

# **Introduction to Statistics and Data Science**

Elizabeth Tipton      Arend M Kuyper      Danielle Sass  
Kaitlyn G. Fitzgerald      Adapted from ModernDive by  
Chester Ismay and Albert Y. Kim

2022-09-19

# Table of contents

<b>Welcome</b>	<b>5</b>
License . . . . .	5
<b>Preface</b>	<b>6</b>
Introduction for students . . . . .	6
What you will learn from this book . . . . .	7
Data/science pipeline . . . . .	9
Reproducible research . . . . .	10
<b>I Getting started</b>	<b>12</b>
<b>1 Getting Started with Data in R</b>	<b>13</b>
1.1 What are R and RStudio? . . . . .	13
1.1.1 Using RStudio Cloud . . . . .	14
1.1.2 Installing R and RStudio on your personal computer . . . . .	15
1.1.3 Using R via RStudio . . . . .	15
1.2 How do I code in R? . . . . .	16
1.2.1 Creating your first R Markdown document . . . . .	17
1.2.2 Basic programming concepts and terminology . . . . .	18
1.2.3 Errors, warnings, and messages . . . . .	19
1.2.4 Tips on learning to code . . . . .	20
1.3 What are R packages? . . . . .	21
1.3.1 Package installation . . . . .	22
1.3.2 Package loading . . . . .	23
1.3.3 Package use . . . . .	24
1.4 Explore your first dataset . . . . .	24
1.4.1 <code>nycflights13</code> package . . . . .	25
1.4.2 <code>flights</code> data frame . . . . .	26
1.4.3 Exploring data frames . . . . .	27
1.4.4 Help files . . . . .	29
1.5 Conclusion . . . . .	30
1.5.1 Additional resources . . . . .	30

<b>II Data Exploration via the tidyverse</b>	<b>31</b>
<b>2 Data Visualization</b>	<b>32</b>
Packages Needed . . . . .	32
2.1 The Grammar of Graphics . . . . .	32
2.1.1 Components of the Grammar . . . . .	33
2.1.2 Gapminder data . . . . .	33
2.1.3 Other components . . . . .	35
2.1.4 ggplot2 package . . . . .	35
2.2 Five Named Graphs - The 5NG . . . . .	36
2.3 5NG#1: Scatterplots . . . . .	36
2.3.1 Scatterplots via geom_point . . . . .	37
2.3.2 Over-plotting . . . . .	40
2.3.3 Summary . . . . .	45
2.4 5NG#2: Linegraphs . . . . .	45
2.4.1 Linegraphs via geom_line . . . . .	46
2.4.2 Summary . . . . .	48
2.5 5NG#3: Histograms . . . . .	48
2.5.1 Histograms via geom_histogram . . . . .	50
2.5.2 Adjusting the bins . . . . .	53
2.5.3 Summary . . . . .	56
2.6 Facets . . . . .	56
2.7 5NG#4: Boxplots . . . . .	59
2.7.1 Boxplots via geom_boxplot . . . . .	62
2.7.2 Summary . . . . .	65
2.8 5NG#5: Barplots . . . . .	66
2.8.1 Barplots via geom_bar or geom_col . . . . .	67
2.8.2 Must avoid pie charts! . . . . .	70
2.8.3 Two categorical variables . . . . .	73
2.8.4 Summary . . . . .	78
2.9 Conclusion . . . . .	78
2.9.1 Summary table . . . . .	78
2.9.2 Argument specification . . . . .	79
2.9.3 Additional resources . . . . .	79
2.9.4 What's to come . . . . .	80
<b>3 Data Wrangling</b>	<b>82</b>
<b>4 Data Importing &amp; “Tidy Data”</b>	<b>83</b>

<b>III Data Modeling</b>	<b>84</b>
<b>5 Basic Regression</b>	<b>85</b>
<b>6 Multiple Regression</b>	<b>86</b>
<b>IV Statistical Theory</b>	<b>87</b>
<b>7 Randomization and Causality</b>	<b>88</b>
<b>8 Populations and Generalizability</b>	<b>89</b>
<b>9 Sampling Distributions</b>	<b>90</b>
<b>V Statistical Inference</b>	<b>91</b>
<b>10 Confidence Intervals</b>	<b>92</b>
<b>11 P-values</b>	<b>93</b>
<b>12 Hypothesis tests</b>	<b>94</b>
<b>13 Putting it all together</b>	<b>95</b>
<b>References</b>	<b>96</b>
<b>Appendices</b>	<b>96</b>
<b>A Statistical Background</b>	<b>97</b>

# Welcome

This is the website for **Introduction to Statistics and Data Science**. This book starts you down the path of learning how to think with data using R. You'll learn the basics of how to engage, explore, and examine many types of data arising from several contexts. Hopefully you'll have fun and see how valuable it is to be able to critically think with data.

## ⚠️ Warning

Please note that this is a “development version” of this book for the new design of STAT 202. Meaning this is a work in progress being edited and updated as we go.

We would appreciate any feedback on typos and errors.

This open textbook is produced with support from [Northwestern University Libraries](#) and [The Alumnae of Northwestern University](#).



Figure 1



Figure 2

## License

This website is (and will always be) **free to use**, and is licensed under the [Creative Commons Zero v1.0 Universal License](#). If you'd like to give back, please consider reporting a typo or leaving a pull request at [github.com/NUstat/intro-stat-data-sci](https://github.com/NUstat/intro-stat-data-sci).

# Preface

Help! I'm new to R and RStudio and I need to learn about them! However, I'm completely new to coding! What do I do?



Figure 3



Figure 4

If you're asking yourself this question, then you've come to the right place! Start with our "Introduction for Students".

## Introduction for students

This book assumes no prerequisites: no algebra, no calculus, and no prior programming/coding experience. This is intended to be a gentle introduction to the practice of analyzing data and answering questions using data the way statisticians, data scientists, data journalists, and other researchers would.

In Figure 5 we present a flowchart of what you'll cover in this book. You'll first get started with data in Chapter 1, where you'll learn about the difference between R and RStudio, start coding in R, understand what R packages are, and explore your first dataset: all domestic departure flights from a New York City airport in 2013. Then

1. **Data Exploration:** You'll assemble your data science toolbox using `tidyverse` packages. In particular:
  - Ch. 2: Visualizing data via the `ggplot2` package.
  - Ch. 3: Wrangling data via the `dplyr` package.
  - Ch. 4: Understanding the concept of “tidy” data as a standardized data input format for all packages in the `tidyverse`
2. **Data Modeling:** Using these data science tools, you'll start performing data modeling. In particular:

- Ch. 5: Constructing basic regression models.
  - Ch. 6: Constructing multiple regression models.

3. **Statistical Theory:** Now you'll learn about the role of randomization in making inferences and the general frameworks used to make inferences in statistics. In particular:

  - Ch. 7: Randomization and causality.
  - Ch. 8: Populations and generalizability.
  - Ch. 9: Sampling distributions.

4. **Statistical Inference:** You'll learn to combine your newly acquired data analysis and modeling skills with statistical theory to make inferences. In particular:

  - Ch. 10: Building confidence intervals.
  - Ch. 11: Calculating p-values.
  - Ch. 12: Conducting hypothesis tests.

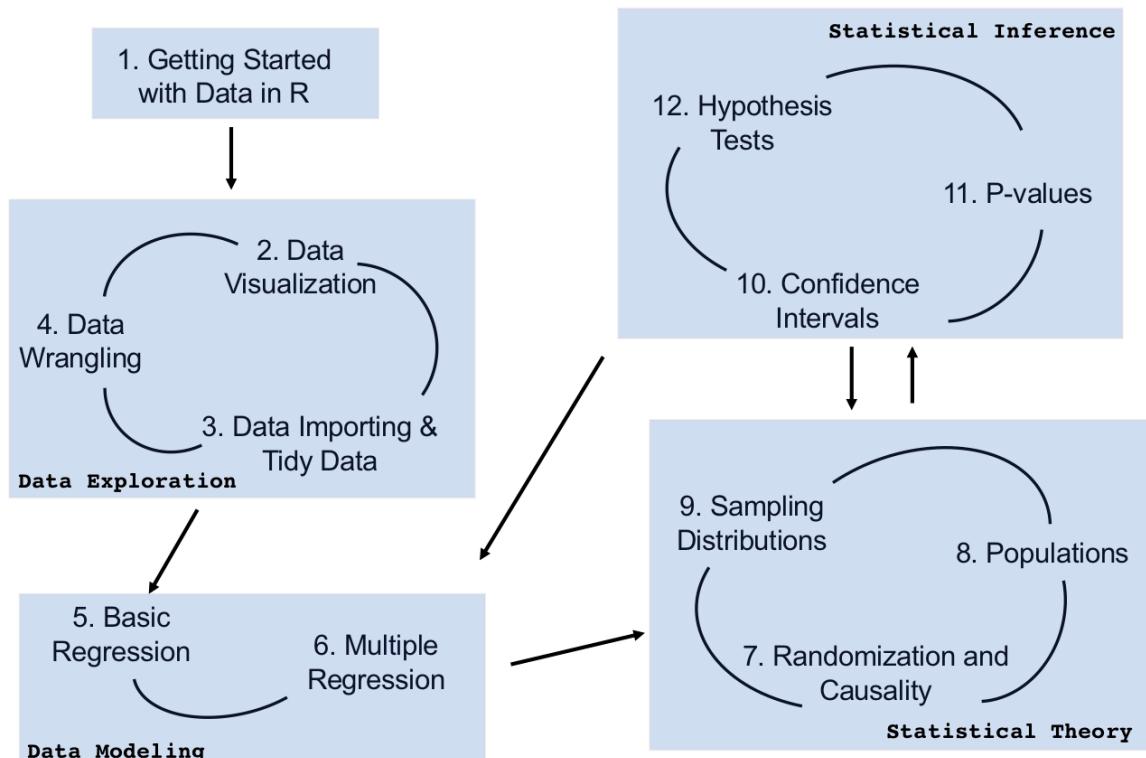


Figure 5: Course Flowchart

## **What you will learn from this book**

We hope that by the end of this book, you'll have learned

1. How to use R to explore data.
2. How to generate research questions and hypotheses.
3. How to think like a statistician and the role of chance in your data.
4. How to answer statistical questions using tools like confidence intervals and hypothesis tests.
5. How to effectively create “data stories” using these tools.

What do we mean by data stories? We mean any analysis involving data that engages the reader in answering questions with careful visuals and thoughtful discussion, such as [How strong is the relationship between per capita income and crime in Chicago neighborhoods?](#) and [How many f\\*\\*ks does Quentin Tarantino give \(as measured by the amount of swearing in his films\)?](#). Further discussions on data stories can be found in this [Think With Google article](#).

For other examples of data stories constructed by students like yourselves, look at the final projects for two courses that have previously used a version of this book:

- Middlebury College [MATH 116 Introduction to Statistical and Data Sciences](#) using student collected data.
- Pacific University [SOC 301 Social Statistics](#) using data from the [fivethirtyeight R package](#).

This book will help you develop your “data science toolbox”, including tools such as data visualization, data formatting, data wrangling, and data modeling using regression. With these tools, you’ll be able to perform the entirety of the “data/science pipeline” while building data communication skills.

In particular, this book will lean heavily on data visualization. In today’s world, we are bombarded with graphics that attempt to convey ideas. We will explore what makes a good graphic and what the standard ways are to convey relationships with data. You’ll also see the use of visualization to introduce concepts like mean, median, standard deviation, distributions, etc. In general, we’ll use visualization as a way of building almost all of the ideas in this book.

To impart the statistical lessons in this book, we have intentionally minimized the number of mathematical formulas used and instead have focused on developing a conceptual understanding via data visualization, statistical computing, and simulations. We hope this is a more intuitive experience than the way statistics has traditionally been taught in the past and how it is commonly perceived.

Finally, you’ll learn the importance of literate programming. By this we mean you’ll learn how to write code that is useful not just for a computer to execute but also for readers to understand exactly what your analysis is doing and how you did it. This is part of a greater effort to encourage reproducible research (see subsection *Reproducible research* for more details). Hal Abelson coined the phrase that we will follow throughout this book:

“Programs must be written for people to read, and only incidentally for machines to execute.”

We understand that there may be challenging moments as you learn to program. We still continue to struggle and find ourselves often using web searches to find answers and reach out to colleagues for help. In the long run though, we all can solve problems faster and more elegantly via programming. We wrote this book as our way to help you get started and you should know that there is a huge community of R users that are always happy to help everyone along as well. This community exists in particular on the internet on various forums and websites such as [stackoverflow.com](https://stackoverflow.com).

## Data/science pipeline

You may think of statistics as just being a bunch of numbers. We commonly hear the phrase “statistician” when listening to broadcasts of sporting events. Statistics (in particular, data analysis), in addition to describing numbers like with baseball batting averages, plays a vital role in all of the sciences. You’ll commonly hear the phrase “statistically significant” thrown around in the media. You’ll see articles that say “Science now shows that chocolate is good for you.” Underpinning these claims is data analysis and a theoretical model relating the data collected in a sample to a larger population. By the end of this book, you’ll be able to better understand whether these claims should be trusted or whether we should be wary. Inside data analysis are many sub-fields that we will discuss throughout this book (though not necessarily in this order):

- data collection
- data wrangling
- data visualization
- data modeling
- statistical inference
- correlation and regression
- interpretation of results
- data communication/storytelling

These sub-fields are summarized in what Grolemund and Wickham term the “[Data/Science Pipeline](#)” in Figure 6.

We will begin by digging into the gray **Understand** portion of the cycle with data visualization, then with a discussion on what is meant by tidy data and data wrangling, and then conclude by talking about interpreting and discussing the results of our models via **Communication**. These steps are vital to any statistical analysis. But why should you care about statistics? “Why did they make me take this class?”

There’s a reason so many fields require a statistics course. Scientific knowledge grows through an understanding of statistical significance and data analysis. You needn’t be intimidated by

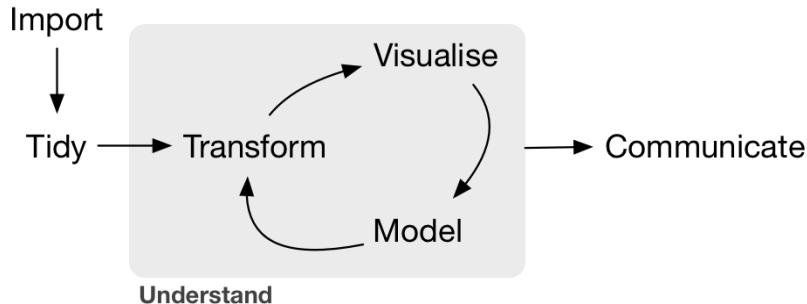


Figure 6: Data/Science Pipeline

statistics. It's not the beast that it used to be and, paired with computation, you'll see how reproducible research in the sciences particularly increases scientific knowledge.

## Reproducible research

“The most important tool is the *mindset*, when starting, that the end product will be reproducible.” – Keith Baggerly

Another goal of this book is to help readers understand the importance of reproducible analyses. The hope is to get readers into the habit of making their analyses reproducible from the very beginning. This means we'll be trying to help you build new habits. This will take practice and be difficult at times. You'll see just why it is so important for you to keep track of your code and well-document it to help yourself later and any potential collaborators as well.

Copying and pasting results from one program into a word processor is not the way that efficient and effective scientific research is conducted. It's much more important for time to be spent on data collection and data analysis and not on copying and pasting plots back and forth across a variety of programs.

In a traditional analysis if an error was made with the original data, we'd need to step through the entire process again: recreate the plots and copy and paste all of the new plots and our statistical analysis into your document. This is error prone and a frustrating use of time. We'll see how to use R Markdown to get away from this tedious activity so that we can spend more time doing science.

“We are talking about *computational* reproducibility.” - Yihui Xie

Reproducibility means a lot of things in terms of different scientific fields. Are experiments conducted in a way that another researcher could follow the steps and get similar results? In this book, we will focus on what is known as **computational reproducibility**. This refers to being able to pass all of one's data analysis, data-sets, and conclusions to someone else and have them get exactly the same results on their machine. This allows for time to be

spent interpreting results and considering assumptions instead of the more error prone way of starting from scratch or following a list of steps that may be different from machine to machine.

# **Part I**

## **Getting started**

# 1 Getting Started with Data in R

Before we can start exploring data in R, there are some key concepts to understand first:

1. What are R and RStudio?
2. How do I code in R?
3. What are R packages?

We'll introduce these concepts in upcoming Sections 1.1 - 1.3 If you are already somewhat familiar with these concepts, feel free to skip to Section 1.4 where we'll introduce our first data set: all domestic flights departing a New York City airport in 2013. This is a dataset we will explore in depth in this book.

## 1.1 What are R and RStudio?

For much of this book, we will assume that you are using R via RStudio. First time users often confuse the two. At its simplest:

- R is like a car's engine.
- RStudio is like a car's dashboard.

R: Engine	RStudio: Dashboard
	

More precisely, R is a programming language that runs computations while RStudio is an *integrated development environment (IDE)* that provides an interface by adding many convenient features and tools. So just as having access to a speedometer, rearview mirrors, and a

navigation system makes driving much easier, using RStudio's interface makes using R much easier as well.

### 1.1.1 Using RStudio Cloud

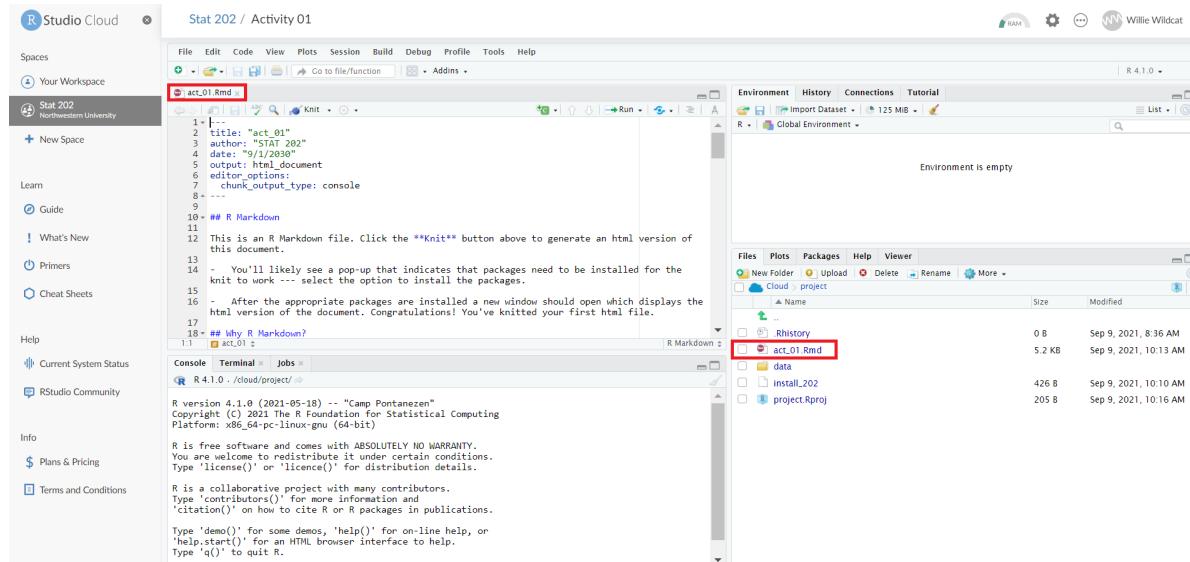
RStudio Cloud (<https://rstudio.cloud>) is a hosted version of RStudio that allows you to begin coding directly from your browser - there is no software to install and nothing to configure on your computer.

To begin using RStudio Cloud use the link provided by your instructor to gain access to the classroom workspace. You will be prompted to create a free account or log in if you have an existing account.

After you open RStudio Cloud, you should now have access to the classroom under 'Spaces' on the left hand side (in this case 'Stat 202').



Throughout class you will be working on various activities. Once the instructor has made an activity available you will click on the classroom Workspace (Stat 202) to access the available projects. To begin working on an activity click 'Start'. Once that activity project is open navigate to the 'File' pane and open the R Markdown 'Rmd' file.



You can use RStudio Cloud for personal use as well by creating projects in ‘Your Workspace’. However, RStudio Cloud limits the number of projects and amount of accessible time so it is recommended that you later install the software on your own computer.

