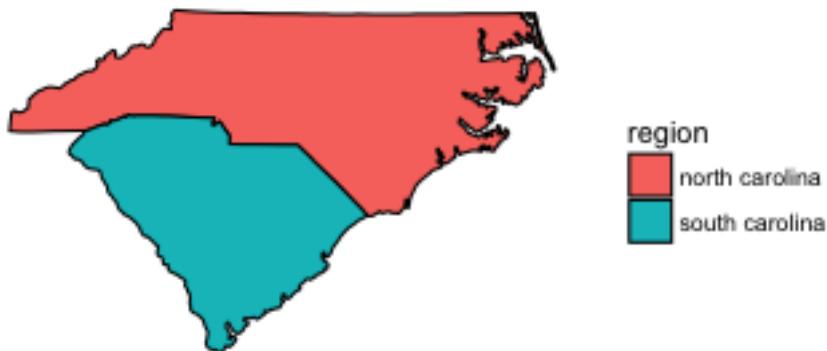
Using `group = group` avoids the extra lines from the earlier map:

```
ggplot(carolinas, aes(x = long, y = lat,
                      group = group)) +
  geom_path()
```

**plot of chunk cay**

To plot filled regions, use `geom_polygon` with `fill = region`. Also, the “void” theme is often useful when mapping:

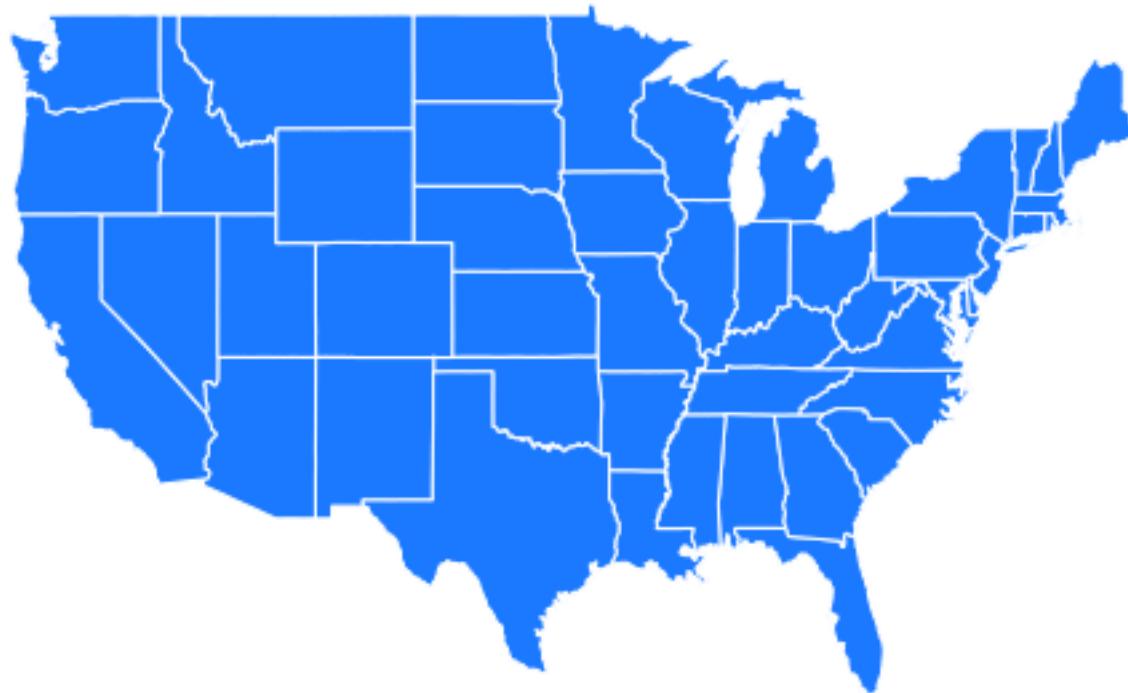
```
ggplot(carolinas, aes(x = long, y = lat,
                      group = group,
                      fill = region)) +
  geom_polygon(color = "black") +
  theme_void()
```

**plot of chunk caz**

Here is an example of plotting all of the US by state:

```
map_1 <- ggplot(us_map, aes(x = long, y = lat,
                             group = group)) +
  geom_polygon(fill = "dodgerblue",
               color = "white") +
  theme_void()
```

```
map_1
```



plot of chunk cbb

To add points to these maps, you can use `geom_point`, again using longitude and latitude to define position. \bigskip

Here I'll use an example of data points related to the story told in last year's “Serial” podcast.

```
serial <- read_csv(paste0("https://raw.githubusercontent.com/dgrtwo/serial-ggvis/",
                           "master/input_data/serial_podcast_data/serial_map_data.csv"))

head(serial, 3)
# A tibble: 3 × 5
  x     y     Type  Name Description
<int> <int> <chr> <chr>      <chr>
1   356   437 cell-site L688       <NA>
2   740   360 cell-site L698       <NA>
3   910   340 cell-site L654       <NA>
```

David Robinson figured out a way to convert the x and y coordinates in this data to latitude and longitude coordinates. I'm also adding a column for whether or not the point is a cell tower.

```

serial <- serial %>%
  mutate(long = -76.8854 + 0.00017022 * x,
        lat = 39.23822 + 1.371014e-04 * y,
        tower = Type == "cell-site")

serial[c(1:2, (nrow(serial) - 1):nrow(serial)),
       c("Type", "Name", "long", "lat", "tower")]
# A tibble: 4 × 5
  Type      Name     long     lat tower
  <chr>    <chr>   <dbl>   <dbl> <lgl>
1 cell-site L688 -76.82480 39.29813 TRUE
2 cell-site L698 -76.75944 39.28758 TRUE
3 base-location Adnan's house -76.76284 39.30622 FALSE
4 base-location Jenn's house -76.72301 39.29443 FALSE

```

Now I can map just Baltimore City and Baltimore County in Maryland and add these points. I used `map_data` to pull the “county” map and specified “region” as “maryland”, to limit the map just to Maryland counties.

```

baltimore <- map_data('county', region = 'maryland')
head(baltimore, 3)
  long      lat group order  region subregion
1 -78.64992 39.53982     1     1 maryland allegany
2 -78.70148 39.55128     1     2 maryland allegany
3 -78.74159 39.57420     1     3 maryland allegany

```

From that, I subset out rows where the `subregion` column was “baltimore city” or “baltimore”.

```

baltimore <- subset(baltimore,
                     subregion %in% c("baltimore city",
                                     "baltimore"))

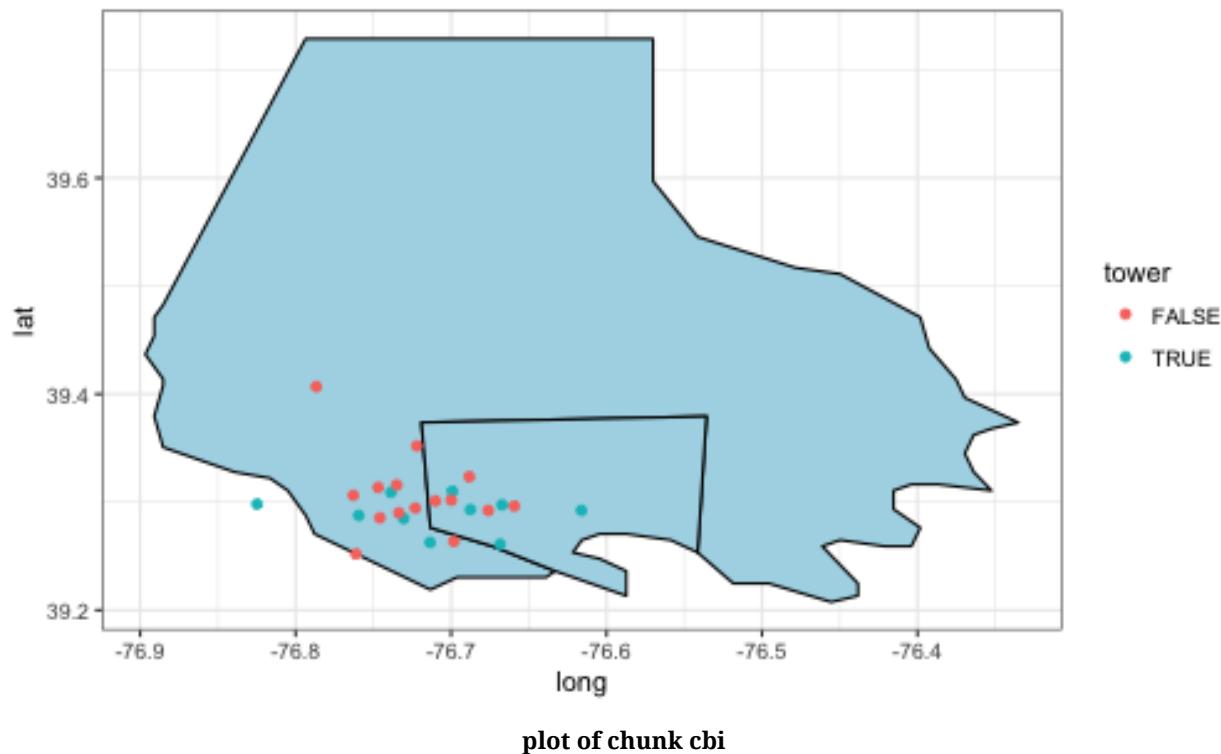
head(baltimore, 3)
  long      lat group order  region subregion
114 -76.88521 39.35074     3    114 maryland baltimore
115 -76.89094 39.37939     3    115 maryland baltimore
116 -76.88521 39.40804     3    116 maryland baltimore

```

I used `geom_point` to plot the points. `ggplot` uses the `group` column to group together counties, but we don’t need that in the points, so I needed to set `group = NA` in the `geom_point` statement. I put `color = tower` inside the `aes` statement so that the points would be one color for cell towers and another color for everything else.

```
balt_plot <- ggplot(baltimore,
                     aes(x = long, y = lat, group = group)) +
  geom_polygon(fill = "lightblue", color = "black") +
  geom_point(data = serial, aes(x = long, y = lat,
                                group = NA,
                                color = tower)) +
  theme_bw()
```

```
balt_plot
```



ggmap, Google Maps API

The `ggmap` package allows you to use tools from Google Maps directly from R.

```
## install.packages("ggmap")
library(ggmap)
```

This package uses the Google Maps API, so you should read their [terms of service](#) and make sure you follow them. In particular, you are limited to just a certain number of queries per time.

You can use the `get_map` function to get maps for different locations. You can either use the longitude and latitude of the center point of the map, along with the `zoom` option to say how

much to zoom in (3: continent to 20: building) or you can use a character string to specify a location.

If you do the second, `get_map` will actually use the Google Maps API to geocode the string to a latitude and longitude and then get the map (you can imagine that this is like searching in Google Maps in the search box for a location).

```
beijing <- get_map("Beijing", zoom = 12)
ggmap(beijing)
```



plot of chunk cby

With this package, you can get maps from the following different sources:

- Google Maps
- OpenStreetMap
- Stamen Maps
- CloudMade Maps (You may need a separate API key for this)

Here are different examples of Beijing using different map sources. (Also, note that I'm using the option `extent = "device"` to fill up the whole plot area with the map, instead of including axis labels and titles.)

```
beijing_a <- get_map("Beijing", zoom = 12,
                      source = "stamen", maptype = "toner")
a <- ggmap(beijing_a, extent = "device")

beijing_b <- get_map("Beijing", zoom = 12,
                      source = "stamen", maptype = "watercolor")
b <- ggmap(beijing_b, extent = "device")

beijing_c <- get_map("Beijing", zoom = 12,
                      source = "google", maptype = "hybrid")
c <- ggmap(beijing_c, extent = "device")

gridExtra::grid.arrange(a, b, c, nrow = 1)
```



plot of chunk ccb

As with the maps from `ggplot2`, you can add points to these maps:

```
serial_phone <- read_csv(paste0("https://raw.githubusercontent.com/dgrtwo/",
                                "serial-ggvis/master/input_data/",
                                "serial_podcast_data/serial_phone_data.csv")) %>%
  mutate(Cell_Site = substring(Cell_Site, 1, 4),
         Call_Time = as.POSIXct(Call_Time, format = "%d/%m/%y %H:%M",
                                tz = "EST")) %>%
  left_join(serial, by = c("Cell_Site" = "Name")) %>%
  select(Person_Called, Call_Time, Duration, long, lat) %>%
  filter(!(Person_Called %in% c("incoming", "# + Adnan cell"))) %>%
  arrange(Call_Time)

serial_map <- get_map(c(-76.7, 39.3), zoom = 12,
                      source = "stamen",
                      maptype = "toner")
serial_map <- ggmap(serial_map, extent = "device") +
  geom_point(data = serial_phone,
             aes(x = long, y = lat),
             color = "red", size = 3,
             alpha = 0.4) +
  geom_point(data = subset(serial,
                          Type != "cell-site"),
             aes(x = long, y = lat),
             color = "darkgoldenrod1",
             size = 2)
```



plot of chunk ccg

You can also use the Google Maps API, through the geocode function, to get the latitude and longitude of specific locations. Basically, if the string would give you the right location if you typed it in Google Maps, geocode should be able to geocode it.

For example, you can get the location of CSU:

```
geocode("Colorado State University")
       lon lat
1     NA  NA
```

You can also get a location by address:

```
geocode("1 First St NE, Washington, DC")
       lon      lat
1 -77.00465 38.89051
```

You can get distances, too, using the mapdist function with two locations. This will give you distance and also time.

```
mapdist("Fort Collins CO",
       "1 First St NE, Washington, DC") %>%
  select(from, miles, hours)
  from     miles    hours
1 Fort Collins CO 1670.349 24.6225
```

Mapping US counties and states

If you need to map US states and counties, the `choroplethr` and `choroplethrMaps` packages offer functions for fast and straightforward mapping.

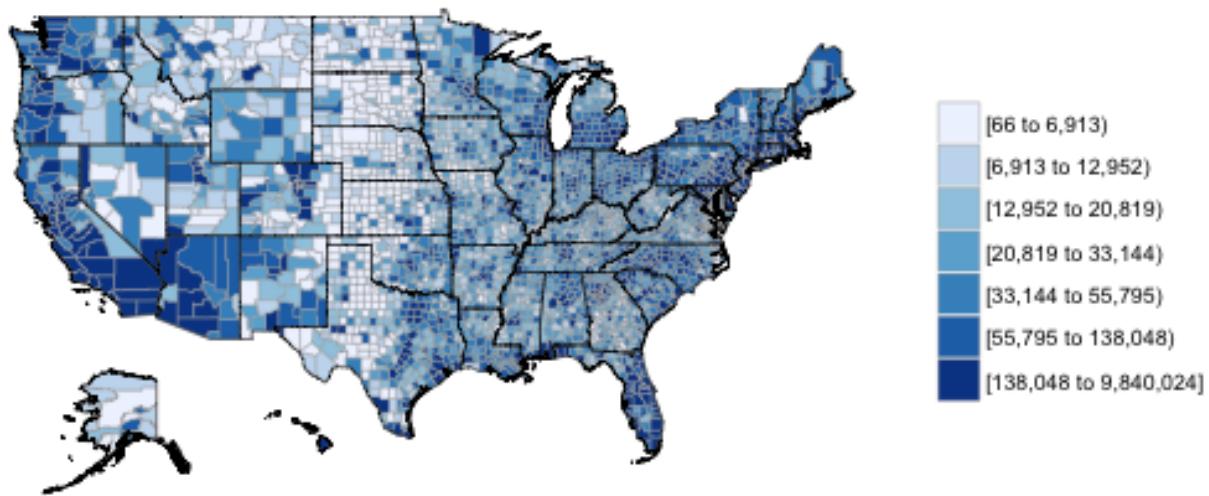
```
library(choroplethr)
library(choroplethrMaps)
```

As an example, I'll use data on county-level population in 2012 that comes as the dataset `df_pop_county` with the `choroplethr` package. This dataset gives the population of each county (`value`) and the county FIPS number (`region`), which is a unique identification number for each US county.

```
data(df_pop_county)
df_pop_county %>% slice(1:3)
  region  value
1 1001 54590
2 1003 183226
3 1005 27469
```

To map population by county, you can use the `countyChoropleth` function.

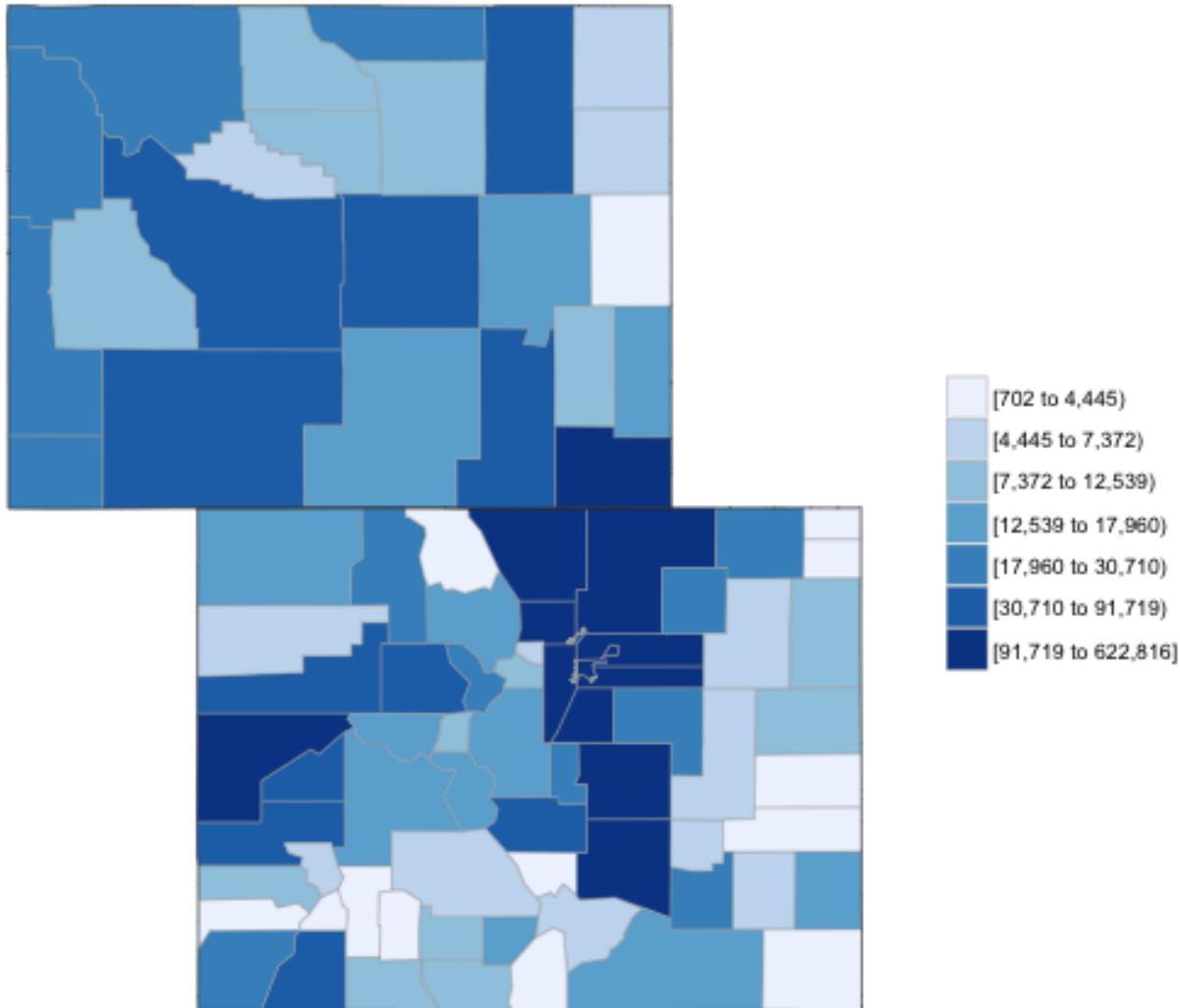
```
county_choropleth(df_pop_county)
```



plot of chunk unnamed-chunk-56

If you want to only plot some of states, you can use the `state_zoom` argument:

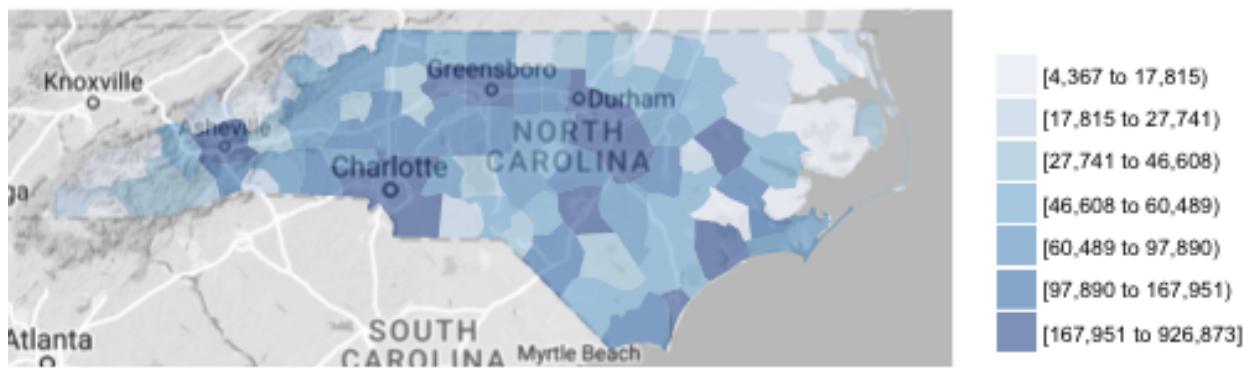
```
county_choropleth(df_pop_county, state_zoom = c("colorado", "wyoming"))
```



plot of chunk unnamed-chunk-57

To plot values over a reference map from Google Maps, you can use the `reference_map` argument:

```
county_choropleth(df_pop_county, state_zoom = c("north carolina"),
                   reference_map = TRUE)
```



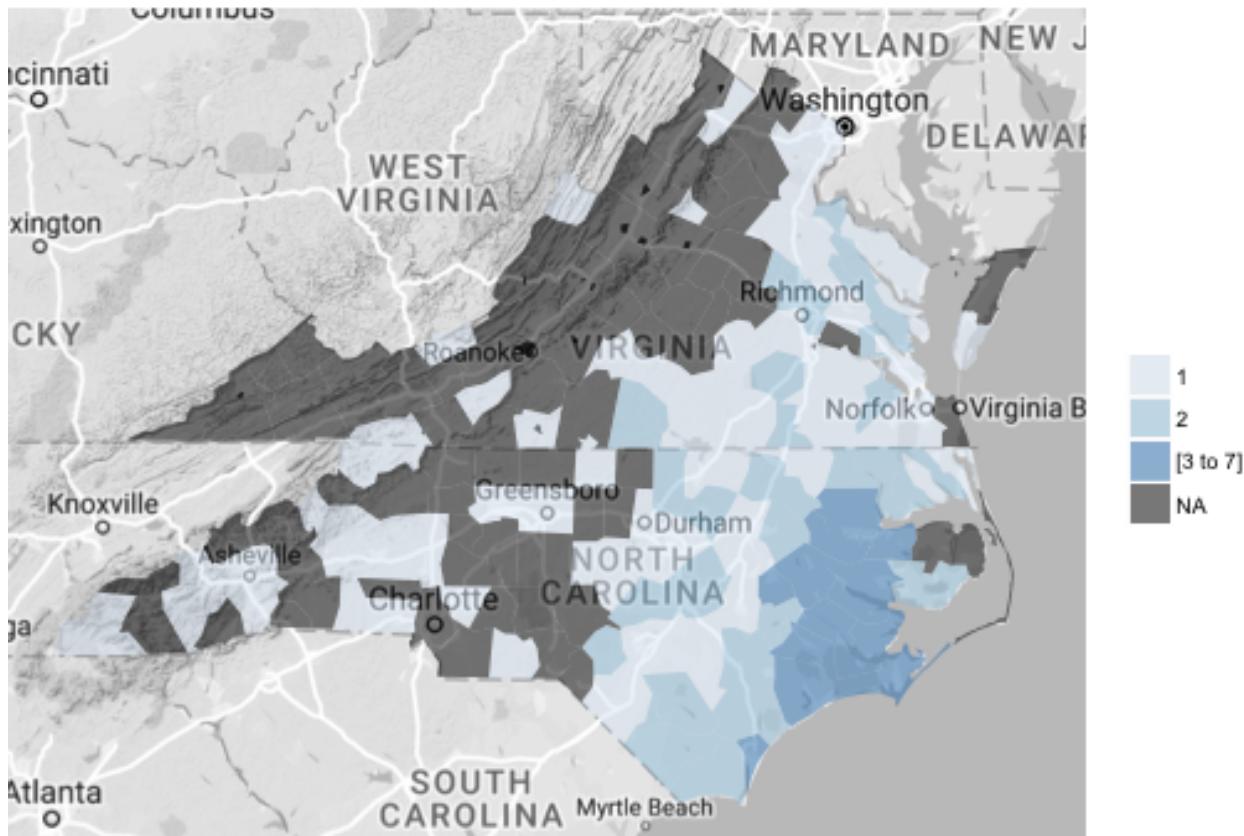
plot of chunk unnamed-chunk-58

This example is using one of the datasets that comes with the `choroplethr` package, but you can map any dataset that includes a column with county FIPS and a column with the value you would like to plot. All you have to do is make sure the county FIPS is in a numeric class and give the columns for FIPS and the value to plot as “region” and “value”, respectively (the `rename` function from `dplyr` is useful here). For example, here is a dataframe giving storm events that were listed in NOAA’s Storm Events database for near Hurricane Floyd’s track in the three days when the storm was closest to each county:

```
library(readr)
floyd_events <- read_csv("data/floyd_events.csv")
floyd_events %>% slice(1:3)
# A tibble: 3 × 4
  begin_date   end_date   fips      type
  <date>       <date>     <chr>    <chr>
1 1999-09-16 1999-09-17 25011 Heavy Rain
2 1999-09-16 1999-09-17 25001 Heavy Rain
3 1999-09-16 1999-09-17 25015 Heavy Rain
```

You can use the following code to plot the number of events listed for each US county:

```
floyd_events %>%
  dplyr::group_by(fips) %>%
  dplyr::summarize(n_events = n()) %>%
  dplyr::mutate(fips = as.numeric(fips)) %>%
  dplyr::rename(region = fips,
               value = n_events) %>%
  county_choropleth(state_zoom = c("north carolina", "virginia"),
                     reference_map = TRUE)
```

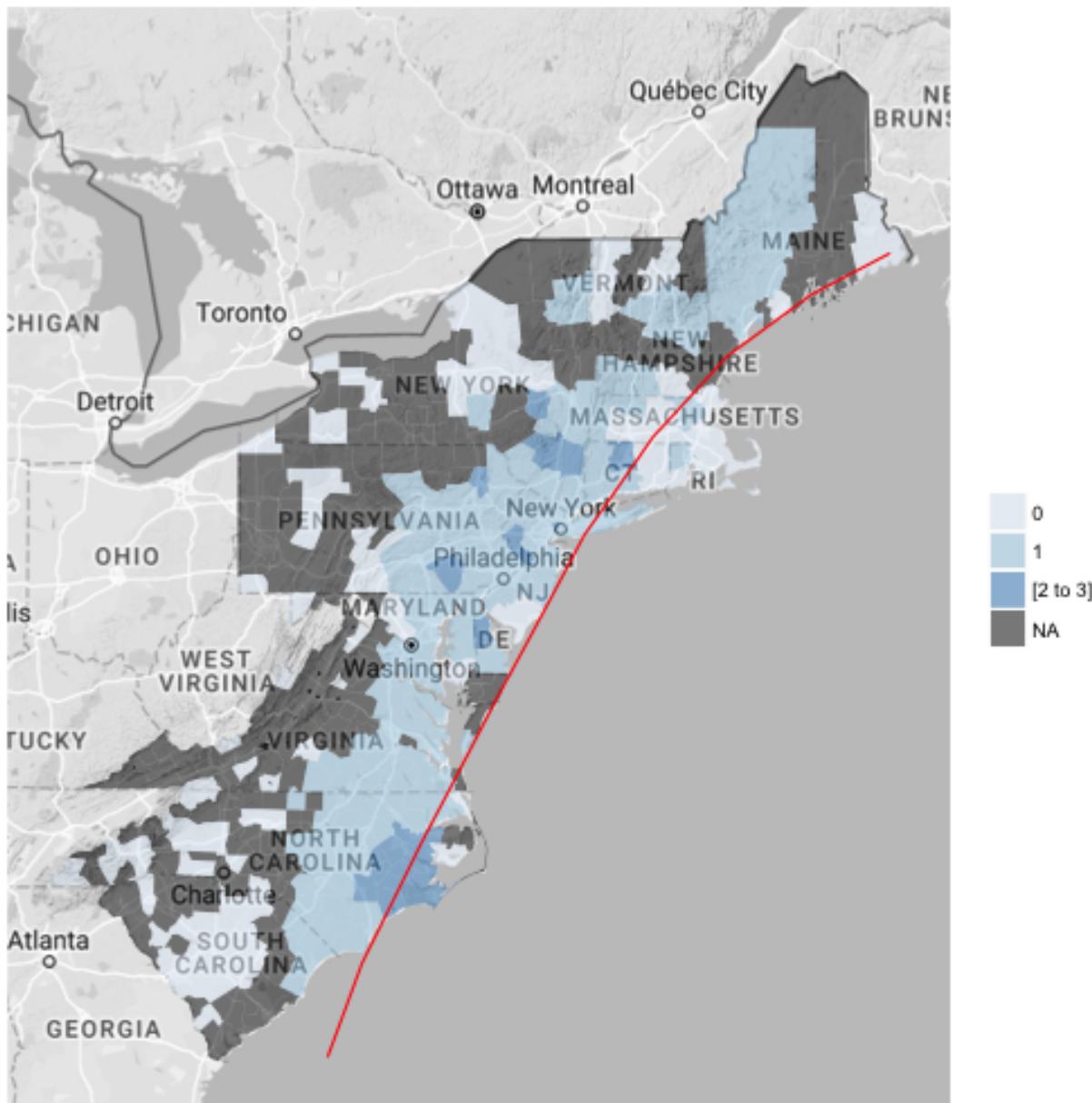


plot of chunk unnamed-chunk-60

The map created by `county_choropleth` (and the other maps created by functions in the `choroplethr` package) is a `ggplot` object, so you can add elements to it. For example, to create a map of flood events that includes the track of Hurricane Floyd on the map, you can run:

```
floyd_track <- read_csv("data/floyd_track.csv")

floyd_events %>%
  dplyr::group_by(fips) %>%
  dplyr::summarize(flood = sum(grep1("Flood", type))) %>%
  dplyr::mutate(fips = as.numeric(fips)) %>%
  dplyr::rename(region = fips,
                value = flood) %>%
  county_choropleth(state_zoom = c("north carolina", "maryland",
                                    "delaware", "new jersey",
                                    "virginia", "south carolina",
                                    "pennsylvania", "new york",
                                    "connecticut", "massachusetts",
                                    "new hampshire", "vermont",
                                    "maine", "rhode island"),
                     reference_map = TRUE) +
  geom_path(data = floyd_track, aes(x = -longitude, y = latitude,
                                    group = NA),
            color = "red")
```



plot of chunk unnamed-chunk-61

To create county choropleths with the `choroplethr` package that are more customized, you can use the package's `CountyChoropleth`, which is an R6 object for creating custom county choropleths. To create an object, you can run `CountyChoropleth$new` with the data you'd like to map. As with `county_choropleth`, this data should have a column named "region" with county FIPS codes in a numeric class and a column named "values" with the values to plot. To map counties in which a flood event was reported around the time of Floyd, you can start by cleaning your data and then an object using `CountyChoropleth$new`:

```
floyd_floods <- floyd_events %>%
  dplyr::filter(grep1("Flood", type)) %>%
  dplyr::mutate(fips = as.numeric(fips)) %>%
  dplyr::group_by(fips) %>%
  dplyr::summarize(value = 1) %>%
  dplyr::ungroup() %>%
  dplyr::rename(region = fips)

floyd_map <- CountyChoropleth$new(floyd_floods)
```

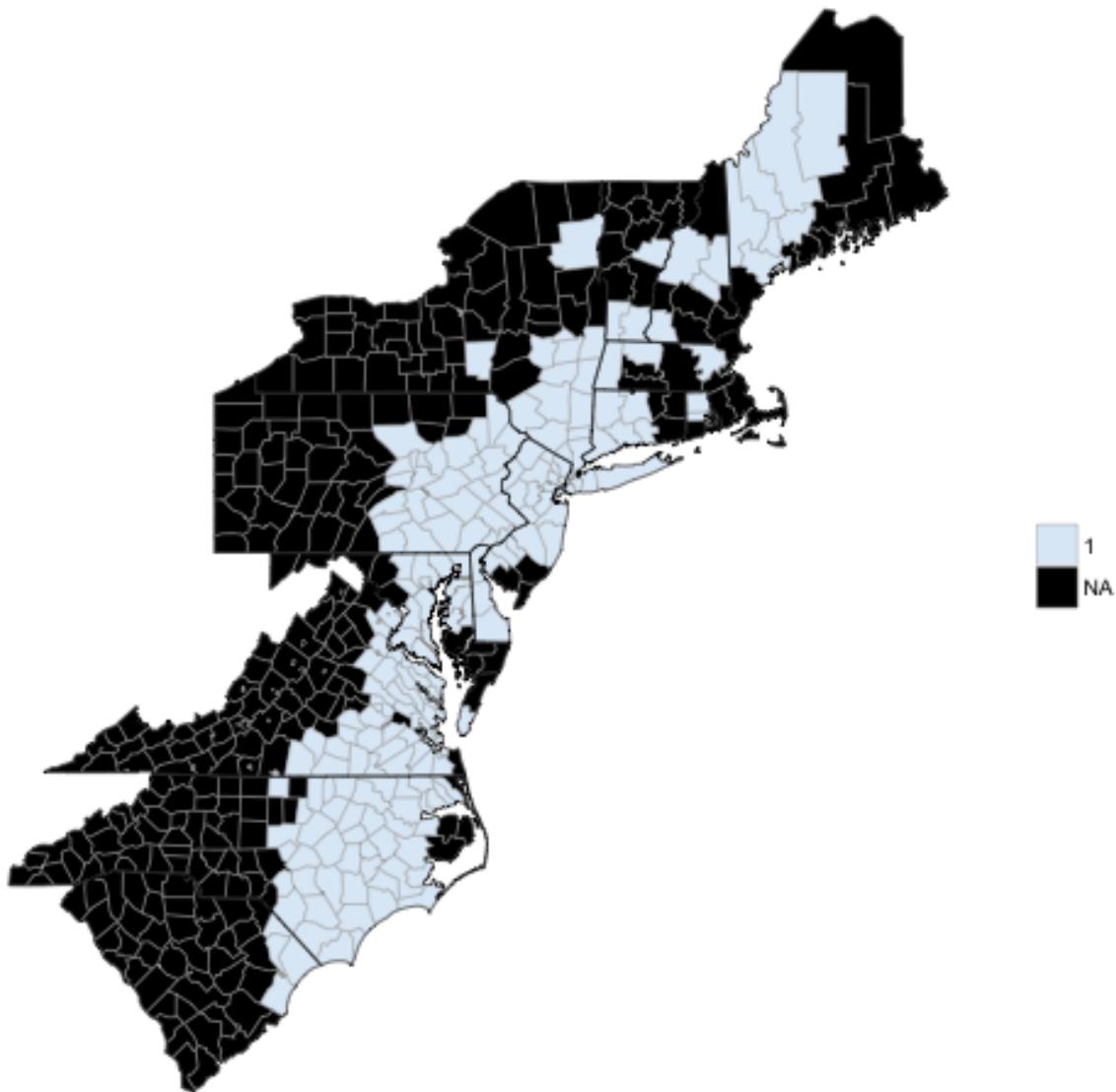
As a note, in cleaning the data here, I wanted to limit the dataset to only observations where the event type included the word “Flood” (this will pull events listed as “Flood” or “Flash Flood”), so I’ve used the `grep1` function to help filter to just those observations. The choropleth mapping functions require that each county is included only once, so I used `group_by` and `summarize` to collapse the dataframe to have only a single observation for each county.

Once you have created a basic object using `CountyChoropleth`, you can use the methods for this type of object to customize the map substantially. For example, you can set the states using the `set_zoom` method:

```
floyd_map$set_zoom(c("north carolina", "maryland",
                     "delaware", "new jersey",
                     "virginia", "south carolina",
                     "pennsylvania", "new york",
                     "connecticut", "massachusetts",
                     "new hampshire", "vermont",
                     "maine", "rhode island"))
```

At any point, you can render the object using the `render` method (or `render_with_reference_map`, to plot the map with the Google reference map added):

```
floyd_map$render()
```



plot of chunk unnamed-chunk-64

To find out what options are available for this object type, in terms of methods you can use or attributes you can change, you can run:

```
names(floyd_map)
[1] ".__enclos_env__"           "add_state_outline"
[3] "ggplot_polygon"           "projection"
[5] "ggplot_scale"              "warn"
[7] "legend"                   "title"
[9] "choropleth.df"             "map.df"
[11] "user.df"                  "clone"
[13] "clip"                     "initialize"
[15] "set_zoom"                 "render_state_outline"
[17] "render_helper"             "render"
[19] "set_num_colors"            "get_zoom"
[21] "format_levels"             "theme_inset"
[23] "theme_clean"               "get_scale"
[25] "prepare_map"                "bind"
[27] "discretize"                "render_with_reference_map"
[29] "get_choropleth_as_polygon" "get_reference_map"
[31] "get_y_scale"                "get_x_scale"
[33] "get_bounding_box"           "get_max_lat"
[35] "get_min_lat"                "get_max_long"
[37] "get_min_long"
```

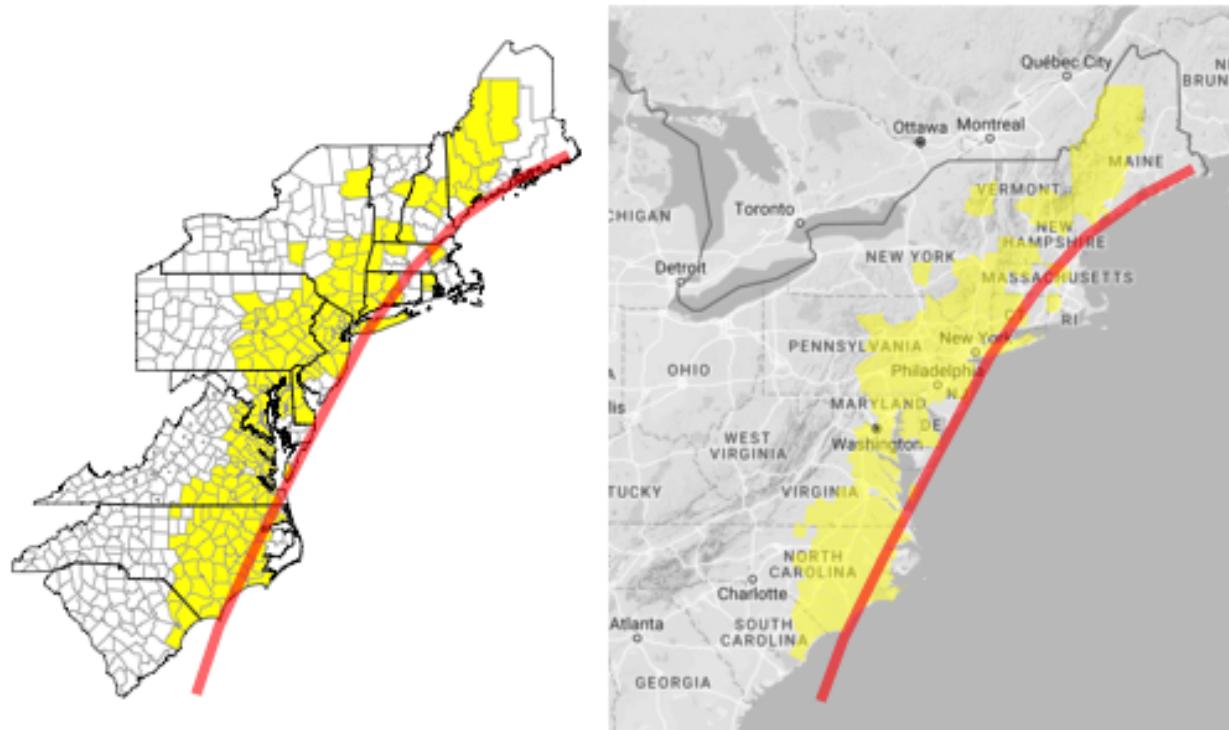
The following code shows an example of customizing this county choropleth map:

```
floyd_map$add_state_outline <- TRUE
floyd_map$ggplot_scale <- scale_fill_manual(values = c("yellow"),
                                              guide = FALSE)
floyd_xlim <- floyd_map$get_bounding_box()[c(1, 3)]
floyd_ylim <- floyd_map$get_bounding_box()[c(2, 4)]

a <- floyd_map$render() +
  geom_path(data = floyd_track, aes(x = -longitude, y = latitude,
                                    group = NA),
            color = "red", size = 2, alpha = 0.6) +
  xlim(floyd_map$get_bounding_box()[c(1, 3)]) +
  ylim(floyd_map$get_bounding_box()[c(2, 4)])

b <- floyd_map$render_with_reference_map() +
  geom_path(data = floyd_track, aes(x = -longitude, y = latitude,
                                    group = NA),
            color = "red", size = 2, alpha = 0.6) +
  xlim(floyd_xlim) +
  ylim(floyd_ylim)

library(gridExtra)
grid.arrange(a, b, ncol = 2)
```



plot of chunk unnamed-chunk-66

Here, I've used `$add_state_outline` to change the object to include state outlines (this only shows up when you render the map without a background reference map). I've also used `ggplot_scale` to customize the colors used for plotting counties with flood events and to remove the map legend. The `get_bounding_box` method pulls the range of latitudes and longitudes for the current map. I'm planning to add the hurricane track to the map, and the hurricane track extends well beyond this range. To later limit the map to these states, I'm using this `get_bounding_box` method to get these boundaries, and then I've used those values for the `xlim` and `ylim` functions when I create the final `ggplot` objects. Finally, the rendered maps are `ggplot` objects, so to include the hurricane track, I can add `ggplot` elements to the

map using `+`, as with any `ggplot` object. I used the `grid.arrange` function from the `gridExtra` package to put the two maps (with and without the background Google map) side-by-side.

More advance mapping- Spatial objects

Spatial objects in R

R has a series of special object types for spatial data. For many mapping / GIS tasks, you will need your data to be in one of these objects.

Spatial objects:

- `SpatialPolygons`
- `SpatialPoints`
- `SpatialLines`

Spatial objects + dataframes:

- `SpatialPolygonsDataFrame`
- `SpatialPointsDataFrame`
- `SpatialLinesDataFrame`

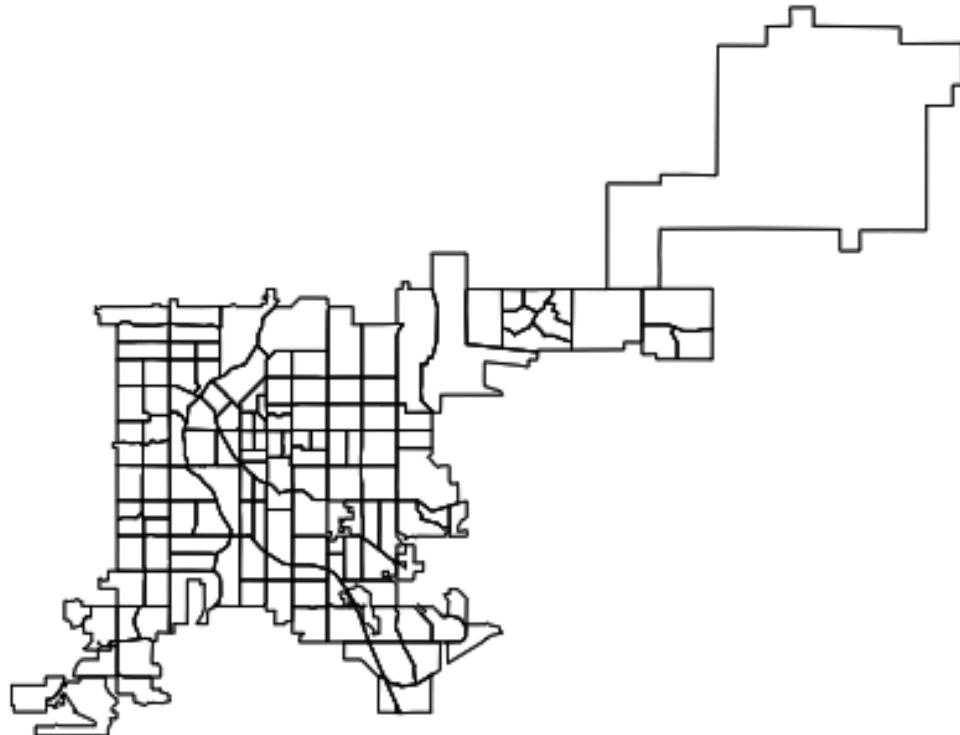
The `tigris` package lets you pull spatial data directly from the US Census. This data comes in directly as a spatial object.

```
library(tigris)
denver_tracts <- tracts(state = "CO", county = 31, cb = TRUE)
class(denver_tracts)
[1] "SpatialPolygonsDataFrame"
attr(,"package")
[1] "sp"
```

For more on this package, see the related article in *The R Journal*: <https://journal.r-project.org/archive/accepted/>

You can plot a spatial object in R just by calling `plot`:

```
plot(denver_tracts)
```



Plotting a spatial object

These spatial objects come with a number of special *methods*, or functions that work for the specific object type. You can list these methods using `name`:

```
names(summary(denver_tracts))
[1] "class"        "bbox"          "is.projected" "proj4string"
[5] "data"
```

For example, `bbox` will print out the *bounding box* of the spatial object (range of latitudes and longitudes included).

```
bbox(denver_tracts)
      min         max
x -105.10993 -104.60030
y   39.61443   39.91425
```

The `is.projected` and `proj4string` functions give you some information about the current Coordinate Reference System of the data.

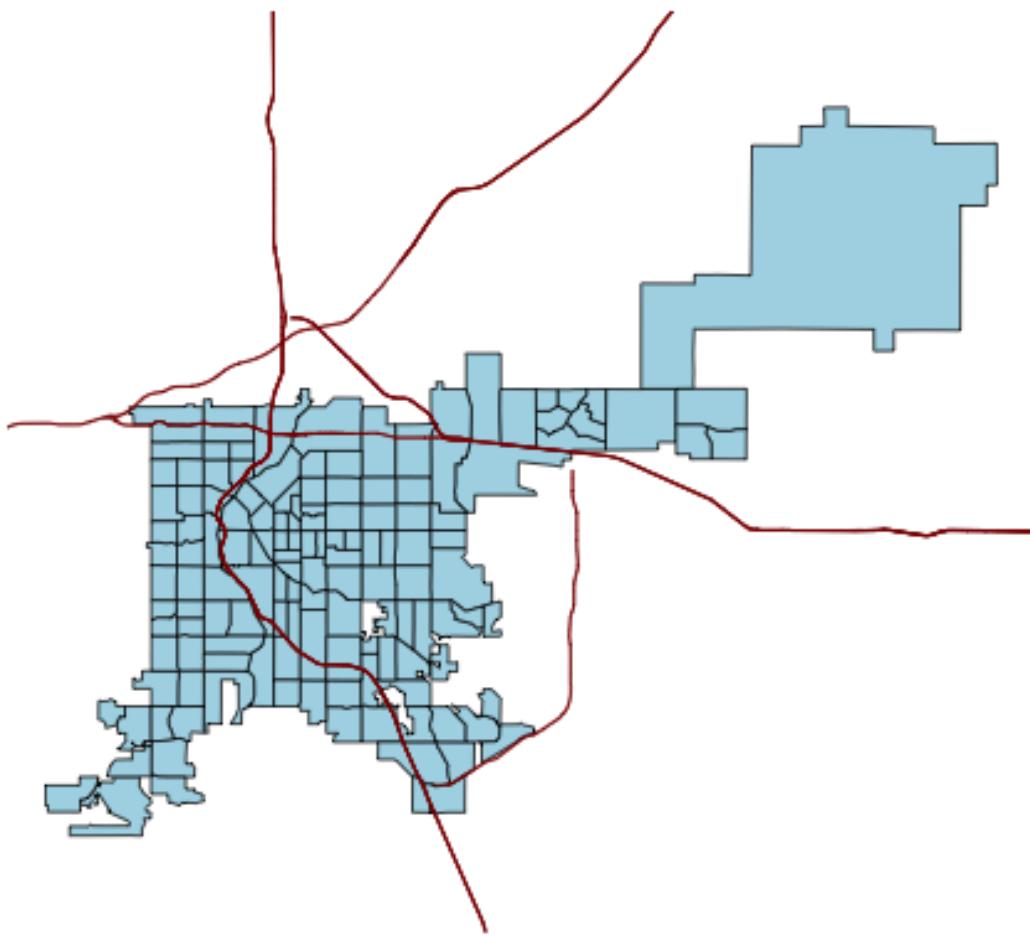
```
is.projected(denver_tracts)
[1] FALSE
proj4string(denver_tracts)
[1] "+proj=longlat +datum=NAD83 +no_defs +ellps=GRS80 +towgs84=0,0,0"
```

You can access a “slot” in a spatial object with a dataframe to pull out the data. This is similar to indexing a list. Just use `@` instead of `$`. For example, here’s the dataframe for the `denver_tracts` spatial object:

```
head(denver_tracts@data[, 1:4])
  STATEFP COUNTYFP TRACTCE      AFFGEOID
25      08        031 000201 1400000US08031000201
26      08        031 000302 1400000US08031000302
27      08        031 001101 1400000US08031001101
28      08        031 002802 1400000US08031002802
29      08        031 003300 1400000US08031003300
30      08        031 004006 1400000US08031004006
```

You can add different layers of spatial objects onto the same plot. To do that, just use `add = TRUE` for added layers. For example, to add primary roads to the Denver census tract map, you could run:

```
denver_roads <- primary_roads()
plot(denver_tracts, col = "lightblue")
plot(denver_roads, add = TRUE, col = "darkred")
```



Adding layers of spatial objects

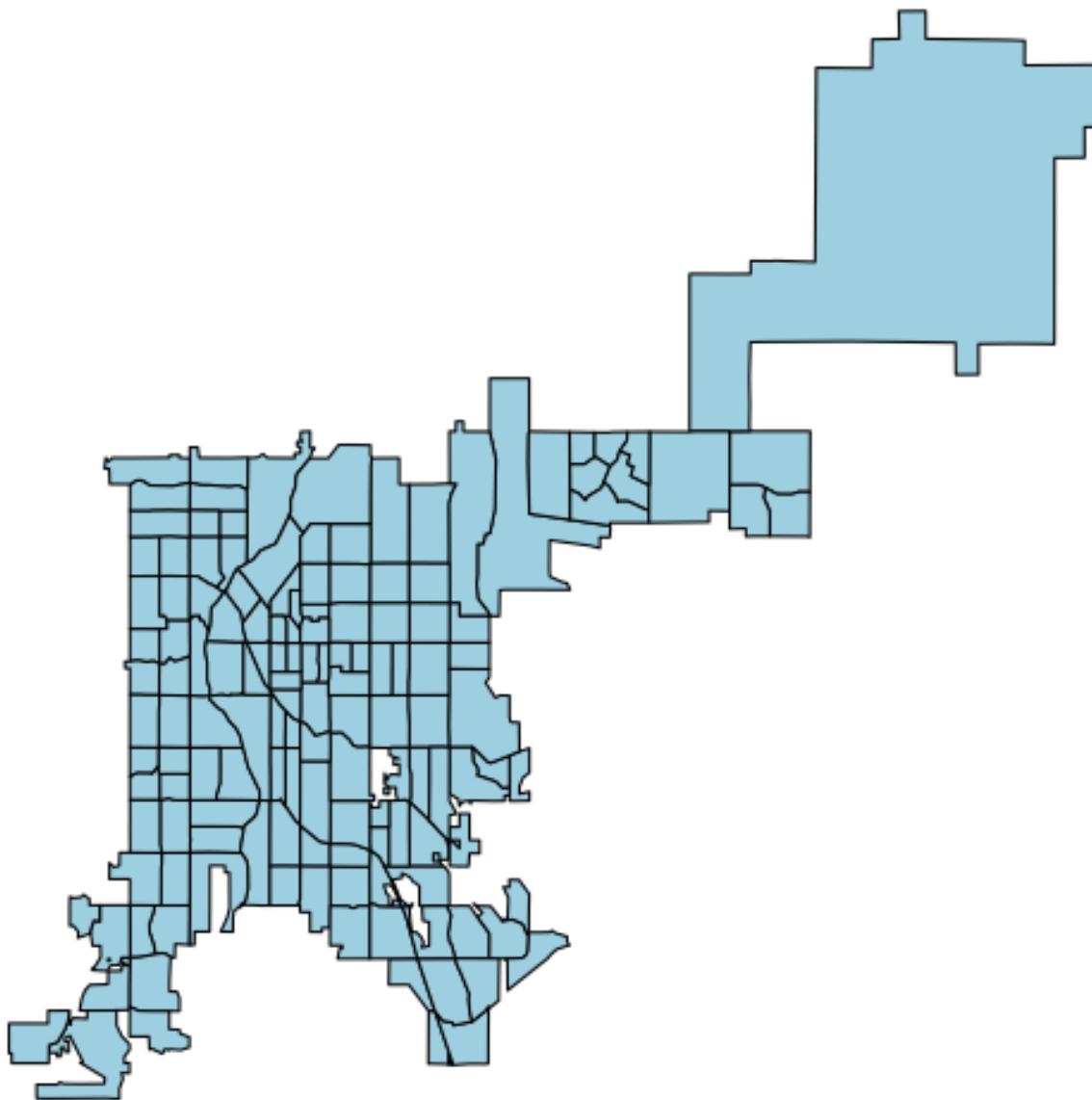
If you read in a shapefile, it will automatically be one of these shape objects. However, you can also convert other data into shape objects.

- Functions from `sp` package convert data into spatial objects
- `fortify` converts from a spatial object to a data frame (useful for `ggplot` plotting)

You can use the `fortify` function from `ggplot2` to convert the spatial object into a data frame, so you can plot it using polygons in `ggplot2`.

```
fortify(denver_tracts) %>%
  dplyr::select(1:4) %>% dplyr::slice(1:5)
  long      lat order  hole
1 -105.0251 39.79400    1 FALSE
2 -105.0213 39.79398    2 FALSE
3 -105.0208 39.79109    3 FALSE
4 -105.0158 39.79107    4 FALSE
5 -105.0064 39.79105    5 FALSE

denver_tracts %>%
  fortify() %>%
  ggplot(aes(x = long, y = lat, group = group)) +
  geom_polygon(fill = "lightblue", color = "black") +
  theme_void()
```



Using `fortify()` to convert spatial object to data frame

Coordinate reference systems

Spatial objects can have different Coordinate Reference Systems (CRSs). CRSs can be *geographic* (e.g., WGS84, for longitude-latitude data) or *projected* (e.g., UTM, NADS83).

There is a website that lists projection strings and can be useful in setting projection information or re-projecting data: <http://www.spatialreference.org>

Here is an excellent resource on projections and maps in R from Melanie Frazier: <https://www.nceas.ucsb.edu/~mfrazier/teaching/RforSpatialData/RforSpatialData.html>

] To tell R the Coordinate Reference System of some data, set this attribute with `proj4string`

(similar to setting column names with `colnames`):

```
## Generic code
proj4string(my_spatial_object) <- "+proj=longlat +datum=NAD83"
```

This does not create a projection. Instead, this is just how you tell R what projection the data already is in.

The `CRS` function creates CRS class objects that can be used to specify projections. You input a character string of projection arguments into this function (for example, `CRS("+proj=longlat +datum=NAD27")`). You can also use, however, use a shorter EPSG code for a projection (for example, `CRS("+init=epsg:28992")`).

```
library(sp)
CRS("+proj=longlat +datum=NAD27")
CRS arguments:
+proj=longlat +datum=NAD27 +ellps=clrk66
+nadgrids=@conus,@alaska,@ntv2_0.gsb,@ntv1_can.dat
CRS("+init=epsg:28992")
CRS arguments:
+init=epsg:28992 +proj=sterea +lat_0=52.15616055555555
+lon_0=5.38763888888889 +k=0.9999079 +x_0=155000 +y_0=463000
+ellps=bessel
+towgs84=565.4171,50.3319,465.5524,-0.398957,0.343988,-1.87740,4.0725
+units=m +no_defs
```

To **change** the projection of a spatial object, you can use the `spTransform` function from the `sp` package.

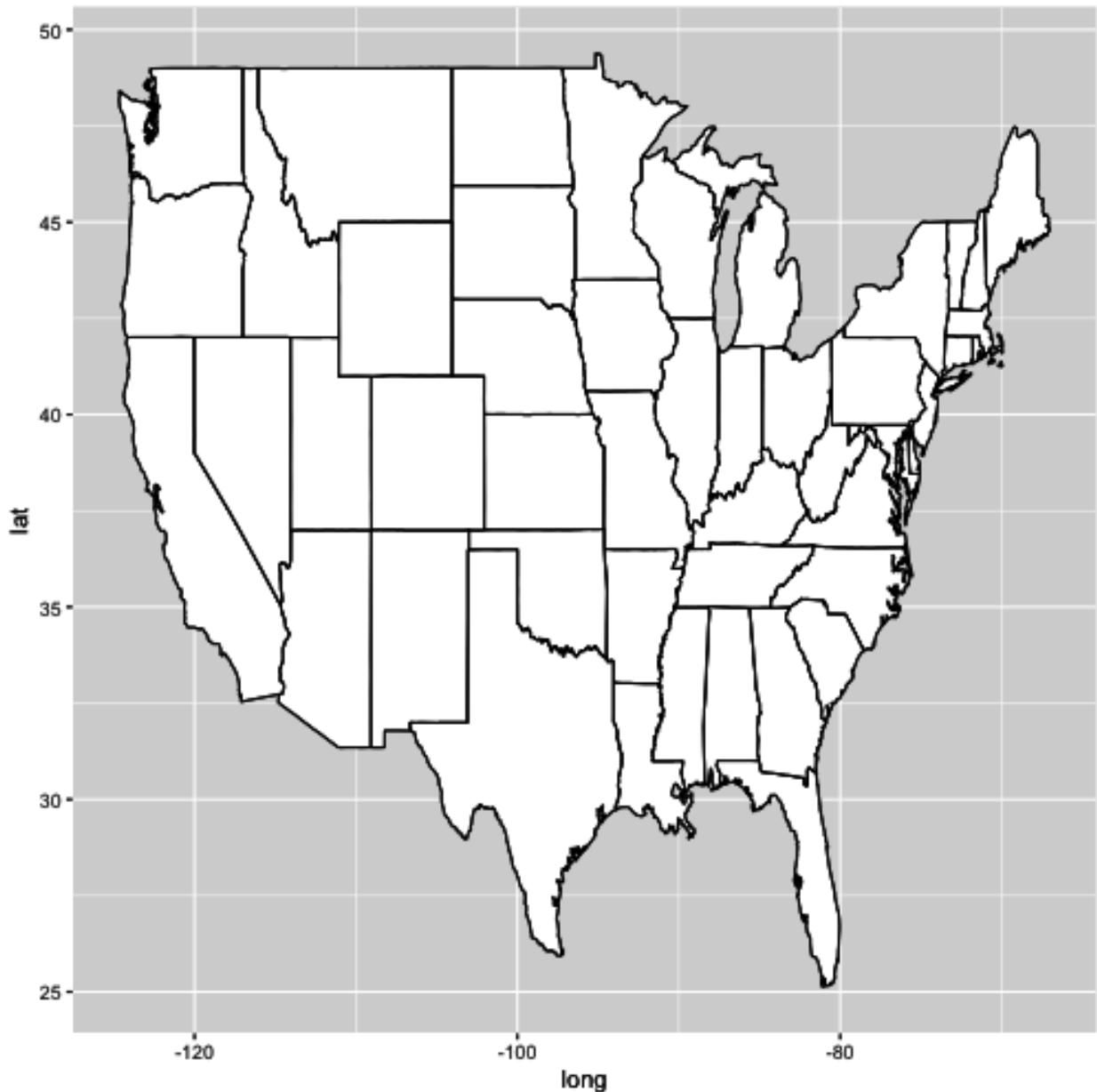
```
## Generic code
my_spatial_object <- spTransform(my_spatial_object,
                                 proj4string(another_sp_object))
```

The `coord_map` function in `ggplot2` can help you with map projections, as well. Here's an example from the help file with a US map.

```
states <- map_data("state")
usamap <- ggplot(states,
                  aes(long, lat, group = group)) +
  geom_polygon(fill = "white", colour = "black")
```

The default is Cartesian coordinates:

```
usamap
```



Mercator projection:

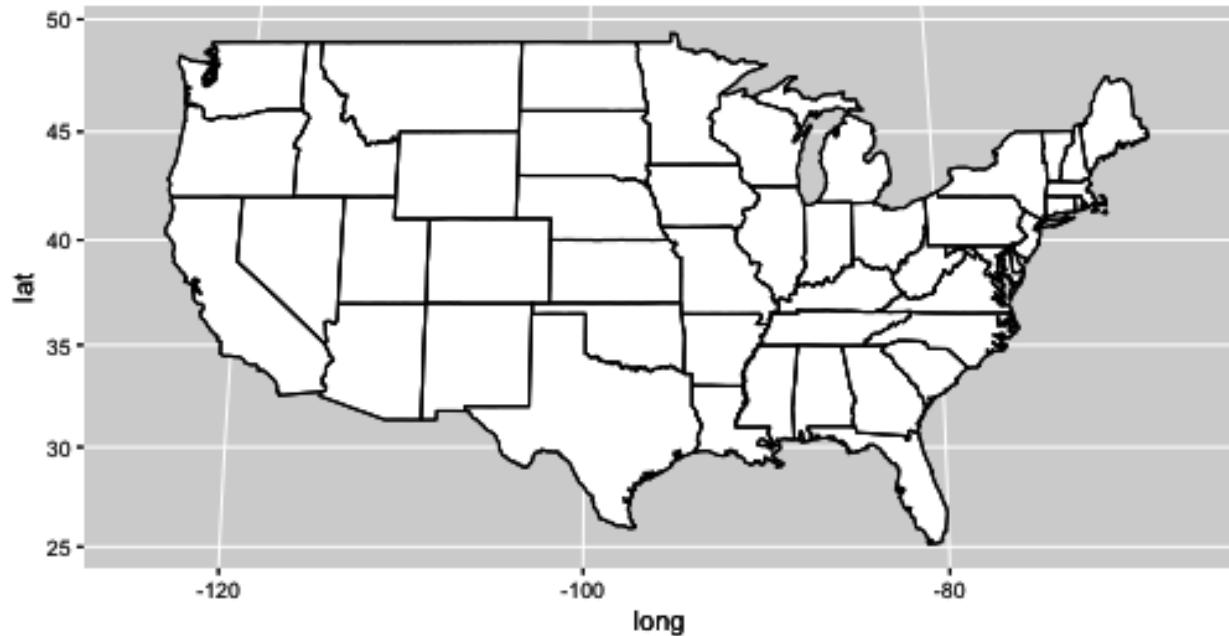
```
usamap + coord_map("mercator")
```



Mercator projection

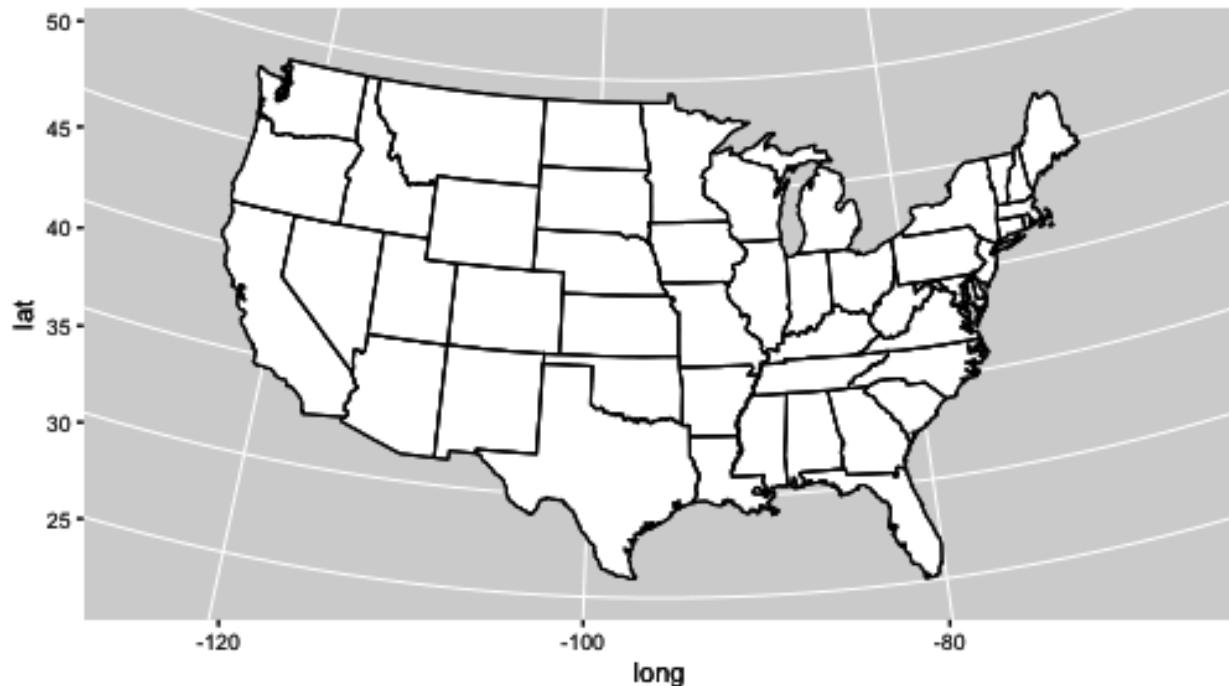
Gilbert projection

```
usamap + coord_map("gilbert")
```

**Gilbert projection**

Conic projection:

```
usamap + coord_map("conic", lat0 = 30)
```



Conic projection

Shapefiles

- File format (ESRI, but usable by other software)
- Not a single file, but rather a directory of related files (e.g., .shp, .shx, .dbf, .prj)
- Typically includes both geographic information (e.g., locations of county boundaries) and attribute information (e.g., median income of each county)
- To read shapefiles into R, use the `readOGR` function from the `rgdal` package
- You can also write out shapefiles you've created or modified in R, using `writeOGR`.