### **1.1.2 Installing R and RStudio on your personal computer**

**Note about RStudio Server or RStudio Cloud:** If your instructor has provided you with a link and access to RStudio Server or RStudio Cloud, then you can skip this section. We do recommend after a few months of working on RStudio Server/Cloud that you return to these instructions to install this software on your own computer though. You will first need to download and install both R and RStudio (Desktop version) on your computer. It is important that you install R first and then install RStudio second.

**1. You must do this first: [Download and install R.](#)**

- If you are a Windows user: Click on “Download R for Windows”, then click on “base”, then click on the Download link.
- If you are macOS user: Click on “Download R for (Mac) OS X”, then under “Latest release:” click on R-X.X.X.pkg, where R-X.X.X is the version number. For example, the latest version of R as of August 10, 2019 was R-3.6.1.

**2. You must do this second: [Download and install RStudio.](#)**

- Scroll down to “Installers for Supported Platforms” near the bottom of the page.
- Click on the download link corresponding to your computer’s operating system.

### **1.1.3 Using R via RStudio**

Recall our car analogy from above. Much as we don’t drive a car by interacting directly with the engine but rather by interacting with elements on the car’s dashboard, we won’t be using R directly but rather we will use RStudio’s interface. After you install R and RStudio on your computer, you’ll have two new programs AKA applications you can open. We will always work in RStudio and not R. In other words:

---

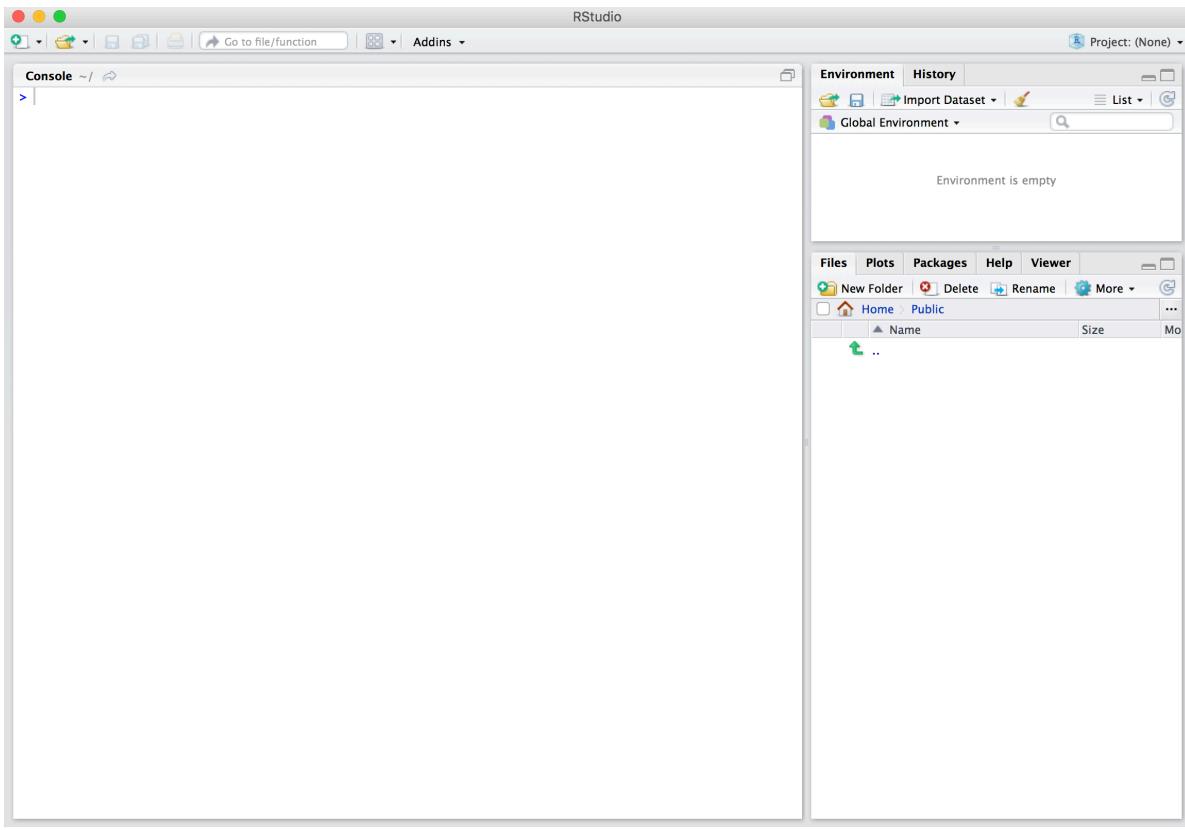
R: Do not open this

RStudio: Open this

---



After you open RStudio, you should see the following:



Note the three panes, which are three panels dividing the screen: The *Console pane*, the *Files pane*, and the *Environment pane*. Over the course of this chapter, you'll come to learn what purpose each of these panes serve.

## 1.2 How do I code in R?

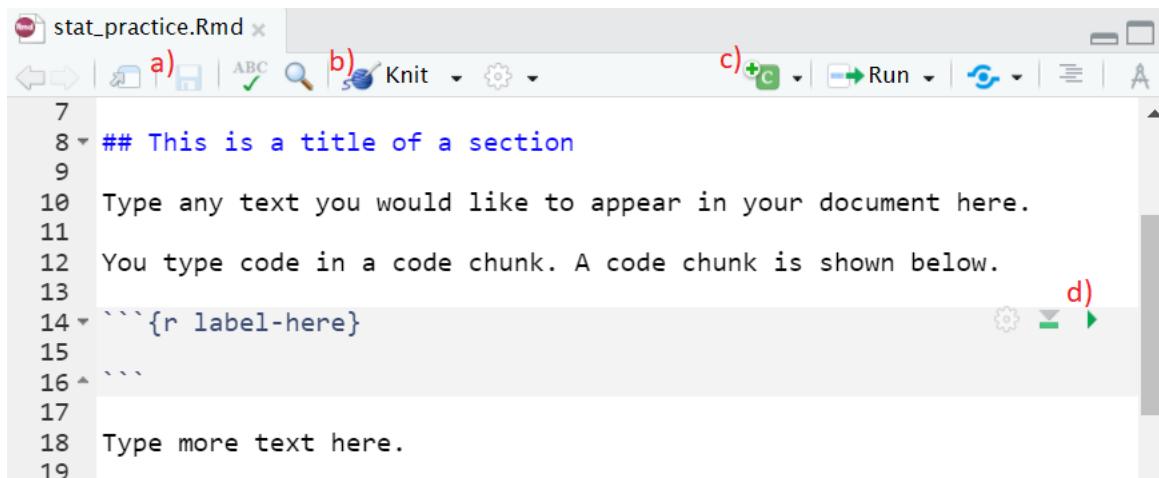
Now that you're set up with R and RStudio, you are probably asking yourself "OK. Now how do I use R?" The first thing to note is that unlike other statistical software programs like Excel, STATA, or SAS that provide [point and click](#) interfaces, R is an [interpreted language](#), meaning you have to enter in R commands written in R code. In other words, you have to code/program in R. Note that we'll use the terms "coding" and "programming" interchangeably in this book.

While it is not required to be a seasoned coder/computer programmer to use R, there is still a set of basic programming concepts that R users need to understand. Consequently, while this book is not a book on programming, you will still learn just enough of these basic programming concepts needed to explore and analyze data effectively.

### 1.2.1 Creating your first R Markdown document

R Markdown allows you to easily create a document which combines your code, the results from your code, as well as any text that accompanies the analysis. To create a new R Markdown file, in R Studio select File>New File>R Markdown. Then, you will see a window pop-up titled *New R Markdown*. Here, you specify the type of file you wish to create. HTML is generally the recommended document type since it does not have traditional *page* separators like PDF and Word do. You can also choose a title and author for your document using their respective fields. Finally, select *Ok* to create your new R Markdown file. You will see it appear as a tab in your R Studio session. Click the *save icon* to save your new document.

The following is an example of an R Markdown document:



```
stat_practice.Rmd x
a) ABC b) Knit c) Run d)
7
8 ## This is a title of a section
9
10 Type any text you would like to appear in your document here.
11
12 You type code in a code chunk. A code chunk is shown below.
13
14 ``{r label=here}
15
16 ...
17
18 Type more text here.
19
```

- a) Save your document.
- b) Click *knit* to compile your R Markdown into the document file type that you specified. The file will be saved in your *Files pane*. This will also save your document.
- c) Insert a new code chunk in your document where the cursor is located. You will often have many code chunks in your document.
- d) Run the current code chunk.

When you create your Markdown file and *knit* it into a document, the chunks are run in order and any output from them is shown in the document, in the order and location that their respective chunk appears. Sometimes you may wish to type code or analyze data without it printing in the document. If that is the case, you type the code in the *Console* rather than in the *.Rmd* file.

While you read through this book, it will be helpful to have an RMarkdown document open so you can copy code provided and paste it into a code chunk to run.

### 1.2.2 Basic programming concepts and terminology

We now introduce some basic programming concepts and terminology. Instead of asking you to learn all these concepts and terminology right now, we'll guide you so that you'll “learn by doing.” Note that in this book we will always use a different font to distinguish regular text from `computer_code`. The best way to master these topics is, in our opinions, “learning by doing” and lots of repetition.

- Basics:
  - *Console*: Where you enter in commands.
  - *Running code*: The act of telling R to perform an action by giving it commands in the console.
  - *Objects*: Where values are saved in R. In order to do useful and interesting things in R, we will want to *assign* a name to an object. For example we could do the following assignments: `x <- 44 - 20` and `three <- 3`. This would allow us to run `x + three` which would return 27.
  - *Data types*: Integers, doubles/numerics, logicals, and characters.

In R Studio try typing the following code into the console or code chunk.

```
x <- 44-20
three <- 3
x+three
```

```
[1] 27
```

You should see `x` and `three` appear as stored objects in the *Environment* pane. Anything you store in the *Environment* pane can be referenced and used later. R can also be used as a calculator, notice how it evaluates `x+three`.

- *Vectors*: A series of values. These are created using the `c()` function, where `c()` stands for “combine” or “concatenate”. For example: `c(6, 11, 13, 31, 90, 92)`.
- *Factors*: *Categorical data* are represented in R as factors.
- *Data frames*: Data frames are like rectangular spreadsheets: they are representations of datasets in R where the rows correspond to *observations* and the columns correspond to *variables* that describe the observations. We'll cover data frames later in Section Section 1.4.
- *Conditionals*:

- Testing for equality in R using `==` (and not `=` which is typically used for assignment). Ex: `2 + 1 == 3` compares `2 + 1` to `3` and is correct R code, while `2 + 1 = 3` will return an error.
- Boolean algebra: TRUE/FALSE statements and mathematical operators such as `<` (less than), `<=` (less than or equal), and `!=` (not equal to).
- Logical operators: `&` representing “and” as well as `|` representing “or.” Ex: `(2 + 1 == 3) & (2 + 1 == 4)` returns FALSE since both clauses are not TRUE (only the first clause is TRUE). On the other hand, `(2 + 1 == 3) | (2 + 1 == 4)` returns TRUE since at least one of the two clauses is TRUE.
- *Functions*, also called *commands*: Functions perform tasks in R. They take in inputs called *arguments* and return outputs. You can either manually specify a function’s arguments or use the function’s *default values*.

This list is by no means an exhaustive list of all the programming concepts and terminology needed to become a savvy R user; such a list would be so large it wouldn’t be very useful, especially for novices. Rather, we feel this is a minimally viable list of programming concepts and terminology you need to know before getting started. We feel that you can learn the rest as you go. Remember that your mastery of all of these concepts and terminology will build as you practice more and more.

### 1.2.3 Errors, warnings, and messages

One thing that intimidates new R and RStudio users is how it reports *errors*, *warnings*, and *messages*. R reports errors, warnings, and messages in a glaring red font, which makes it seem like it is scolding you. However, seeing red text in the console is not always bad.

R will show red text in the console pane in three different situations:

- **Errors:** When the red text is a legitimate error, it will be prefaced with “Error in...” and try to explain what went wrong. Generally when there’s an error, the code will not run. For example, we’ll see in Subsection 1.3.3 if you see `Error in ggplot(...)` : `could not find function "ggplot"`, it means that the `ggplot()` function is not accessible because the package that contains the function (`ggplot2`) was not loaded with `library(ggplot2)`. Thus you cannot use the `ggplot()` function without the `ggplot2` package being loaded first.
- **Warnings:** When the red text is a warning, it will be prefaced with “Warning:” and R will try to explain why there’s a warning. Generally your code will still work, but with some caveats. For example, you will see in Chapter 2 if you create a scatterplot based on a dataset where one of the values is missing, you will see this warning: `Warning: Removed 1 rows containing missing values (geom_point)`. R will still produce the scatterplot with all the remaining values, but it is warning you that one of the points isn’t there.

- **Messages:** When the red text doesn't start with either "Error" or "Warning", it's *just a friendly message*. You'll see these messages when you load *R packages* in the upcoming Subsection 1.3.2 or when you read data saved in spreadsheet files with the `read_csv()` function as you'll see in Chapter 4. These are helpful diagnostic messages and they don't stop your code from working. Additionally, you'll see these messages when you install packages too using `install.packages()`.

Remember, when you see red text in the console, *don't panic*. It doesn't necessarily mean anything is wrong. Rather:

- If the text starts with "Error", figure out what's causing it. Think of errors as a red traffic light: something is wrong!
- If the text starts with "Warning", figure out if it's something to worry about. For instance, if you get a warning about missing values in a scatterplot and you know there are missing values, you're fine. If that's surprising, look at your data and see what's missing. Think of warnings as a yellow traffic light: everything is working fine, but watch out/pay attention.
- Otherwise the text is just a message. Read it, wave back at R, and thank it for talking to you. Think of messages as a green traffic light: everything is working fine.

#### 1.2.4 Tips on learning to code

Learning to code/program is very much like learning a foreign language, it can be very daunting and frustrating at first. Such frustrations are very common and it is very normal to feel discouraged as you learn. However just as with learning a foreign language, if you put in the effort and are not afraid to make mistakes, anybody can learn.

Here are a few useful tips to keep in mind as you learn to program:

- **Remember that computers are not actually that smart:** You may think your computer or smartphone are "smart," but really people spent a lot of time and energy designing them to appear "smart." Rather you have to tell a computer everything it needs to do. Furthermore the instructions you give your computer can't have any mistakes in them, nor can they be ambiguous in any way.
- **Take the "copy, paste, and tweak" approach:** Especially when learning your first programming language, it is often much easier to take existing code that you know works and modify it to suit your ends, rather than trying to write new code from scratch. We call this the *copy, paste, and tweak* approach. So early on, we suggest not trying to write code from memory, but rather take existing examples we have provided you, then copy, paste, and tweak them to suit your goals. Don't be afraid to play around!

- **The best way to learn to code is by doing:** Rather than learning to code for its own sake, we feel that learning to code goes much smoother when you have a goal in mind or when you are working on a particular project, like analyzing data that you are interested in.
- **Practice is key:** Just as the only method to improving your foreign language skills is through practice, practice, and practice; so also the only method to improving your coding is through practice, practice, and practice. Don't worry however; we'll give you plenty of opportunities to do so!

## 1.3 What are R packages?

Another point of confusion with many new R users is the idea of an R package. R packages extend the functionality of R by providing additional functions, data, and documentation. They are written by a world-wide community of R users and can be downloaded for free from the internet. For example, among the many packages we will use in this book are the `ggplot2` package for data visualization in Chapter 2, the `dplyr` package for data wrangling in Chapter 3, and the `moderndive` package that accompanies this book.

A good analogy for R packages is they are like apps you can download onto a mobile phone:

R: A new phone	R Packages: Apps you can download
	 

So R is like a new mobile phone: while it has a certain amount of features when you use it for the first time, it doesn't have everything. R packages are like the apps you can download onto your phone from Apple's App Store or Android's Google Play.

Let's continue this analogy by considering the Instagram app for editing and sharing pictures. Say you have purchased a new phone and you would like to share a recent photo you have taken on Instagram. You need to:

1. *Install the app:* Since your phone is new and does not include the Instagram app, you need to download the app from either the App Store or Google Play. You do this once

and you're set. You might do this again in the future any time there is an update to the app.

2. *Open the app:* After you've installed Instagram, you need to open the app.

Once Instagram is open on your phone, you can then proceed to share your photo with your friends and family. The process is very similar for using an R package. You need to:

1. *Install the package:* This is like installing an app on your phone. Most packages are not installed by default when you install R and RStudio. Thus if you want to use a package for the first time, you need to install it first. Once you've installed a package, you likely won't install it again unless you want to update it to a newer version.
2. *"Load" the package:* "Loading" a package is like opening an app on your phone. Packages are not "loaded" by default when you start RStudio on your computer; you need to "load" each package you want to use every time you start RStudio.

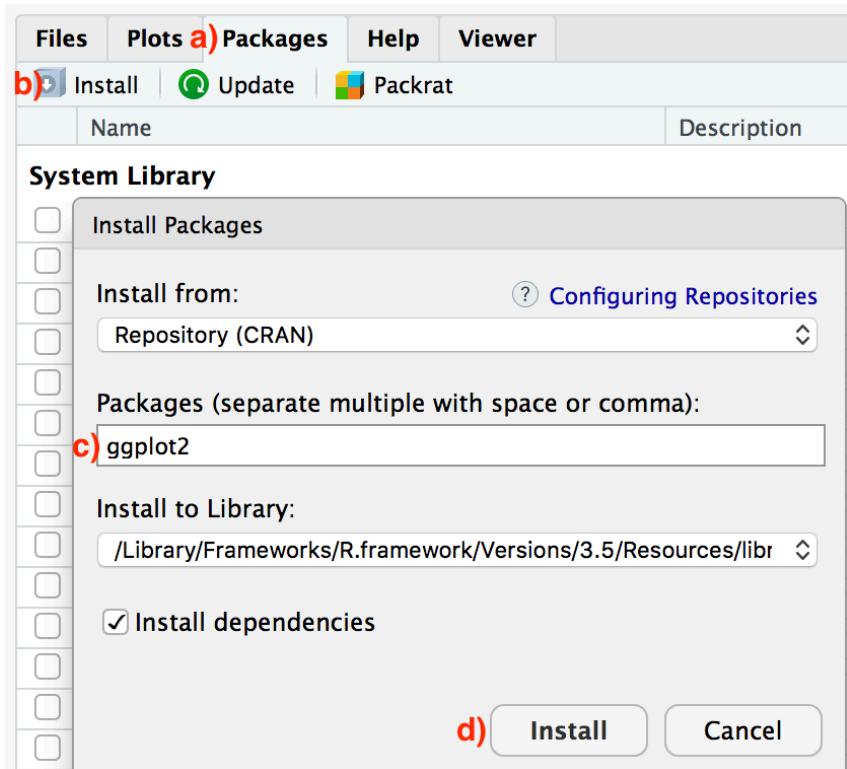
Let's now show you how to perform these two steps for the `ggplot2` package for data visualization.

### 1.3.1 Package installation

**Note about RStudio Server:** If your instructor has provided you with a link and access to RStudio Server, you probably will not need to install packages, as they have likely been pre-installed for you by your instructor. That being said, it is still a good idea to know this process for later on when you are not using RStudio Server, but rather RStudio Desktop on your own computer.

There are two ways to install an R package. For example, to install the `ggplot2` package:

1. **Easy way:** In the Files pane of RStudio:
  - a) Click on the "Packages" tab
  - b) Click on "Install"
  - c) Type the name of the package under "Packages (separate multiple with space or comma):" In this case, type `ggplot2`
  - d) Click "Install"



2. **Slightly harder way:** An alternative but slightly less convenient way to install a package is by typing `install.packages("ggplot2")` in the Console pane of RStudio and hitting enter. Note you must include the quotation marks.

Much like an app on your phone, you only have to install a package once. However, if you want to update an already installed package to a newer verions, you need to re-install it by repeating the above steps.

#### Learning Check 1.1

Repeat the above installing steps for the `dplyr`, `nycflights13`, and `knitr` packages. This will install the earlier mentioned `dplyr` package, the `nycflights13` package containing data on all domestic flights leaving a NYC airport in 2013, and the `knitr` package for writing reports in R.

### 1.3.2 Package loading

Recall that after you've installed a package, you need to “load” it, in other words open it. We do this by using the `library()` command. For example, to load the `ggplot2` package, run the

following code in the Console pane. What do we mean by “run the following code”? Either type or copy & paste the following code into the Console pane and then hit the enter key.

```
library(ggplot2)
```

If after running the above code, a blinking cursor returns next to the > “prompt” sign, it means you were successful and the `ggplot2` package is now loaded and ready to use. If however, you get a red “error message” that reads...

```
Error in library(ggplot2) : there is no package called ‘ggplot2’
```

... it means that you didn’t successfully install it. In that case, go back to the previous subsection “Package installation” and install it.

#### Learning Check 1.2

“Load” the `dplyr`, `nycflights13`, and `knitr` packages as well by repeating the above steps.

### 1.3.3 Package use

One extremely common mistake new R users make when wanting to use particular packages is that they forget to “load” them first by using the `library()` command we just saw. Remember: *you have to load each package you want to use every time you start RStudio*. If you don’t first “load” a package, but attempt to use one of its features, you’ll see an error message similar to:

```
Error: could not find function
```

R is telling you that you are trying to use a function in a package that has not yet been “loaded.” Almost all new users forget do this when starting out, and it is a little annoying to get used to. However, you’ll remember with practice.

## 1.4 Explore your first dataset

Let’s put everything we’ve learned so far into practice and start exploring some real data! Data comes to us in a variety of formats, from pictures to text to numbers. Throughout this book, we’ll focus on datasets that are saved in “spreadsheet”-type format; this is probably the most common way data are collected and saved in many fields. Remember from Subsection 1.2.2

that these “spreadsheet”-type datasets are called *data frames* in R; we will focus on working with data saved as data frames throughout this book.

Let’s first load all the packages needed for this chapter, assuming you’ve already installed them. Read Section 1.3 for information on how to install and load R packages if you haven’t already.

```
library(nycflights13)
library(dplyr)
library(knitr)
```

At the beginning of all subsequent chapters in this text, we’ll always have a list of packages that you should have installed and loaded to work with that chapter’s R code.

### 1.4.1 nycflights13 package

Many of us have flown on airplanes or know someone who has. Air travel has become an ever-present aspect in many people’s lives. If you live in or are visiting a relatively large city and you walk around that city’s airport, you see gates showing flight information from many different airlines. And you will frequently see that some flights are delayed because of a variety of conditions. Are there ways that we can avoid having to deal with these flight delays?

We’d all like to arrive at our destinations on time whenever possible. (Unless you secretly love hanging out at airports. If you are one of these people, pretend for the moment that you are very much anticipating being at your final destination.) Throughout this book, we’re going to analyze data related to flights contained in the `nycflights13` package (Wickham 2021). Specifically, this package contains five data sets saved in five separate data frames with information about all domestic flights departing from New York City in 2013. These include Newark Liberty International (EWR), John F. Kennedy International (JFK), and LaGuardia (LGA) airports:

- `flights`: Information on all 336,776 flights
- `airlines`: A table matching airline names and their two letter IATA airline codes (also known as carrier codes) for 16 airline companies
- `planes`: Information about each of 3,322 physical aircraft used.
- `weather`: Hourly meteorological data for each of the three NYC airports. This data frame has 26,115 rows, roughly corresponding to the  $365 \times 24 \times 3 = 26,280$  possible hourly measurements one can observe at three locations over the course of a year.
- `airports`: Airport names, codes, and locations for 1,458 destination airports.

### 1.4.2 flights data frame

We will begin by exploring the `flights` data frame that is included in the `nycflights13` package and getting an idea of its structure. Run the following code in your console (either by typing it or cutting & pasting it): it loads in the `flights` dataset into your Console. Note depending on the size of your monitor, the output may vary slightly.

```
flights

# A tibble: 336,776 x 19
  year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time arr_delay carrier
  <int> <int> <int>    <int>      <dbl>    <int>    <int>      <dbl>    <chr>
1 2013     1     1      517       515      2     830      819      11  UA
2 2013     1     1      533       529      4     850      830      20  UA
3 2013     1     1      542       540      2     923      850      33  AA
4 2013     1     1      544       545     -1    1004     1022     -18  B6
5 2013     1     1      554       600      -6     812      837     -25  DL
6 2013     1     1      554       558      -4     740      728      12  UA
7 2013     1     1      555       600      -5     913      854      19  B6
8 2013     1     1      557       600      -3     709      723     -14  EV
9 2013     1     1      557       600      -3     838      846      -8  B6
10 2013    1     1      558       600      -2     753      745      8  AA
# ... with 336,766 more rows, 9 more variables: flight <int>, tailnum <chr>,
#   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
#   minute <dbl>, time_hour <dttm>, and abbreviated variable names
#   1: sched_dep_time, 2: dep_delay, 3: arr_time, 4: sched_arr_time,
#   5: arr_delay
```

Let's unpack this output:

- `A tibble: 336,776 x 19`: A `tibble` is a kind of data frame used in R. This particular data frame has
  - 336,776 rows
  - 19 columns corresponding to 19 variables describing each observation
- `year month day dep_time sched_dep_time dep_delay arr_time sched_arr_time arr_delay carrier` are different columns, in other words variables, of this data frame.
- We then have the first 10 rows of observations corresponding to 10 flights.
- `... with 336,766 more rows, and 11 more variables:` indicating to us that 336,766 more rows of data and 11 more variables could not fit in this screen.

Unfortunately, this output does not allow us to explore the data very well. Let's look at different tools to explore data frames.

### 1.4.3 Exploring data frames

Among the many ways of getting a feel for the data contained in a data frame such as `flights`, we present three functions that take as their “argument”, in other words their input, the data frame in question. We also include a fourth method for exploring one particular column of a data frame:

1. Using the `View()` function built for use in RStudio. We will use this the most.
2. Using the `glimpse()` function, which is included in the `dplyr` package.
3. Using the `kable()` function, which is included in the `knitr` package.
4. Using the `$` operator to view a single variable in a data frame.

#### 1. `View()`:

Run `View(flights)` in your Console in RStudio, either by typing it or cutting & pasting it into the Console pane, and explore this data frame in the resulting pop-up viewer. You should get into the habit of always `Viewing` any data frames that come your way. Note the capital “V” in `View`. R is case-sensitive so you’ll receive an error if you run `view(flights)` instead of `View(flights)`.

#### Learning Check 1.3

What does any *ONE* row in this `flights` dataset refer to?

- a. Data on an airline
- b. Data on a flight
- c. Data on an airport
- d. Data on multiple flights

By running `View(flights)`, we see the different *variables* listed in the columns and we see that there are different types of variables. Some of the variables like `distance`, `day`, and `arr_delay` are what we will call *quantitative* variables. These variables are numerical in nature. Other variables here are *categorical*.

Note that if you look in the leftmost column of the `View(flights)` output, you will see a column of numbers. These are the row numbers of the dataset. If you glance across a row with the same number, say row 5, you can get an idea of what each row corresponds to. In other words, this will allow you to identify what object is being referred to in a given row. This is often called the *observational unit*. The *observational unit* in this example is an individual flight departing New York City in 2013. You can identify the observational unit by determining what “thing” is being measured or described by each of the variables.

## 2. `glimpse()`:

The second way to explore a data frame is using the `glimpse()` function included in the `dplyr` package. Thus, you can only use the `glimpse()` function after you've loaded the `dplyr` package. This function provides us with an alternative method for exploring a data frame:

```
glimpse(flights)
```

```
Rows: 336,776
Columns: 19
$ year           <int> 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2-
$ month          <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1-
$ day            <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1-
$ dep_time        <int> 517, 533, 542, 544, 554, 554, 555, 557, 557, 558, 558, ~
$ sched_dep_time <int> 515, 529, 540, 545, 600, 558, 600, 600, 600, 600, ~
$ dep_delay       <dbl> 2, 4, 2, -1, -6, -4, -5, -3, -3, -2, -2, -2, -2, -1-
$ arr_time        <int> 830, 850, 923, 1004, 812, 740, 913, 709, 838, 753, 849, ~
$ sched_arr_time  <int> 819, 830, 850, 1022, 837, 728, 854, 723, 846, 745, 851, ~
$ arr_delay       <dbl> 11, 20, 33, -18, -25, 12, 19, -14, -8, 8, -2, -3, 7, -1-
$ carrier         <chr> "UA", "UA", "AA", "B6", "DL", "UA", "B6", "EV", "B6", "-
$ flight          <int> 1545, 1714, 1141, 725, 461, 1696, 507, 5708, 79, 301, 4-
$ tailnum         <chr> "N14228", "N24211", "N619AA", "N804JB", "N668DN", "N394-
$ origin          <chr> "EWR", "LGA", "JFK", "JFK", "LGA", "EWR", "EWR", "LGA", ~
$ dest            <chr> "IAH", "IAH", "MIA", "BQN", "ATL", "ORD", "FLL", "IAD", ~
$ air_time         <dbl> 227, 227, 160, 183, 116, 150, 158, 53, 140, 138, 149, 1-
$ distance         <dbl> 1400, 1416, 1089, 1576, 762, 719, 1065, 229, 944, 733, ~
$ hour             <dbl> 5, 5, 5, 5, 6, 5, 6, 6, 6, 6, 6, 6, 6, 5, 6, 6, 6-
$ minute           <dbl> 15, 29, 40, 45, 0, 58, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 59, 0-
$ time_hour        <dttm> 2013-01-01 05:00:00, 2013-01-01 05:00:00, 2013-01-01 0-
```

We see that `glimpse()` will give you the first few entries of each variable in a row after the variable. In addition, the *data type* (see Subsection 1.2.2) of the variable is given immediately after each variable's name inside `< >`. Here, `int` and `dbl` refer to “integer” and “double”, which are computer coding terminology for quantitative/numerical variables. In contrast, `chr` refers to “character”, which is computer terminology for text data. Text data, such as the `carrier` or `origin` of a flight, are categorical variables. The `time_hour` variable is an example of one more type of data type: `dttm`. As you may suspect, this variable corresponds to a specific date and time of day. However, we won't work with dates in this class and leave it to a more advanced book on data science.

### Learning Check 1.4

What are some examples in this dataset of **categorical** variables? What makes them different than **quantitative** variables?

#### 3. `kable()`:

The another way to explore the entirety of a data frame is using the `kable()` function from the `knitr` package. Let's explore the different carrier codes for all the airlines in our dataset two ways. Run both of these lines of code in your Console:

```
airlines  
kable(airlines)
```

At first glance, it may not appear that there is much difference in the outputs. However when using tools for document production such as [R Markdown](#), the latter code produces output that is much more legible and reader-friendly.

#### 4. `$ operator`

Lastly, the `$` operator allows us to explore a single variable within a data frame. For example, run the following in your console

```
airlines  
airlines$name
```

We used the `$` operator to extract only the `name` variable and return it as a vector of length 16. We will only be occasionally exploring data frames using this operator, instead favoring the `View()` and `glimpse()` functions.

#### 1.4.4 Help files

Another nice feature of R is the help system. You can get help in R by entering a `?` before the name of a function or data frame in question and you will be presented with a page showing the documentation. For example, let's look at the help file for the `flights` data frame:

```
?flights
```

A help file should pop-up in the Help pane of RStudio. If you have questions about a function or data frame included in an R package, you should get in the habit of consulting the help file right away.

## 1.5 Conclusion

We've given you what we feel are the most essential concepts to know before you can start exploring data in R. Is this chapter exhaustive? Absolutely not. To try to include everything in this chapter would make the chapter so large it wouldn't be useful!

### 1.5.1 Additional resources

If you are completely new to the world of coding, R, and RStudio and feel you could benefit from a more detailed introduction, we suggest you check out Chester Ismay's short book [Getting used to R, RStudio, and R Markdown](#) (Ismay 2016), which includes screencast recordings that you can follow along and pause as you learn. Furthermore, there is an introduction to R Markdown, a tool used for reproducible research in R.

The screenshot shows the front page of the book. On the left is a sidebar with a table of contents:

<a href="#">1 Introduction</a>
<a href="#">2 Why R?</a>
<a href="#">3 R and RStudio Basics</a>
<a href="#">4 R Markdown</a>
<a href="#">5 Intro to R using R Markdown</a>
<a href="#">6 Deciphering Common R Errors</a>
<a href="#">7 Concluding Remarks</a>
<a href="#">8 References</a>

Below the sidebar, there are links: "Published with bookdown", "Create a GitHub Issue", and "Email Chester".

The main content area has the title **Getting used to R, RStudio, and R Markdown**. Below the title are the authors' names: *Chester Ismay* and *Patrick C. Kennedy*, and the publication date: *2018-05-23*.

## 1 Introduction

This book was written to give people who are new to [R](#), [RStudio](#), and [R Markdown](#) the tools they need to begin making their own research reproducible. R is an open-source programming language that has seen its popularity grow tremendously in recent years, with developers adding new functionality via packages on a daily basis. RStudio is a graphical development environment that makes it easier to write and view the results of R code, and R Markdown provides an easy way to produce rich, fully-documented, reproducible analyses.

## **Part II**

# **Data Exploration via the tidyverse**

# 2 Data Visualization

We begin the development of your data science toolbox with data visualization. By visualizing our data, we gain valuable insights that we couldn't initially see from just looking at the raw data in spreadsheet form. We will use the `ggplot2` package as it provides an easy way to customize your plots. `ggplot2` is rooted in the data visualization theory known as *The Grammar of Graphics* (Wilkinson 2005).

At the most basic level, graphics/plots/charts (we use these terms interchangeably in this book) provide a nice way for us to get a sense for how quantitative variables compare in terms of their center (where the values tend to be located) and their spread (how they vary around the center). Graphics should be designed to emphasize the findings and insight you want your audience to understand. This does however require a balancing act. On the one hand, you want to highlight as many meaningful relationships and interesting findings as possible; on the other you don't want to include so many as to overwhelm your audience.

As we will see, plots/graphics also help us to identify patterns and outliers in our data. We will see that a common extension of these ideas is to compare the *distribution* of one quantitative variable (i.e., what the spread of a variable looks like or how the variable is *distributed* in terms of its values) as we go across the levels of a different categorical variable.

## Packages Needed

Let's load all the packages needed for this chapter (this assumes you've already installed them). Read Section 1.3 for information on how to install and load R packages.

```
library(nycflights13)
library(ggplot2)
library(dplyr)
```

## 2.1 The Grammar of Graphics

We begin with a discussion of a theoretical framework for data visualization known as “The Grammar of Graphics,” which serves as the foundation for the `ggplot2` package. Think of how we construct sentences in English to form sentences by combining different elements,

like nouns, verbs, particles, subjects, objects, etc. However, we can't just combine these elements in any arbitrary order; we must do so following a set of rules known as a linguistic grammar. Similarly to a linguistic grammar, "The Grammar of Graphics" define a set of rules for constructing *statistical graphics* by combining different types of *layers*. This grammar was created by Leland Wilkinson (Wilkinson 2005) and has been implemented in a variety of data visualization software including R.

### 2.1.1 Components of the Grammar

In short, the grammar tells us that:

**A statistical graphic is a mapping of data variables to aesthetic attributes of geometric objects.**

Specifically, we can break a graphic into three essential components:

1. **data**: the data set composed of variables that we map.
2. **geom**: the geometric object in question. This refers to the type of object we can observe in a plot. For example: points, lines, and bars.
3. **aes**: aesthetic attributes of the geometric object. For example, x-position, y-position, color, shape, and size. Each assigned aesthetic attribute can be mapped to a variable in our data set.

You might be wondering why we wrote the terms **data**, **geom**, and **aes** in a computer code type font. We'll see very shortly that we'll specify the elements of the grammar in R using these terms. However, let's first break down the grammar with an example.

### 2.1.2 Gapminder data

In February 2006, a statistician named Hans Rosling gave a TED talk titled "[The best stats you've ever seen](#)" where he presented global economic, health, and development data from the website [gapminder.org](#). For example, for the 142 countries included from 2007, let's consider only the first 6 countries when listed alphabetically in Table 2.1.

Table 2.1: Gapminder 2007 Data: First 6 of 142 countries

Country	Continent	Life Expectancy	Population	GDP per Capita
Afghanistan	Asia	43.8	31889923	975
Albania	Europe	76.4	3600523	5937
Algeria	Africa	72.3	33333216	6223
Angola	Africa	42.7	12420476	4797
Argentina	Americas	75.3	40301927	12779

Country	Continent	Life Expectancy	Population	GDP per Capita
Australia	Oceania	81.2	20434176	34435

Each row in this table corresponds to a country in 2007. For each row, we have 5 columns:

1. **Country**: Name of country.
2. **Continent**: Which of the five continents the country is part of. (Note that “Americas” includes countries in both North and South America and that Antarctica is excluded.)
3. **Life Expectancy**: Life expectancy in years.
4. **Population**: Number of people living in the country.
5. **GDP per Capita**: Gross domestic product (in US dollars).

Now consider Figure 2.1, which plots this data for all 142 countries in the data.

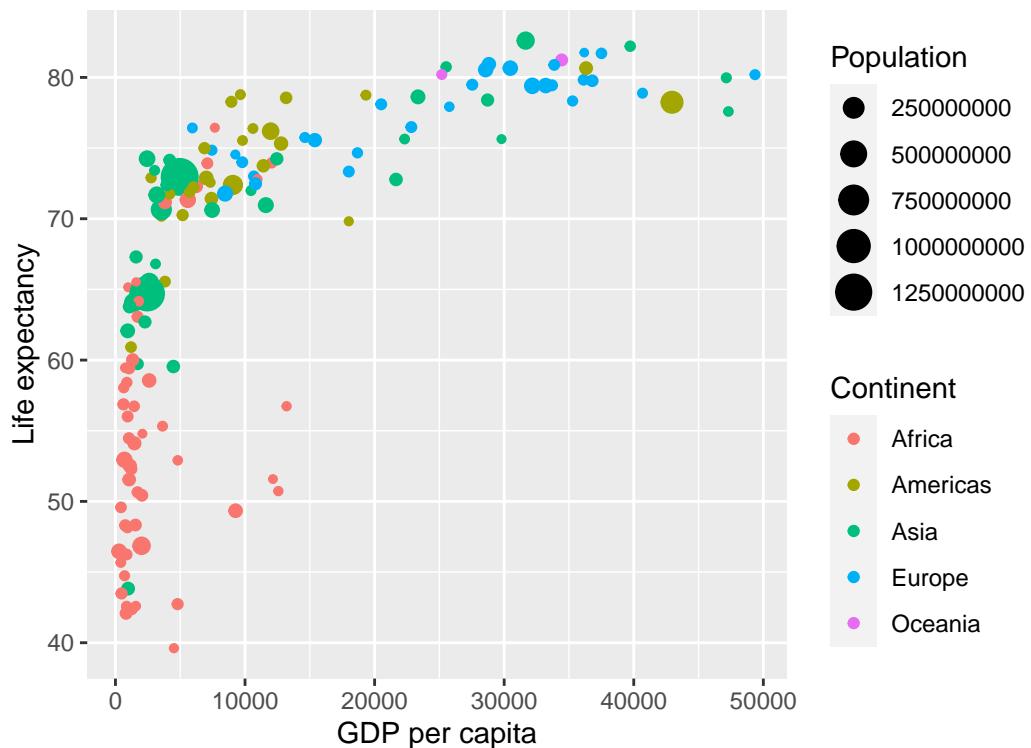


Figure 2.1: Life Expectancy over GDP per Capita in 2007

Let's view this plot through the grammar of graphics:

1. The **data** variable **GDP per Capita** gets mapped to the **x-position** **aesthetic** of the points.

2. The **data** variable **Life Expectancy** gets mapped to the **y**-position **aesthetic** of the points.
3. The **data** variable **Population** gets mapped to the **size** aesthetic of the points.
4. The **data** variable **Continent** gets mapped to the **color** aesthetic of the points.

We'll see shortly that **data** corresponds to the particular data frame where our data is saved and a "data variable" corresponds to a particular column in the data frame. Furthermore, the type of **geometric object** considered in this plot are points. That being said, while in this example we are considering points, graphics are not limited to just points. Other plots involve lines while others involve bars.