R as GIS

You can use R for a number of GIS-style tasks:

- Clipping
- Creating buffers
- Measuring area of a polygon
- Counting points in polygon

There are some advantages to using R for this:

- R is free
- You can write all code in a script, so research is more reproducible
- You save time and effort by staying in one software system, not going between different software

There are some advantages to GIS, too, though:

- More user-friendly at the start (point-and-click)
- R spatial functionality is still spread over lots of packages, with different syntax and conventions.

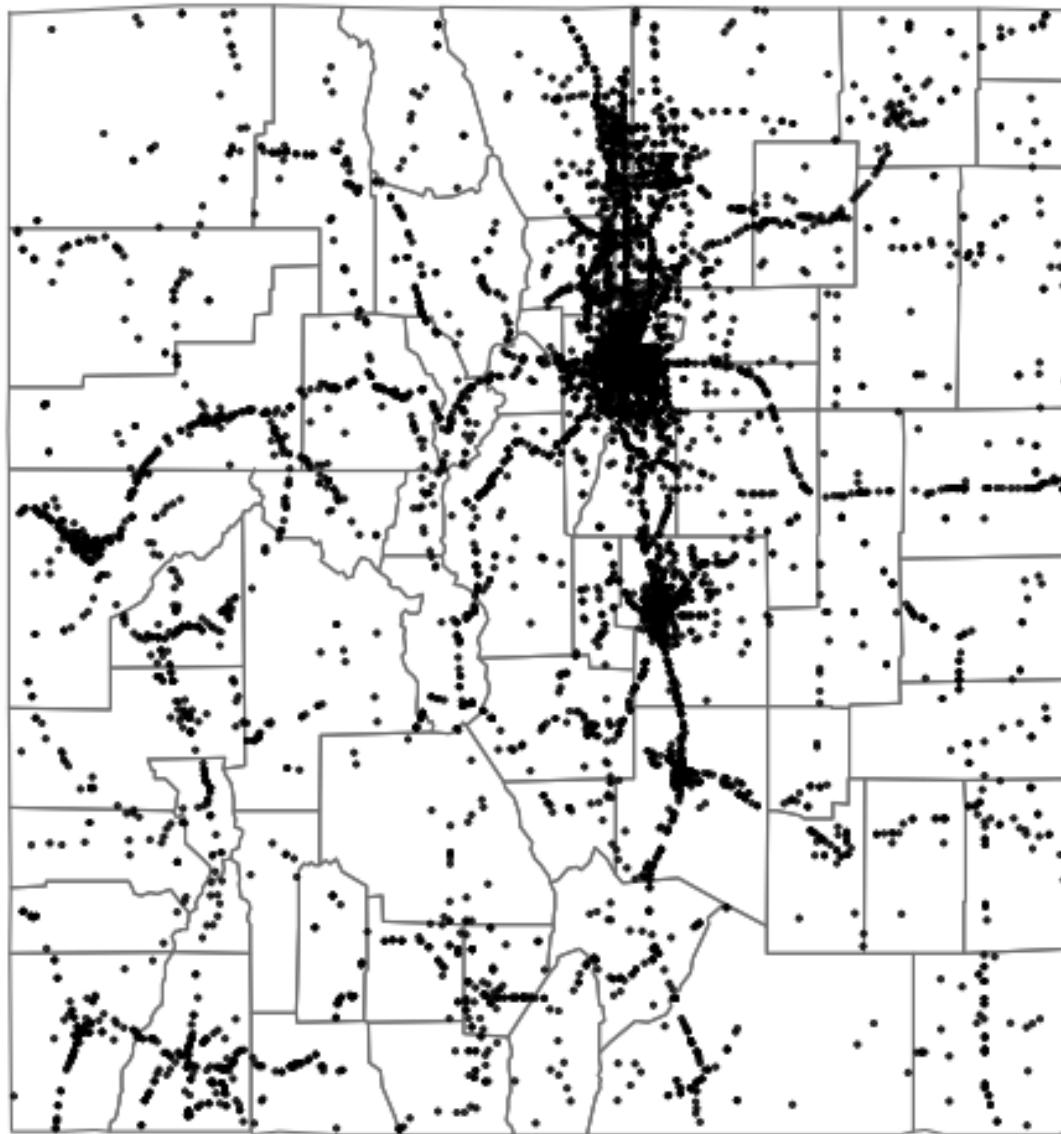
For an example, I've cleaned up some FARs data at the driver level for 2001–2010:

```
load("data/fars_colorado.RData")
driver_data %>%
  dplyr::select(1:5) %>% dplyr::slice(1:5)
  state st_case county          date latitude
1     8   80001      51 2001-01-01 10:00:00 39.10972
2     8   80002      31 2001-01-04 19:00:00 39.68215
3     8   80003      31 2001-01-03 07:00:00 39.63500
4     8   80004      31 2001-01-05 20:00:00 39.71304
5     8   80005      29 2001-01-05 10:00:00 39.09733
```

Here is how you would plot fatal accidents (by driver) in Colorado without using spatial objects:

```
ggplot2::map_data("county", region = "Colorado") %>%
  ggplot2::ggplot(ggplot2::aes(x = long, y = lat,
                                group = subregion)) +
  ggplot2::geom_polygon(color = "gray", fill = NA) +
  ggplot2::theme_void() +
  ggplot2::geom_point(data = driver_data,
                      ggplot2::aes(x = longitud,
                                   y = latitude,
                                   group = NULL),
                      alpha = 0.5, size = 0.7 )
```

Fatal accidents (by driver) in Colorado, 2001–2010:



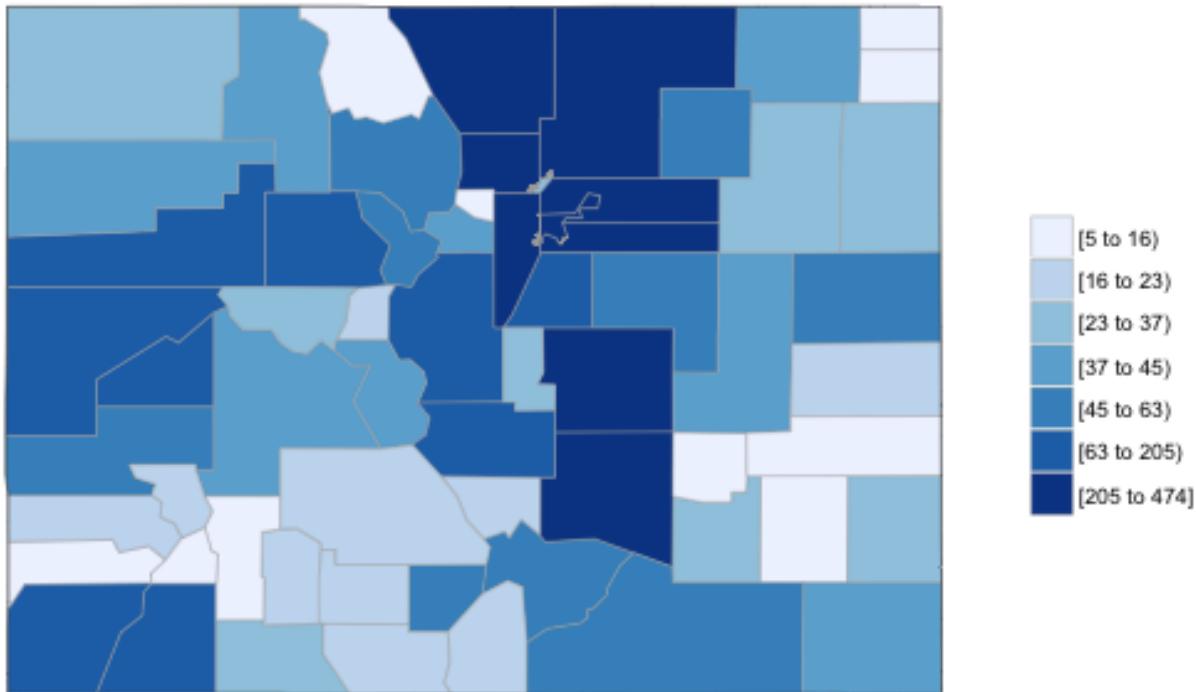
Fatal accidents in Colorado, 2001–2010

You can also make a choropleth of county accident counts without using spatial data, by using the `choroplethr` package.

To do this, you'll first need to use `dplyr` functions to limit to unique accidents (rather than drivers) and add up the number of accidents in each county. In this case, it's possible to add up accidents by county because `county` is included as a column in our data.

```
county_accidents <- driver_data %>%
  dplyr::select(state, st_case, county, latitude, longitud) %>%
  dplyr::distinct() %>%
  dplyr::mutate(county = str_pad(county, width = 3,
                                 side = "left", pad = "0")) %>%
  tidyverse::unite(region, state, county, sep = "") %>%
  dplyr::group_by(region) %>%
  dplyr::summarize(value = n()) %>%
  dplyr::mutate(region = as.numeric(region))
county_accidents %>% slice(1:4)
# A tibble: 4 × 2
  region value
  <dbl> <int>
1     8001    372
2     8003     47
3     8005    305
4     8007     31

county_choropleth(county_accidents, state_zoom = "colorado")
```



plot of chunk unnamed-chunk-89

This “out-of-the-box” solution let us looks at accident counts by county, but what if we want to look at a geographical unit for which we don’t have an identifying column?

For example, we might want to look at accident counts by census tract in Denver. To do this, we’ll need to link each accident (point) to a census tract (polygon), and then we can count up the number of points linked to each polygon.

First, I’ve created a dataframe with only accidents in Denver (based on the `county` column in the accident data):

```
denver_fars <- driver_data %>%
  filter(county == 31)
```

To do this, both the census tracts and the accident data need to be in spatial objects.

```
library(sp)
denver_fars_sp <- denver_fars %>%
  dplyr::rename(longitude = longitud)
coordinates(denver_fars_sp) <- c("longitude", "latitude")
proj4string(denver_fars_sp) <- CRS("+init=epsg:4326")
```

Note that the dataframe is changed into a spatial object by changing its `coordinates` attribute, and that the CRS was set uas the `proj4string` attribute.

```
summary(denver_fars_sp)
Object of class SpatialPointsDataFrame
Coordinates:
min          max
longitude -105.10973 -104.0122
latitude   39.61715  39.8381
Is projected: FALSE
proj4string :
[+init=epsg:4326 +proj=longlat +datum=WGS84 +no_defs +ellps=WGS84
+towgs84=0,0,0]
Number of points: 695
Data attributes:
      state    st_case    county       date
Min.   :8   Min.   :80001   Min.   :31   Min.   :2001-01-03 07:00:00
1st Qu.:8   1st Qu.:80121   1st Qu.:31   1st Qu.:2003-01-06 15:00:00
Median :8   Median :80268   Median :31   Median :2005-01-29 01:00:00
Mean   :8   Mean   :80264   Mean   :31   Mean   :2005-07-05 15:58:42
3rd Qu.:8   3rd Qu.:80390   3rd Qu.:31   3rd Qu.:2007-12-04 19:30:00
Max.   :8   Max.   :80669   Max.   :31   Max.   :2010-12-12 17:00:00

      fatalities    drunk_dr       age       alc_res
Min.   :1.000   Min.   :0.0000   Min.   :13.00   Min.   :0.0000
1st Qu.:1.000   1st Qu.:0.0000   1st Qu.:24.00   1st Qu.:0.0000
Median :1.000   Median :1.0000   Median :32.00   Median :0.1000
Mean   :1.047   Mean   :0.6374   Mean   :36.72   Mean   :0.1077
3rd Qu.:1.000   3rd Qu.:1.0000   3rd Qu.:45.50   3rd Qu.:0.2000
Max.   :3.000   Max.   :4.0000   Max.   :91.00   Max.   :0.4000
NA's   :48        NA's   :48        NA's   :454
```

To be able to pair up polygons and points, their spatial objects need to have the same CRS. To help later with calculating the area of each polygon, I'll use a projected CRS that is reasonable for Colorado.

```
proj4string(denver_tracts)
[1] "+proj=longlat +datum=NAD83 +no_defs +ellps=GRS80 +towgs84=0,0,0"
CRS(proj4string(denver_tracts))
CRS arguments:
+proj=longlat +datum=NAD83 +no_defs +ellps=GRS80 +towgs84=0,0,0
```

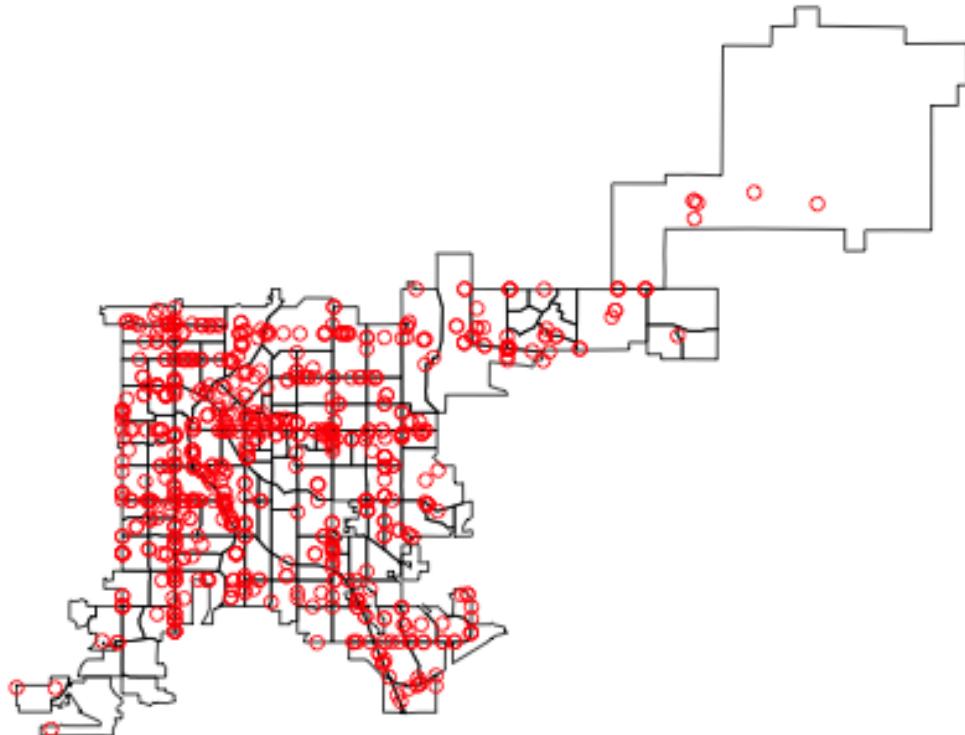
To reproject spatial data, you can use the `spTransform` function:

```
denver_tracts_proj <- spTransform(denver_tracts,
                                    CRS("+init=epsg:26954"))
denver_fars_proj <- spTransform(denver_fars_sp,
                                CRS(proj4string(denver_tracts)))
```

The `spTransform` function transforms the coordinates in a spatial object into a new coordinate reference system.

Here is a map of the tracts with the accidents overlaid:

```
plot(denver_tracts)
plot(denver_fars_proj, add = TRUE, col = "red", pch = 1)
```



Denver tracts with accidents overlaid

Now, the data's in a format where we can link spatial points to spatial polygons.

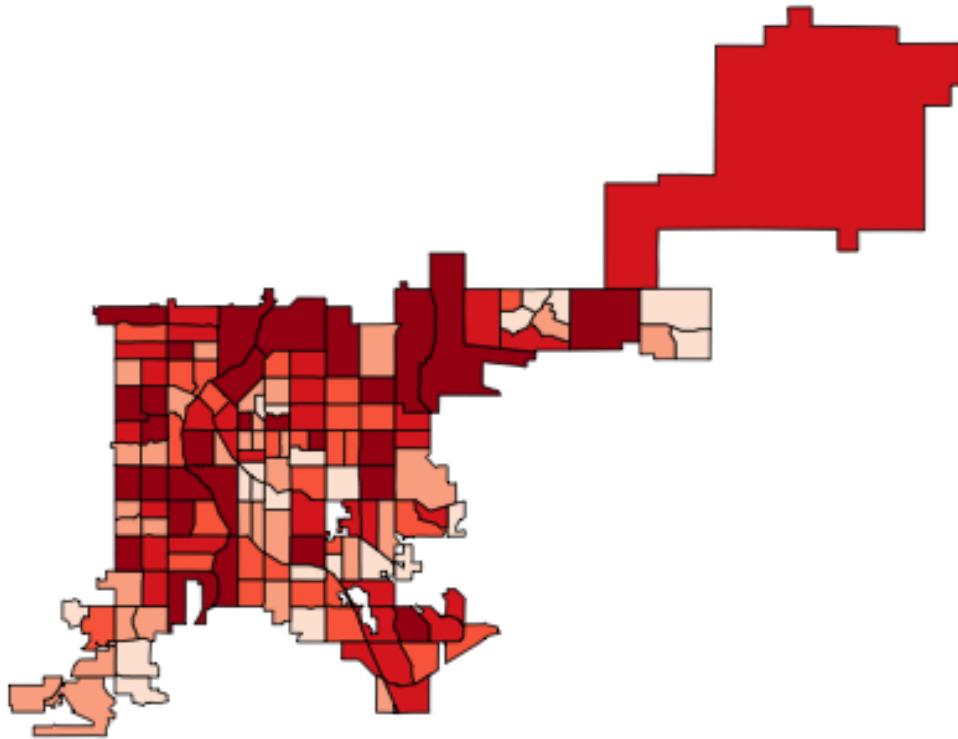
The `poly.counts` function in the `GISTools` package will measure the number of points that fall within each polygon.

It results in a vector with one element for each polygon (census tract in our example), where the element name identifies the polygon and the cell value gives the count within that polygon.

```
tract_counts <- poly.counts(denver_fars_proj, denver_tracts)
head(tract_counts)
25 26 27 28 29 30
7 2 2 0 0 4
```

You can use a choropleth to show these accident counts. The quickest way to do this is probably to use the `choropleth` function in the `GISTools` package.

```
choropleth(denver_tracts, tract_counts)
```

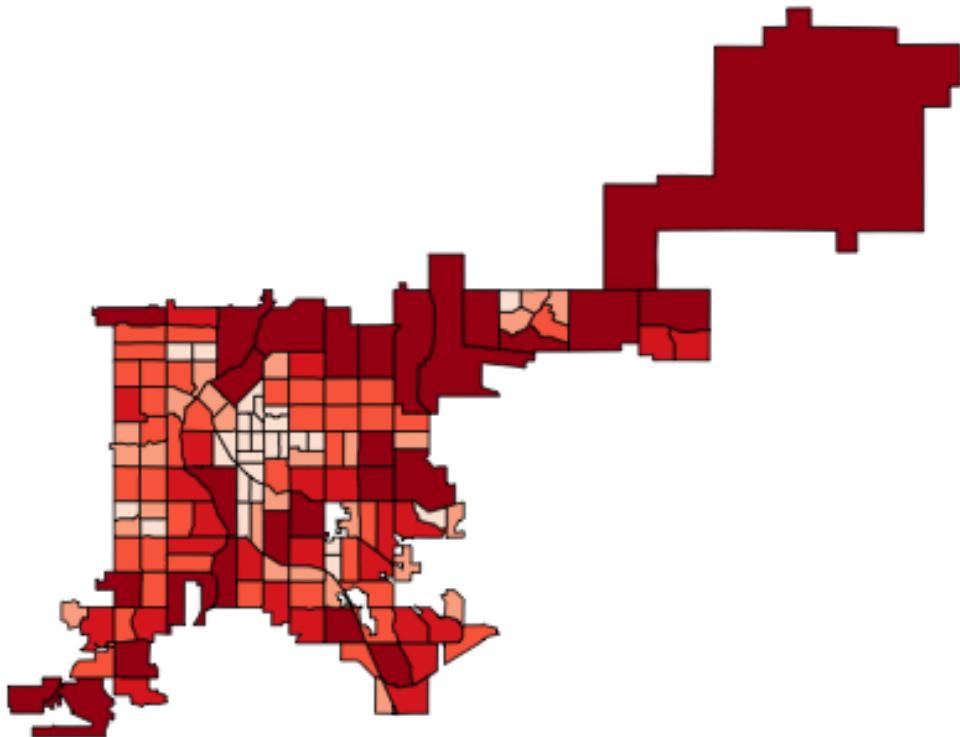


Choropleth map of accident counts

There is another function in the package that calculates the area of each polygon.

```
library(GISTools)
head(poly.areas(denver_tracts_proj))
 25      26      27      28      29      30
2100172.2 1442824.1 897886.3 881530.5 1282812.2 1948187.1
```

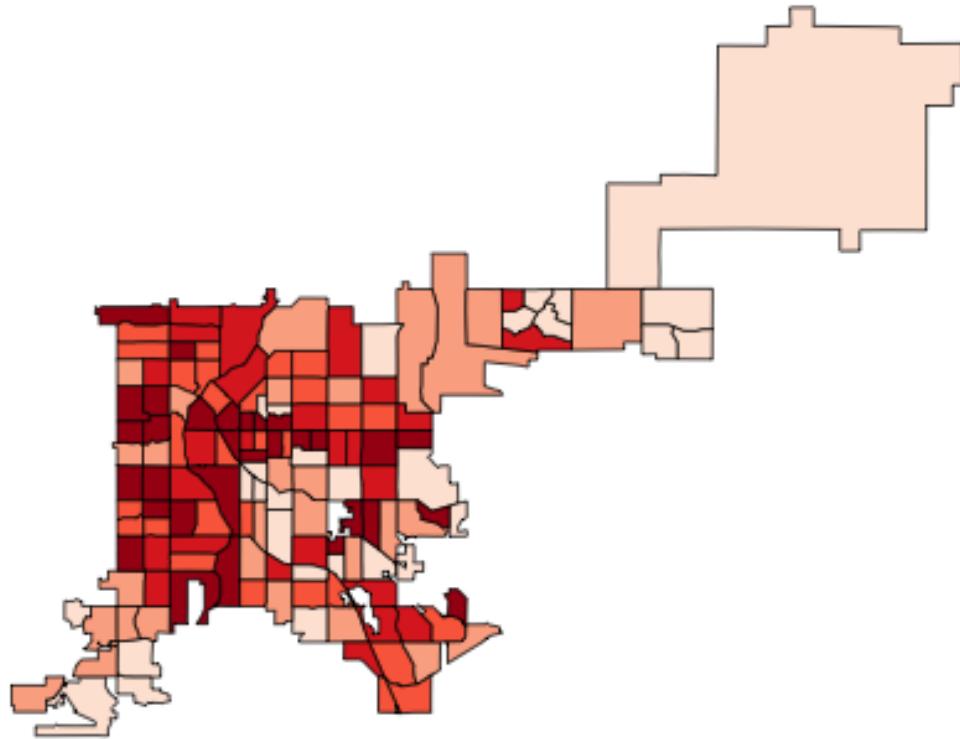
```
choropleth(denver_tracts, poly.areas(denver_tracts_proj))
```



Polygon areas

You can combine these ideas to create a choropleth of the rate of fatal accidents per population in Denver census tracts.

```
choropleth(denver_tracts, tract_counts /  
          poly.areas(denver_tracts_proj))
```



Accidents per population in Denver

Raster data

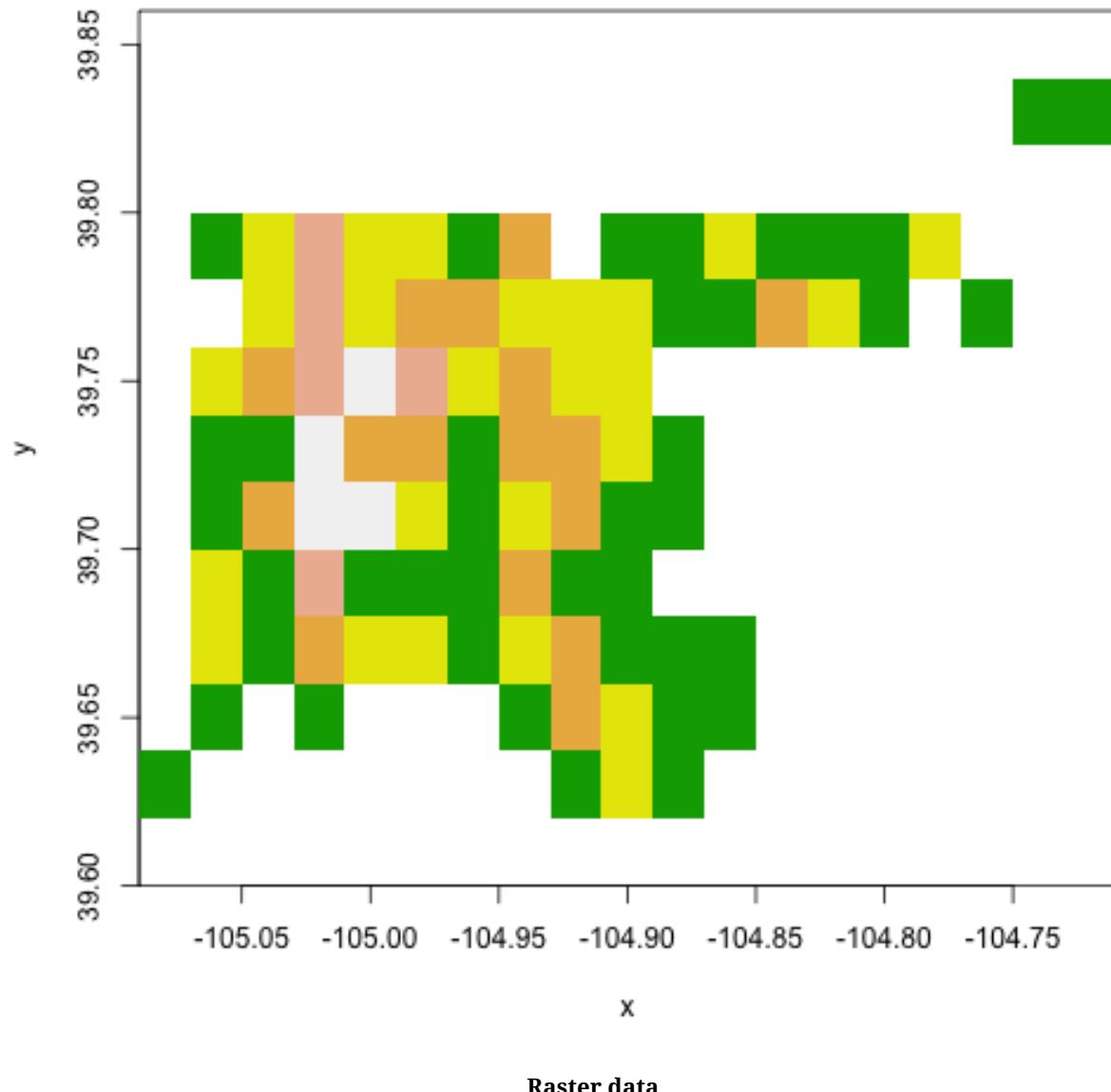
When mapping in R, you may also need to map *raster data*.

You can think of this as pixels—the graphing region is divided into even squares, and color is constant within each square.

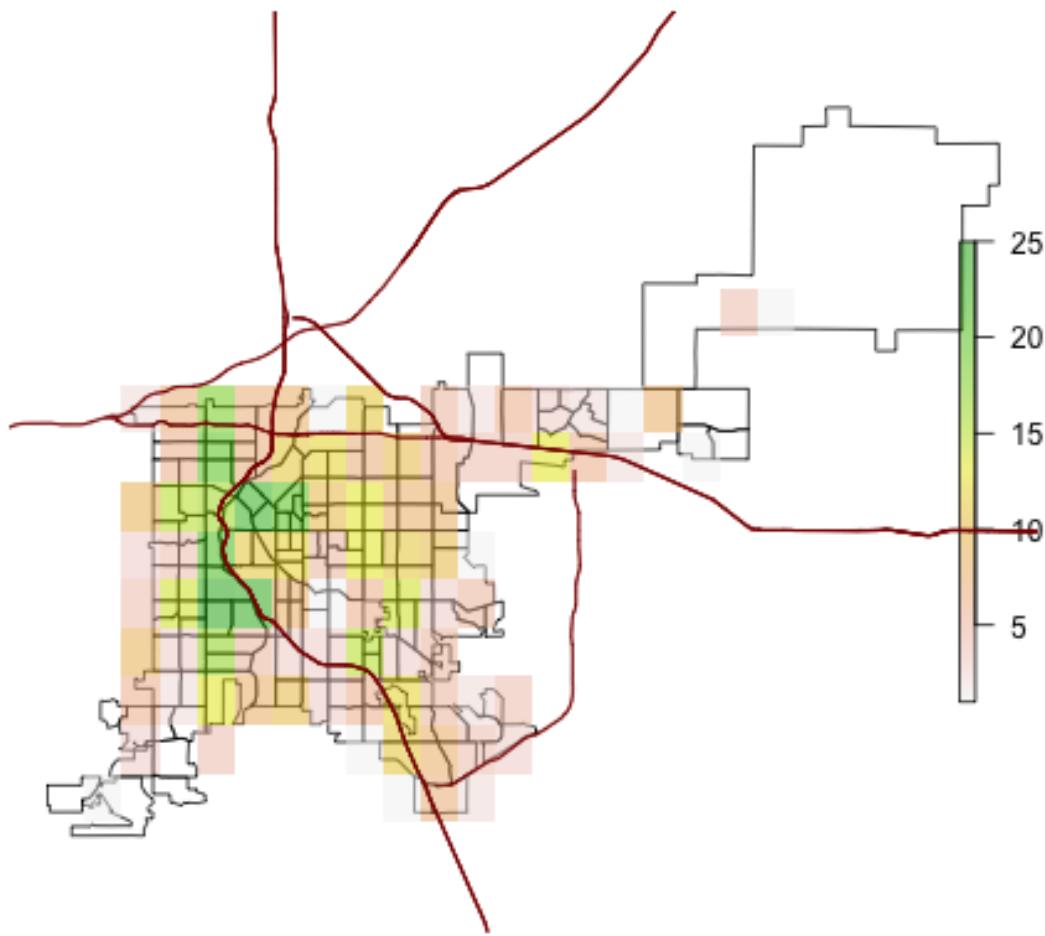
There are functions that allow you to “rasterize” data. That is, you take spatial points data, divide the region into squares, and count the number of points (or other summary) within each square.

```
bbox(denver_fars_sp)
      min      max
longitude -105.10973 -104.0122
latitude   39.61715  39.8381
library(raster)
denver_raster <- raster(xmn = -105.09, ymn = 39.60,
                        xmx = -104.71, ymx = 39.86,
                        res = 0.02)
den_acc_raster <- rasterize(geometry(denver_fars_sp),
                            denver_raster,
                            fun = "count")

image(den_acc_raster, col = terrain.colors(5))
```



```
plot(denver_tracts)
plot(den_acc_raster, add = TRUE, alpha = 0.5)
plot(denver_roads, add = TRUE, col = "darkred")
```



Denver tracts with raster data

Other capabilities

You can also use R for other spatial tasks:

- Kernel density estimation
- Identifying clusters
- Kriging
- Measuring spatial autocorrelation

Where to find more on mapping with R

- *Applied Spatial Data Analysis with R* by Roger Bivand (available online through CSU library)
- *An Introduction to R for Spatial Analysis and Mapping* by Chris Brunsdon and Lex Comber
- [CRAN Spatial Data Task View](#)
- [R Spatial Cheatsheet](#)
- [Great blog post \(among many\) by Zev Ross](#)

4.4 htmlWidgets

Overview of htmlWidgets

Very smart people have been working on creating interactive graphics in R for a long time. So far, nothing coded in R has taken off in a big way.