Let's summarize the three essential components of the Grammar in Table 2.2.

Table 2.2: Summary of Grammar of Graphics for this plot

data variable	aes	geom
GDP per Capita	x	point
Life Expectancy	y	point
Population	size	point
Continent	color	point

### 2.1.3 Other components

There are other components of the Grammar of Graphics we can control as well. As you start to delve deeper into the Grammar of Graphics, you'll start to encounter these topics more frequently. In this book however, we'll keep things simple and only work with the two additional components listed below:

- **faceting** breaks up a plot into small multiples corresponding to the levels of another variable (Section 2.6)
- **position** adjustments for barplots (Section 2.8)

Other more complex components like **scales** and coordinate systems are left for a more advanced text such as [R for Data Science](#) (Grolmund and Wickham 2016). Generally speaking, the Grammar of Graphics allows for a high degree of customization of plots and also a consistent framework for easily updating and modifying them.

### 2.1.4 ggplot2 package

In this book, we will be using the **ggplot2** package for data visualization, which is an implementation of the Grammar of Graphics for R (Wickham et al. 2022). As we noted earlier, a lot of the previous section was written in a computer code type font. This is because the various

components of the Grammar of Graphics are specified in the `ggplot()` function included in the `ggplot2` package, which expects at a minimum as arguments (i.e. inputs):

- The data frame where the variables exist: the `data` argument.
- The mapping of the variables to aesthetic attributes: the `mapping` argument which specifies the `aesthetic` attributes involved.

After we've specified these components, we then add *layers* to the plot using the `+` sign. The most essential layer to add to a plot is the layer that specifies which type of `geometric` object we want the plot to involve: points, lines, bars, and others. Other layers we can add to a plot include layers specifying the plot title, axes labels, visual themes for the plots, and facets (which we'll see in Section 2.6).

Let's now put the theory of the Grammar of Graphics into practice.

## 2.2 Five Named Graphs - The 5NG

In order to keep things simple, we will only focus on five types of graphics in this book, each with a commonly given name. We term these “five named graphs” the **5NG**:

1. scatterplots
2. linegraphs
3. boxplots
4. histograms
5. barplots

We will discuss some variations of these plots, but with this basic repertoire of graphics in your toolbox you can visualize a wide array of different variable types. Note that certain plots are only appropriate for categorical variables and while others are only appropriate for quantitative variables. You'll want to quiz yourself often as we go along on which plot makes sense a given a particular problem or data set.

## 2.3 5NG#1: Scatterplots

The simplest of the 5NG are *scatterplots*, also called bivariate plots. They allow you to visualize the relationship between two numerical variables. While you may already be familiar with scatterplots, let's view them through the lens of the Grammar of Graphics. Specifically, we will visualize the relationship between the following two numerical variables in the `flights` data frame included in the `nycflights13` package:

1. `dep_delay`: departure delay on the horizontal “x” axis and
2. `arr_delay`: arrival delay on the vertical “y” axis

for Alaska Airlines flights leaving NYC in 2013. This requires paring down the data from all 336,776 flights that left NYC in 2013, to only the 714 *Alaska Airlines* flights that left NYC in 2013.

What this means computationally is: we'll take the `flights` data frame, extract only the 714 rows corresponding to Alaska Airlines flights, and save this in a new data frame called `alaska_flights`. Run the code below to do this:

```
alaska_flights <- flights %>%
  filter(carrier == "AS")
```

For now we suggest you ignore how this code works; we'll explain this in detail in Chapter 3 when we cover data wrangling. However, convince yourself that this code does what it is supposed to by running `View(alaska_flights)`: it creates a new data frame `alaska_flights` consisting of only the 714 Alaska Airlines flights.

We'll see later in Chapter 3 on data wrangling that this code uses the `dplyr` package for data wrangling to achieve our goal: it takes the `flights` data frame and filters it to only return the rows where `carrier` is equal to "AS", Alaska Airlines' carrier code. Other examples of carrier codes include "AA" for American Airlines and "UA" for United Airlines. Recall from Section 1.2 that testing for equality is specified with `==` and not `=`. Fasten your seat belts and sit tight for now however, we'll introduce these ideas more fully in Chapter 3.

#### Learning Check 2.1

Take a look at both the `flights` and `alaska_flights` data frames by running `View(flights)` and `View(alaska_flights)`. In what respect do these data frames differ?

### 2.3.1 Scatterplots via `geom_point`

Let's now go over the code that will create the desired scatterplot, keeping in mind our discussion on the Grammar of Graphics in Section 2.1. We'll be using the `ggplot()` function included in the `ggplot2` package.

```
ggplot(data = alaska_flights, mapping = aes(x = dep_delay, y = arr_delay)) +
  geom_point()
```

Let's break this down piece-by-piece:

- Within the `ggplot()` function, we specify two of the components of the Grammar of Graphics as arguments (i.e. inputs):
  1. The data frame to be `alaska_flights` by setting `data = alaska_flights`.

2. The aesthetic mapping by setting `aes(x = dep_delay, y = arr_delay)`. Specifically:
  - the variable `dep_delay` maps to the `x` position aesthetic
  - the variable `arr_delay` maps to the `y` position aesthetic
- We add a layer to the `ggplot()` function call using the `+` sign. The layer in question specifies the third component of the grammar: the geometric object. In this case the geometric object are points, set by specifying `geom_point()`.

After running the above code, you'll notice two outputs: a warning message and the graphic shown in Figure 2.2. Let's first unpack the warning message:

```
Warning: Removed 5 rows containing missing values (geom_point).
```

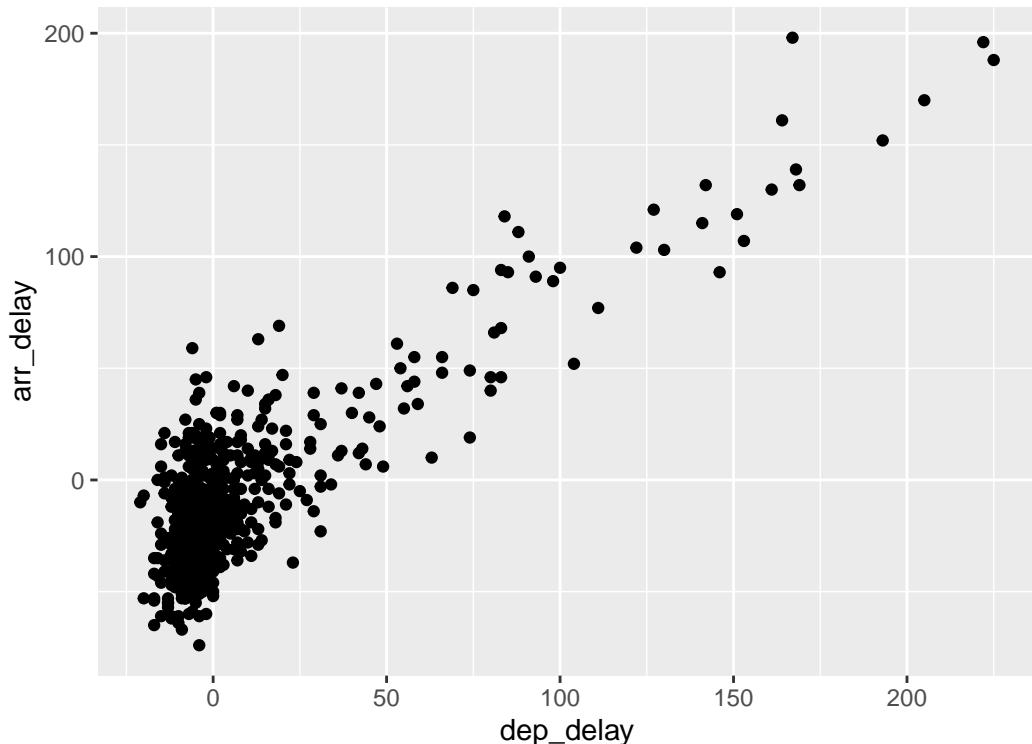


Figure 2.2: Arrival Delays vs Departure Delays for Alaska Airlines flights from NYC in 2013

After running the above code, R returns a warning message alerting us to the fact that 5 rows were ignored due to them being missing. For 5 rows either the value for `dep_delay` or `arr_delay` or both were missing (recorded in R as `NA`), and thus these rows were ignored in our plot. Turning our attention to the resulting scatterplot in Figure 2.2, we see that a positive

relationship exists between `dep_delay` and `arr_delay`: as departure delays increase, arrival delays tend to also increase. We also note the large mass of points clustered near  $(0, 0)$ .

Before we continue, let's consider a few more notes on the layers in the above code that generated the scatterplot:

- Note that the `+` sign comes at the end of lines, and not at the beginning. You'll get an error in R if you put it at the beginning.
- When adding layers to a plot, you are encouraged to start a new line after the `+` so that the code for each layer is on a new line. As we add more and more layers to plots, you'll see this will greatly improve the legibility of your code.
- To stress the importance of adding layers in particular the layer specifying the geometric object, consider Figure 2.3 where no layers are added. A not very useful plot!

```
ggplot(data = alaska_flights, mapping = aes(x = dep_delay, y = arr_delay))
```

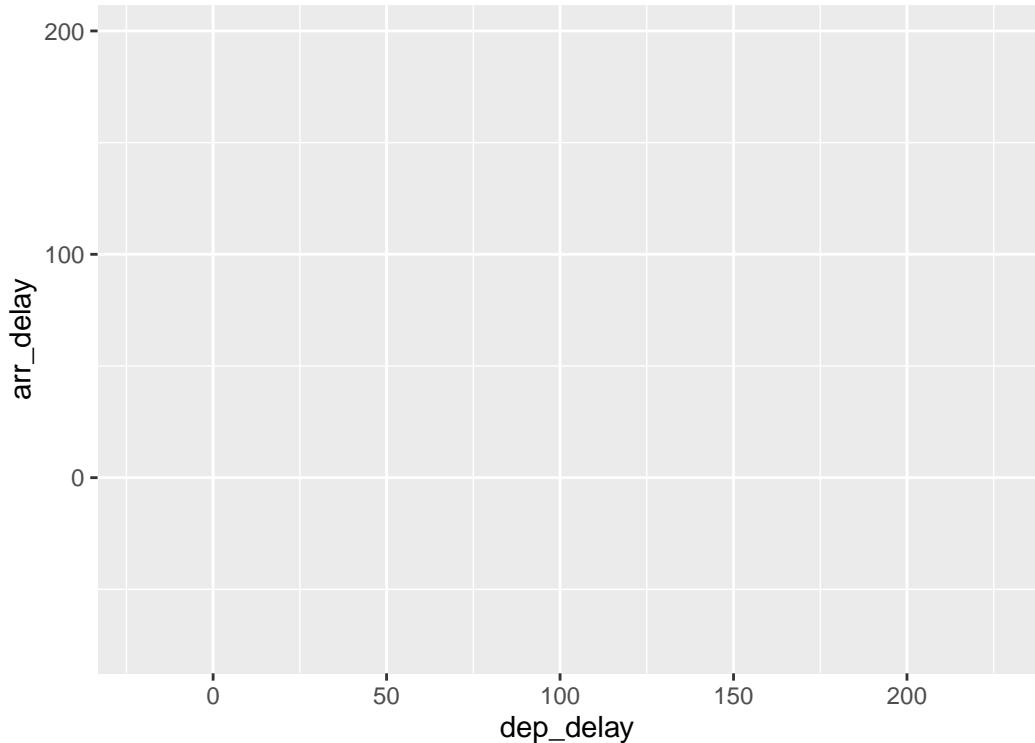


Figure 2.3: Plot with no layers

#### Learning Check 2.2

What are some practical reasons why `dep_delay` and `arr_delay` have a positive relationship?

#### Learning Check 2.3

What variables (not necessarily in the `flights` data frame) would you expect to have a negative correlation (i.e. a negative relationship) with `dep_delay`? Why? Remember that we are focusing on numerical variables here.

#### Learning Check 2.4

Why do you believe there is a cluster of points near (0, 0)? What does (0, 0) correspond to in terms of the Alaskan flights?

#### Learning Check 2.5

What are some other features of the plot that stand out to you?

#### Learning Check 2.6

Create a new scatterplot using different variables in the `alaska_flights` data frame by modifying the example above.

### 2.3.2 Over-plotting

The large mass of points near (0, 0) in Figure 2.2 can cause some confusion as it is hard to tell the true number of points that are plotted. This is the result of a phenomenon called *overplotting*. As one may guess, this corresponds to values being plotted on top of each other *over* and *over* again. It is often difficult to know just how many values are plotted in this way when looking at a basic scatterplot as we have here. There are two methods to address the issue of overplotting:

1. By adjusting the transparency of the points.
2. By adding a little random “jitter”, or random “nudges”, to each of the points.

#### Method 1: Changing the transparency

The first way of addressing overplotting is by changing the transparency of the points by using the `alpha` argument in `geom_point()`. By default, this value is set to 1. We can change this to any value between 0 and 1, where 0 sets the points to be 100% transparent and 1

sets the points to be 100% opaque. Note how the following code is identical to the code in Section 2.3 that created the scatterplot with overplotting, but with `alpha = 0.2` added to the `geom_point()`:

```
ggplot(data = alaska_flights, mapping = aes(x = dep_delay, y = arr_delay)) +
  geom_point(alpha = 0.2)
```

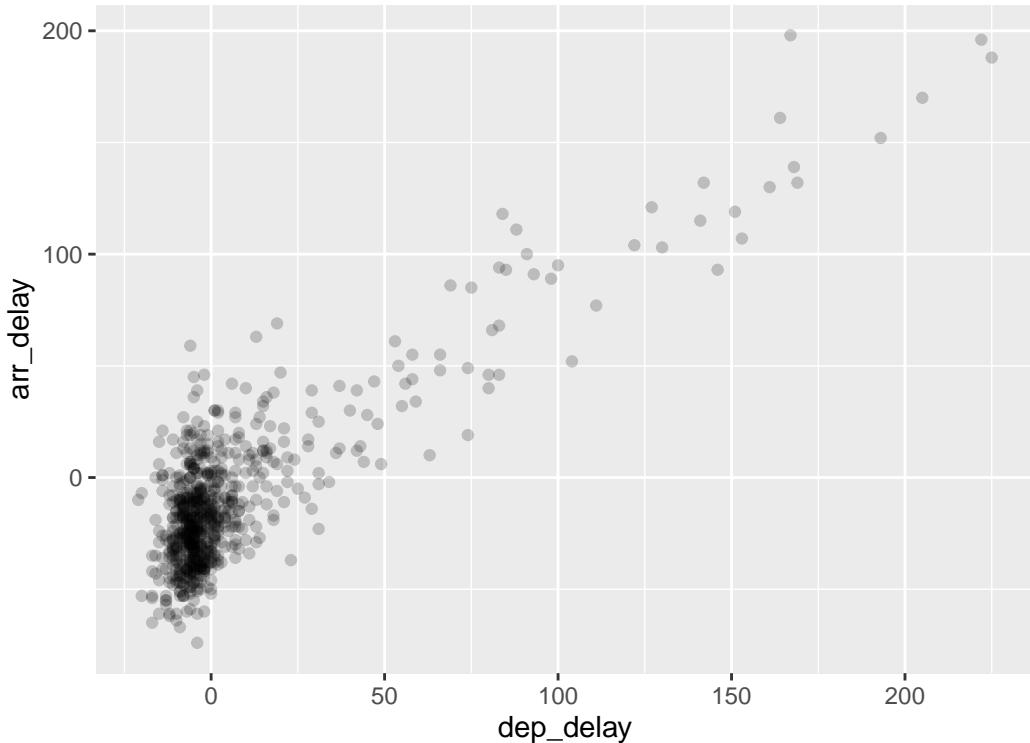


Figure 2.4: Delay scatterplot with  $\alpha = 0.2$

The key feature to note in Figure 2.4 is that the transparency of the points is cumulative: areas with a high-degree of overplotting are darker, whereas areas with a lower degree are less dark. Note furthermore that there is no `aes()` surrounding `alpha = 0.2`. This is because we are not mapping a variable to an aesthetic attribute, but rather merely changing the default setting of `alpha`. In fact, you'll receive an error if you try to change the second line above to read `geom_point(aes(alpha = 0.2))`.

### Method 2: Jittering the points

The second way of addressing overplotting is by *jittering* all the points, in other words give each point a small nudge in a random direction. You can think of “jittering” as shaking the

points around a bit on the plot. Let's illustrate using a simple example first. Say we have a data frame `jitter_example` with 4 rows of identical value 0 for both `x` and `y`:

```
# A tibble: 4 x 2
  x     y
  <dbl> <dbl>
1 0     0
2 0     0
3 0     0
4 0     0
```

We display the resulting scatterplot in Figure 2.5; observe that the 4 points are superimposed on top of each other. While we know there are 4 values being plotted, this fact might not be apparent to others.

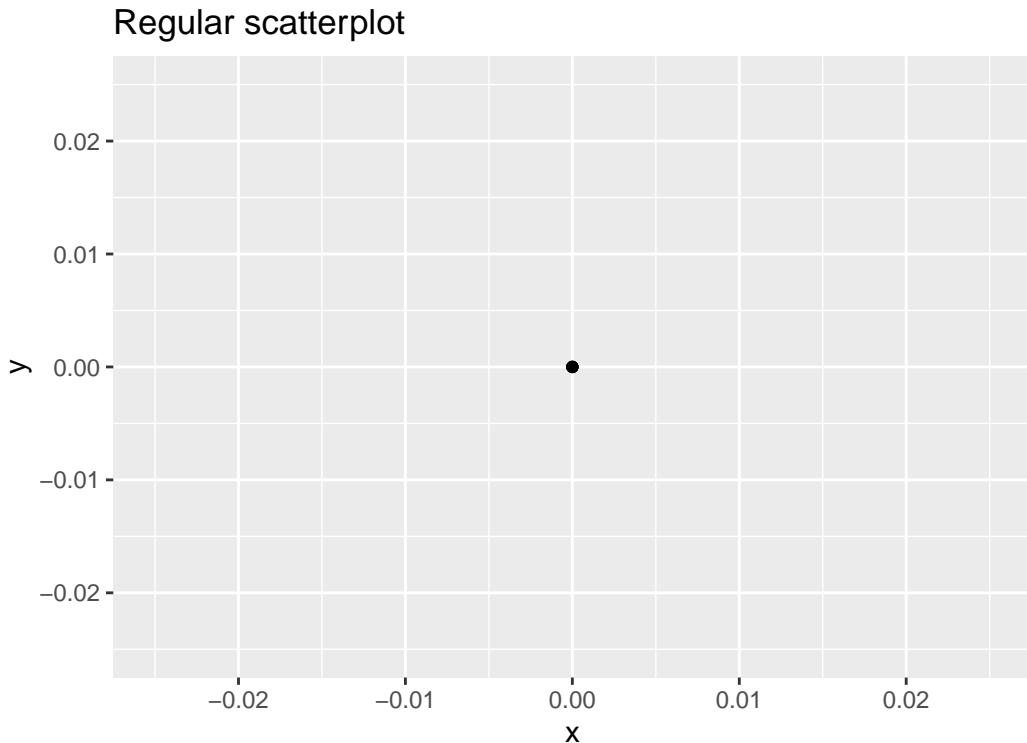


Figure 2.5: Regular scatterplot of jitter example data

In Figure 2.6 we instead display a *jittered scatterplot* where each point is given a random “nudge.” It is now plainly evident that this plot involves four points. Keep in mind that jittering is strictly a visualization tool; even after creating a jittered scatterplot, the original values saved in `jitter_example` remain unchanged.