JavaScript has developed a number of interactive graphics libraries that can be used for documents viewed in a web browser. There is now a series of R packages that allow you to create plots from these JavaScript libraries from within R.

There is a website with much more on these `htmlWidgets` at <http://www.htmlwidgets.org>.

Some of the packages available to help you create interactive graphics from R using JavaScript graphics libraries:

- `leaflet`: Mapping (we'll cover this in the next section)
- `dygraphs`: Time series
- `plotly`: A variety of plots, including maps
- `rbokeh`: A variety of plots, including maps
- `networkD3`: Network data
- `d3heatmap`: Heatmaps
- `DT`: Data tables
- `DiagrammeR`: Diagrams and flowcharts

These packages can be used to make some pretty cool interactive visualizations for HTML output from R Markdown or Shiny (you can also render any of them in RStudio).

There are, however, a few limitations:

- Written by different people. The different packages have different styles as well as different interfaces. Learning how to use one package may not help you much with the other of these packages.
- Many are still in development, often in early development.

plotly package

From the package documentation:

“Easily translate `ggplot2` graphs to an interactive web-based version and / or create custom web-based visualizations directly from R.”

- Like many of the packages today, draws on functionality external to R, but within a package that allows you to work exclusively within R.
- Allows you to create interactive graphs from R. Some of the functions extend the `ggplot2` code you’ve learned.
- Interactivity will only work within RStudio or on documents rendered to HTML.

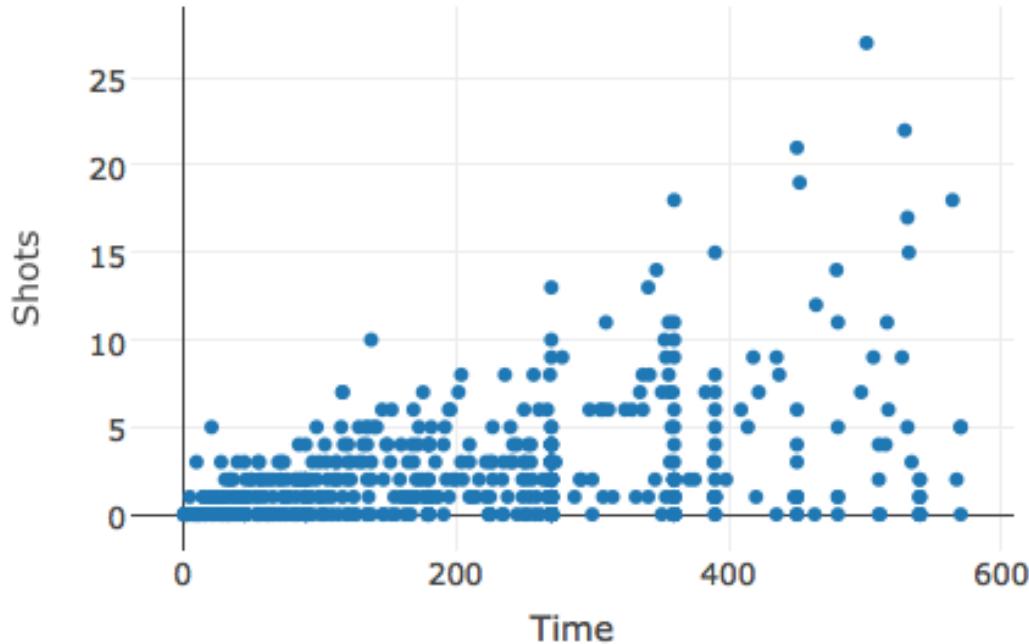
The `plotly` package allows an interface to let you work with `plotly.js` code directly using R code.

`plotly.js` is an open source library for creating interactive graphs in JavaScript. This JavaScript library is built on `d3.js` (Data-Driven Documents), which is a key driver in interactive web-based data graphics today.

There are two main ways of create plots within `plotly`:

- Use one of the functions to create a customized interactive graphic:
 - `plot_ly`: Workhorse of `plotly`, renders most non-map types of graphs
 - `plot_geo, plot_mapbox`: Specific functions for creating `plotly` maps.
- Create a `ggplot` object and then convert it to a `plotly` object using the `ggplotly` function.

```
library(faraway)
data(worldcup)
library(plotly)
plot_ly(worldcup, type = "scatter", x = ~ Time, y = ~ Shots)
```

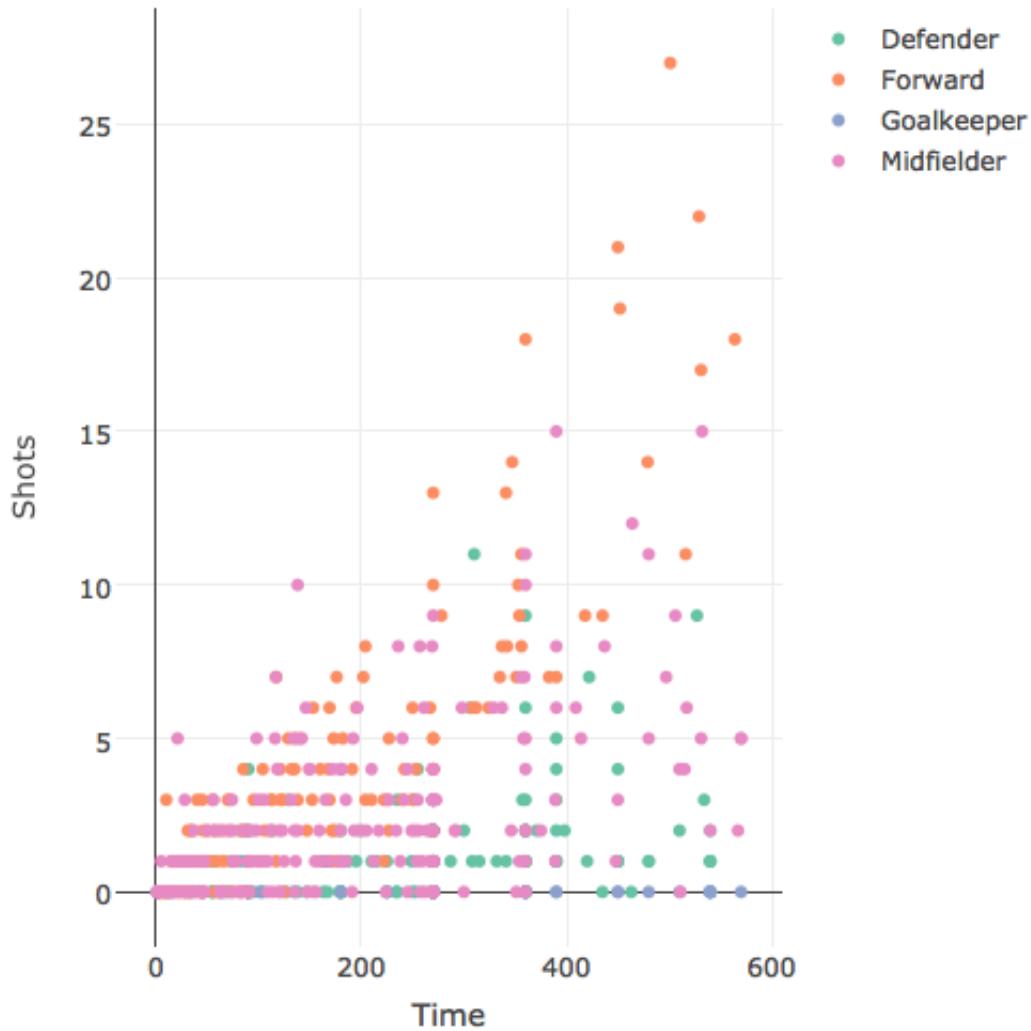


plot of chunk unnamed-chunk-104

Just like with `ggplot2`, the mappings you need depend on the type of plot you are creating. For example, scatterplots (`type = "scatter"`) need `x` and `y` defined, while a surface plot (`type = "surface"`) can be created with a single vector of elevation (we'll see an example in a few slides).

The help file for `plot_ly` includes a link with more documentation on the types of plots that can be made and the required mappings for each.

```
plot_ly(worldcup, type = "scatter", x = ~ Time, y = ~ Shots,
        color = ~ Position)
```



plot of chunk unnamed-chunk-105

The `plotly` package is designed so you can pipe data into `plot_ly` and add elements by piping into `add_*` functions (this idea is similar to adding elements to a `ggplot` object with `+`).

```
worldcup %>%
  plot_ly(x = ~ Time, y = ~ Shots, color = ~ Position) %>%
  add_markers()
```

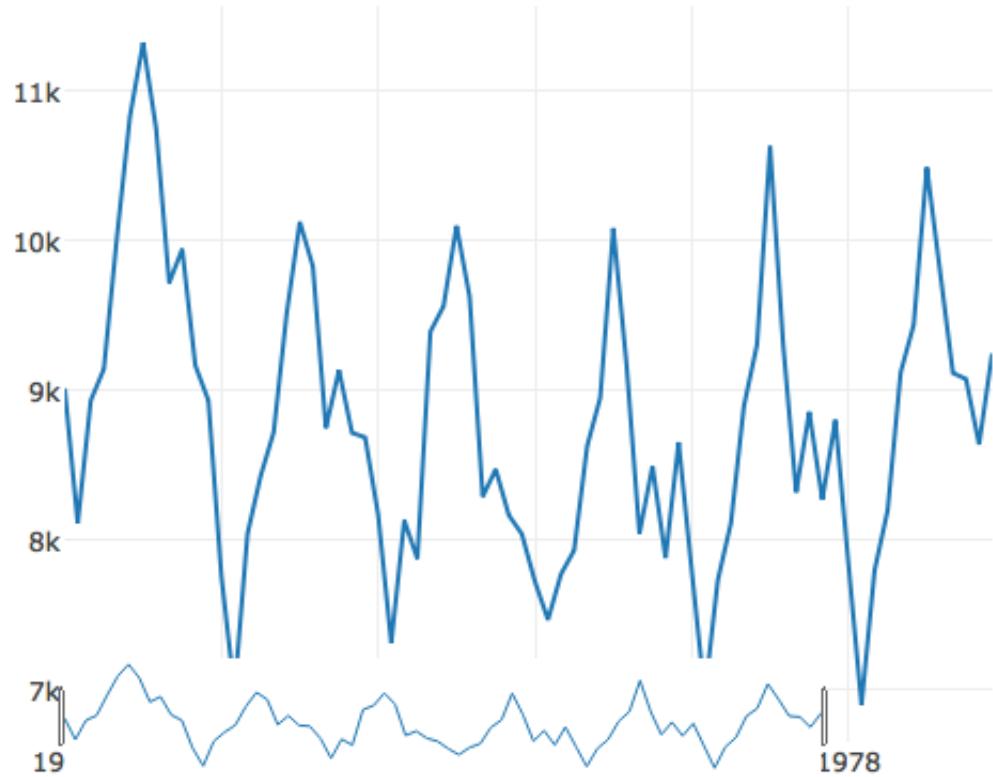
Some of the `add_*` functions include:

- add_markers
- add_lines
- add_paths
- add_polygons
- add_segments
- add_histogram

If you pipe to the `rangeslider` function, it allows the viewer to zoom in on part of the x range. (This can be particularly nice for time series.) \bigbreak

You should have a dataset available through your R session named `USAccDeaths`. This gives a monthly count of accidental deaths in the US for 1973 to 1978. This code will plot it and add a range slider on the lower x-axis.

```
plot_ly(x = time(USAccDeaths), y = USAccDeaths) %>%  
  add_lines() %>%  
  rangeslider()
```



plot of chunk unnamed-chunk-107

For a 3-D scatterplot, add a mapping to the `z` variable:

```
worldcup %>%
  plot_ly(x = ~ Time, y = ~ Shots, z = ~ Passes,
          color = ~ Position, size = I(3)) %>%
  add_markers()
```

- Defender
- Forward
- Goalkeeper
- Midfielder

Webgl is not supported by your browser - visit
<http://get.webgl.org> for more info

plot of chunk unnamed-chunk-108

The `volcano` data comes with R and is in a matrix format. Each value gives the elevation for a particular pair of x- and y-coordinates.

```
dim(volcano)
[1] 87 61
volcano[1:4, 1:4]
 [,1] [,2] [,3] [,4]
[1,] 100 100 101 101
[2,] 101 101 102 102
[3,] 102 102 103 103
[4,] 103 103 104 104

plot_ly(z = ~ volcano, type = "surface")
```

Webgl is not supported by your browser - visit <http://get.webgl.org> for more info



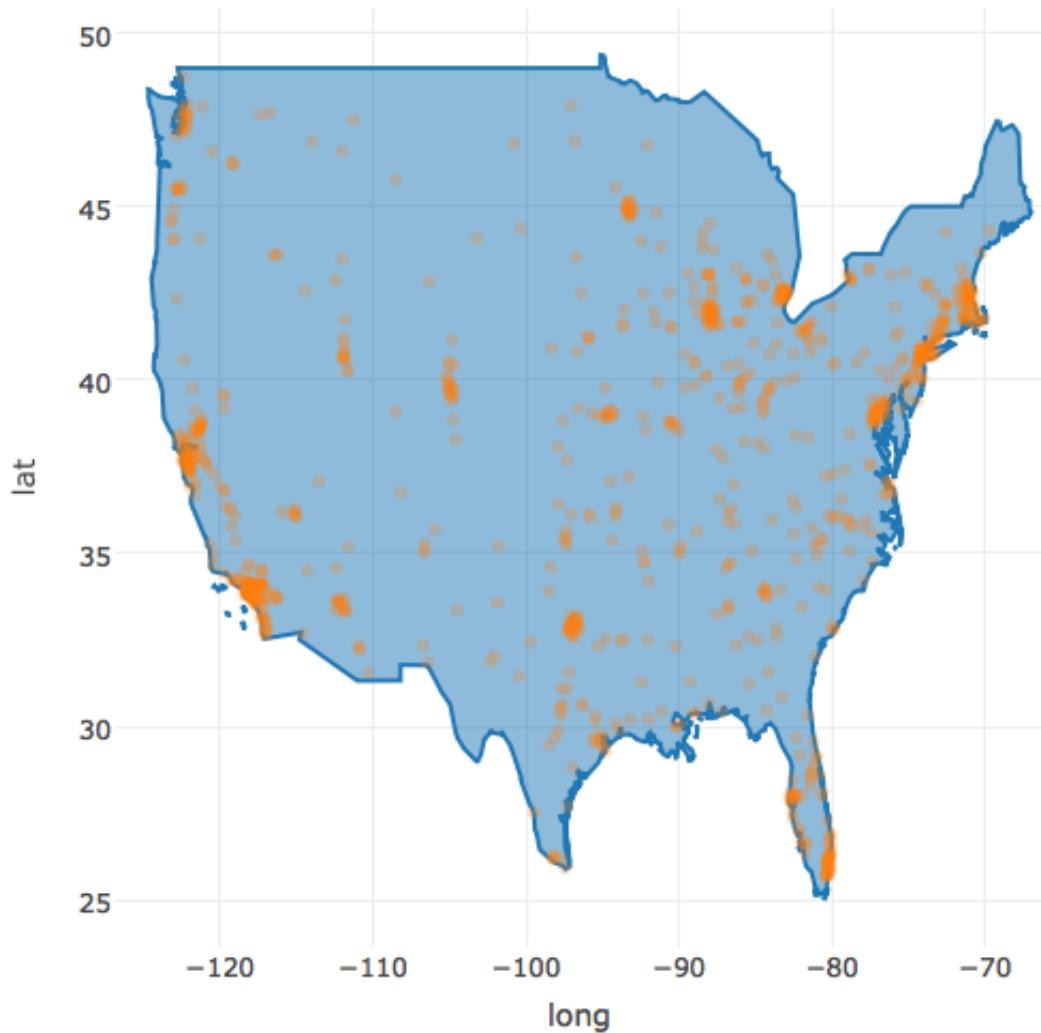
plot of chunk unnamed-chunk-110

Mapping with `plotly` can build on some data that comes with base R or other packages you've likely added (or can add easily, as with the `map_data` function from `ggplot2`). For example, we can map state capitals and cities with > 40,000 people using data in the `us.cities` data frame in the `maps` package:

```
head(maps::us.cities, 3)
  name country.etc    pop      lat      long capital
1 Abilene TX 113888 32.45 -99.74      0
2 Akron OH 206634 41.08 -81.52      0
3 Alameda CA  70069 37.77 -122.26      0
```

Here is code you can use to map all of these cities on a US map:

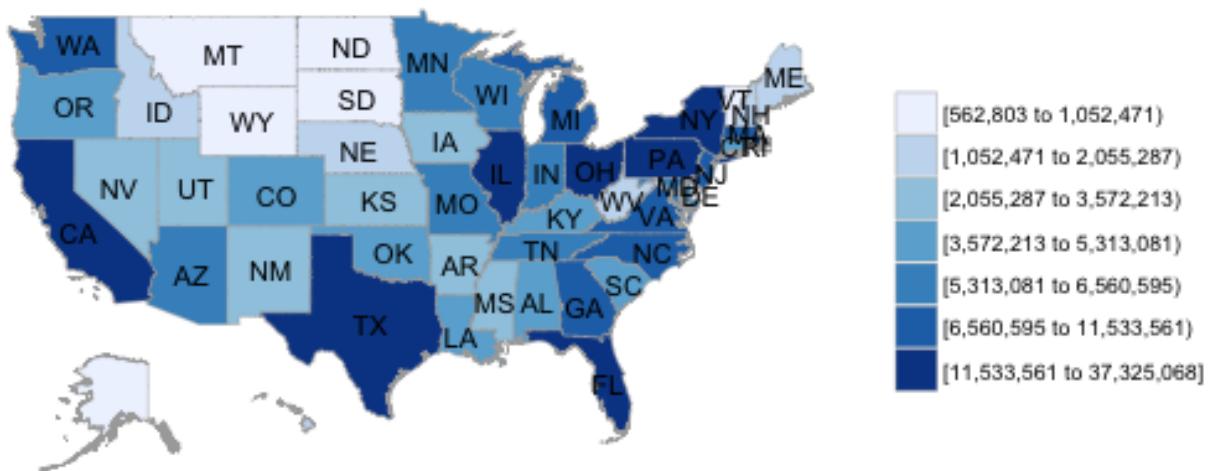
```
ggplot2::map_data("world", "usa") %>%
  group_by(group) %>% filter(-125 < long & long < -60 &
                           25 < lat & lat < 52) %>%
  plot_ly(x = ~long, y = ~lat) %>%
  add_polygons(hoverinfo = "none") %>%
  add_markers(text = ~paste(name, "<br />", pop), hoverinfo = "text",
              alpha = 0.25,
  data = filter(maps::us.cities, -125 < long & long < -60 &
                25 < lat & lat < 52)) %>%
  layout(showlegend = FALSE)
```



plot of chunk unnamed-chunk-112

You can also make choropleths interactive. Remember that we earlier created a choropleth of US state populations with the following code:

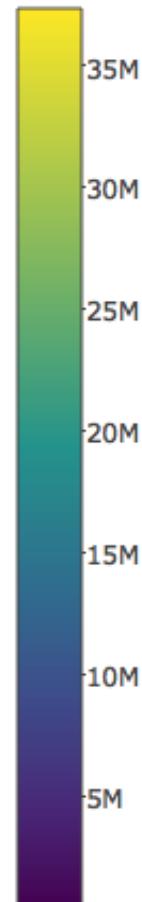
```
library(choroplethr)
data(df_pop_state)
state_choropleth(df_pop_state)
```



plot of chunk unnamed-chunk-113

You can use the following code with `plotly` to make an interactive choropleth instead:

```
us_map <- list(scope = 'usa',
  projection = list(type = 'albers usa'),
  lakecolor = toRGB('white'))
plot_geo() %>%
  add_trace(z = df_pop_state$value[df_pop_state$region !=
    "district of columbia"],
  text = state.name, locations = state.abb,
  locationmode = 'USA-states') %>%
add_markers(x = state.center[["x"]], y = state.center[["y"]],
  size = I(2), symbol = I(8), color = I("white"), hoverinfo = "none") %>%
layout(geo = us_map)
```



plot of chunk unnamed-chunk-114

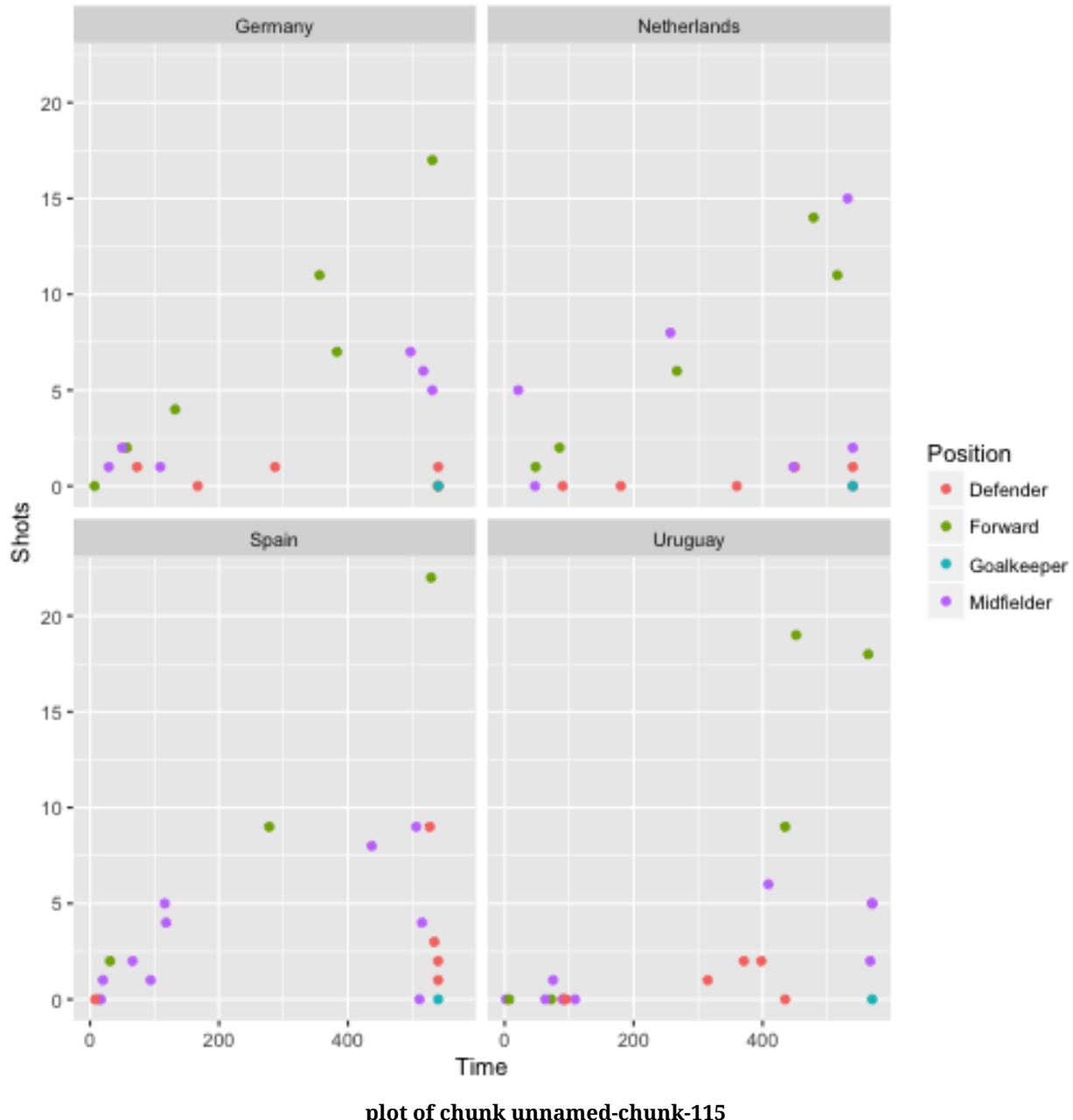
The other way to create a `plotly` graph is to first create a `ggplot` object and then transform it into an interactive graphic using the `ggplotly` function.

The following code can be used to plot Time versus Shots for the World Cup data in a regular, non-interactive plot:

```

shots_vs_time <- worldcup %>%
  mutate(Name = rownames(worldcup)) %>%
  filter(Team %in% c("Netherlands", "Germany", "Spain", "Uruguay")) %>%
  ggplot(aes(x = Time, y = Shots, color = Position, group = Name)) +
  geom_point() +
  facet_wrap(~ Team)
shots_vs_time

```



To make the plot interactive, just pass the `ggplot` object to `ggplotly`:

```
ggplotly(shots_vs_time)
```

With R, not only can you pull things from another website using an API, you can also upload or submit things.

There is a function in the `plotly` library, `plotly_POST`, that lets you post a plot you create in R to <https://plot.ly>.

You need a `plot.ly` account to do that, but there are free accounts available.

The creator of the R `plotly` package has written a bookdown book on the package that you can read [here](#). It provides extensive details and examples for using `plotly`.

Getting Started with D3 by Mike Dewar (a short book on D3 in JavaScript) is available for free [here](#).

Leaflet

Example data:

```
library(tigris)
denver_tracts <- tracts(state = "CO", county = 31, cb = TRUE)
load("data/fars_colorado.RData")
accident_data <- driver_data %>%
  dplyr::select(state, st_case, county, latitude, longitud,
                date, fatalities, drunk_dr) %>%
  dplyr::filter(county == 31 & longitud < -104.1) %>%
  dplyr::distinct()
```

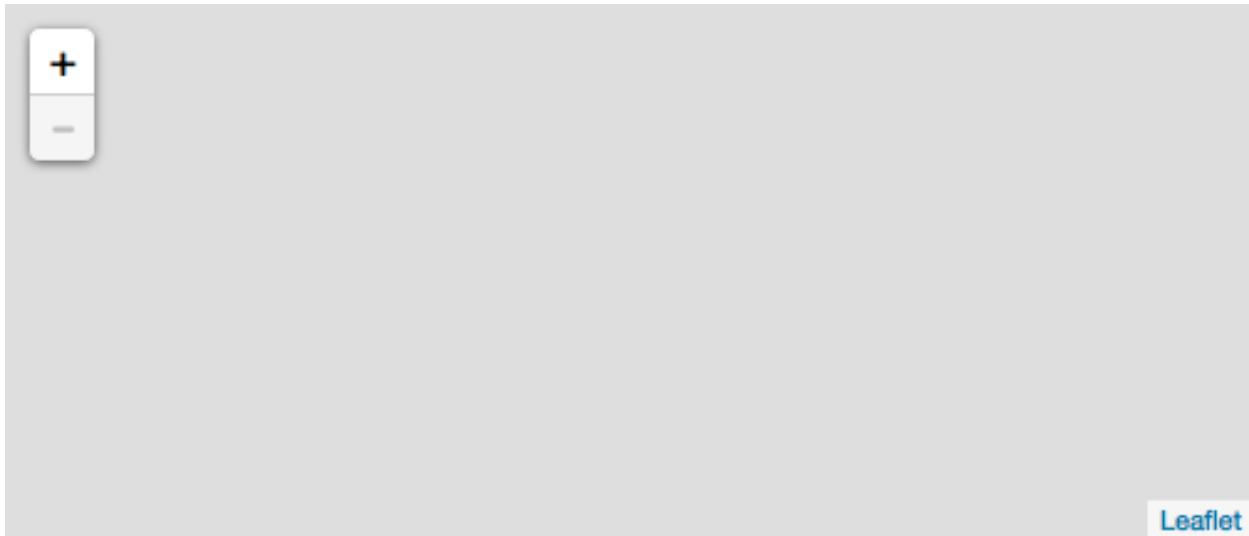
“Leaflet” is a JavaScript library for making interactive maps. You can find out more about the JavaScript version here: <http://leafletjs.com>

The `leaflet` package brings this functionality to R. The R Studio group has created a website on `leaflet`: <http://rstudio.github.io/leaflet/>. This website walks you through different options available with `leaflet`.

```
library(leaflet)
```

If you just run `leaflet()`, you just get a blank leaflet area:

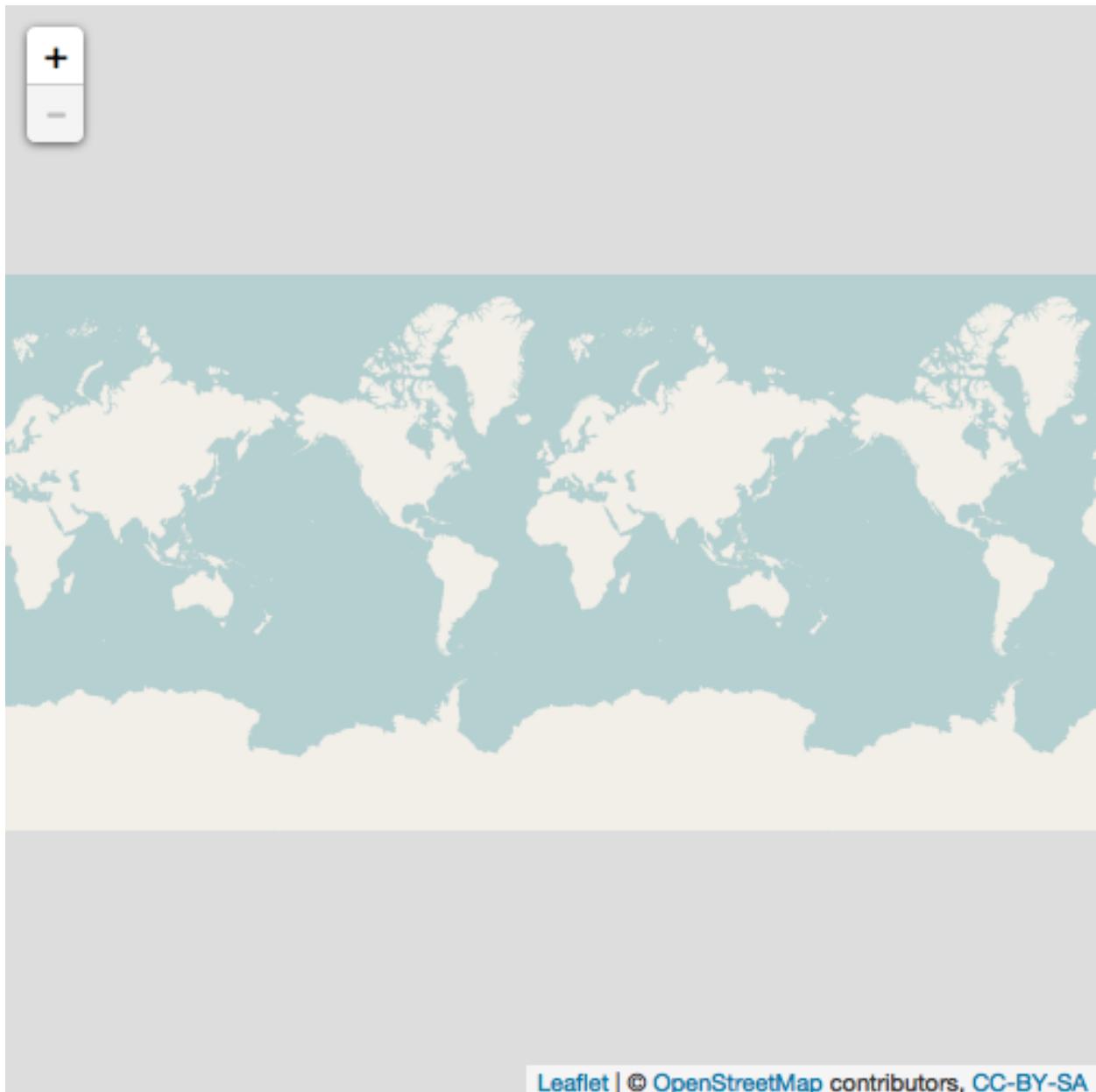
```
leaflet()
```



plot of chunk unnamed-chunk-119

In `leaflet`, the map background is composed of *tiles*. To get something more interesting, you'll need to add tiles to your leaflet map. If you don't include any other data, the leaflet map will include the world:

```
leaflet() %>%
  addTiles()
```

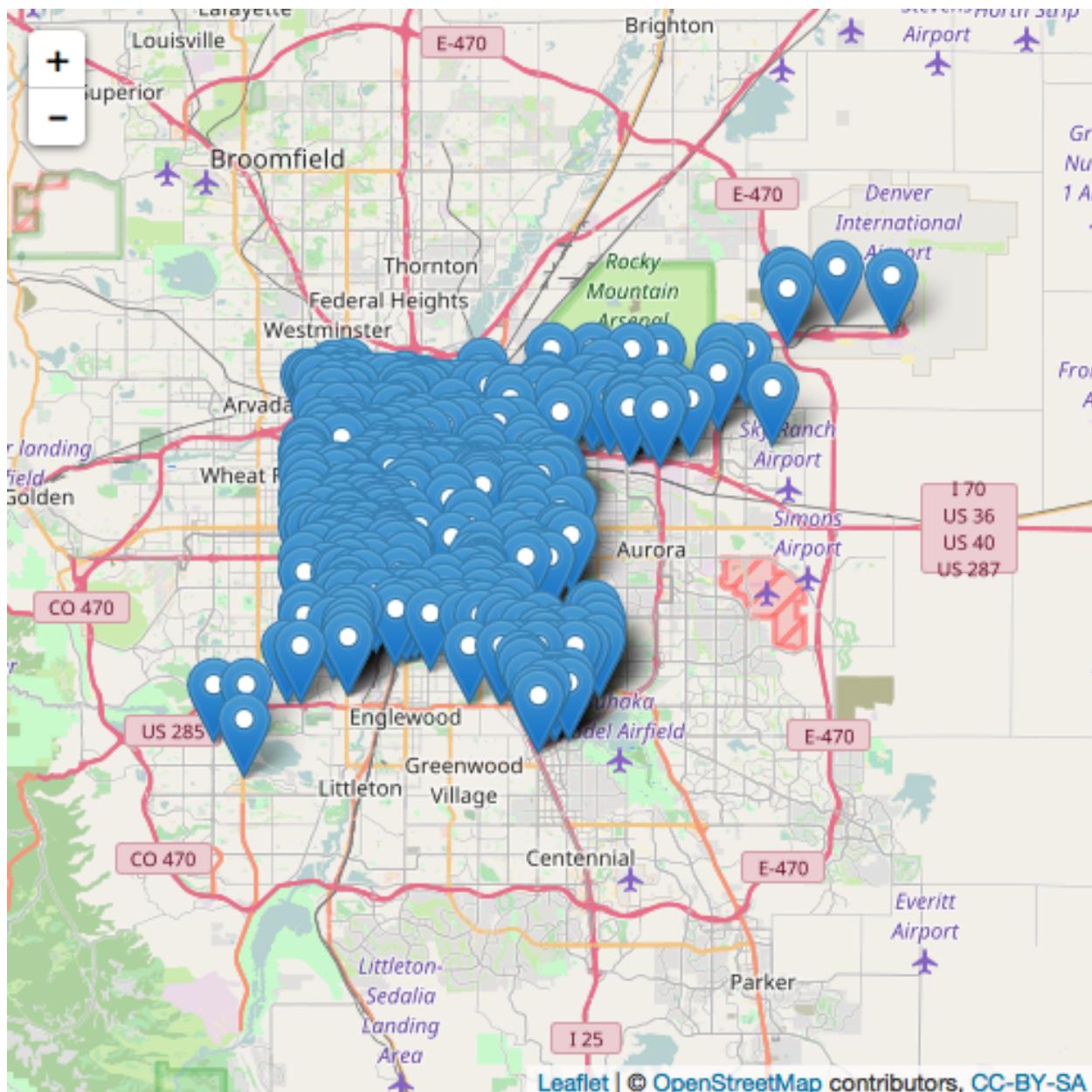


plot of chunk unnamed-chunk-120

For `htmlWidgets`, points are often referred to as *markers*.