Jittered scatterplot

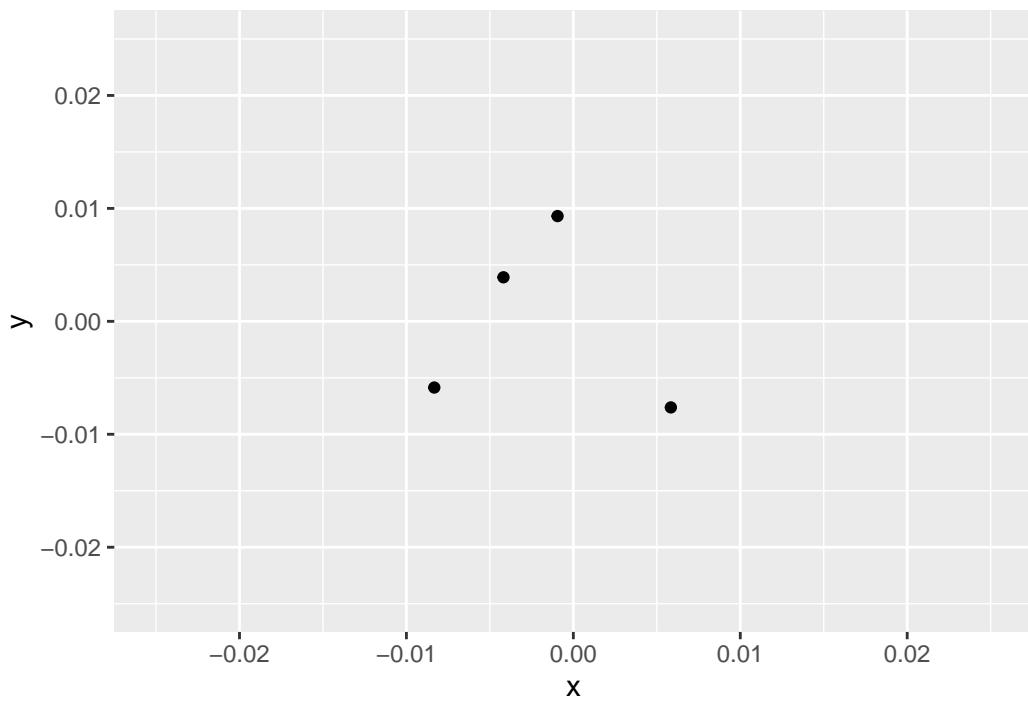


Figure 2.6: Jittered scatterplot of jitter example data

To create a jittered scatterplot, instead of using `geom_point()`, we use `geom_jitter()`. To specify how much jitter to add, we adjust the `width` and `height` arguments. This corresponds to how hard you'd like to shake the plot in units corresponding to those for both the horizontal and vertical variables (in this case minutes).

```
ggplot(data = alaska_flights, mapping = aes(x = dep_delay, y = arr_delay)) +
  geom_jitter(width = 30, height = 30)
```

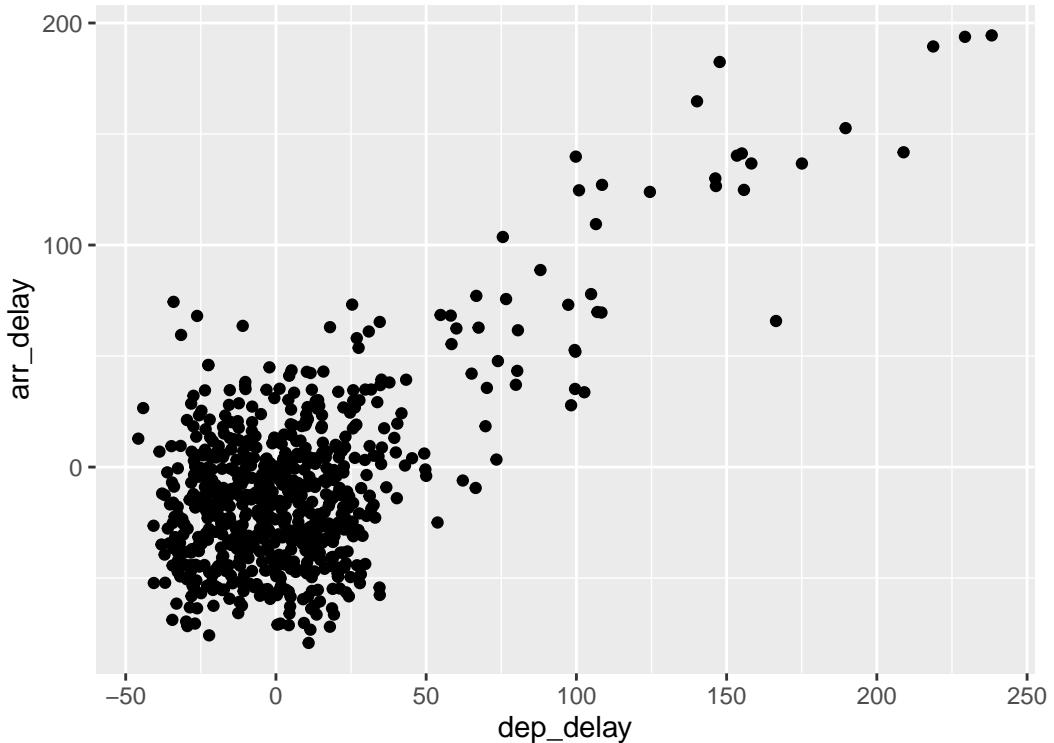


Figure 2.7: Jittered delay scatterplot

Observe how the above code is identical to the code that created the scatterplot with overplotting in Subsection 2.3.1, but with `geom_point()` replaced with `geom_jitter()`.

The resulting plot in Figure 2.7 helps us a little bit in getting a sense for the overplotting, but with a relatively large data set like this one (714 flights), it can be argued that changing the transparency of the points by setting `alpha` proved more effective. In terms of how much jitter one should add using the `width` and `height` arguments, it is important to add just enough jitter to break any overlap in points, but not so much that we completely alter the overall pattern in points.

#### Learning Check 2.7

Why is setting the `alpha` argument value useful with scatterplots? What further information does it give you that a regular scatterplot cannot?

#### Learning Check 2.8

After viewing the Figure 2.4 above, give an approximate range of arrival delays and departure delays that occur the most frequently. How has that region changed compared to when you observed the same plot without the `alpha = 0.2` set in Figure 2.2?

### 2.3.3 Summary

Scatterplots display the relationship between two numerical variables. They are among the most commonly used plots because they can provide an immediate way to see the trend in one variable versus another. However, if you try to create a scatterplot where either one of the two variables is not numerical, you might get strange results. Be careful!

With medium to large data sets, you may need to play around with the different modifications one can make to a scatterplot. This tweaking is often a fun part of data visualization, since you'll have the chance to see different relationships come about as you make subtle changes to your plots.

## 2.4 5NG#2: Linegraphs

The next of the five named graphs are linegraphs. Linegraphs show the relationship between two numerical variables when the variable on the x-axis, also called the *explanatory* variable, is of a sequential nature; in other words there is an inherent ordering to the variable. The most common example of linegraphs have some notion of time on the x-axis: hours, days, weeks, years, etc. Since time is sequential, we connect consecutive observations of the variable on the y-axis with a line. Linegraphs that have some notion of time on the x-axis are also called *time series* plots. Linegraphs should be avoided when there is not a clear sequential ordering to the variable on the x-axis. Let's illustrate linegraphs using another data set in the `nycflights13` package: the `weather` data frame.

Let's get a sense for the `weather` data frame:

- Explore the `weather` data by running `View(weather)`.
- Run `?weather` to bring up the help file.

We can see that there is a variable called `temp` of hourly temperature recordings in Fahrenheit at weather stations near all three airports in New York City: Newark (`origin` code `EWR`), JFK, and La Guardia (`LGA`). Instead of considering hourly temperatures for all days in 2013 for all three airports however, for simplicity let's only consider hourly temperatures at only Newark airport for the first 15 days in January.

Recall in Section 2.3 we used the `filter()` function to only choose the subset of rows of `flights` corresponding to Alaska Airlines flights. We similarly use `filter()` here, but by using the `&` operator we only choose the subset of rows of `weather` where

1. The `origin` is "EWR" and
2. the `month` is January and
3. the `day` is between 1 and 15

```
early_january_weather <- weather %>%
  filter(origin == "EWR" & month == 1 & day <= 15)
```

#### Learning Check 2.9

Take a look at both the `weather` and `early_january_weather` data frames by running `View(weather)` and `View(early_january_weather)`. In what respect do these data frames differ?

#### Learning Check 2.10

`View()` the `flights` data frame again. Why does the `time_hour` variable uniquely identify the hour of the measurement whereas the `hour` variable does not?

### 2.4.1 Linegraphs via `geom_line`

Let's plot a linegraph of hourly temperatures in `early_january_weather` by using `geom_line()` instead of `geom_point()` like we did for scatterplots:

```
ggplot(data = early_january_weather, mapping = aes(x = time_hour, y = temp)) +
  geom_line()
```

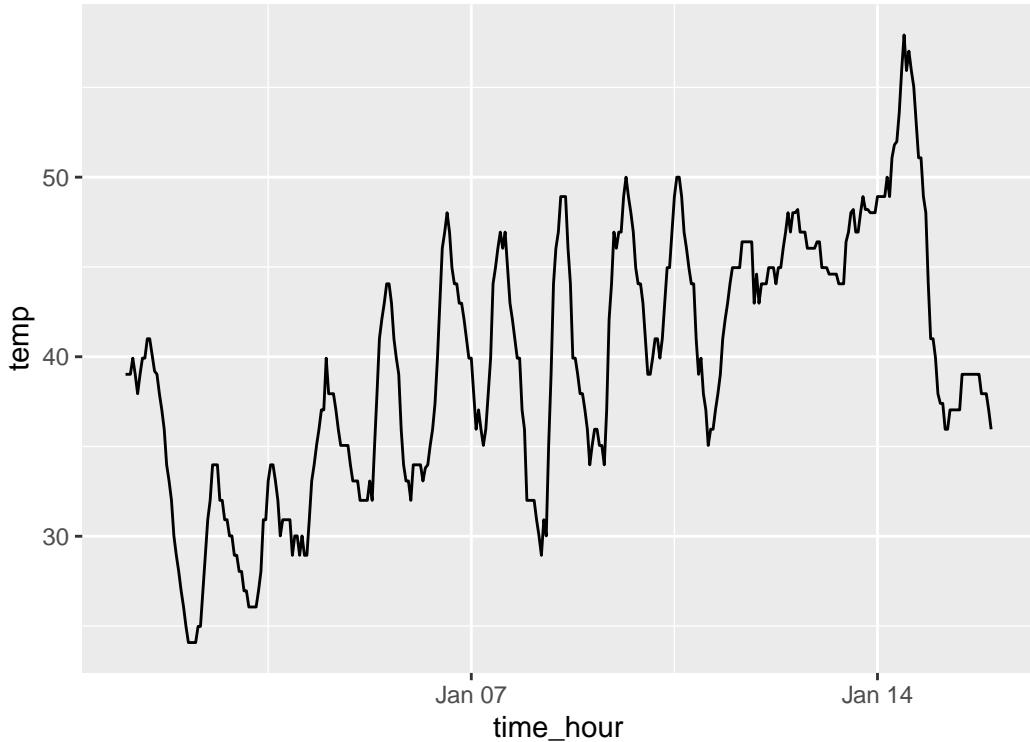


Figure 2.8: Hourly Temperature in Newark for January 1-15, 2013

Much as with the `ggplot()` code that created the scatterplot of departure and arrival delays for Alaska Airlines flights in Figure 2.2, let's break down the above code piece-by-piece in terms of the Grammar of Graphics:

- Within the `ggplot()` function call, we specify two of the components of the Grammar of Graphics as arguments:
  1. The `data` frame to be `early_january_weather` by setting `data = early_january_weather`
  2. The `aesthetic` mapping by setting `aes(x = time_hour, y = temp)`. Specifically:
    - the variable `time_hour` maps to the `x` position aesthetic.
    - the variable `temp` maps to the `y` position aesthetic
- We add a layer to the `ggplot()` function call using the `+` sign. The layer in question specifies the third component of the grammar: the `geometric object` in question. In this case the geometric object is a `line`, set by specifying `geom_line()`.

#### Learning Check 2.11

Why should linegraphs be avoided when there is not a clear ordering of the horizontal axis?

#### Learning Check 2.12

Why are linegraphs frequently used when time is the explanatory variable on the x-axis?

#### Learning Check 2.13

Plot a time series of a variable other than `temp` for Newark Airport in the first 15 days of January 2013.

### 2.4.2 Summary

Linegraphs, just like scatterplots, display the relationship between two numerical variables. However it is preferred to use linegraphs over scatterplots when the variable on the x-axis (i.e. the explanatory variable) has an inherent ordering, like some notion of time.

## 2.5 5NG#3: Histograms

Let's consider the `temp` variable in the `weather` data frame once again, but unlike with the linegraphs in Section 2.4, let's say we don't care about the relationship of temperature to time, but rather we only care about how the values of `temp` *distribute*. In other words:

1. What are the smallest and largest values?
2. What is the “center” value?
3. How do the values spread out?
4. What are frequent and infrequent values?

One way to visualize this *distribution* of this single variable `temp` is to plot them on a horizontal line as we do in Figure 2.9:



Figure 2.9: Plot of Hourly Temperature Recordings from NYC in 2013

This gives us a general idea of how the values of `temp` distribute: observe that temperatures vary from around 11°F up to 100°F. Furthermore, there appear to be more recorded temperatures between 40°F and 60°F than outside this range. However, because of the high degree of overlap in the points, it's hard to get a sense of exactly how many values are between, say, 50°F and 55°F.

What is commonly produced instead of the above plot is known as a *histogram*. A histogram is a plot that visualizes the *distribution* of a numerical value as follows:

1. We first cut up the x-axis into a series of *bins*, where each bin represents a range of values.
2. For each bin, we count the number of observations that fall in the range corresponding to that bin.
3. Then for each bin, we draw a bar whose height marks the corresponding count.

Let's drill-down on an example of a histogram, shown in @fig-histogramexample.

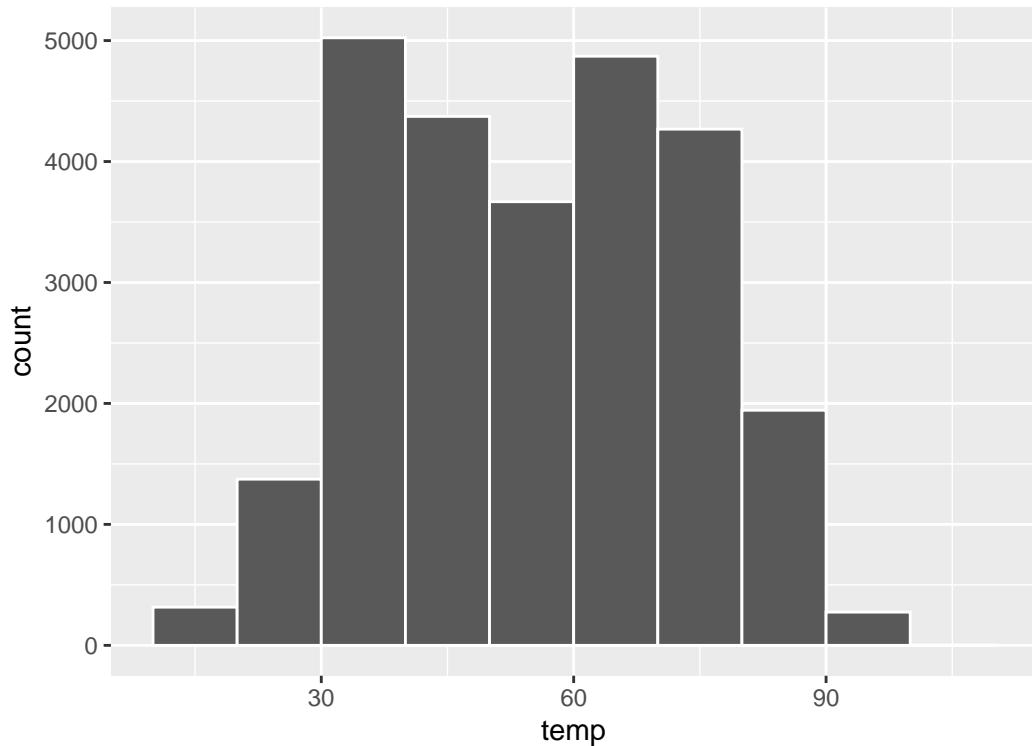


Figure 2.10: Example histogram

Observe that there are three bins of equal width between 30°F and 60°F, thus we have three bins of width 10°F each: one bin for the 30-40°F range, another bin for the 40-50°F range, and another bin for the 50-60°F range. Since:

1. The bin for the 30-40°F range has a height of around 5000, this histogram is telling us that around 5000 of the hourly temperature recordings are between 30°F and 40°F.
2. The bin for the 40-50°F range has a height of around 4300, this histogram is telling us that around 4300 of the hourly temperature recordings are between 40°F and 50°F.
3. The bin for the 50-60°F range has a height of around 3500, this histogram is telling us that around 3500 of the hourly temperature recordings are between 50°F and 60°F.

The remaining bins all have a similar interpretation.

### 2.5.1 Histograms via geom\_histogram

Let's now present the `ggplot()` code to plot your first histogram! Unlike with scatterplots and linegraphs, there is now only one variable being mapped in `aes()`: the single numerical variable `temp`. The y-aesthetic of a histogram gets computed for you automatically. Furthermore, the geometric object layer is now a `geom_histogram()`

```
ggplot(data = weather, mapping = aes(x = temp)) +
  geom_histogram()

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

Warning: Removed 1 rows containing non-finite values (stat_bin).
```

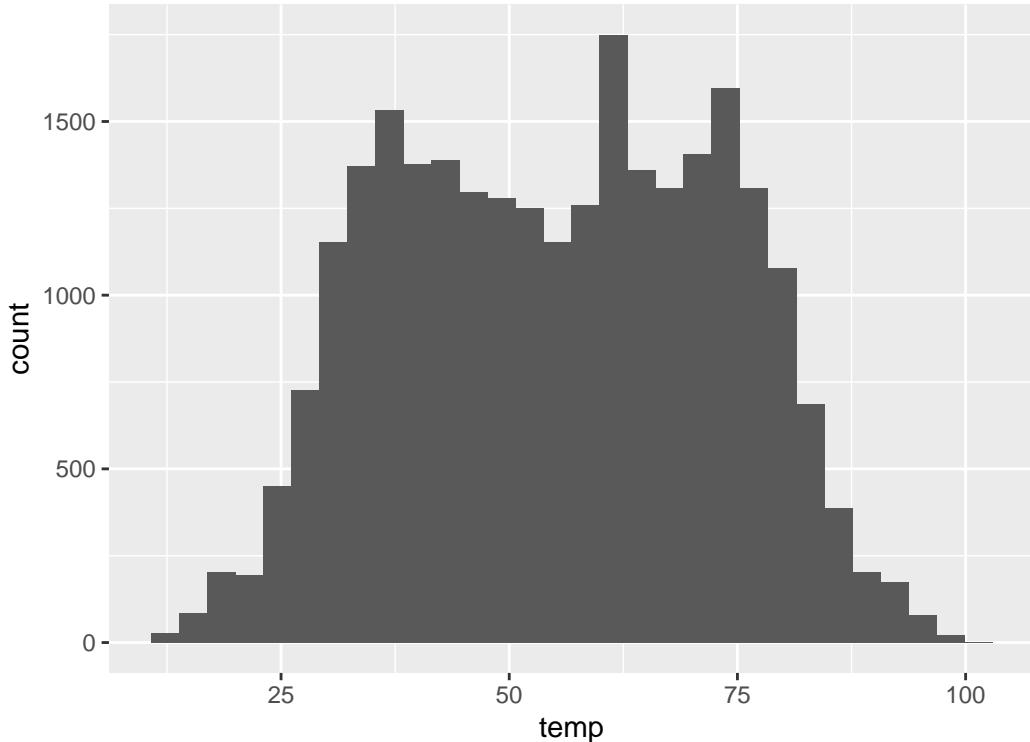


Figure 2.11: Histogram of hourly temperatures at three NYC airports

Let's unpack the messages R sent us first. The first message is telling us that the histogram was constructed using `bins = 30`, in other words 30 equally spaced bins. This is known in computer programming as a default value; unless you override this default number of bins with a number you specify, R will choose 30 by default. We'll see in the next section how to change this default number of bins. The second message is telling us something similar to the warning message we received when we ran the code to create a scatterplot of departure and arrival delays for Alaska Airlines flights in Figure 2.2: that because one row has a missing `NA` value for `temp`, it was omitted from the histogram. R is just giving us a friendly heads up that this was the case.