Once you add these markers, the map will automatically scale to a reasonable size for their bounding box.

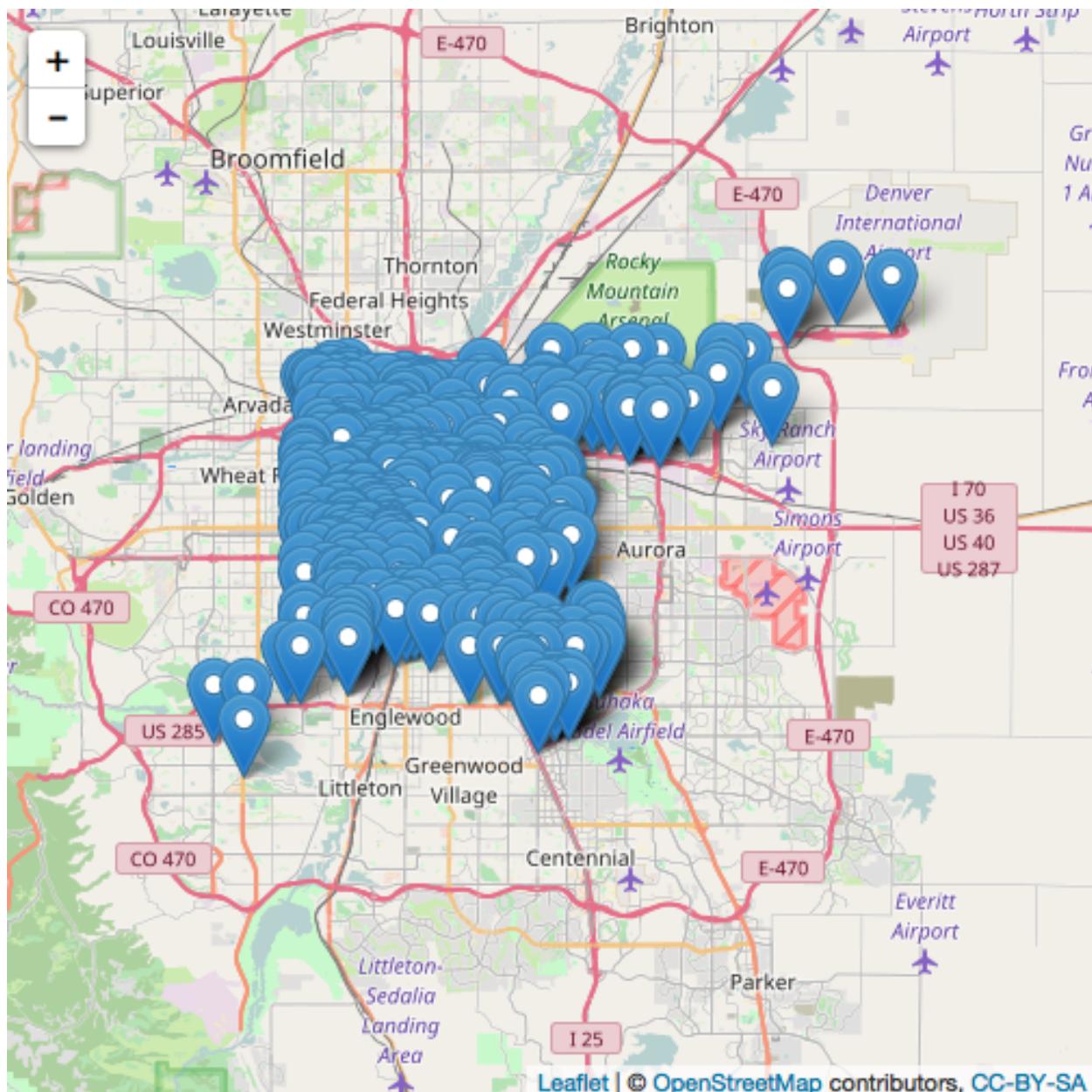
```
leaflet() %>%
  addTiles() %>%
  addMarkers(data = accident_data, lng = ~ longitud, lat = ~ latitude)
```



plot of chunk unnamed-chunk-121

Use `lng` and `lat` to tell R which columns contain data on longitude and latitude for each point. This is not needed if you are using a spatial object (e.g., `SpatialPointsDataFrame`). Further, R will try to guess the columns in a regular dataframe.

```
leaflet() %>%
  addTiles() %>%
  addMarkers(data = accident_data, lng = ~ longitude, lat = ~ latitude)
```



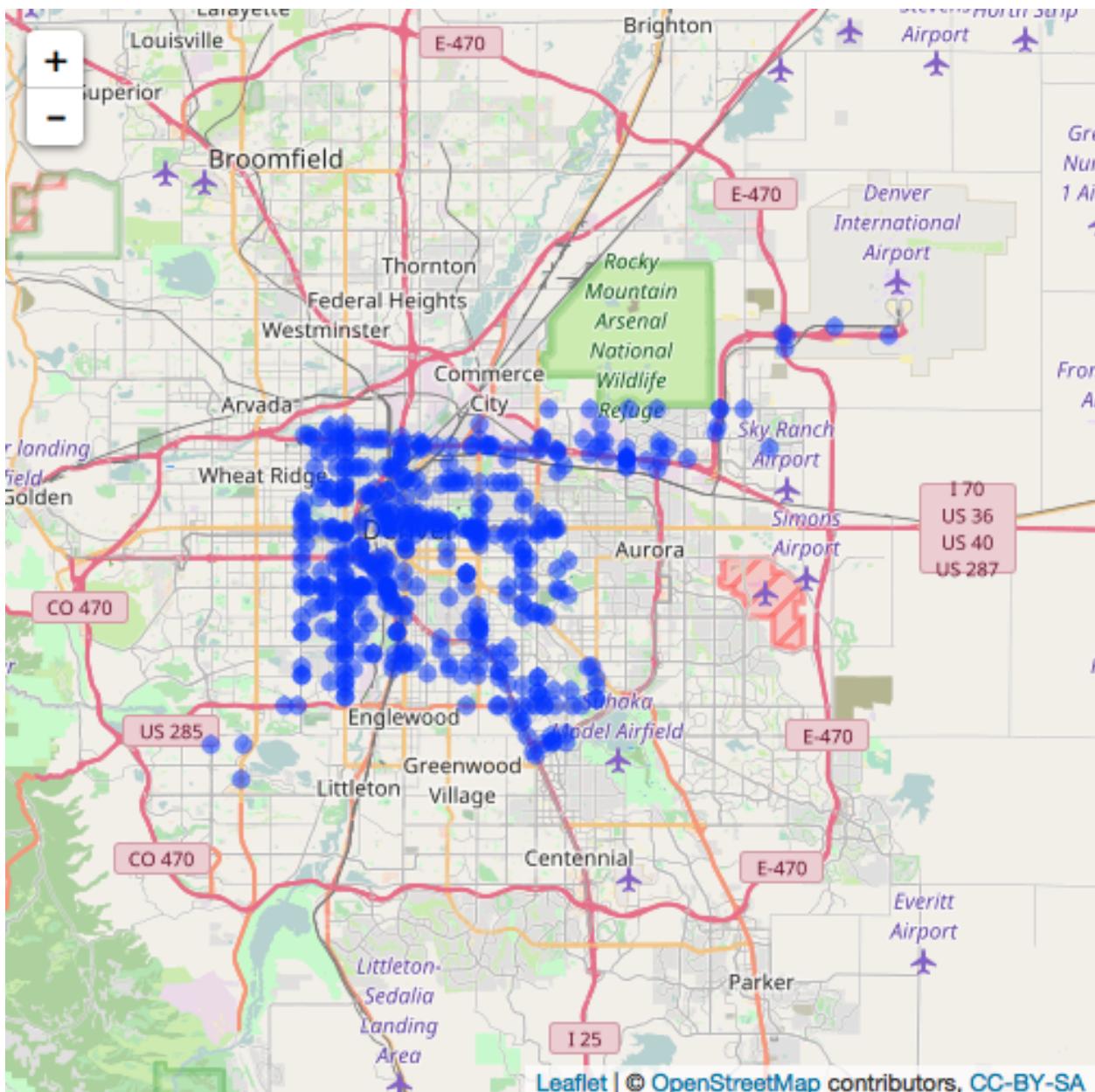
plot of chunk unnamed-chunk-122

You can use several types of R objects for your data for leaflet:

- Dataframe with columns for latitude and longitude
- Spatial objects (SpatialPoints, SpatialLines, etc.)
- Latitude-longitude matrix

You can choose circles for your markers instead by using addCircleMarkers. You can adjust the circle size with radius.

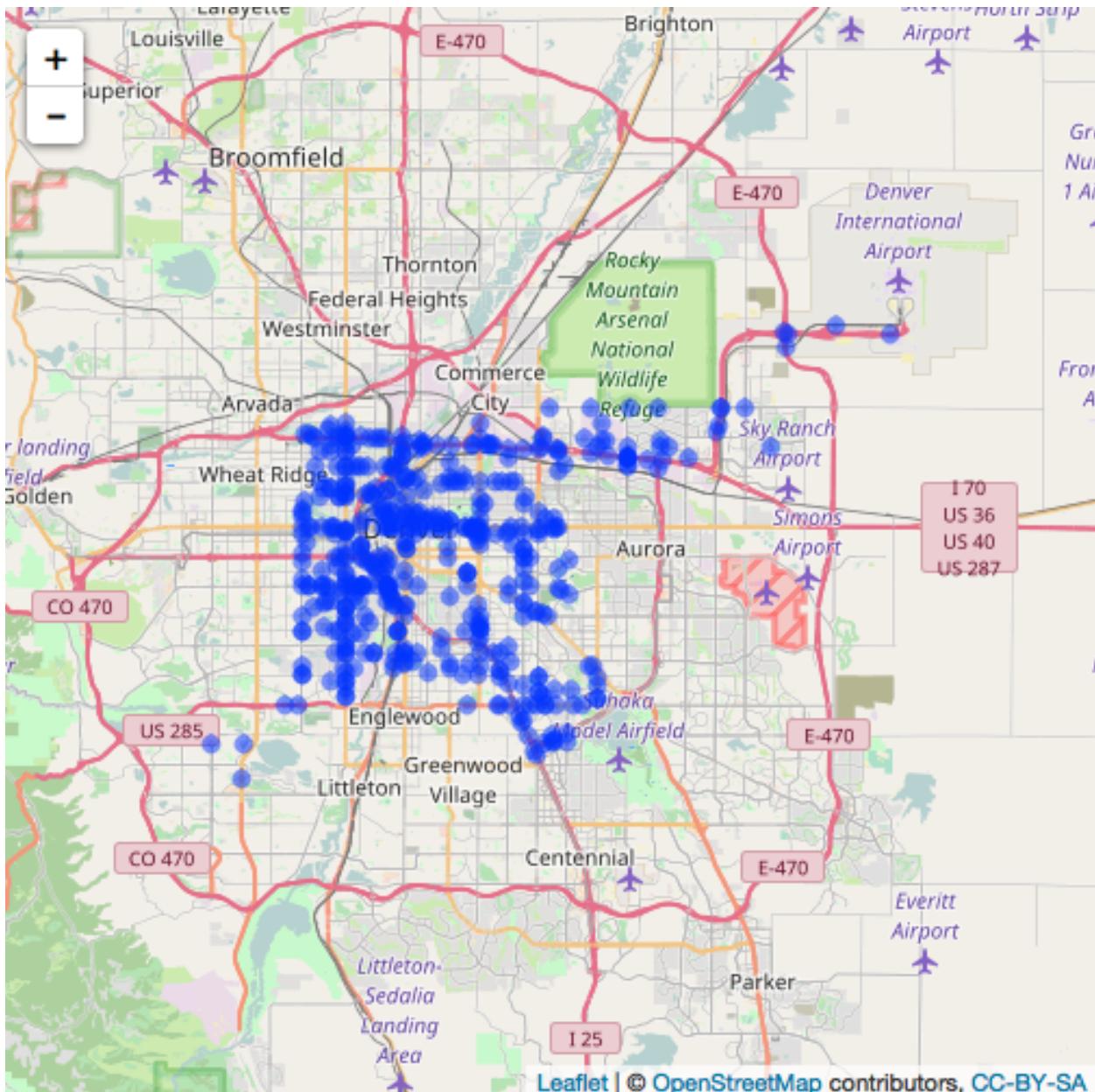
```
leaflet() %>%  
  addTiles() %>%  
  addCircleMarkers(data = accident_data, radius = 2,  
    lng = ~ longitude, lat = ~ latitude)
```



plot of chunk unnamed-chunk-123

The `radius` argument specifies the size of the circle. For `CircleMarkers`, the size will reset as you zoom in and out. If you want something with a constant radius (e.g., in meters), you can add `Circles`.

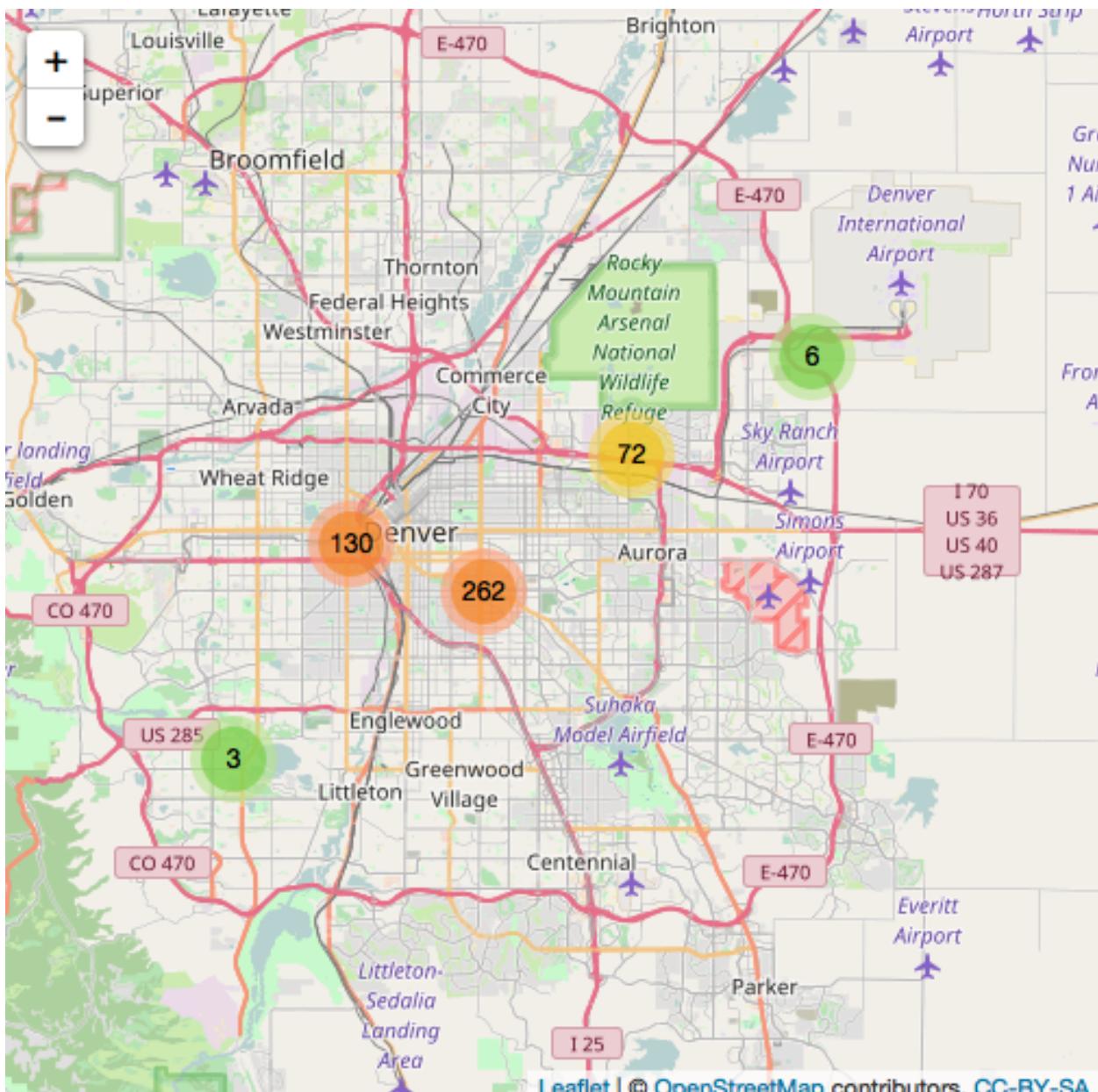
```
leaflet() %>%
  addTiles() %>%
  addCircleMarkers(data = accident_data, radius = 2,
    lng = ~ longitude, lat = ~ latitude)
```



plot of chunk unnamed-chunk-124

If you have a lot of overlapping data, you can also use the `clusterOptions` argument to show the markers as clusters that group together when you zoom out but split up when you zoom in:

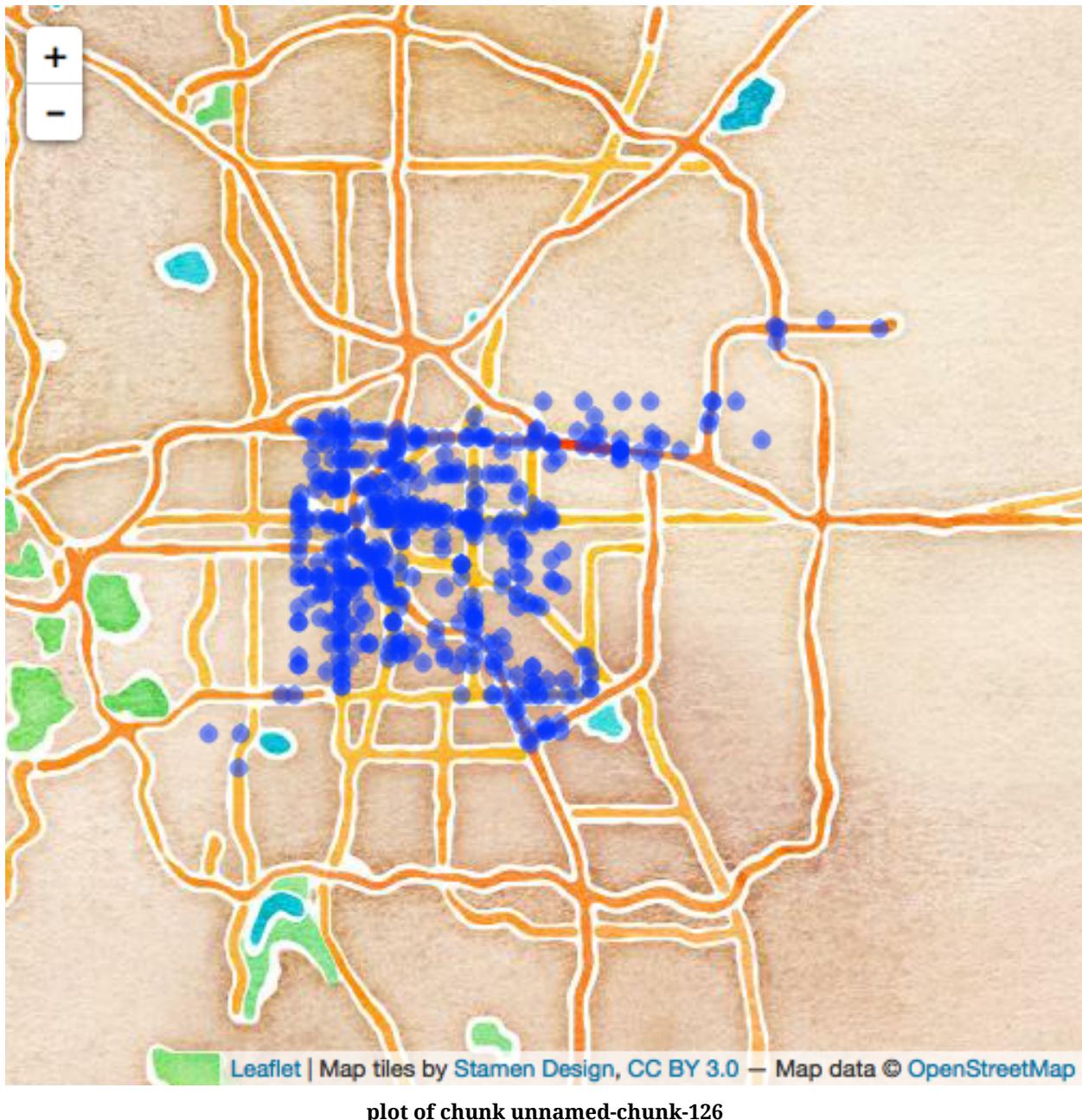
```
leaflet() %>%
  addTiles() %>%
  addMarkers(data = accident_data,
             lng = ~ longitude, lat = ~ latitude,
             clusterOptions = markerClusterOptions())
```



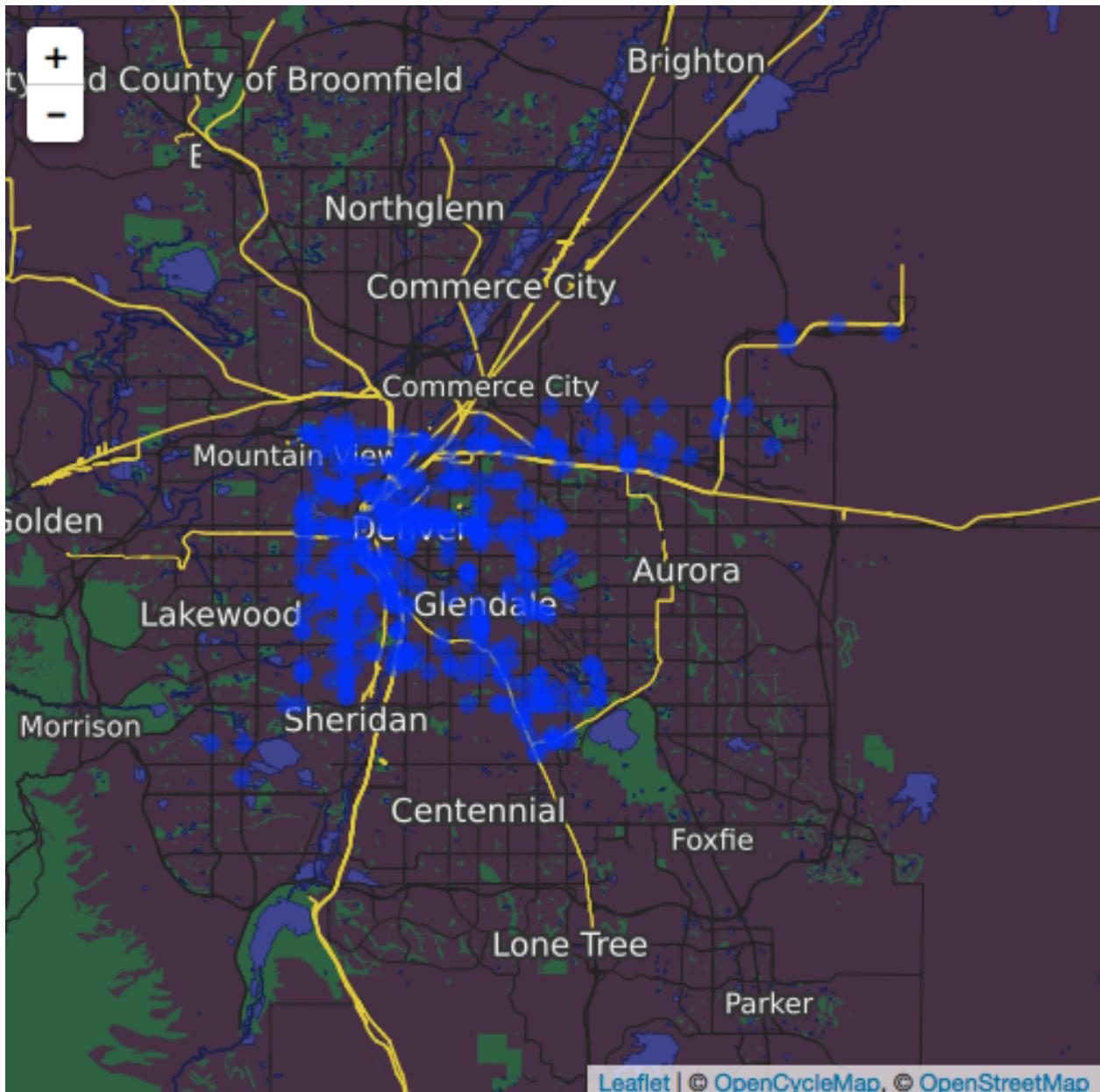
plot of chunk unnamed-chunk-125

For the background, the default is to use map tiles from OpenStreetMap. However, you can change the source of the tiles by using `addProviderTiles`. For example, to use Stamen Watercolor, you can call:

```
leaflet() %>%
  addProviderTiles("Stamen.Watercolor") %>%
  addCircleMarkers(data = accident_data, radius = 2,
    lng = ~ longitud, lat = ~ latitude)
```

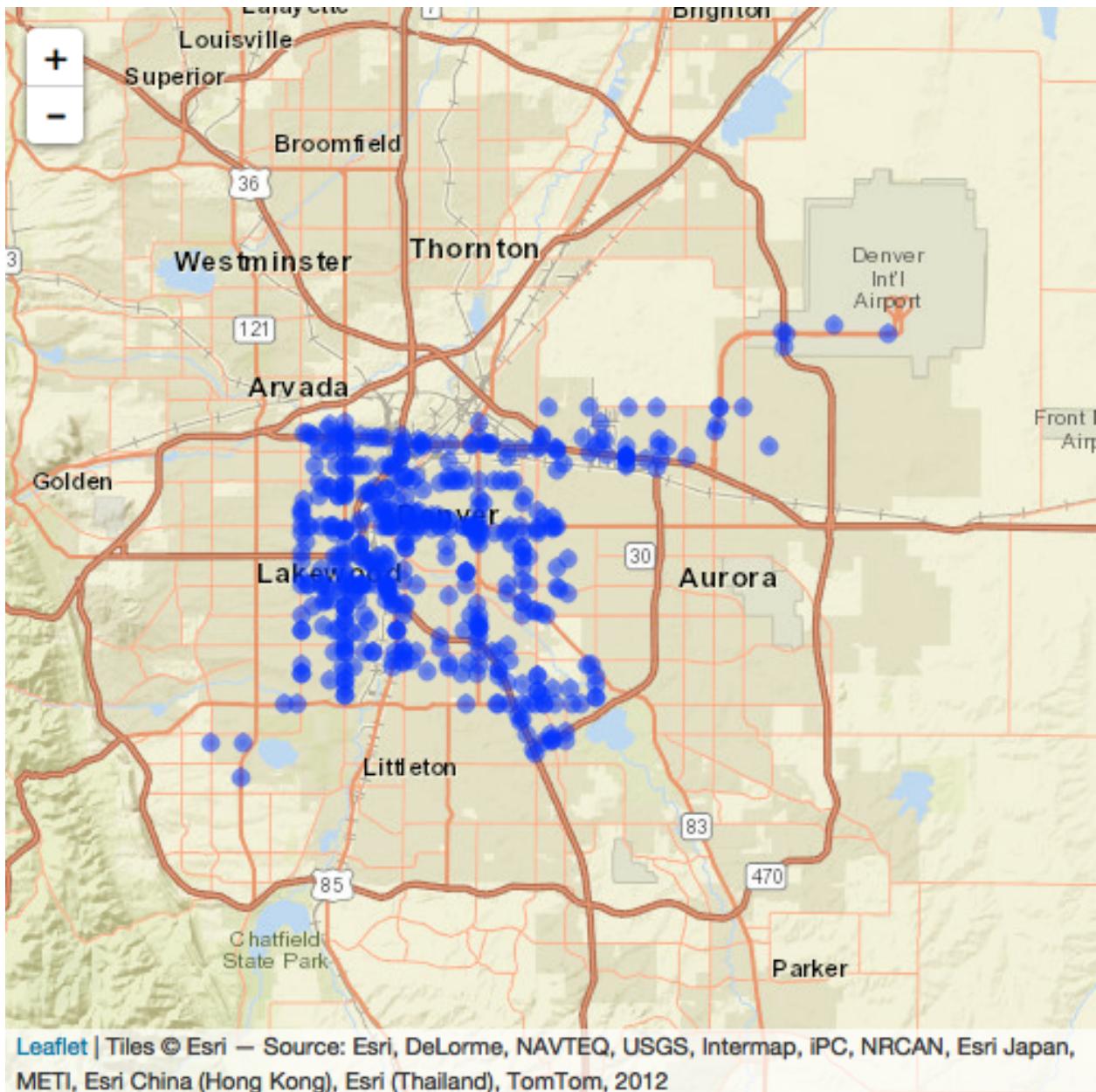


```
leaflet() %>%  
  addProviderTiles("Thunderforest.TransportDark") %>%  
  addCircleMarkers(data = accident_data, radius = 2,  
    lng = ~ longitud, lat = ~ latitude)
```



plot of chunk unnamed-chunk-127

```
leaflet() %>%
  addProviderTiles("Esri.WorldStreetMap") %>%
  addCircleMarkers(data = accident_data, radius = 2,
    lng = ~ longitud, lat = ~ latitude)
```



plot of chunk unnamed-chunk-128

You can see previews of provider choices here: <http://leaflet-extras.github.io/leaflet-provider-s/preview/index.html>.

You can use the `popup` option to show information when the user clicks on a marker.