Now's let's unpack the resulting histogram in Figure 2.11. Observe that values less than 25°F as well as values above 80°F are rather rare. However, because of the large number of bins, its hard to get a sense for which range of temperatures is covered by each bin; everything is one giant amorphous blob. So let's add white vertical borders demarcating the bins by adding a `color = "white"` argument to `geom_histogram()`:

```
ggplot(data = weather, mapping = aes(x = temp)) +
  geom_histogram(color = "white")
```

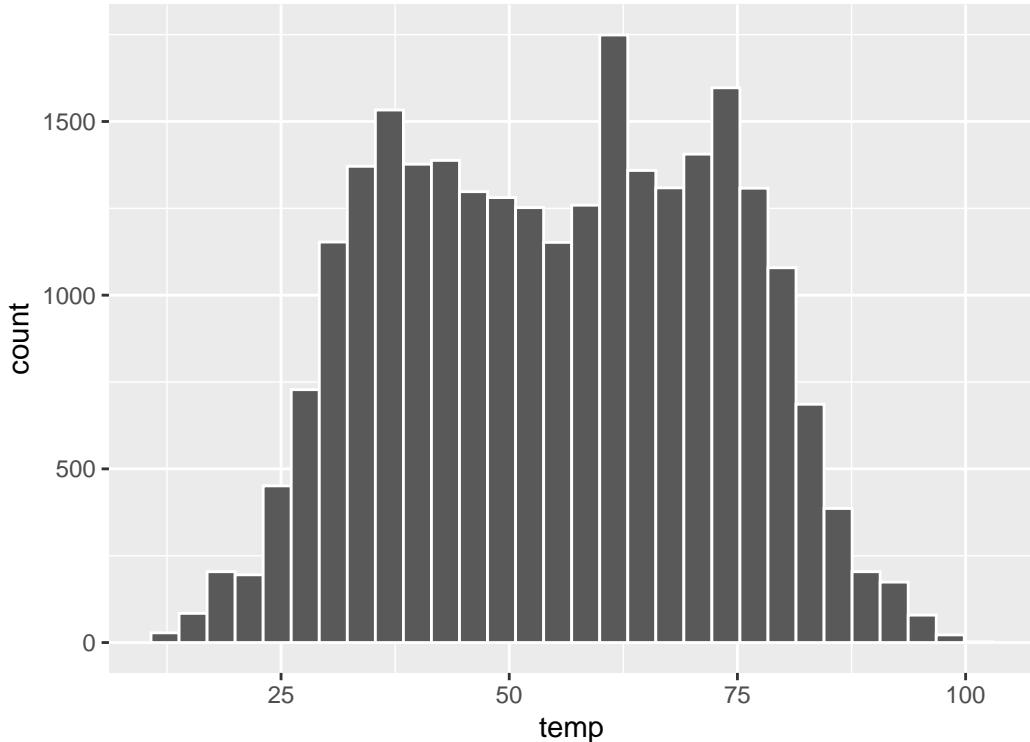


Figure 2.12: Histogram of hourly temperatures at three NYC airports with white borders

We can now better associate ranges of temperatures to each of the bins. We can also vary the color of the bars by setting the `fill` argument. Run `colors()` to see all 657 possible choice of colors!

```
ggplot(data = weather, mapping = aes(x = temp)) +  
  geom_histogram(color = "white", fill = "steelblue")
```

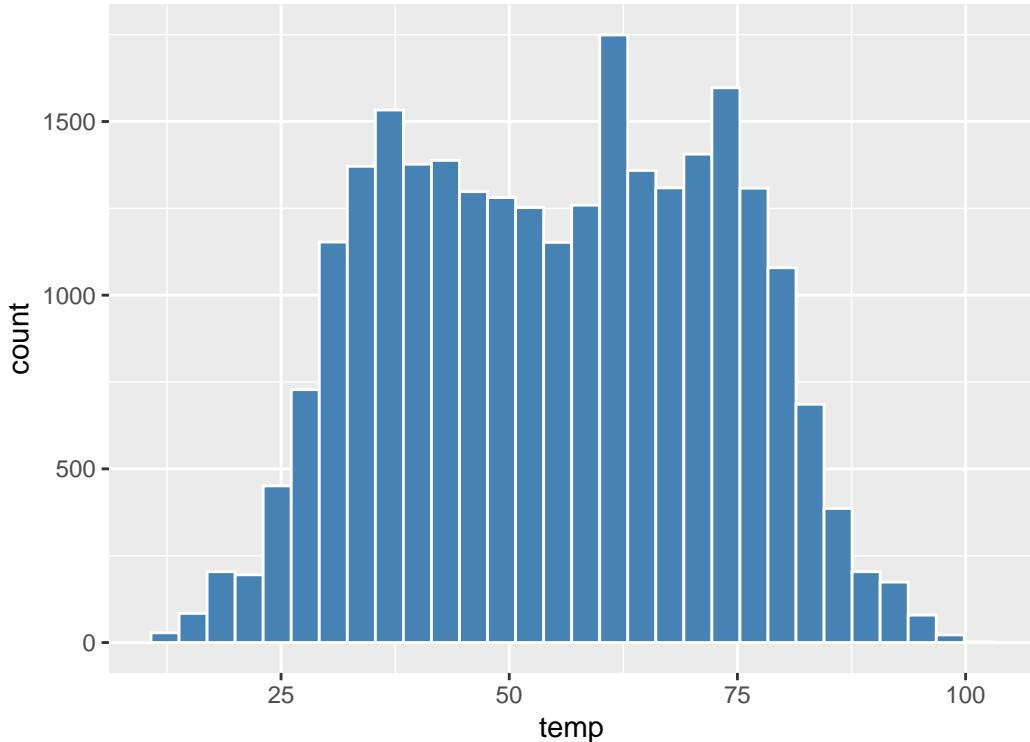


Figure 2.13: Histogram of hourly temperatures at three NYC airports with white borders

### 2.5.2 Adjusting the bins

Observe in both Figure 2.12 and Figure 2.13 that in the 50-75°F range there appear to be roughly 8 bins. Thus each bin has width 25 divided by 8, or roughly 3.12°F which is not a very easily interpretable range to work with. Let's now adjust the number of bins in our histogram in one of two methods:

1. By adjusting the number of bins via the `bins` argument to `geom_histogram()`.
2. By adjusting the width of the bins via the `binwidth` argument to `geom_histogram()`.

Using the first method, we have the power to specify how many bins we would like to cut the x-axis up in. As mentioned in the previous section, the default number of bins is 30. We can override this default, to say 40 bins, as follows:

```
ggplot(data = weather, mapping = aes(x = temp)) +
  geom_histogram(bins = 40, color = "white")
```

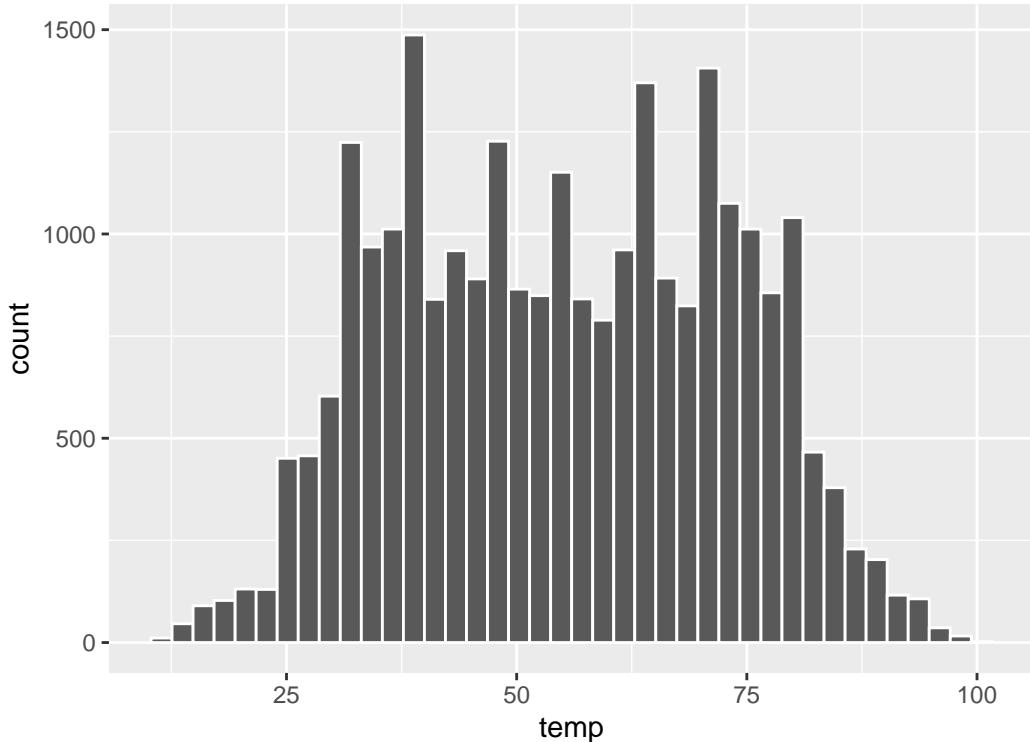


Figure 2.14: Histogram with 40 bins

Using the second method, instead of specifying the number of bins, we specify the width of the bins by using the `binwidth` argument in the `geom_histogram()` layer. For example, let's set the width of each bin to be 10°F.

```
ggplot(data = weather, mapping = aes(x = temp)) +  
  geom_histogram(binwidth = 10, color = "white")
```

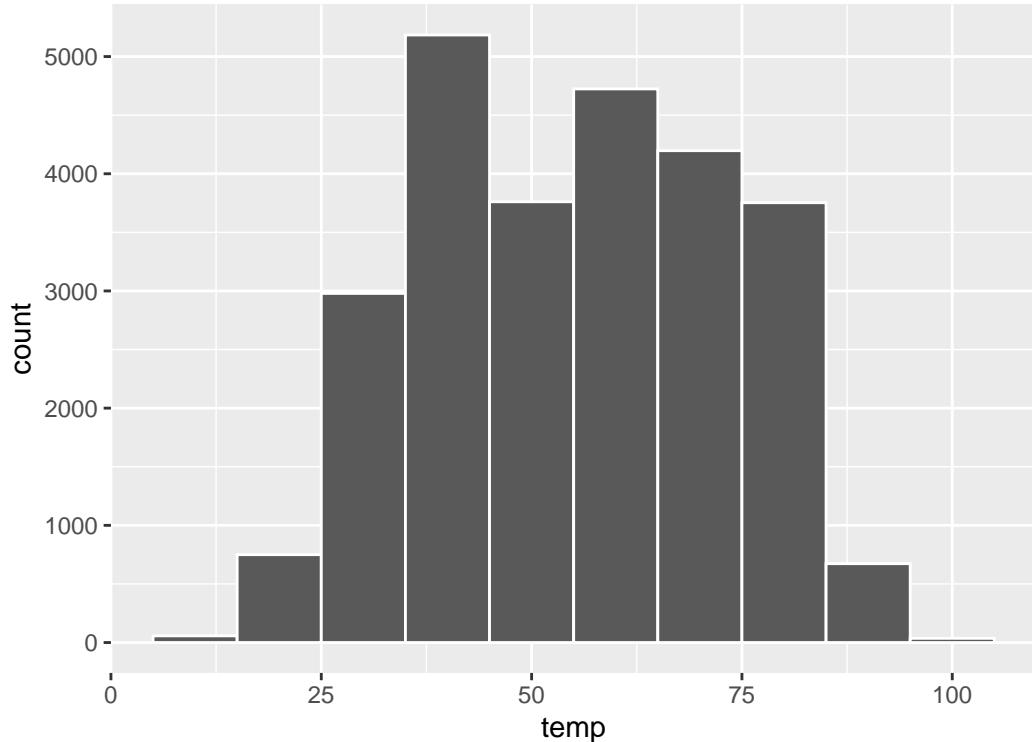


Figure 2.15: Histogram with binwidth 10

Learning Check 2.14

What does changing the number of bins from 30 to 40 tell us about the distribution of temperatures?

Learning Check 2.15

Would you classify the distribution of temperatures as symmetric or skewed?

Learning Check 2.16

What would you guess is the “center” value in this distribution? Why did you make that choice?

Learning Check 2.17

Is this data spread out greatly from the center or is it close? Why?

### 2.5.3 Summary

Histograms, unlike scatterplots and linegraphs, present information on only a single numerical variable. Specifically, they are visualizations of the distribution of the numerical variable in question.

## 2.6 Facets

Before continuing the 5NG, let's briefly introduce a new concept called *faceting*. Faceting is used when we'd like to split a particular visualization of variables by another variable. This will create multiple copies of the same type of plot with matching x and y axes, but whose content will differ.

For example, suppose we were interested in looking at how the histogram of hourly temperature recordings at the three NYC airports we saw in Section 2.5 differed by month. We would "split" this histogram by the 12 possible months in a given year, in other words plot histograms of `temp` for each `month`. We do this by adding `facet_wrap(~ month)` layer.

```
ggplot(data = weather, mapping = aes(x = temp)) +  
  geom_histogram(binwidth = 5, color = "white") +  
  facet_wrap(~ month)
```

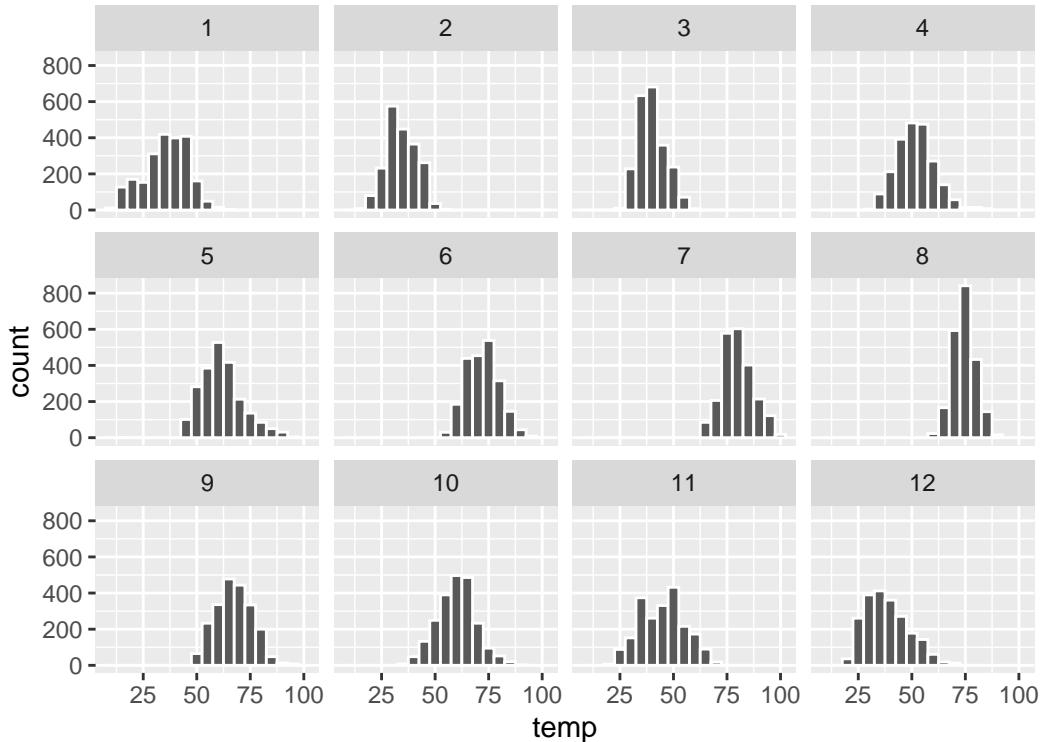


Figure 2.16: Faceted histogram

Note the use of the tilde ~ before month in `facet_wrap()`. The tilde is required and you'll receive the error `Error in as.quoted(facets) : object 'month' not found` if you don't include it before month here. We can also specify the number of rows and columns in the grid by using the `nrow` and `ncol` arguments inside of `facet_wrap()`. For example, say we would like our faceted plot to have 4 rows instead of 3. Add the `nrow = 4` argument to `facet_wrap(~ month)`

```
ggplot(data = weather, mapping = aes(x = temp)) +
  geom_histogram(binwidth = 5, color = "white") +
  facet_wrap(~ month, nrow = 4)
```

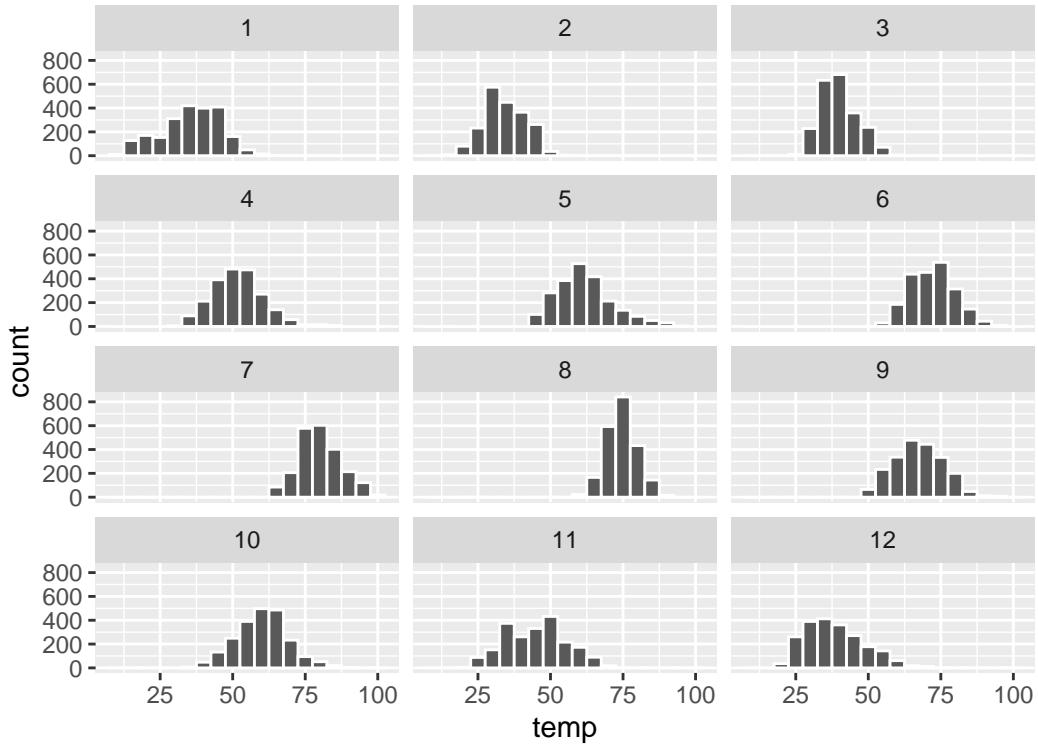


Figure 2.17: Faceted histogram with 4 instead of 3 rows

Observe in both Figure 2.16 and Figure 2.17 that as we might expect in the Northern Hemisphere, temperatures tend to be higher in the summer months, while they tend to be lower in the winter.

#### Learning Check 2.18

What other things do you notice about the faceted plot above? How does a faceted plot help us see relationships between two variables?

#### Learning Check 2.19

What do the numbers 1-12 correspond to in the plot above? What about 25, 50, 75, 100?

#### Learning Check 2.20

For which types of data sets would these types of faceted plots not work well in comparing relationships between variables? Give an example describing the nature of these variables

and other important characteristics.

#### Learning Check 2.21

Does the `temp` variable in the `weather` data set have a lot of variability? Why do you say that?

## 2.7 5NG#4: Boxplots

While faceted histograms are one visualization that allows us to compare distributions of a numerical variable split by another variable, another visualization that achieves this same goal are *side-by-side boxplots*. A boxplot is constructed from the information provided in the *five-number summary* of a numerical variable (see Appendix A). To keep things simple for now, let's only consider hourly temperature recordings for the month of November in Figure 2.18.