It's easiest to do this if you have the information you want to show in the dataframe with the location data. For example, we have date-time, number of fatalities, and number of drunk drivers in this data:

```
accident_data %>%
  dplyr::select(date, fatals, drunk_dr) %>%
  dplyr::slice(1:3)
#> #> #>
#> #> #>
```

| | date | fatals | drunk_dr |
|---|---------------------|--------|----------|
| 1 | 2001-01-04 19:00:00 | 1 | 1 |
| 2 | 2001-01-03 07:00:00 | 1 | 1 |
| 3 | 2001-01-05 20:00:00 | 1 | 1 |

If we want to show day of the week, month, hour, and number of fatalities, go ahead and calculate any value not already in the dataset:

```
library(lubridate)
accident_data <- accident_data %>%
  mutate(weekday = wday(date, label = TRUE, abbr = FALSE),
        month = month(date, label = TRUE, abbr = FALSE),
        hour = format(date, format = "%H:%M"))
```

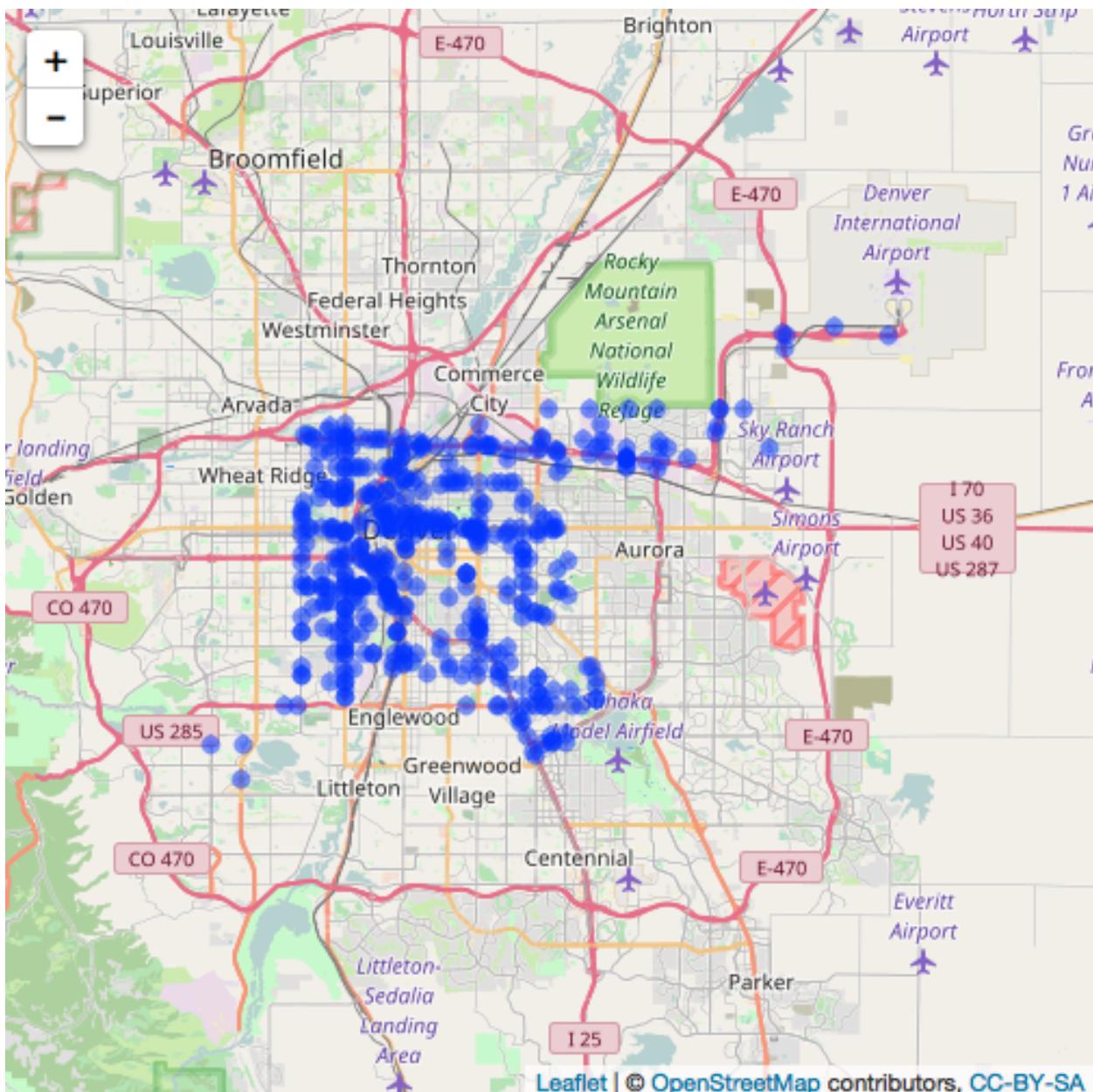
The popup text needs to be a character vector, written in HTML syntax. You can create that vector first, and then pass it to the `popup` argument.

```
popup_info <- paste0("<b>Weekday:</b>  ",
                      accident_data$weekday, "<br/>",
                      "<b>Month:</b>  ",
                      accident_data$month, "<br/>",
                      "<b>Hour:</b>  ",
                      accident_data$hour, "<br/>",
                      "<b>Fatalities:</b>  ",
                      accident_data$fatals)

popup_info[1:3]
[1] "<b>Weekday:</b> Thursday<br/><b>Month:</b> January<br/><b>Hour:</b> 19:00<br/><b>Fatalities:</b> 1"
[2] "<b>Weekday:</b> Wednesday<br/><b>Month:</b> January<br/><b>Hour:</b> 07:00<br/><b>Fatalities:</b> 1"
[3] "<b>Weekday:</b> Friday<br/><b>Month:</b> January<br/><b>Hour:</b> 20:00<br/><b>Fatalities:</b> 1"
```

Now pass that vector to the `popup` argument for the layer you want to pair it with:

```
leaflet() %>%
  addTiles() %>%
  addCircleMarkers(data = accident_data, radius = 2,
    lng = ~ longitud, lat = ~ latitude,
    popup = popup_info)
```



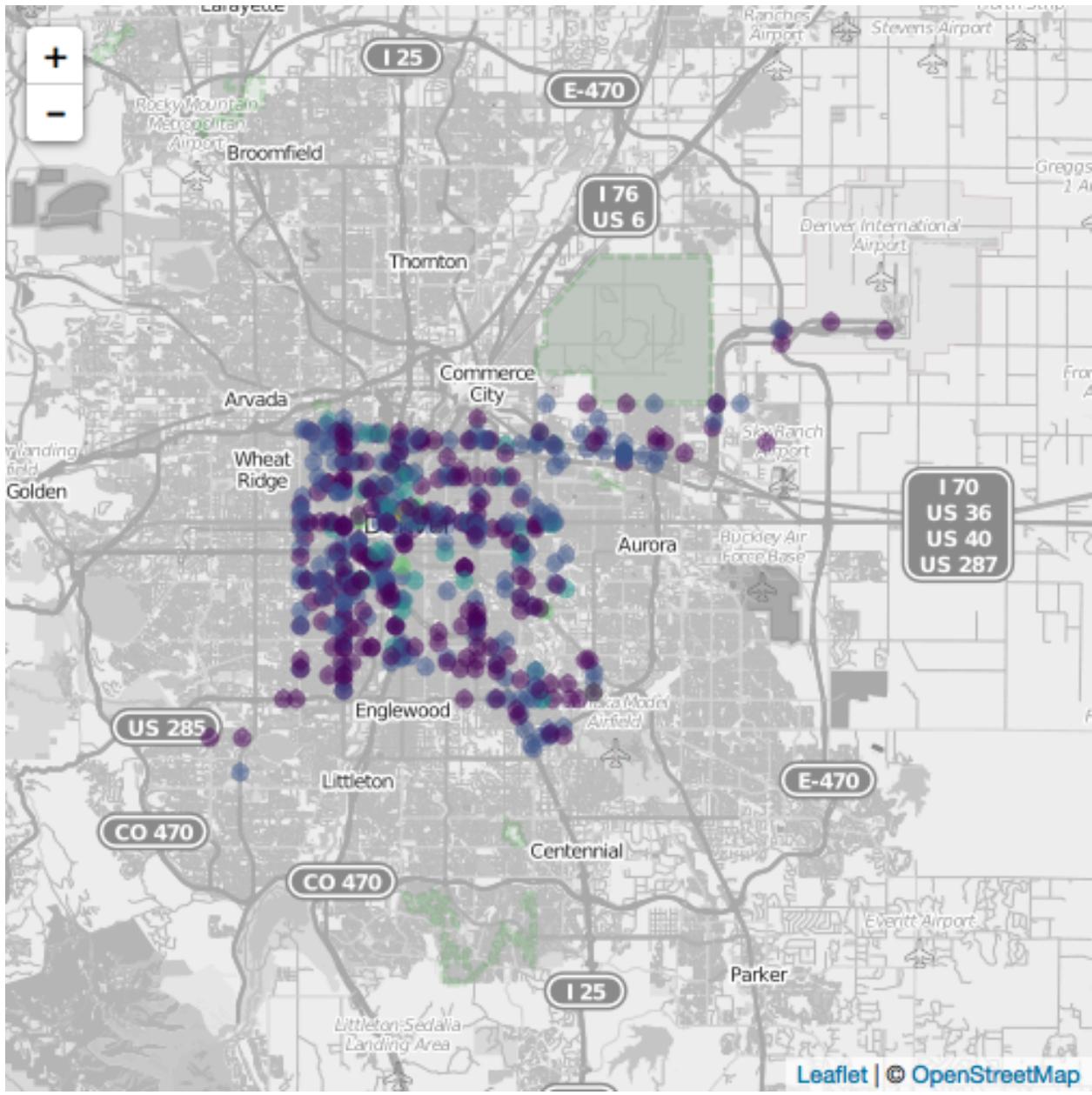
plot of chunk unnamed-chunk-132

If you aren't familiar with HTML syntax, here's one cheatsheet: <http://web.stanford.edu/group/csp/cs21/html天上人間>

In the popups, you can use HTML to format things like color, typeface, and size. You can also add links.

To use color to show a value, you need to do a few things. First, you need to the the `colorFactor` function (or another in its family) to create a function for mapping from values to colors. Then, you need to use this within the call to add the markers.

```
library(viridisLite)
pal <- colorFactor(viridis(5), accident_data$drunk_dr)
leaflet() %>%
  addProviderTiles("OpenStreetMap.BlackAndWhite") %>%
  addCircleMarkers(data = accident_data, radius = 2,
    lng = ~ longitud, lat = ~ latitude,
    popup = popup_info,
    color = pal(accident_data$drunk_dr))
```



plot of chunk unnamed-chunk-133

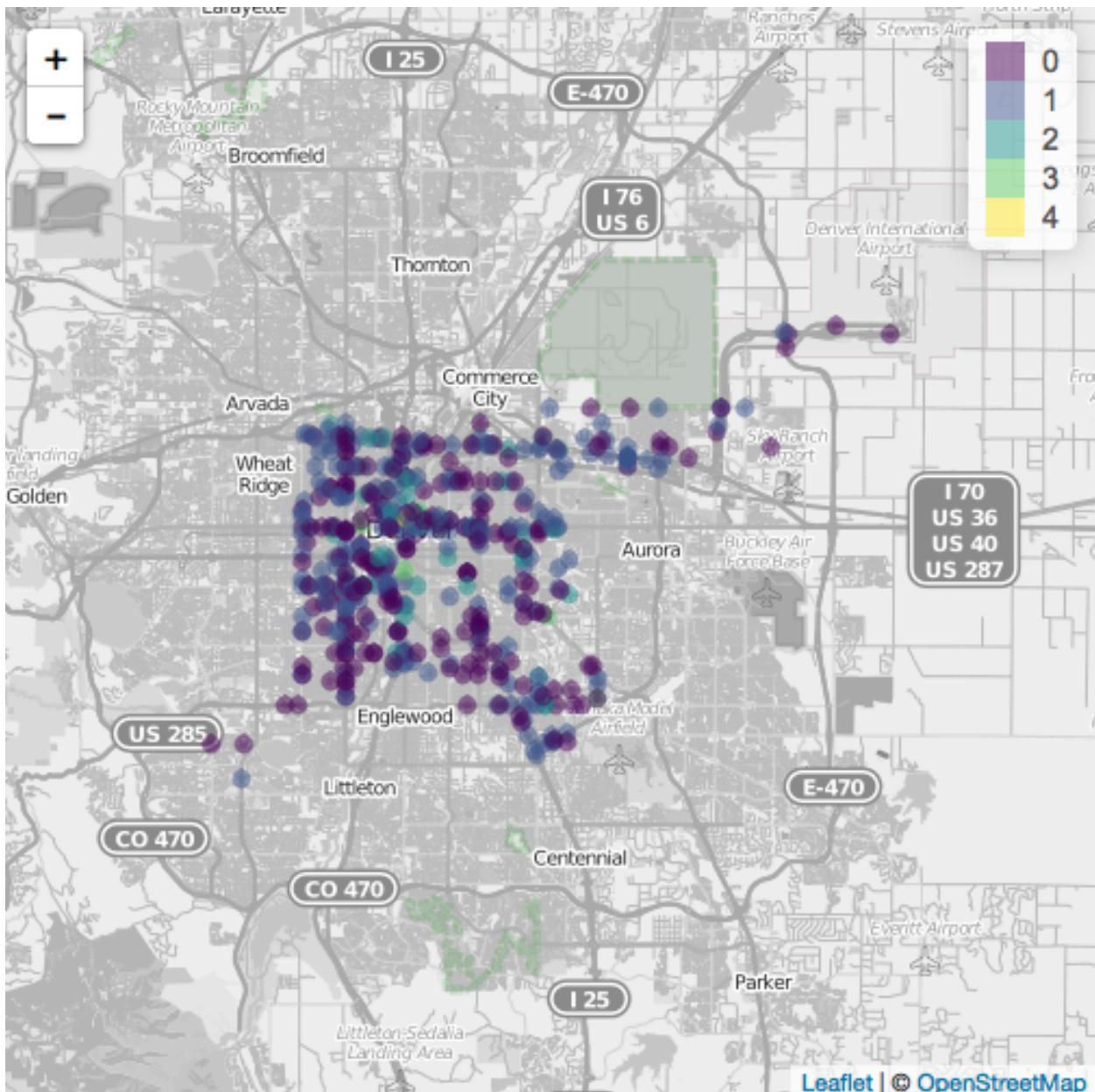
The `colorFactor` function (and friends) are a pretty cool type of function that actually creates a new function:

```
pal <- colorFactor(viridis(5), accident_data$drunk_dr)
class(pal)
[1] "function"
head(pal)

1 structure(function (x)
2 {
3   if (length(x) == 0 || all(is.na(x))) {
4     return(pf(x))
5   }
6   lvs = getLevels(domain, x, lvs, ordered)
```

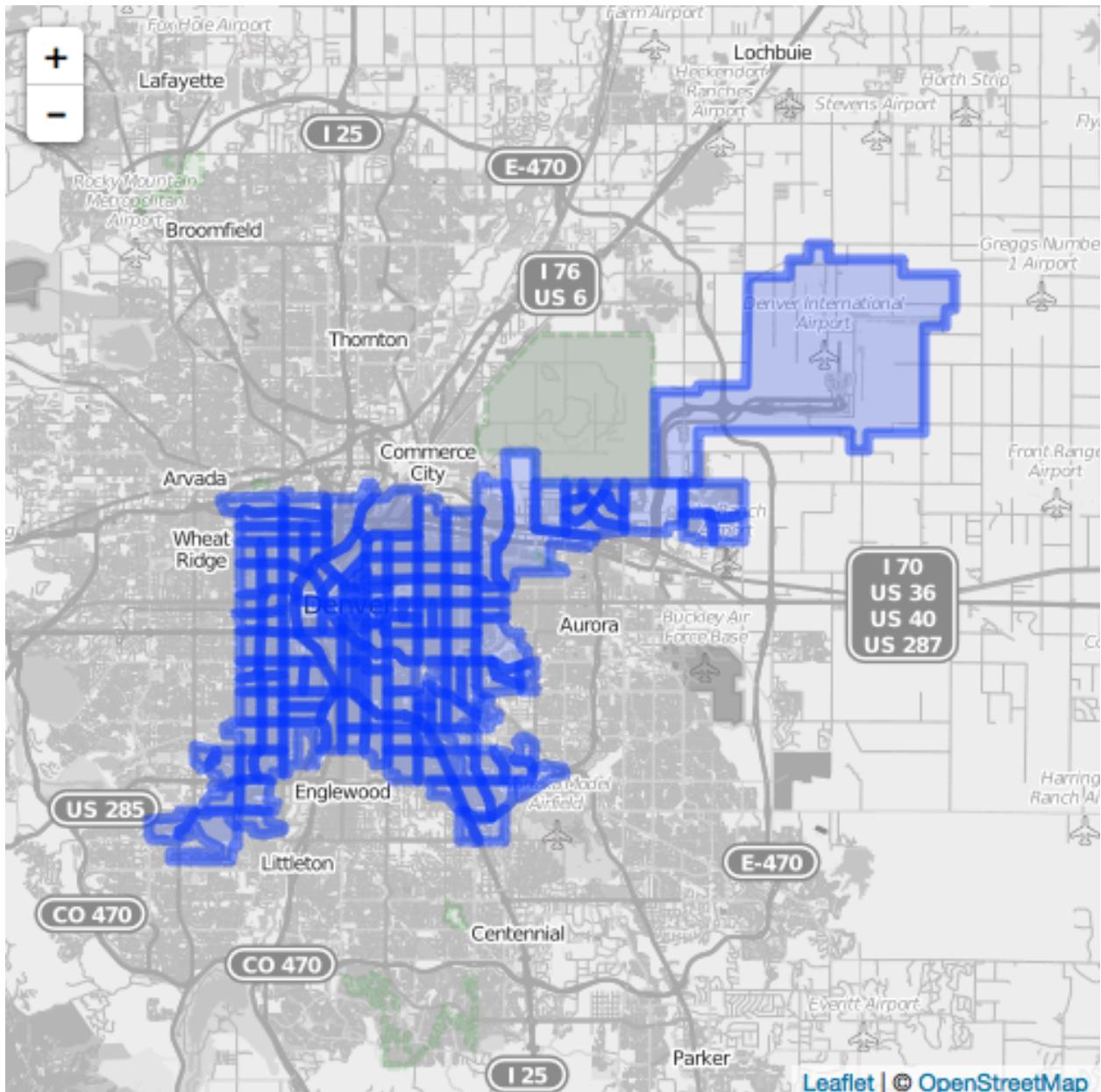
Once you are showing something with color, you can add a legend to explain it. You can do that with the `addLegend` function, which must include values for the color palette and values for each point from this color palette.

```
library(viridisLite)
pal <- colorFactor(viridis(5), accident_data$drunk_dr)
leaflet() %>%
  addProviderTiles("OpenStreetMap.BlackAndWhite") %>%
  addCircleMarkers(data = accident_data, radius = 2,
    lng = ~ longitud, lat = ~ latitude,
    popup = popup_info,
    color = pal(accident_data$drunk_dr)) %>%
  addLegend(pal = pal, values = accident_data$drunk_dr)
```



You can add polygons with the `addPolygons` function.

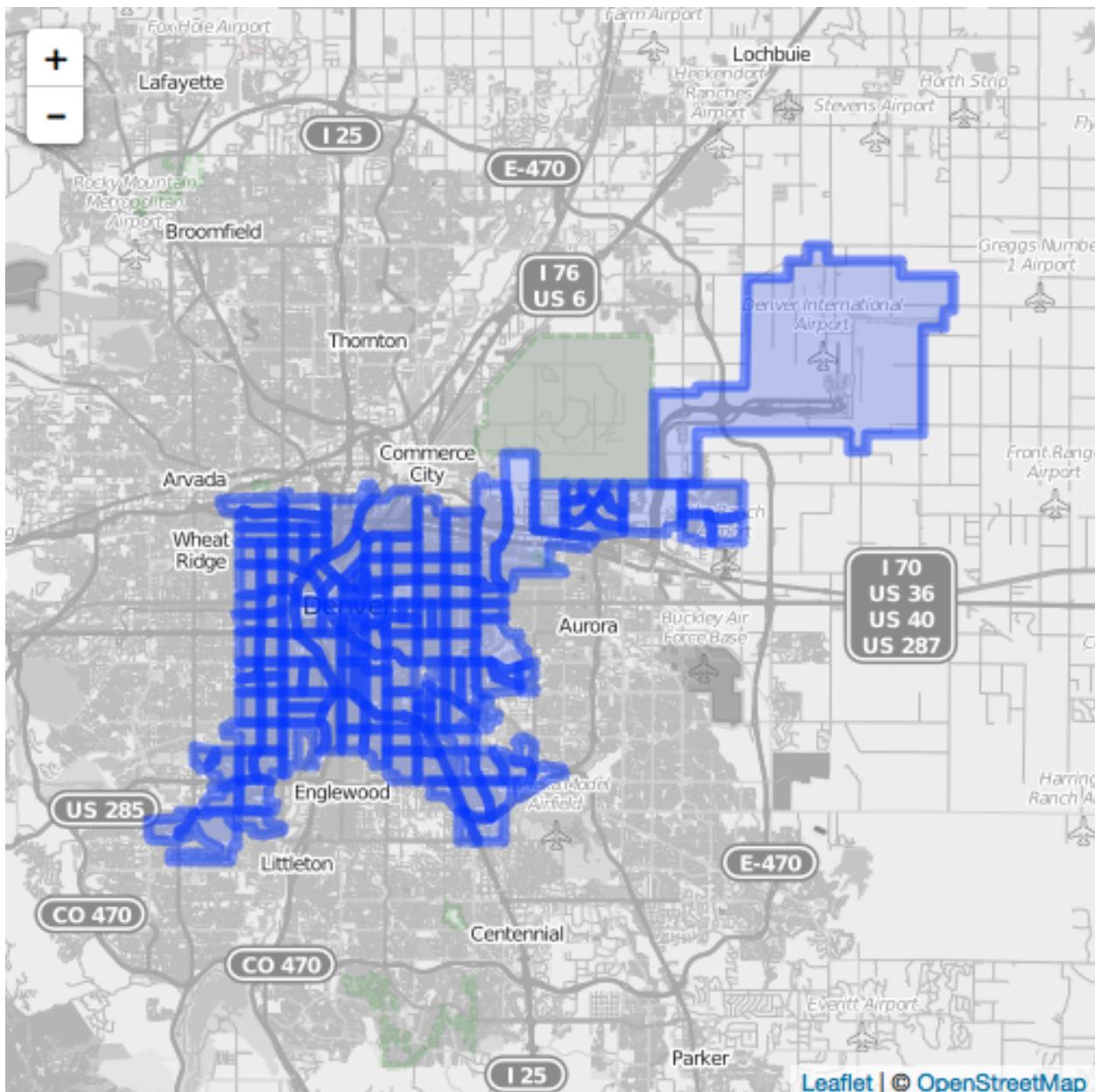
```
leaflet() %>%  
  addProviderTiles("OpenStreetMap.BlackAndWhite") %>%  
  addPolygons(data = denver_tracts)
```



plot of chunk unnamed-chunk-136

You can add popups for polygons, as well.

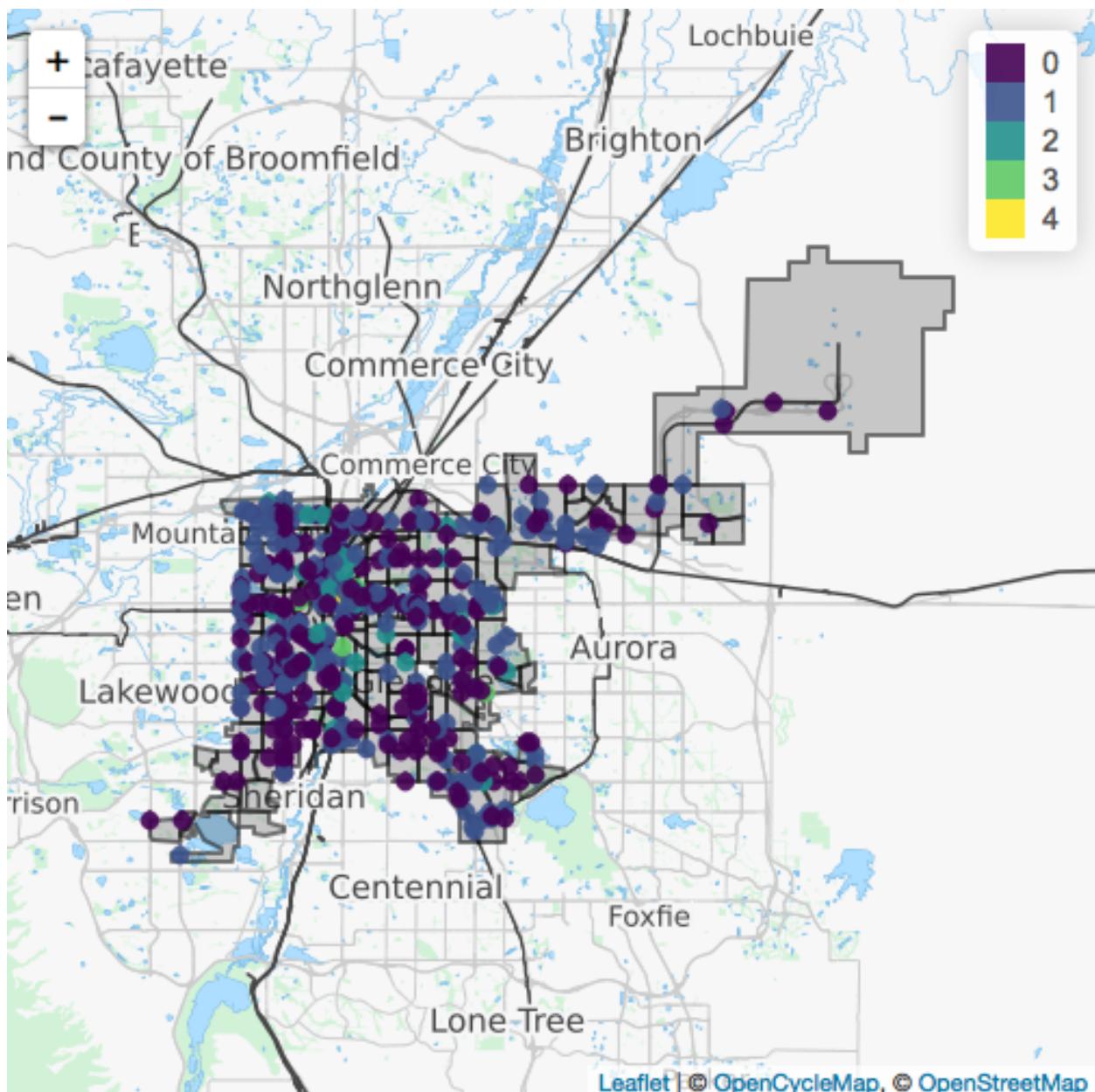
```
polygon_popup <- paste0("Tract ID: ",  
                         denver_tracts@data$NAME)  
leaflet() %>%  
  addProviderTiles("OpenStreetMap.BlackAndWhite") %>%  
  addPolygons(data = denver_tracts, popup = polygon_popup)
```



plot of chunk unnamed-chunk-137

You can overlay different elements. For example, you can show both accidents and tracts:

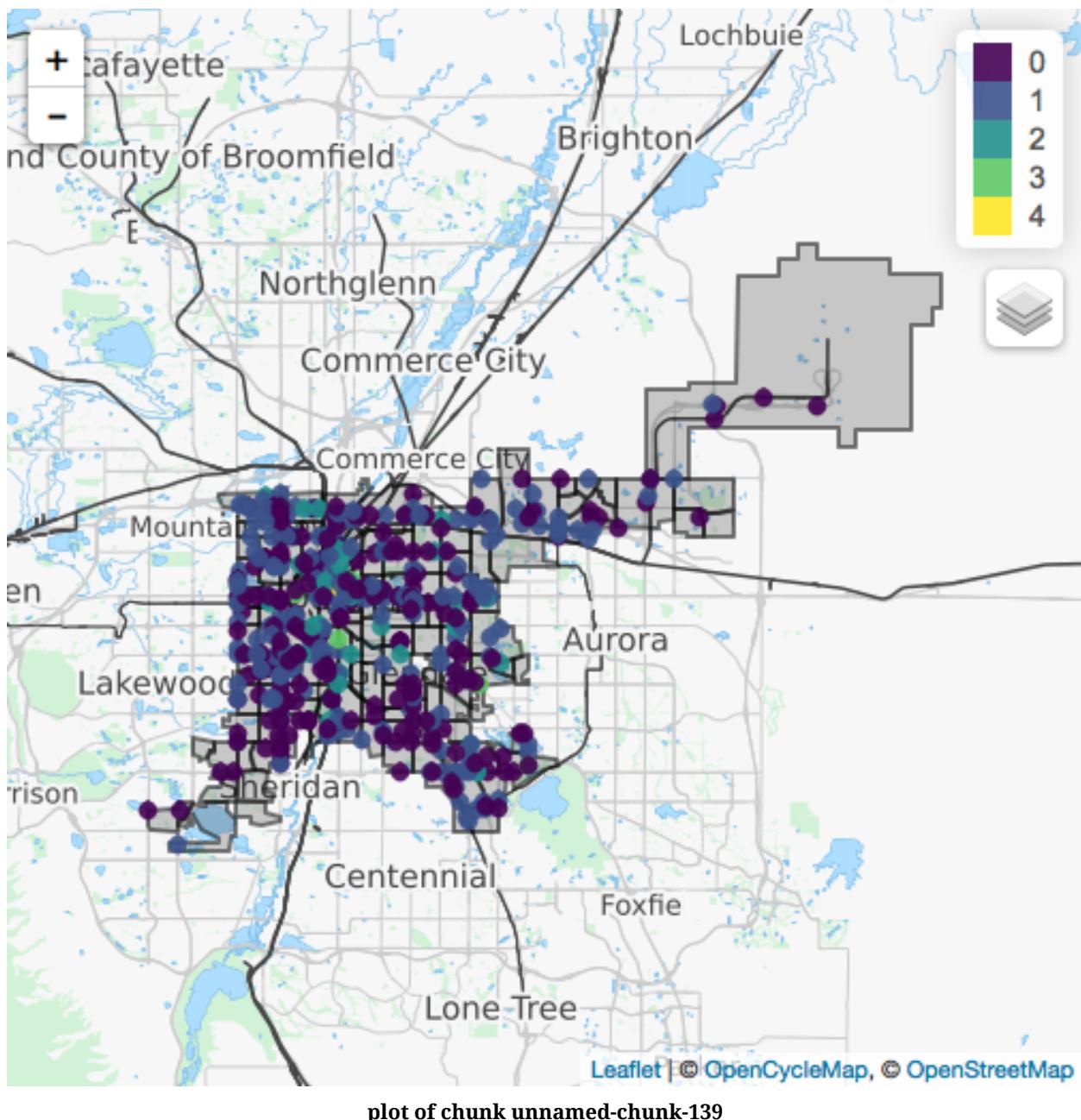
```
polygon_popup <- paste0("Tract ID: ",  
                           denver_tracts@data$NAME)  
accident_data_sp <- accident_data  
coordinates(accident_data_sp) <- c("longitude", "latitude")  
proj4string(accident_data_sp) <- CRS(proj4string(denver_tracts))  
leaflet() %>%  
  addProviderTiles("Thunderforest.Transport") %>%  
  addPolygons(data = denver_tracts, popup = polygon_popup,  
               color = "#000000", fillColor = "969696",  
               weight = 2) %>%  
  addCircleMarkers(data = accident_data_sp, radius = 2,  
                   popup = popup_info, opacity = 0.9,  
                   color = pal(accident_data$drunk_dr)) %>%  
  addLegend(pal = pal, values = accident_data$drunk_dr, opacity = 0.9)
```



plot of chunk unnamed-chunk-138

You can add the ability for the user to pick which layers to see using `addLayersControls`.

```
leaflet() %>%
  addProviderTiles("Thunderforest.Transport", group = "base map") %>%
  addPolygons(data = denver_tracts, popup = polygon_popup,
              color = "#000000", fillColor = "969696",
              weight = 2, group = "tracts") %>%
  addCircleMarkers(data = accident_data_sp, radius = 2,
                   popup = popup_info, opacity = 0.9,
                   color = pal(accident_data$drunk_dr),
                   group = "accidents") %>%
  addLegend(pal = pal, values = accident_data$drunk_dr, opacity = 0.9) %>%
  addLayersControl(baseGroups = c("base map"),
                   overlayGroups = c("tracts", "accidents"))
```



Find out more

Here are some good tutorials for trying out other examples of `leaflet` in R:

<http://zevross.com/blog/2015/10/14/manipulating-and-mapping-us-census-data-in-r-using-the-acs-tigris-and-leaflet-packages-3/>

<http://robinlovelace.net/r/2015/02/01/leaflet-r-package.html>

<http://trendct.org/2015/06/26/tutorial-how-to-put-dots-on-a-leaflet-map-with-r/>

Creating your own widget

If you find a JavaScript visualization library and would like to create bindings to R, you can create your own package for a new `htmlWidget`.

There is advice on creating your own widget for R available at http://www.htmlwidgets.org/develop_intro.html.

4.5 The grid Package

The `grid` package in R implements the primitive graphical functions that underly the `ggplot2` plotting system. While one typically does not interact directly with the `grid` package (it is imported by the `ggplot2` package), it is necessary to understand some aspects of the `grid` package in order to build new geoms and graphical elements for `ggplot2`. In this section we will discuss key elements of the `grid` package that can be used in extending `ggplot2`.

While the `grid` package can be used to produce graphical output directly, it is seldom used for that purpose. Rather, the `grid` package provides a set of functions and classes that represent graphical objects or grobs, that can be manipulated like any other R object. With grobs, we can manipulate graphical elements (“edit” them) using standard R functions.

Overview of grid graphics

Grid graphics is a system for plotting within R, and should be thought of as a separate system than base R graphics. The `ggplot2` package is built on top of grid graphics, so grid graphics functions can be used to customize and manipulate `ggplot2` objects. The grid graphics functions interact less well with plots created using the base R graphics system. While we have focused on plotting using `ggplot2` in this course, we have covered a few plots created using base R, specifically the maps created by running a `plot` call on a spatial object, like a `SpatialPoints` object.

While it is not straightforward to use grid graphics functions to manipulate plots created by base graphics, Paul Murrell has developed a package called `gridBase` to help interface the two R graphing systems. However, even with this useful tool, base R graphics cannot be manipulated with ease by grid graphics in the same way that objects built using the grid graphics system (include `ggplot` objects) can.



The `grid` package is now a base package, which means it is intalled automatically when you install R. This means you won't have to install it using `install.packages()` before you use it. You will, however, have to load the package with `library()` when you want to use it in an R session.

Grobs

The most critical concept of grid graphics to understand for extending `ggplot2` it the concept of grobs. Grobs are graphical objects that you can make and change with grid graphics

functions. For example, you may create a circle grob or points grobs. Once you have created one or more of these grobs, you can add them to or take them away from larger grid graphics objects, including ggplot objects. These grobs are the actual objects that get printed to a graphics device when you print a grid graphics plot; if you tried to create a grid graphics plot without any grobs, you would get a blank plot.

The grid package has a Grob family of functions that either make or change grobs. If you want to build a custom geom for ggplot that is unusual enough that you cannot rely on inheriting from an existing geom, you will need to use functions from the Grob family of functions to code your geom. Here is a list of functions from the Grob family that can be used to create new grobs:

- `arcCurvature`
- `bezierGrob`: Bezier curve
- `circleGrob`
- `frameGrob`
- `functionGrob`
- `legendGrob`
- `linesGrob`
- `nullGrob`
- `polygonGrob`
- `polylineGrob`
- `rasterGrob`
- `rectGrob`
- `roundrectGrob`
- `segmentsGrob`
- `textGrob`
- `xaxisGrob`
- `xsplineGrob`
- `yaxisGrob`

Here is a list of functions that can be used to change existing grobs:

- `addGrob`
- `clipGrob`
- `delayGrob`
- `editGrob`
- `forceGrob`
- `getGrob`
- `grobName`
- `packGrob`
- `pathGrob`
- `placeGrob`

- recordGrob
- removeGrob
- reorderGrob
- setGrob
- showGrob

Once you have created a grob, you can print it using the `grid.draw` function. Functions that create grobs typically include parameters to specify the location where the grobs should be placed. For example, the `pointsGrob` function includes `x` and `y` parameters, while the `segmentsGrob` includes parameters for the starting and ending location of each segment (`x0, x1, y0, y1`).

The grob family of functions also includes a parameter called `gp` for setting graphical parameters like color, fill, line type, line width, point size, etc., for grob objects. The input to this function must be a `gpar` object, which can be created using the `gpar` function. For example, to create and draw a gray circle grob, you could run:

```
my_circle <- circleGrob(x = 0.5, y = 0.5, r = 0.5,
                        gp = gpar(col = "gray", lty = 3))
grid.draw(my_circle)
```

Multiple grobs can be combined into a single grob using the `gTree` function. A `gTree` grob contains one or more “children” grobs. It can be very useful for creating grobs that need to contain multiple elements, like an axis, which needs to include the axis line, axis ticks, axis labels, etc.

Viewports

Much of the power of grid graphics comes from the ability to move in and out of working spaces around the full graph area. As an example, say you would like to create a map of the states of the US with a small pie chart added at the centroid of each state showing the distribution of population in that state by education level. This kind of plot is where grid graphics shines (although it appears that you now can create such a plot directly in ggplot2). In this case, you want to zoom in at the coordinates of a state centroid, have your own smaller working space at that location, add a pie chart showing data specific to that state, then zoom out and do the process again for a different state centroid.

In grid graphics, these smaller working spaces within the larger plot are called *viewports*. You can navigate to one of the viewports, make some changes, and then pop back up and navigate deeply into another viewport in the plot.

Using grid graphics, you can:

- Make a new viewport
- Remove an old viewport
- Navigate down to a specific viewport

- Navigate up to the full plotting area

You can only operate in one viewport at a time. Once you are in that viewport, you can write grobs within the viewport. If you want to place the next grob in a different viewport, you will need to navigate the that viewport before you can do so.

A grid graphics object can end up with a complex tree of viewports and grobs. Any of these elements can be customized, as long as you know can navigate back down to the specific element you want to change.

Grid graphics coordinate systems

Once you have created a grob and moved into the viewport in which you want to plot it, you need a way to specify where in the viewport to write the grob. The numbers you use to specify x- and y-placements for a grob will depend on the coordinate system you use. In grid graphics, you have a variety of options for the units to use in this coordinate system, and picking the right units for this coordinate system will make it much easier to create the plot you want.

There are several units that can be used for coordinate systems, and you typically will use different units to place objects. For example, you may want to add points to a plot based on the current x- and y-scales in that plot region, in which case you can use *native* units. The *native* unit is often the most useful when creating extensions for ggplot2, for example. The *npc* units are also often useful in designing new plots—these set the x- and y-ranges to go from 0 to 1, so you can use these units if you need to place an object in, for example, the exact center of a viewport (`c(0.5, 0.5)` in *npc* units), or create a viewport in the top right quarter of the plot region. Grid graphics also allows the use of some units with absolute values, including inches (`inches`), centimeters (`cm`), and millimeters (`mm`).

You can specify the coordinate system you would like to use when placing an object by with the `unit` function (`unit([numeric vector], units = "native")`).

The `gridExtra` package

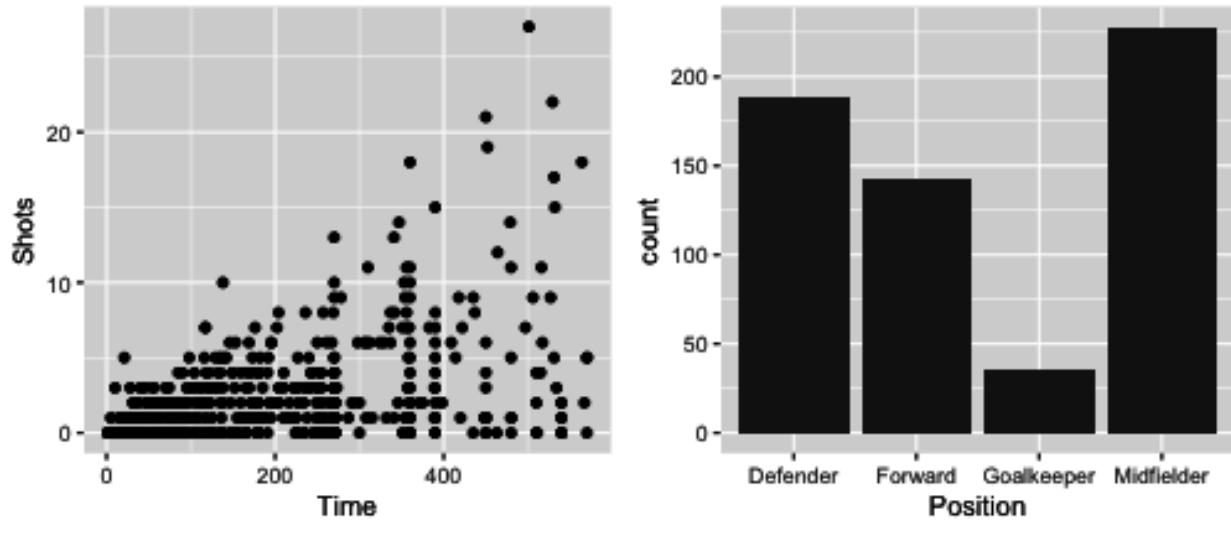
The `gridExtra` package provides useful extensions to the grid system, with an emphasis on higher-level functions to work with grid graphic objects, rather than the lower-level utilities in the `grid` package that are used to create and edit specific lower-level elements of a plot. This package has particularly useful functions that allow you to arrange and write multiple grobs to a graphics device and to include tables in grid graphics objects.

The `grid.arrange` function from the `gridExtra` package makes it easy to create a plot with multiple grid objects plotting to it. Because ggplot2 was built on grid graphics, you can also use this function to plot multiple ggplot objects to a graphics device. For example, say you wanted to create a plot that has two plots based on the World Cup data side-by-side. To create this plot, you can assign the ggplot objects for each separate graph to objects in your R global environment (`time_vs_shots` and `player_positions` in this example), and then input these objects to a `grid.arrange` call:

```
library(gridExtra)

time_vs_shots <- ggplot(worldcup, aes(x = Time, y = Shots)) +
  geom_point()
player_positions <- ggplot(worldcup, aes(x = Position)) +
  geom_bar()

grid.arrange(time_vs_shots, player_positions, ncol = 2)
```



plot of chunk unnamed-chunk-141

The `gridExtra` also has a function, `tableGrob`, that facilitates in adding tables to grid graphic objects.

Find out more about grid graphics

Grid graphics provides an extensive graphics system that can allow you to create almost any plot you can imagine in R. It takes quite a bit of work to fully understand all elements of the grid graphics system, but you might find it worthwhile to study grid graphics in greater depth if you often need to create very tailored, unusual graphs.

There are a number of resources you can use to learn more about grid graphics. The most comprehensive is the **R Graphics** book by Paul Murrell, the creator of grid graphics. This book is now in its second edition, and its first edition was written before ggplot2 became so popular. It is worth try to get the second edition, which includes some content specifically on ggplot2 and how that package relates to grid graphics. The vignettes that go along with the grid package are also by Paul Murrell and give a useful introduction to grid graphics, and the vignettes for the gridExtra package are also a useful next step for finding out more.

- Links to pdfs of vignettes for the grid graphics package are available at <https://stat.ethz.ch/R-manual/R-devel/library/grid/doc/index.html>

- Links to pdfs of vignettes for the gridGraphics package are available on the package’s CRAN page: <https://cran.r-project.org/web/packages/gridExtra/index.html>

4.6 Building a New Theme

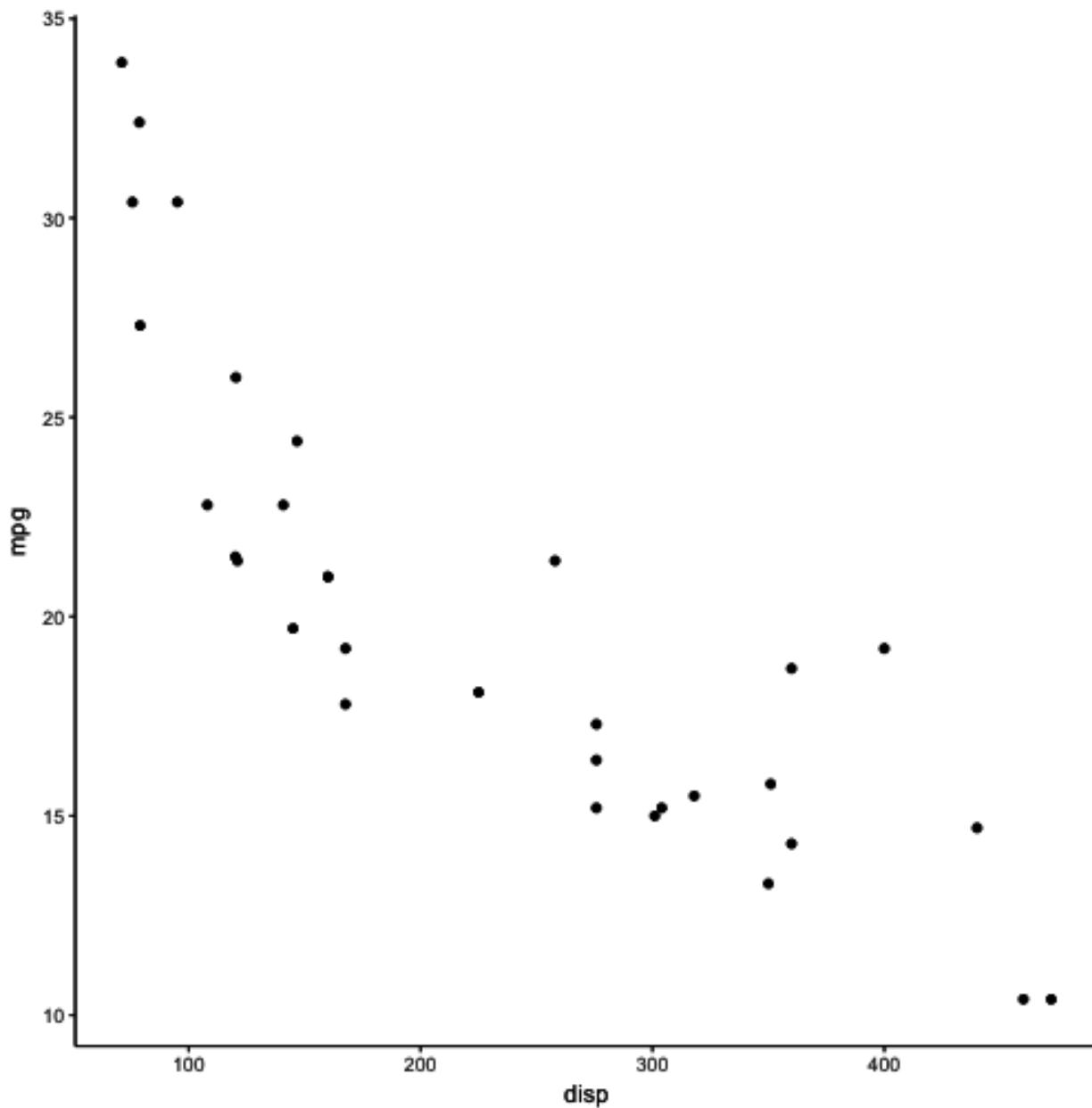
Building and modifying a theme in `ggplot2` is a key feature of the `ggplot2` package and system for building data graphics. The original base graphics system in R did not have a notion of a “theme” for how graphical elements are presented—users were left to individually customize each graphic without any clear way to programmatically implement shared elements across plots.

The `ggplot2` package implements the notion of a theme for its plots by allowing you to modify many different elements of a plot and to store all those modifications as a special “theme” object. Those elements that can be modified are documented in the help page `?theme`, which documents the `theme()` function.

The default theme for `ggplot2` is encapsulated by the `theme_gray()` function. Like other elements in the `ggplot2` universe, themes can be “added” using the `+` operator to plot commands in order to change the look and feel of a plot. Adding a theme (either existing or custom built by you) will override elements of any default theme.

For example, here is a plot that uses the `theme_classic()` function:

```
library(ggplot2)
ggplot(data = mtcars, aes(x = disp, y = mpg)) +
  geom_point() +
  theme_classic()
```



plot of chunk unnamed-chunk-142

Notice how the look and the feel of the plot is substantially different from the default gray theme of `ggplot2`. The key differences in the `theme_classic()` setup are the background color (white instead of gray), the colors of the grid lines (none instead of white), and the presence of solid black x- and y-axes. Other elements are the same, like the plotting character (solid circle) and fonts.



Note that themes in `ggplot2` only allow you to modify the **non-data** elements of a plot. Things like the title, axis labels, background, etc. can be modified with a theme. If you want to change data elements, like the plotting symbol or colors, you can modify those things separately in their respective `geom_*` functions.

Why Build a New Theme?

Why would one want to build a new theme? For many people, it is a matter of personal preference with respect to colors, shapes, fonts, positioning of labels, etc. Because plots, much like writing, are an expression of your ideas, it is often desirable to customize those plots so that they accurately represent your vision.

In corporate or institutional settings, developing themes can be a powerful branding tool. Plots that are distributed on the web or through marketing materials that have a common theme can be useful for reinforcing a brand. For example, plots made by the [FiveThirtyEight.com](#) web site have a distinct look and feel (see [this article](#) by Walt Hickey for one of many examples). When you see one of those plots you instinctively know that it is a “FiveThirtyEight” plot. Developing a theme for your organization can help to get others to better understand what your organization is about when it produces data graphics.

Another advantage of having a pre-programmed theme is that it removes the need for you to think about it later! One key reason why news organizations like FiveThirtyEight or the New York Times have common themes for their data graphics is because they are constantly producing those graphics on a daily basis. If every plot required a custom look and feel with a separate palette of colors, the entire process would grind to a halt. If you are in an environment where there is a need for reproducible graphics with a consistent feel, then developing a custom theme is probably a good idea. While using the default `ggplot2` theme is perfectly fine from a data presentation standpoint, why not try to stand out from the crowd?

Default Theme

As noted above, `ggplot2` has a default theme, which is `theme_gray()`. This theme produces the familiar gray-background-white-grid-lines plot. You can obtain the default theme using the `theme_get()` function.

```
x <- theme_get()
class(x)
[1] "theme" "gg"
```

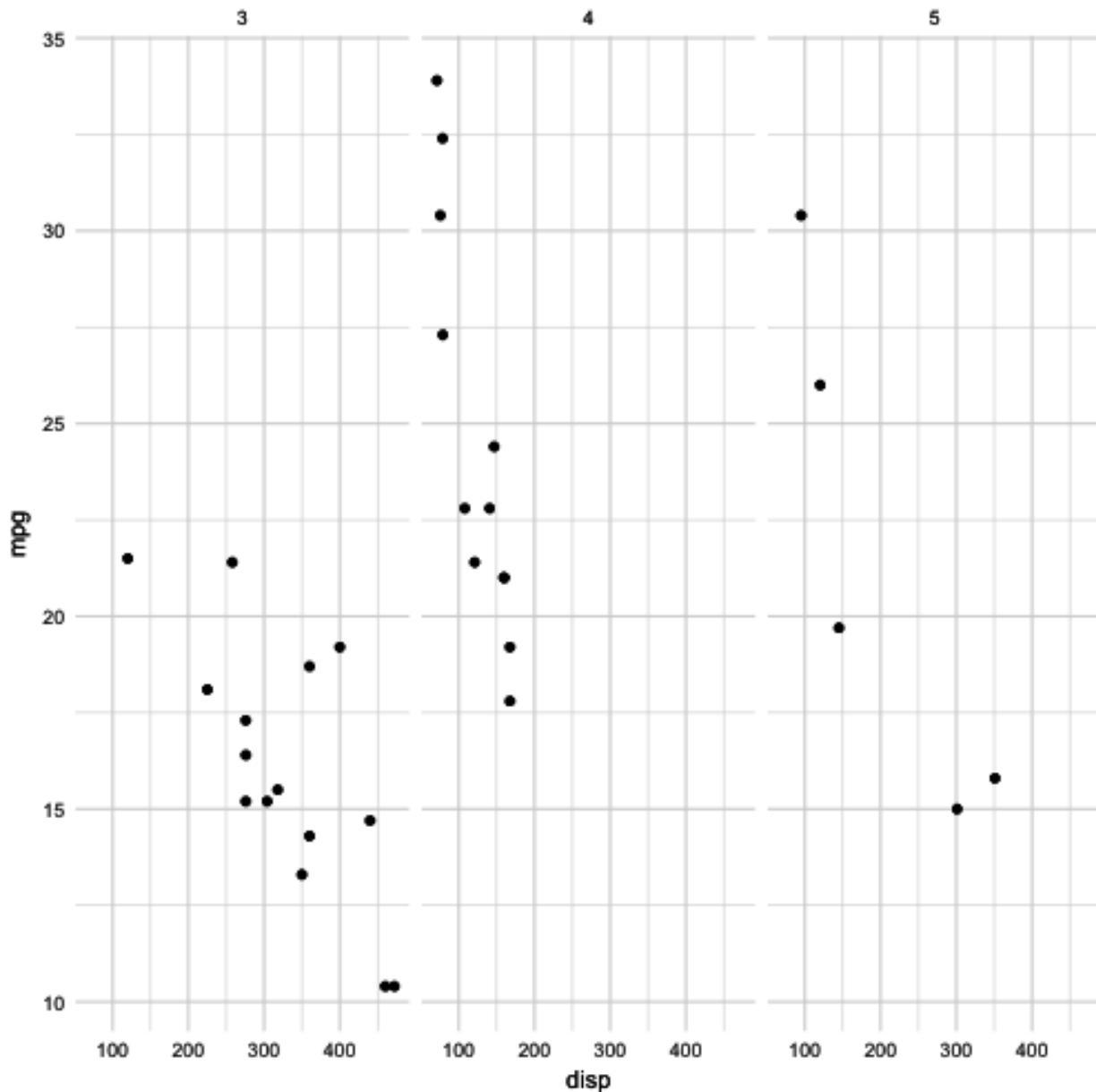
The object returned by `theme_get()` is rather large so it's not recommended to print it to the console. Notice that the object returned by `theme_get()` is an S3 object of class "theme" and "gg". This is the kind of object you will need to create or modify in order to customize your theme.

You can modify the default theme by using the `theme_set()` function and passing it a `theme` object. For example, if I want all my plots to use the `theme_minimal()` theme, I could do

```
new_theme <- theme_minimal()
theme_set(new_theme)
```

Now your plots will use the `theme_minimal()` theme without you having to specify it.

```
ggplot(data = mtcars, aes(disp, mpg)) +  
  geom_point() +  
  facet_grid( . ~ gear)
```



plot of chunk unnamed-chunk-145

Quitting R will erase the default theme setting. If you load `ggplot2` in a future session it will revert to the default gray theme. If you'd like for `ggplot2` to always use a different theme (either yours or one of the built-in ones), you can set a load hook and put it in your `.Rprofile` file. For example, the following hook sets the default theme to be `theme_minimal()` every time the `ggplot2` package is loaded.

```
setHook(packageEvent("ggplot2", "onLoad"),
       function(...) ggplot2::theme_set(ggplot2::theme_minimal()))
```

Of course, you can always override this default theme by adding a `theme` object to any of your plots that you construct in `ggplot2`.

Creating a New Theme

Perhaps the easiest thing to start with when customizing your own theme is to modify an existing theme (i.e. one that comes built-in to `ggplot2`). In case you are interested in thoroughly exploring this area and learning from others, there is also the [ggthemes package](#) on CRAN which provides a number of additional themes for `ggplot2`.

Looking at the help page `?theme` you'll see that there are many things to modify. We will start simple here by illustrating the general approach to making theme modifications. We will begin with the `theme_bw()` theme. This theme is a simple black and white theme that has little ornamentation and few features.

Modifying theme attributes

Suppose we want to make the default color for plot titles to be dark red. We can change just that element by adding a `theme()` modification to the existing theme.

```
newtheme <- theme_bw() + theme(plot.title = element_text(color = "darkred"))
```

Note that in our call to `theme()`, when we modify the `plot.title` attribute, we cannot simply say `color = "darkred"`. This must be wrapped in a call to the `element_text()` function so that the elements of `plot.title` are appropriately modified. In the help page for `theme()`, you will see that each attribute of a theme is modified by using one of four `element_*` functions:

- `element_text()`: specify the display of text elements
- `element_line()`: specify the display of lines (i.e. axis lines)
- `element_rect()`: specify the display of borders and backgrounds
- `element_blank()`: draw nothing

All of these functions work in the same way (although they contain different elements) and each of them returns a list of values inheriting from the class “element”. The `ggplot2` functions know how to handle objects of this class and will modify the theme of a plot accordingly.

Let's change a few more things about our new theme. We can make the box surrounding the plot to look a little different by modifying the `panel.border` element of the theme. First let's take a look at what the value is by default.

```
newtheme$panel.border
List of 5
$ fill      : logi NA
$ colour    : chr "grey20"
$ size      : NULL
$ linetype   : NULL
$ inherit.blank: logi TRUE
- attr(*, "class")= chr [1:2] "element_rect" "element"
```

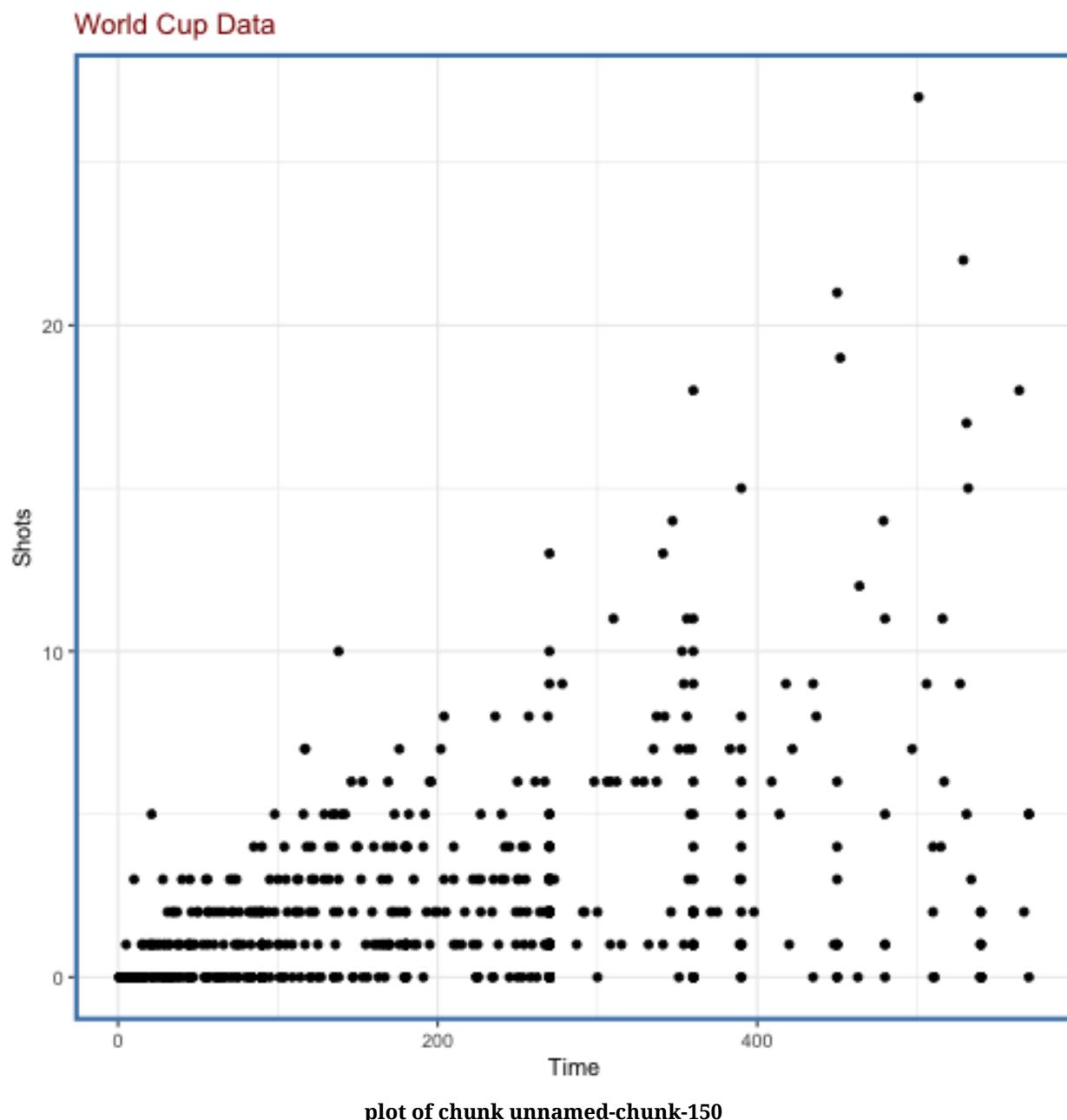
You can see that this is an object of class `element_rect` and there are 5 elements in this list, including the `fill`, `colour` (or `color`), `size`, and `linetype`. These attributes have the same meaning as they do in the usual `ggplot2` context.

We can modify the `color` attribute to make it “steelblue” and modify the `size` attribute to make it a little bigger.

```
newtheme <- newtheme +
  theme(panel.border = element_rect(color = "steelblue", size = 2))
```

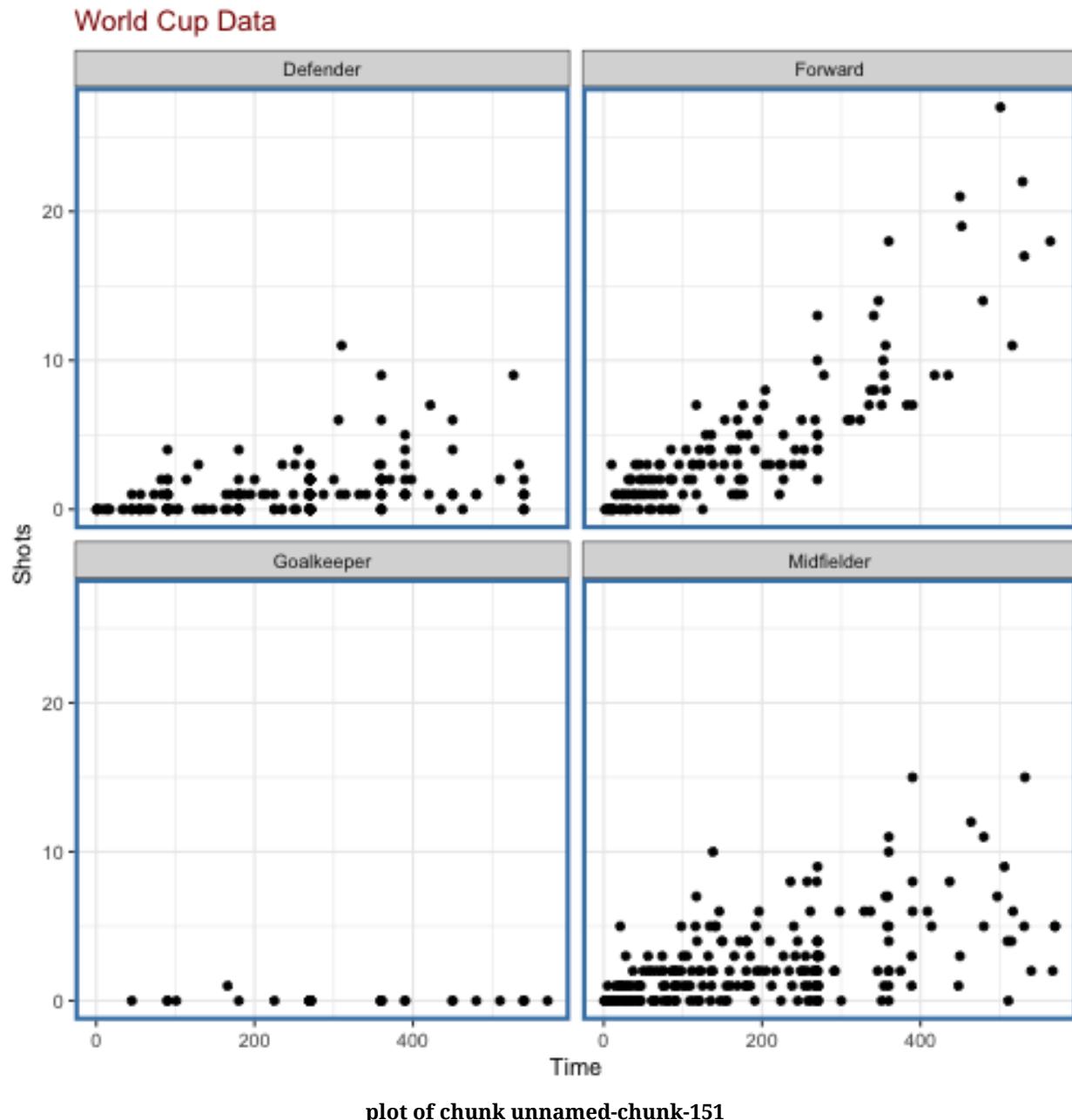
Now let’s see what a typical plot might look like. The following is a plot of minutes played an shots attempted from the `worldcup` dataset in the `faraway` package.

```
library(faraway)
ggplot(data = worldcup, aes(Time, Shots)) +
  geom_point() +
  ggtitle("World Cup Data") +
  newtheme
```



This may not be your idea of a great-looking theme, but it is certainly different! If we were to facet the data by `Position` the theme attributes would extend to the individual panels, which is a nice thing to get for free.

```
ggplot(data = worldcup, aes(Time, Shots)) +  
  geom_point() +  
  facet_wrap(facets = ~ Position, ncol = 2) +  
  ggtitle("World Cup Data") +  
  newtheme
```



Complete themes

When using the `theme()` function to modify specific elements of an existing theme, the default value for the argument `complete` is `FALSE`. This simply indicates, that the `theme()` function is *not* returning a complete theme where every element is appropriately specified. Rather, it is just modifying the theme element that you specified.