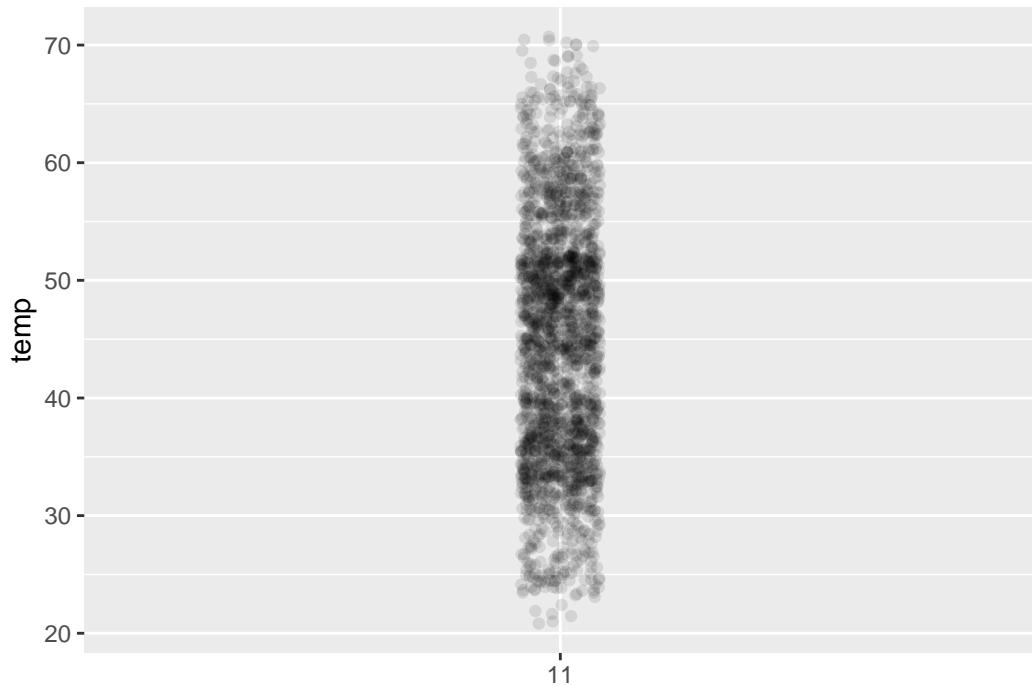


Figure 2.18: November temperatures

These 2141 observations have the following five-number summary:

1. Minimum: 21.02°F

2. First quartile AKA 25<sup>th</sup> percentile: 35.96°F
3. Median AKA second quartile AKA 50<sup>th</sup> percentile: 44.96°F
4. Third quartile AKA 75<sup>th</sup> percentile: 51.98°F
5. Maximum: 71.06°F

Let's mark these 5 values with dashed horizontal lines in Figure 2.19.

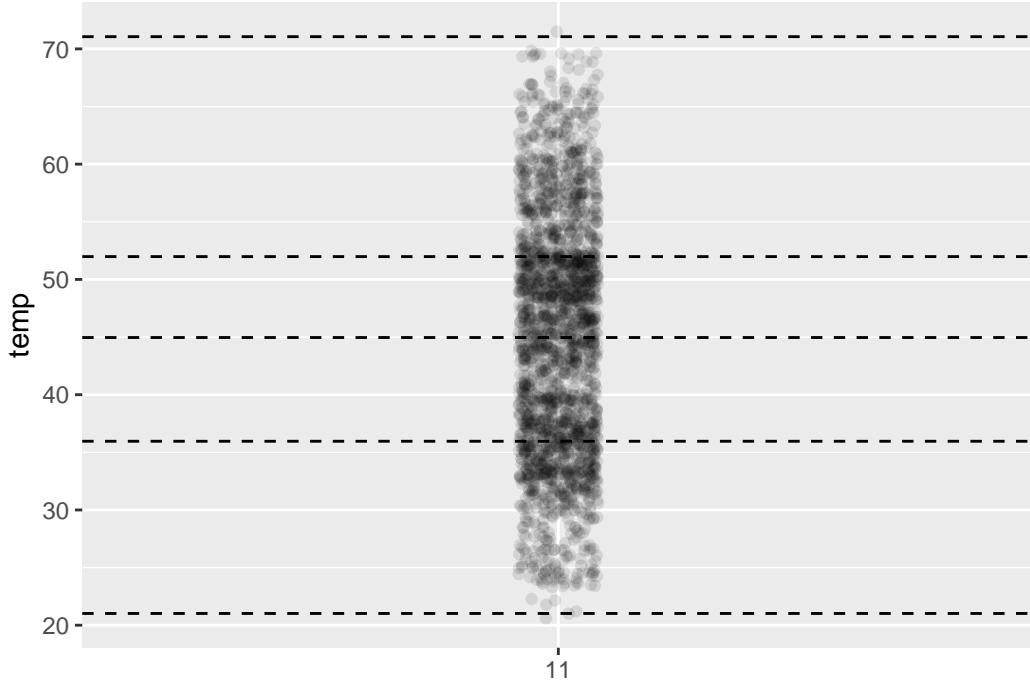


Figure 2.19: November temperatures

Let's add the boxplot underneath these points and dashed horizontal lines in Figure 2.20.

What the boxplot does summarize the 2141 points by emphasizing that:

1. 25% of points (about 534 observations) fall below the bottom edge of the box, which is the first quartile of 35.96°F. In other words 25% of observations were colder than 35.96°F.
2. 25% of points fall between the bottom edge of the box and the solid middle line, which is the median of 44.96°F. In other words 25% of observations were between 35.96 and 44.96°F and 50% of observations were colder than 44.96°F.
3. 25% of points fall between the solid middle line and the top edge of the box, which is the third quartile of 51.98°F. In other words 25% of observations were between 44.96 and 51.98°F and 75% of observations were colder than 51.98°F.

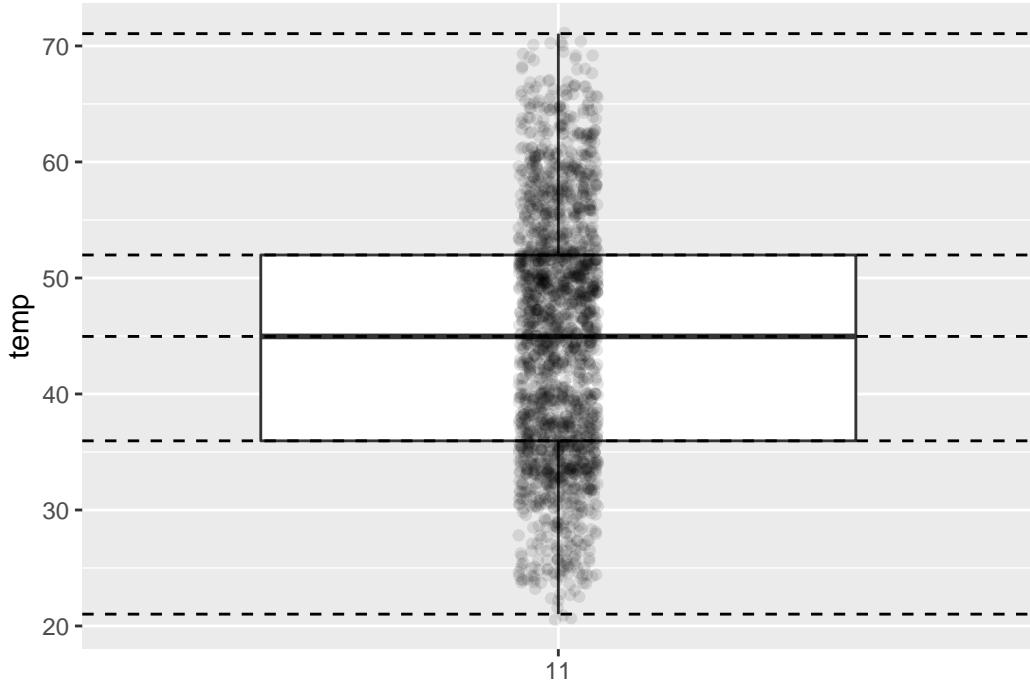


Figure 2.20: November temperatures

4. 25% of points fall over the top edge of the box. In other words 25% of observations were warmer than  $51.98^{\circ}\text{F}$ .
5. The middle 50% of points lie within the *interquartile range* between the first and third quartile of  $51.98 - 35.96 = 16.02^{\circ}\text{F}$ .

Lastly, for clarity's sake let's remove the points but keep the dashed horizontal lines in Figure 2.21.

We can now better see the *whiskers* of the boxplot. They stick out from either end of the box all the way to the minimum and maximum observed temperatures of  $21.02^{\circ}\text{F}$  and  $71.06^{\circ}\text{F}$  respectively. However, the whiskers don't always extend to the smallest and largest observed values. They in fact can extend no more than  $1.5 \times$  the interquartile range from either end of the box, in this case  $1.5 \times 16.02^{\circ}\text{F} = 24.03^{\circ}\text{F}$  from either end of the box. Any observed values outside this whiskers get marked with points called *outliers*, which we'll see in the next section.

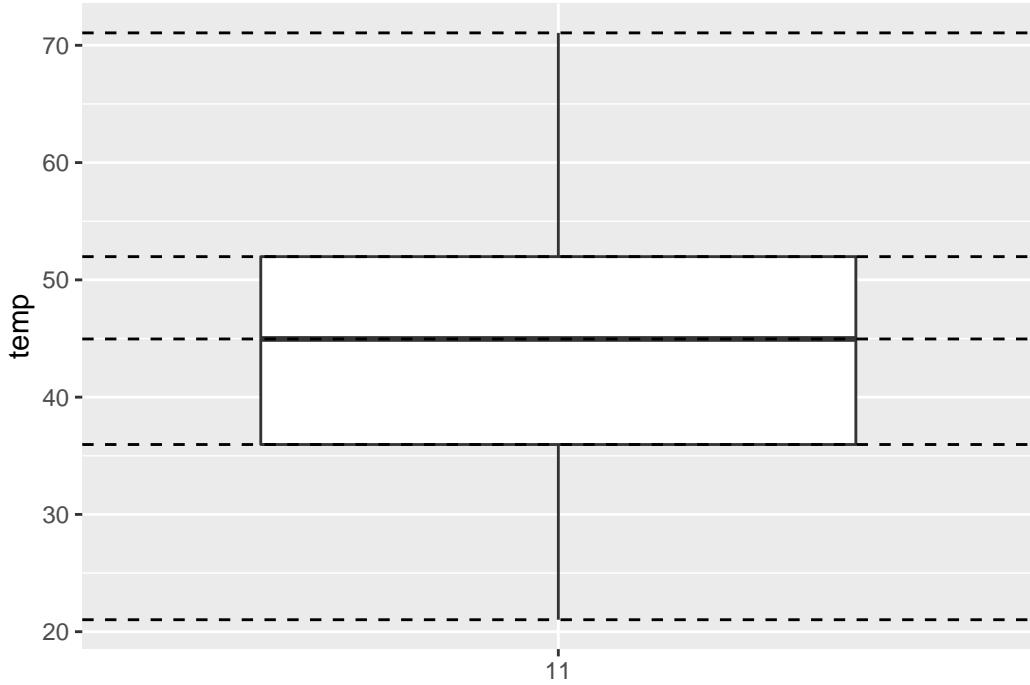


Figure 2.21: November temperatures

### 2.7.1 Boxplots via `geom_boxplot`

Let's now create a side-by-side boxplot of hourly temperatures split by the 12 months as we did above with the faceted histograms. We do this by mapping the `month` variable to the x-position aesthetic, the `temp` variable to the y-position aesthetic, and by adding a `geom_boxplot()` layer:

```
ggplot(data = weather, mapping = aes(x = month, y = temp)) +  
  geom_boxplot()
```

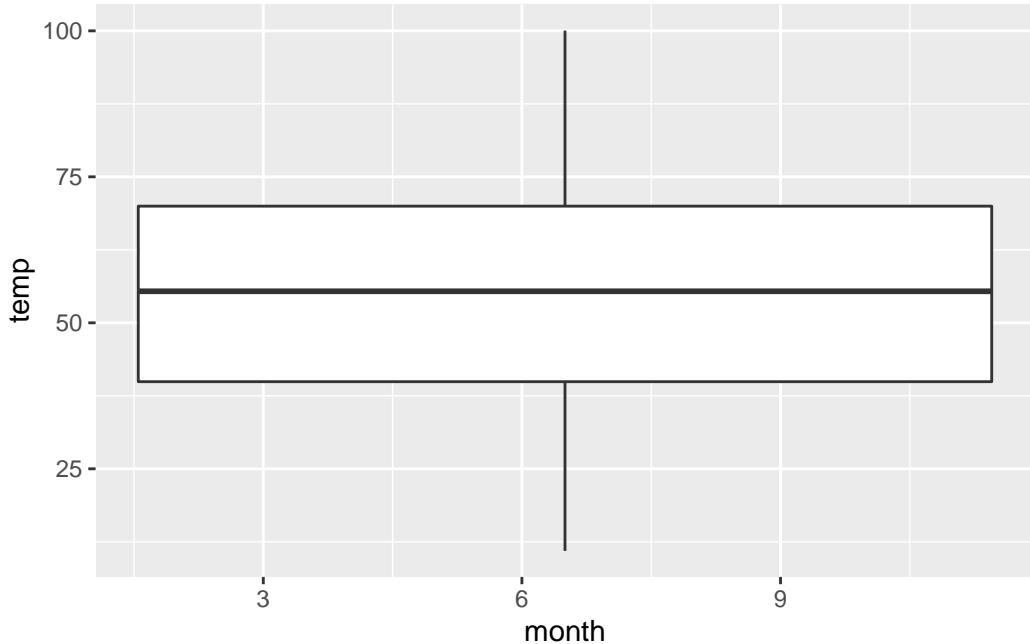


Figure 2.22: Invalid boxplot specification

**Warning messages:**

```
1: Continuous x aesthetic -- did you forget aes(group=...)?  
2: Removed 1 rows containing non-finite values (stat_boxplot).
```

Observe in Figure 2.22 that this plot does not provide information about temperature separated by month. The warning messages clue us in as to why. The second warning message is identical to the warning message when plotting a histogram of hourly temperatures: that one of the values was recorded as NA missing. However, the first warning message is telling us that we have a “continuous”, or numerical variable, on the x-position aesthetic. Boxplots however require a categorical variable on the x-axis.

We can convert the numerical variable `month` into a categorical variable by using the `factor()` function. So after applying `factor(month)`, `month` goes from having numerical values 1, 2, ..., 12 to having labels “1”, “2”, ..., “12.”

```
ggplot(data = weather, mapping = aes(x = factor(month), y = temp)) +  
  geom_boxplot()
```

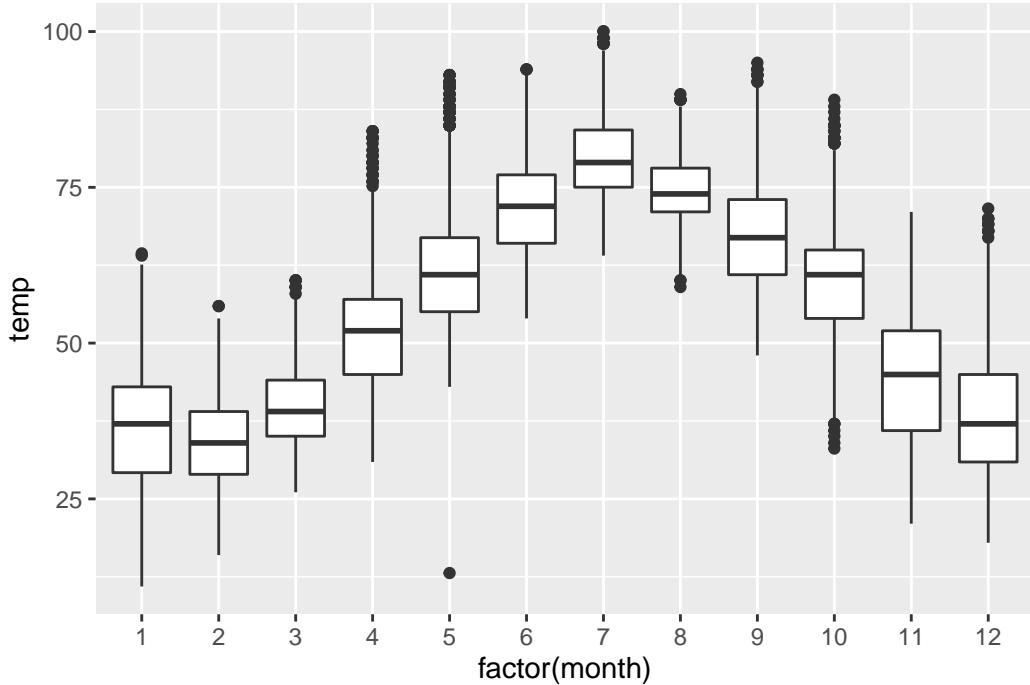


Figure 2.23: Temp by month boxplot

The resulting Figure 2.23 shows 12 separate “box and whiskers” plots with the features we saw earlier focusing only on November:

- The “box” portions of this visualization represent the 1<sup>st</sup> quartile, the median AKA the 2<sup>nd</sup> quartile, and the 3<sup>rd</sup> quartile.
- The “length” of each box, i.e. the value of the 3<sup>rd</sup> quartile minus the value of the 1<sup>st</sup> quartile, is the *interquartile range*. It is a measure of spread of the middle 50% of values, with longer boxes indicating more variability.
- The “whisker” portions of these plots extend out from the bottoms and tops of the boxes and represent points less than the 25<sup>th</sup> percentile and greater than the 75<sup>th</sup> percentiles respectively. They’re set to extend out no more than  $1.5 \times IQR$  units away from either end of the boxes. We say “no more than” because the ends of the whiskers have to correspond to observed temperatures. The length of these whiskers show how the data outside the middle 50% of values vary, with longer whiskers indicating more variability.
- The dots representing values falling outside the whiskers are called *outliers*. These can be thought of as anomalous values.

It is important to keep in mind that the definition of an outlier is somewhat arbitrary and not absolute. In this case, they are defined by the length of the whiskers, which are no

more than  $1.5 \times IQR$  units long. Looking at this plot we can see, as expected, that summer months (6 through 8) have higher median temperatures as evidenced by the higher solid lines in the middle of the boxes. We can easily compare temperatures across months by drawing imaginary horizontal lines across the plot. Furthermore, the height of the 12 boxes as quantified by the interquartile ranges are informative too; they tell us about variability, or spread, of temperatures recorded in a given month.

#### Learning Check 2.22

What does the dot at the bottom of the plot for May correspond to? Explain what might have occurred in May to produce this point.

#### Learning Check 2.23

Which months have the highest variability in temperature? What reasons can you give for this?

#### Learning Check 2.24

We looked at the distribution of the numerical variable `temp` split by the numerical variable `month` that we converted to a categorical variable using the `factor()` function. Why would a boxplot of `temp` split by the numerical variable `pressure` similarly converted to a categorical variable using the `factor()` not be informative?

#### Learning Check 2.25

Boxplots provide a simple way to identify outliers. Why may outliers be easier to identify when looking at a boxplot instead of a faceted histogram?

### 2.7.2 Summary

Side-by-side boxplots provide us with a way to compare and contrast the distribution of a quantitative variable across multiple levels of another categorical variable. One can see where the median falls across the different groups by looking at the center line in the boxes. To see how spread out the variable is across the different groups, look at both the width of the box and also how far the whiskers stretch out away from the box. Outliers are even more easily identified when looking at a boxplot than when looking at a histogram as they are marked with points.

## 2.8 5NG#5: Barplots

Both histograms and boxplots are tools to visualize the distribution of numerical variables. Another common task is visualize the distribution of a categorical variable. This is a simpler task, as we are simply counting different categories, also known as *levels*, of a categorical variable. Often the best way to visualize these different counts, also known as *frequencies*, is with a barplot (also known as a barchart). One complication, however, is how your data is represented: is the categorical variable of interest “pre-counted” or not? For example, run the following code that manually creates two data frames representing a collection of fruit: 3 apples and 2 oranges.

```
fruits <- tibble(  
  fruit = c("apple", "apple", "orange", "apple", "orange")  
)  
  
fruits_counted <- tibble(  
  fruit = c("apple", "orange"),  
  number = c(3, 2)  
)
```

We see both the `fruits` and `fruits_counted` data frames represent the same collection of fruit. Whereas `fruits` just lists the fruit individually...

```
# A tibble: 5 x 1  
fruit  
<chr>  
1 apple  
2 apple  
3 orange  
4 apple  
5 orange  
  
... fruits_counted has a variable number which represents pre-counted values of each fruit.  
  
# A tibble: 2 x 2  
fruit  number  
<chr>   <dbl>  
1 apple      3  
2 orange     2
```

Depending on how your categorical data is represented, you’ll need to use add a different `geom` layer to your `ggplot()` to create a barplot, as we now explore.