Setting `complete = TRUE` in the call to `theme()` tells `ggplot2` that the `theme()` function is returning a complete theme along the lines of `theme_gray()` or `theme_bw()`. In particular, all of the theme elements will inherit from the `blank` element, meaning that there will be no values to “fall back” on in the event that you do not specify them. Setting `complete = TRUE` only if you plan to specify every single theme element that is meaningful to you. If you are creating a brand new theme specific to you, then you may well be doing that. But if you are simply tweaking an existing theme, it’s appropriate to set `complete = FALSE`.

Summary

Building a new theme allows you to customize the look and feel of a plot to match your personal preferences. It also allows you to define a consistent “branded” presentation of your data graphics that can be clearly identified with your organization or company.

4.7 Building New Graphical Elements

Some of the key elements of a data graphic made with `ggplot2` are geoms and stats. The fact is, the `ggplot2` package comes with tremendous capabilities that allow users to make a wide range of interesting and rich data graphics. These graphics can be made through a combination of calls to various `geom_*` and `stat_*` functions (as well as other classes of functions).

So why would one want to build a new geom or stat on top of all that `ggplot2` already provides?

There are two key reasons for building new geoms and stats for `ggplot2`:

1. **Implement a new feature.** There may be something very specific to your application that is not yet implemented—a new statistical modeling approach or a novel plotting symbol. In this case you don’t have much choice and need to extend the functionality of `ggplot2`.
2. **Simplify a complex workflow.** With certain types of analyses you may find yourself producing the same kind of plot elements repeatedly. These elements may involve a combination of points, lines, facets, or text and essentially encapsulate a single idea. In that case it may make sense to develop a new geom to literally encapsulate the collection of plot elements and to make it simpler to include these things in your future plots.

Building new stats and geoms is the plotting equivalent of writing functions (that may sound a little weird because stats and geoms *are* functions, but they are thought of a little differently from generic functions). While the action taken by a function can typically be executed using separate expressions outside of a function context, it is often convenient for the user to encapsulate those actions into a clean function. In addition, writing a function allows you to easily parameterize certain elements of that code. Creating new geoms and stats similarly allows for a simplification of code and for allowing users to easily tweak certain elements of a plot without having to wade through an entire mess of code every time.

Building a Geom

New geoms in `ggplot2` inherit from a top level class called `Geom` and are constructed using a two step process.

1. The `ggproto()` function is used to construct a new class corresponding to your new geom. This new class specifies a number of attributes and functions that describe how data should be drawn on a plot.
2. The `geom_*` function is constructed as a regular function. This function returns a layer to that can be added to a plot created with the `ggplot()` function.

The basic setup for a new geom class will look something like the following.

```
GeomNEW <- ggproto("GeomNEW", Geom,
  required_aes = <a character vector of required aesthetics>,
  default_aes = aes(<default values for certain aesthetics>),
  draw_key = <a function used to draw the key in the legend>,
  draw_panel = function(data, panel_scales, coord) {
    ## Function that returns a grid grob that will
    ## be plotted (this is where the real work occurs)
  }
)
```

The `ggproto` function is used to create the new class. Here, “NEW” will be replaced by whatever name you come up with that best describes what your new geom is adding to a plot. The four things listed inside the class are required of all geoms and must be specified.

The required aesthetics should be straightforward—if your new geom makes a special kind of scatterplot, for example, you will likely need `x` and `y` aesthetics. Default values for aesthetics can include things like the plot symbol (i.e. `shape`) or the color.

Implementing the `draw_panel` function is the hard part of creating a new geom. Here you must have some knowledge of the `grid` package in order to access the underlying elements of a `ggplot2` plot, which is based on the `grid` system. However, you can implement a reasonable amount of things with knowledge of just a few elements of `grid`.

The `draw_panel` function has three arguments to it. The `data` element is a data frame containing one column for each aesthetic specified, `panel_scales` is a list containing information about

the `x` and `y` scales for the current panel, and `coord` is an object that describes the coordinate system of your plot.

The `coord` and the `panel_scales` objects are not of much use except that they transform the data so that you can plot them.

```
library(grid)
GeomMyPoint <- ggproto("GeomMyPoint", Geom,
  required_aes = c("x", "y"),
  default_aes = aes(shape = 1),
  draw_key = draw_key_point,
  draw_panel = function(data, panel_scales, coord) {
    ## Transform the data first
    coords <- coord$transform(data, panel_scales)

    ## Let's print out the structure of the 'coords' object
    str(coords)

    ## Construct a grid grob
    pointsGrob(
      x = coords$x,
      y = coords$y,
      pch = coords$shape
    )
  })
}
```



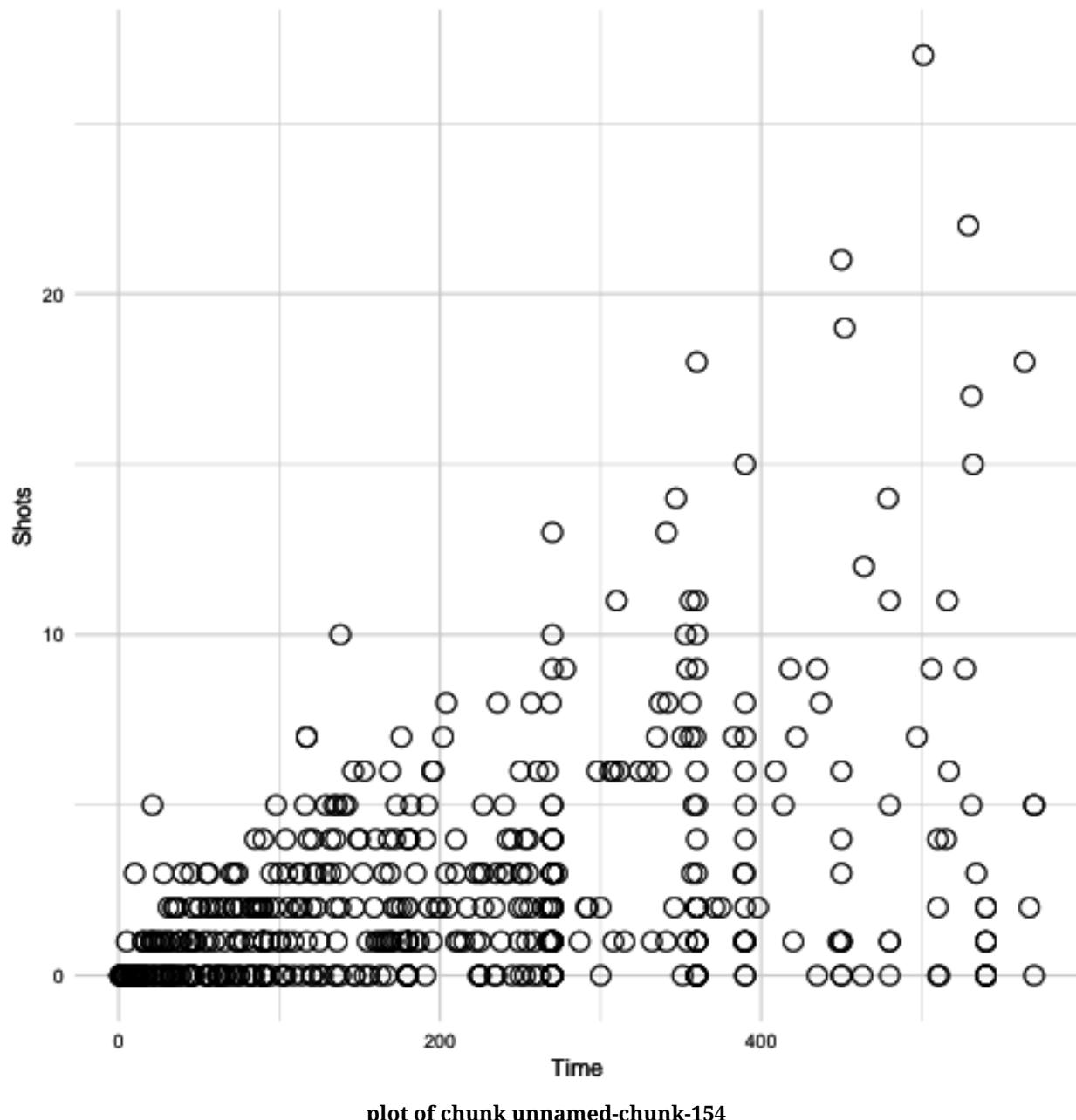
In this example we print out the structure of the `coords` object with the `str()` function just so you can see what is in it. Normally, when building a new geom you wouldn't do this.

In addition to creating a new Geom class, you need to create the `actually` function that will build a layer based on your geom specification. Here, we call that new function `geom_mypoint()`, which is modeled after the built in `geom_point()` function.

```
geom_mypoint <- function(mapping = NULL, data = NULL, stat = "identity",
  position = "identity", na.rm = FALSE,
  show.legend = NA, inherit.aes = TRUE, ...) {
  ggplot2::layer(
    geom = GeomMyPoint, mapping = mapping,
    data = data, stat = stat, position = position,
    show.legend = show.legend, inherit.aes = inherit.aes,
    params = list(na.rm = na.rm, ...))
}
```

Now we can use our new geom on the `worldcup` dataset.

```
ggplot(data = worldcup, aes(Time, Shots)) + geom_mypoint()
```



```
'data.frame':      595 obs. of  5 variables:
 $ x     : num  0.0694 0.6046 0.3314 0.4752 0.1174 ...
 $ y     : num  0.0455 0.0455 0.0455 0.0791 0.1128 ...
 $ PANEL: int  1 1 1 1 1 1 1 1 1 ...
 $ group: int  -1 -1 -1 -1 -1 -1 -1 -1 -1 ...
 $ shape: num  1 1 1 1 1 1 1 1 1 ...
 - attr(*, "vars")= chr "PANEL"
```

From the `str()` output we can see that the `coords` object contains the `x` and `y` aesthetics, as well as the `shape` aesthetic that we specified as the default. Note that both `x` and `y` have been rescaled to be between 0 and 1. This is the normalized parent coordinate system.

Example: An Automatic Transparency Geom

One problem when making scatterplots of large amounts of data is *overplotting*. In particular, with `ggplot2`'s default solid circle as the plotting shape, if there are many overlapping points all you will see is a solid mass of black.

One solution to this problem of overplotting is to make the individual points *transparent* by setting the alpha channel. The alpha channel is a number between 0 and 1 where 0 is totally transparent and 1 is completely opaque. With transparency, if two points overlap each other, they will be darker than a single point sitting by itself. Therefore, you can see more of the “density” of the data when the points are transparent.

The one requirement for using transparency in scatterplots is computing the amount of transparency, or the the alpha channel. Often this will depend on the number of points in the plot. For a simple plot with a few points, no transparency is needed. For a plot with hundreds or thousands of points, transparency is required. Computing the exact amount of transparency may require some experimentation.

The following example creates a geom that computes the alpha channel based on the number of points that are being plotted. First we create the Geom class, which we call `GeomAutoTransparent`. This class sets the `alpha` aesthetic to be 0.3 if the number of data points is between 100 and 200 and 0.15 if the number of data points is over 200. If the number of data points is 100 or less, no transparency is used.

```
GeomAutoTransparent <- ggproto("GeomAutoTransparent", Geom,
  required_aes = c("x", "y"),
  default_aes = aes(shape = 19),
  draw_key = draw_key_point,
  draw_panel = function(data, panel_scales, coord) {
    ## Transform the data first
    coords <- coord$transform(data, panel_scales)

    ## Compute the alpha transparency factor based on the
    ## number of data points being plotted
    n <- nrow(data)
    if(n > 100 && n <= 200)
      coords$alpha <- 0.3
    else if(n > 200)
      coords$alpha <- 0.15
  }
}
```

```

    else if(n > 200)
      coords$alpha <- 0.15
    else
      coords$alpha <- 1
    ## Construct a grid grob
    grid::pointsGrob(
      x = coords$x,
      y = coords$y,
      pch = coords$shape,
      gp = grid::gpar(alpha = coords$alpha)
    )
  })
}

```

Now we need to create the corresponding geom function, which we slightly modify from `geom_point()`. Note that the `geom` argument to the `layer()` function takes our new `GeomAutoTransparent` class as its argument.

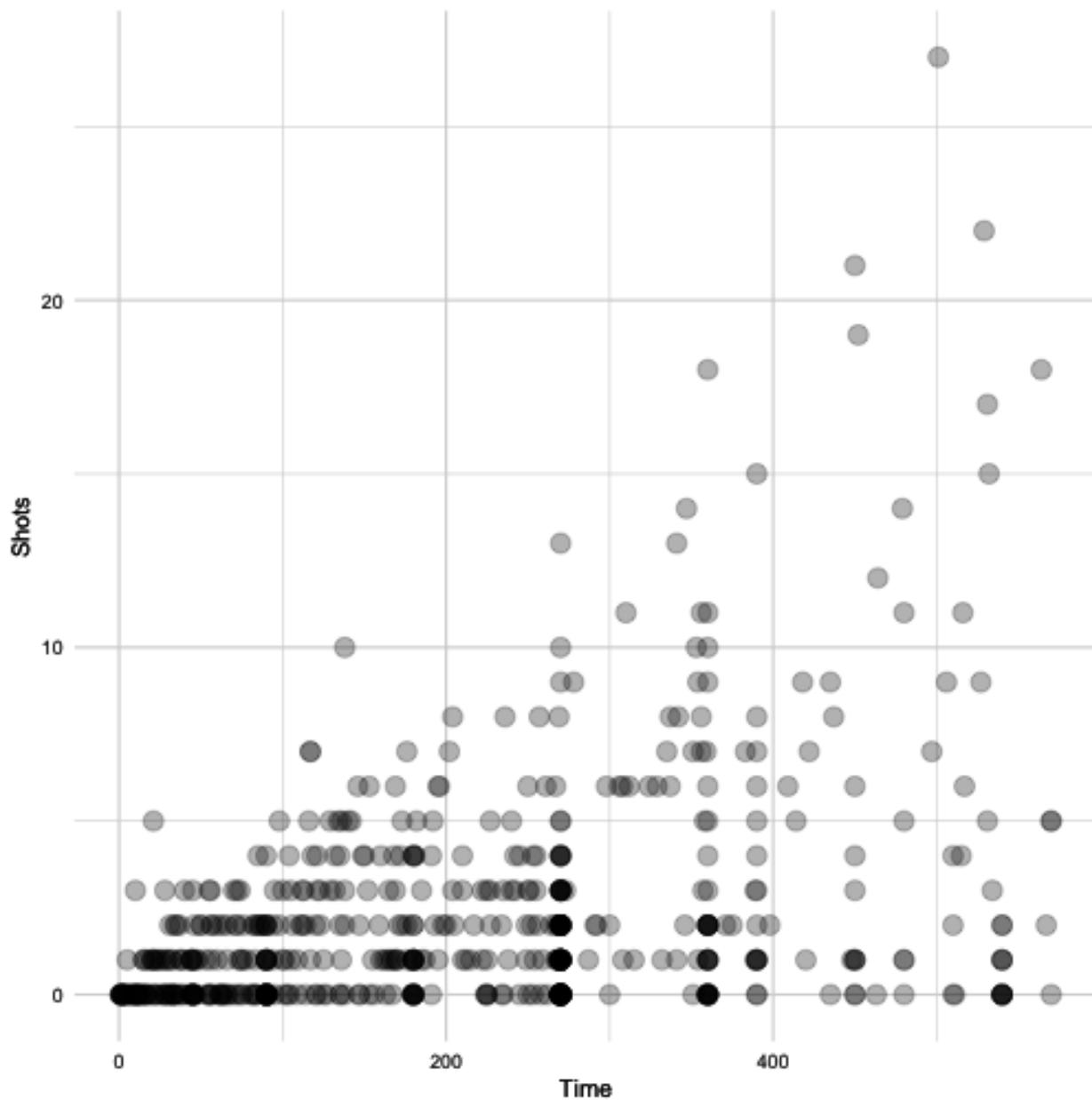
```

geom_transparent <- function(mapping = NULL, data = NULL, stat = "identity",
                           position = "identity", na.rm = FALSE,
                           show.legend = NA, inherit.aes = TRUE, ...) {
  ggplot2::layer(
    geom = GeomAutoTransparent, mapping = mapping,
    data = data, stat = stat, position = position,
    show.legend = show.legend, inherit.aes = inherit.aes,
    params = list(na.rm = na.rm, ...))
}

```

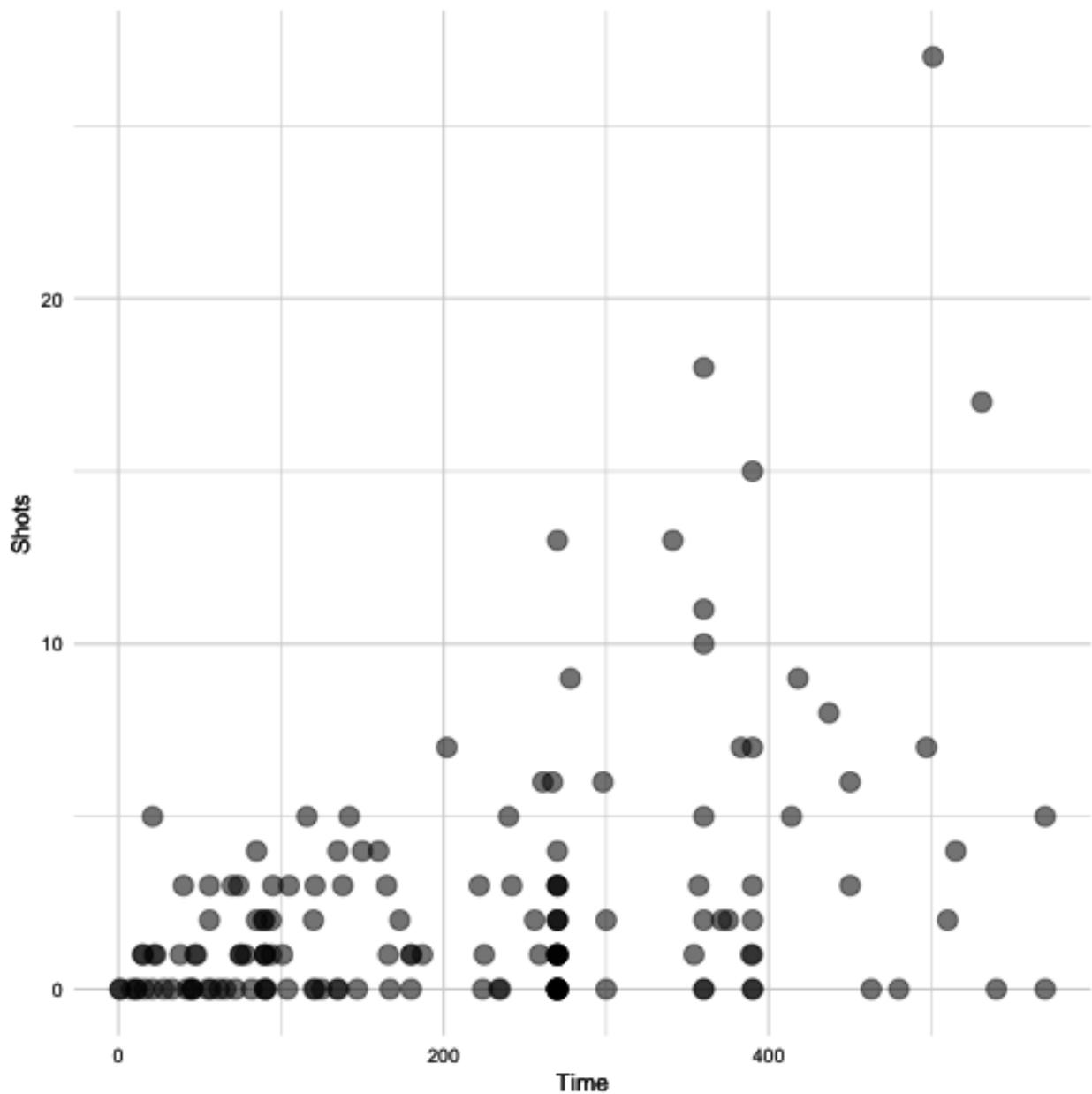
Now we can try out our new `geom_transparent()` function with differing amounts of data to see how the transparency works. Here is the entire `worldcup` dataset, which has 595 observations.

```
ggplot(data = worldcup, aes(Time, Shots)) + geom_transparent()
```



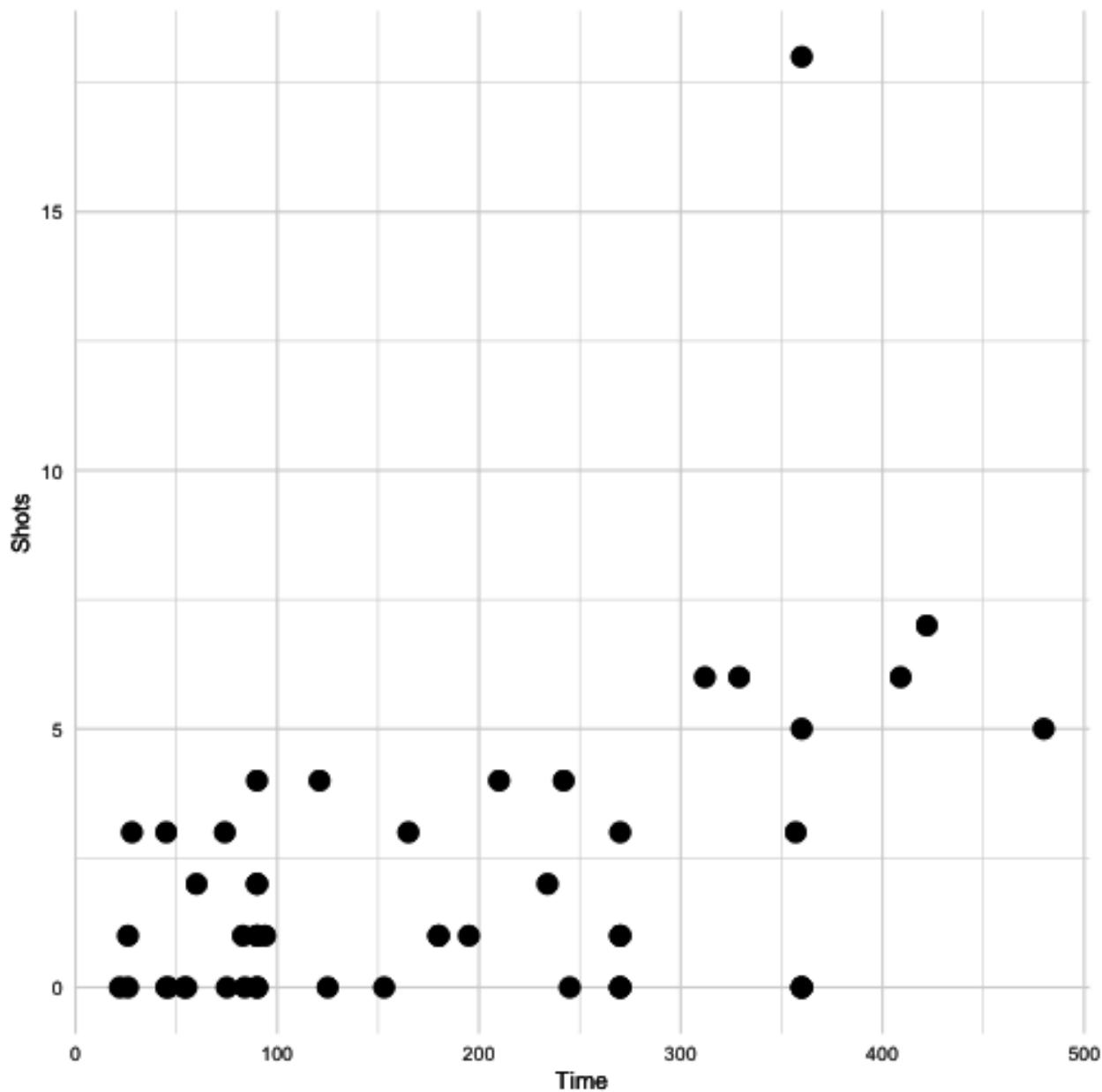
Here we take a random sample of 150 observations. The transparency should be a little less in this plot.

```
library(dplyr)
ggplot(data = sample_n(worldcup, 150), aes(Time, Shots)) +
  geom_transparent()
```



Here we take a random sample of 50 observations. There should be no transparency used in this plot.

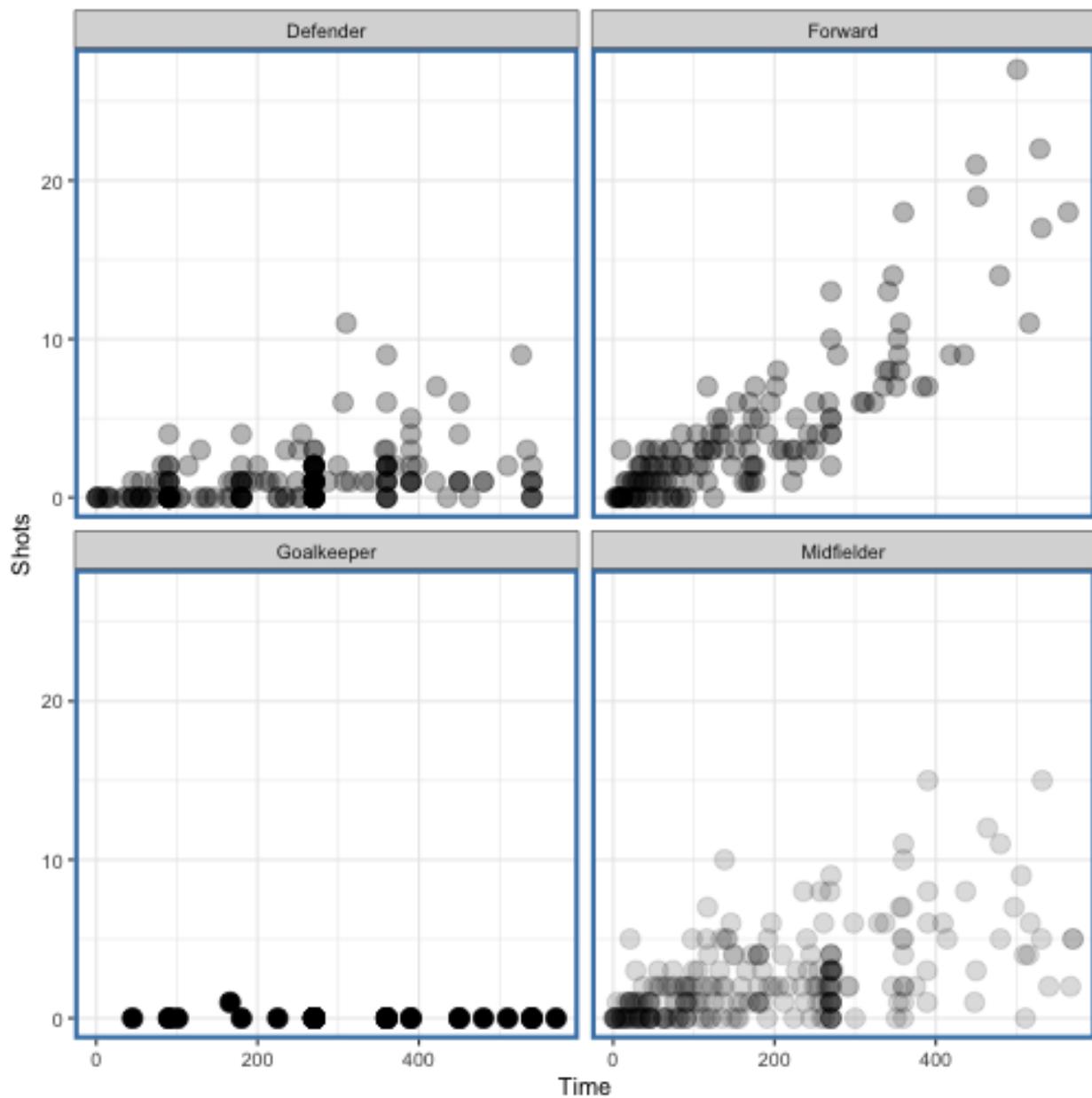
```
ggplot(data = sample_n(worldcup, 50), aes(Time, Shots)) +  
  geom_transparent()
```



plot of chunk unnamed-chunk-160

We can also reproduce a faceted plot from the previous section with our new geom and the features of the geom will propagate to the panels.

```
ggplot(data = worldcup, aes(Time, Shots)) +  
  geom_transparent() +  
  facet_wrap(~ Position, ncol = 2) +  
  newtheme
```



plot of chunk unnamed-chunk-161

Notice that the data for the “Midfielder”, “Defender”, and “Forward” panels have some transparency because there are more points there but the “Goalkeeper” panel has no transparency because it has relatively few points.

It’s worth noting that in this example, a different approach might have been to *not* create a new geom, but rather to compute an “alpha” column in the dataset that was a function of the number of data points (or the number of data points in each subgroup). Then you could have just set the `alpha` aesthetic to be equal to that column and `ggplot2` would have naturally mapped the appropriate alpha value to the the right subgroup. However, there a few issues

with that approach:

1. It involves adding a column to the data that isn't fundamentally related to the data (it is related to *presenting* the data); and
2. Some version of that alpha computation would need to be done every time you plotted the data in a different way. For example if you faceted on a different grouping variable, you'd need to compute the alpha value based on the number of points in the new subgroups.

The advantage of creating a geom in this case is that it abstracts the computation, removes the need to modify the data each time, and allows for a simpler communication of what is trying to be done in this plotting code.

Summary

Building new geoms can be a useful way to implement a completely new graphical procedure or to simplify a complex graphical task that must be used repeatedly in many plots. Building a new geom requires defining a new Geom class via `ggproto()` and defining a new `geom_*` function that builds a layer based on the new Geom class.

Some further resources that are worth investigating if you are interested in building new graphical elements are

- [R Graphics](#) by Paul Murrell, describes the grid graphical system on which `ggplot2` is based.
- [Extending ggplot2 vignette](#), provides further details about how to build new geoms and stats.
- [ggplot2 Extensions](#) web site, provides numerous examples of `ggplot2` extensions that members of the community have developed.

About the Authors

Roger D. Peng is a Professor of Biostatistics at the Johns Hopkins Bloomberg School of Public Health. He is also a Co-Founder of the [Johns Hopkins Data Science Specialization](#), which has enrolled over 1.5 million students, the [Johns Hopkins Executive Data Science Specialization](#), the [Simply Statistics blog](#) where he writes about statistics and data science for the general public, and the [Not So Standard Deviations](#) podcast. Roger can be found on Twitter and GitHub under the user name [rdpeng](#).

Sean Kross is a software developer in the Department of Biostatistics at the Johns Hopkins Bloomberg School of Public Health. Sean's professional interests range between metagenomics, cybersecurity, and human-computer interaction. He is also the lead developer of [swirl](#), a software package designed to teach programming, statistics, and data science in an authentic programming environment. You can find Sean on Twitter and GitHub at [seankross](#).

Brooke Anderson is an Assistant Professor at Colorado State University in the Department of Environmental & Radiological Health Sciences, as well as a Faculty Associate in the Department of Statistics. She is also a member of the universityâ€™s Partnership of Air Quality, Climate, and Health and is a member of the editorial boards of *Epidemiology* and *Environmental Health Perspectives*. Previously, she completed a postdoctoral appointment in Biostatistics at Johns Hopkins Bloomberg School of Public and a PhD in Engineering at Yale University. Her research focuses on the health risks associated with climate-related exposures, including heat waves and air pollution, for which she has conducted several national-level studies. As part of her research, she has also published a number of open source R software packages to facilitate environmental epidemiologic research.