### 2.8.1 Barplots via `geom_bar` or `geom_col`

Let's generate barplots using these two different representations of the same basket of fruit: 3 apples and 2 oranges. Using the `fruits` data frame where all 5 fruits are listed individually in 5 rows, we map the `fruit` variable to the x-position aesthetic and add a `geom_bar()` layer.

```
ggplot(data = fruits, mapping = aes(x = fruit)) +  
  geom_bar()
```

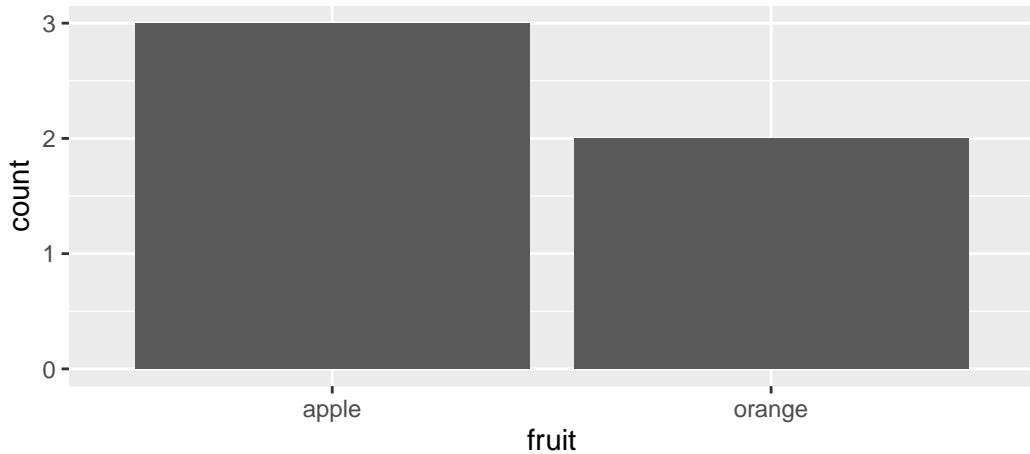


Figure 2.24: Barplot when counts are not pre-counted

However, using the `fruits_counted` data frame where the fruit have been “pre-counted”, we map the `fruit` variable to the x-position aesthetic as with `geom_bar()`, but we also map the `count` variable to the y-position aesthetic, and add a `geom_col()` layer.

```
ggplot(data = fruits_counted, mapping = aes(x = fruit, y = number)) +  
  geom_col()
```

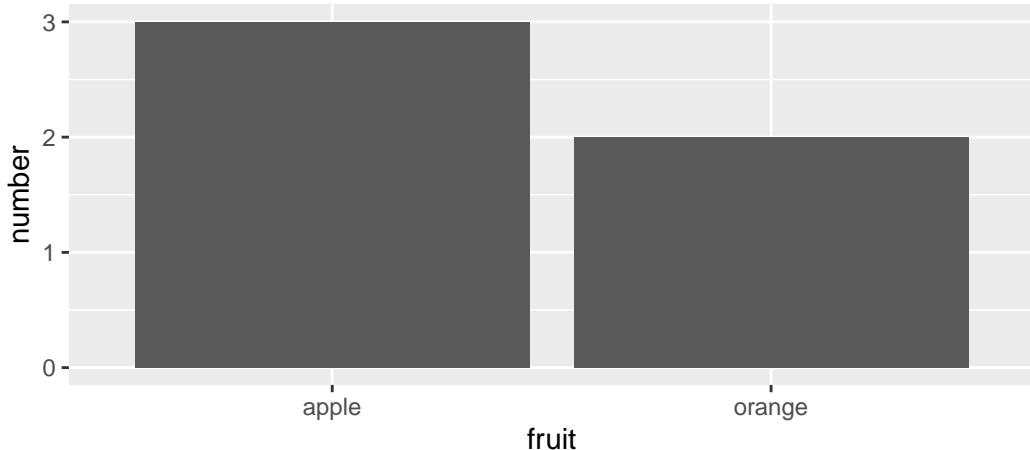


Figure 2.25: Barplot when counts are pre-counted

Compare the barplots in Figure 2.24 and Figure 2.25. They are identical because they reflect count of the same 5 fruit. However depending on how our data is saved, either pre-counted or not, we must add a different `geom` layer. When the categorical variable whose distribution you want to visualize is:

- Is not pre-counted in your data frame: use `geom_bar()`.
- Is pre-counted in your data frame, use `geom_col()` with the y-position aesthetic mapped to the variable that has the counts.

Let's now go back to the `flights` data frame in the `nycflights13` package and visualize the distribution of the categorical variable `carrier`. In other words, let's visualize the number of domestic flights out of the three New York City airports each airline company flew in 2013. Recall from Section 1.4.3 when you first explored the `flights` data frame you saw that each row corresponds to a flight. In other words the `flights` data frame is more like the `fruits` data frame than the `fruits_counted` data frame above, and thus we should use `geom_bar()` instead of `geom_col()` to create a barplot. Much like a `geom_histogram()`, there is only one variable in the `aes()` aesthetic mapping: the variable `carrier` gets mapped to the x-position.

```
ggplot(data = flights, mapping = aes(x = carrier)) +
  geom_bar()
```

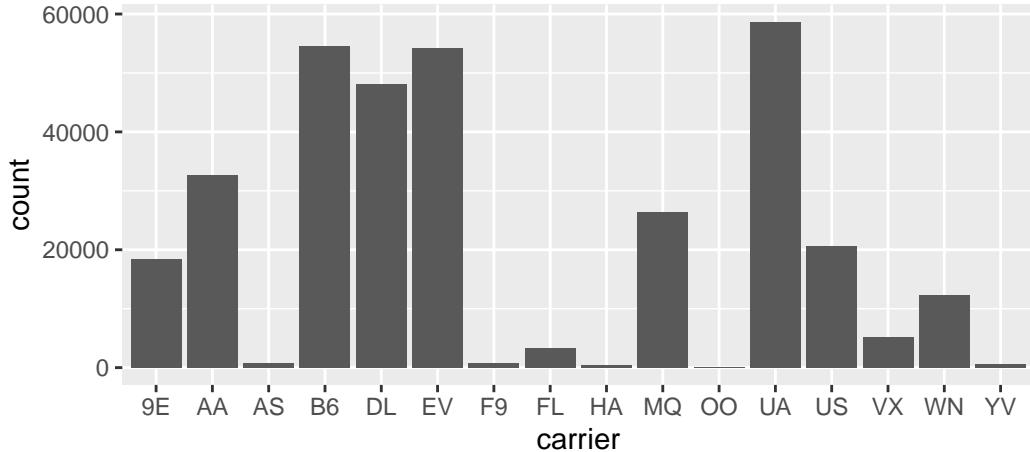


Figure 2.26: Number of flights departing NYC in 2013 by airline using `geom_bar()`

Observe in Figure 2.26 that United Air Lines (UA), JetBlue Airways (B6), and ExpressJet Airlines (EV) had the most flights depart New York City in 2013. If you don't know which airlines correspond to which carrier codes, then run `View(airlines)` to see a directory of airlines. For example: AA is American Airlines; B6 is JetBlue Airways; DL is Delta Airlines; EV is ExpressJet Airlines; MQ is Envoy Air; while UA is United Airlines.

Alternatively, say you had a data frame `flights_counted` where the number of flights for each `carrier` was pre-counted like in Table 2.3.

Table 2.3: Number of flights pre-counted for each carrier

carrier	number
UA	58665
B6	54635
EV	54173
DL	48110
AA	32729
MQ	26397
US	20536
9E	18460
WN	12275
VX	5162
FL	3260
AS	714
F9	685
YV	601

carrier	number
HA	342
OO	32

In order to create a barplot visualizing the distribution of the categorical variable `carrier` in this case, we would use `geom_col()` instead with `x` mapped to `carrier` and `y` mapped to `number` as seen below. The resulting barplot would be identical to Figure 2.26.

```
ggplot(data = flights_table, mapping = aes(x = carrier, y = number)) +
  geom_col()
```

#### Learning Check 2.26

Why are histograms inappropriate for visualizing categorical variables?

#### Learning Check 2.27

What is the difference between histograms and barplots?

#### Learning Check 2.28

How many Envoy Air flights departed NYC in 2013?

#### Learning Check 2.29

What was the seventh highest airline in terms of departed flights from NYC in 2013?  
How could we better present the table to get this answer quickly?

### 2.8.2 Must avoid pie charts!

Unfortunately, one of the most common plots seen today for categorical data is the pie chart. While they may seem harmless enough, they actually present a problem in that humans are unable to judge angles well. As Naomi Robbins describes in her book “Creating More Effective Graphs” (Robbins 2013), we overestimate angles greater than 90 degrees and we underestimate angles less than 90 degrees. In other words, it is difficult for us to determine relative size of one piece of the pie compared to another.

Let’s examine the same data used in our previous barplot of the number of flights departing NYC by airline in Figure 2.26, but this time we will use a pie chart in Figure 2.27.

Try to answer the following questions:

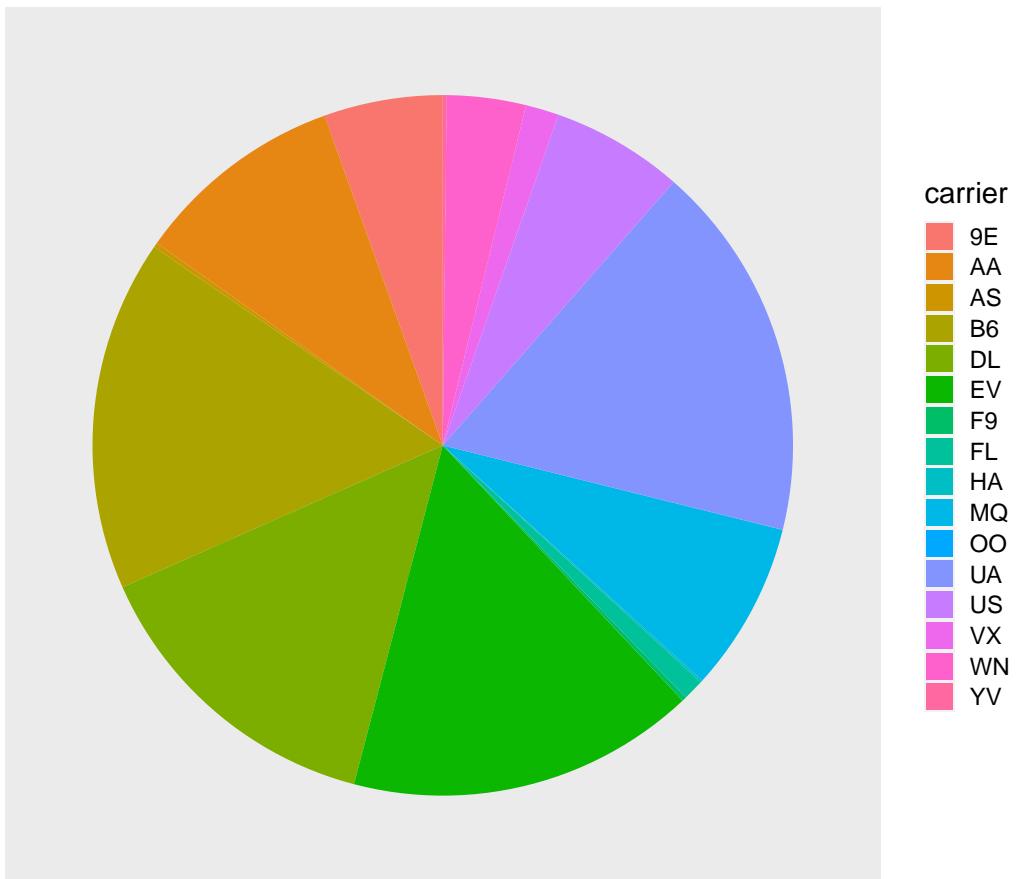
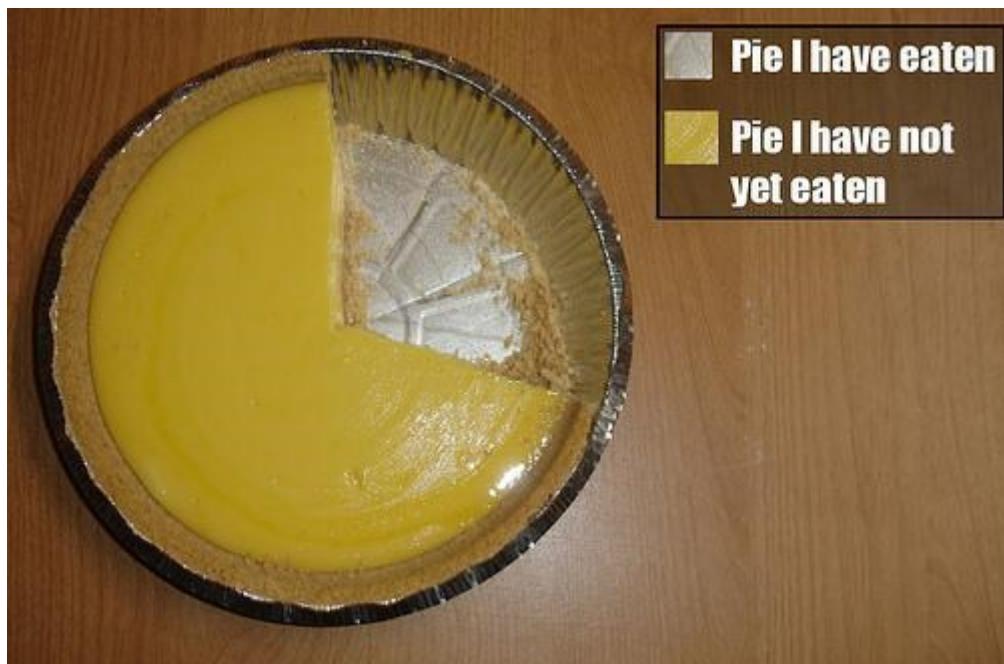


Figure 2.27: The dreaded pie chart

- How much larger the portion of the pie is for ExpressJet Airlines (EV) compared to US Airways (US),
- What the third largest carrier is in terms of departing flights, and
- How many carriers have fewer flights than United Airlines (UA)?

While it is quite difficult to answer these questions when looking at the pie chart in Figure 2.27, we can much more easily answer these questions using the barchart in Figure 2.26. This is true since barplots present the information in a way such that comparisons between categories can be made with single horizontal lines, whereas pie charts present the information in a way such that comparisons between categories must be made by comparing angles.

There may be one exception of a pie chart not to avoid courtesy Nathan Yau at [Flowing-Data.com](#), but we will leave this for the reader to decide:



#### Learning Check 2.30

Why should pie charts be avoided and replaced by barplots?

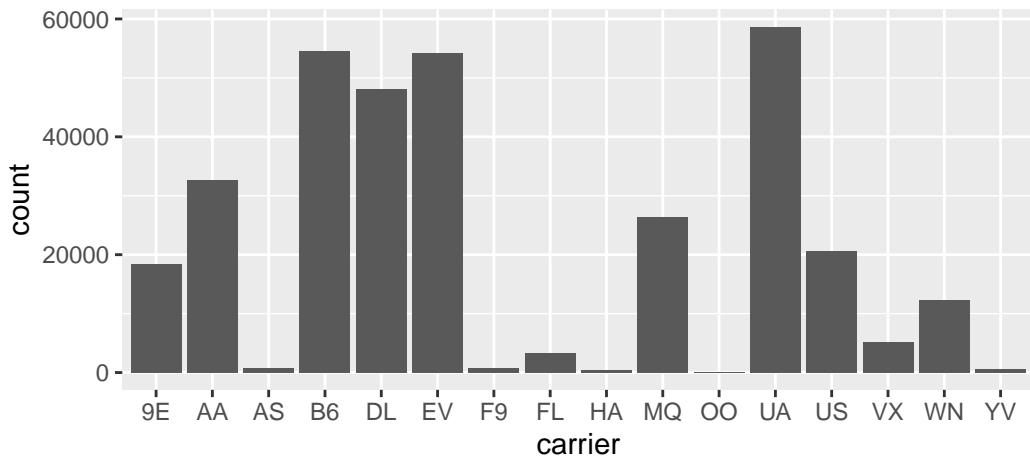
#### Learning Check 2.31

Why do you think people continue to use pie charts?

### 2.8.3 Two categorical variables

Barplots are the go-to way to visualize the frequency of different categories, or levels, of a single categorical variable. Another use of barplots is to visualize the *joint* distribution of two categorical variables at the same time. Let's examine the *joint* distribution of outgoing domestic flights from NYC by `carrier` and `origin`, or in other words the number of flights for each `carrier` and `origin` combination. For example, the number of WestJet flights from JFK, the number of WestJet flights from LGA, the number of WestJet flights from EWR, the number of American Airlines flights from JFK, and so on. Recall the `ggplot()` code that created the barplot of `carrier` frequency in Figure 2.26:

```
ggplot(data = flights, mapping = aes(x = carrier)) +  
  geom_bar()
```



We can now map the additional variable `origin` by adding a `fill = origin` inside the `aes()` aesthetic mapping; the `fill` aesthetic of any bar corresponds to the color used to fill the bars.

```
ggplot(data = flights, mapping = aes(x = carrier, fill = origin)) +  
  geom_bar()
```

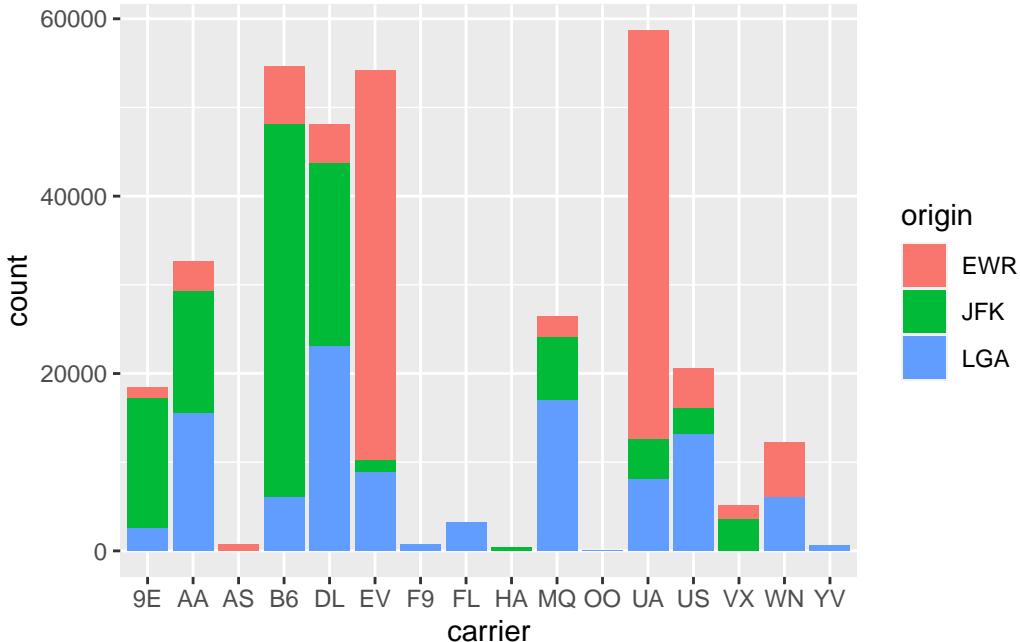


Figure 2.28: Stacked barplot comparing the number of flights by carrier and origin

Figure 2.28 is an example of a *stacked barplot*. While simple to make, in certain aspects it is not ideal. For example, it is difficult to compare the heights of the different colors between the bars, corresponding to comparing the number of flights from each `origin` airport between the carriers.

Before we continue, let's address some common points of confusion amongst new R users. First, note that `fill` is another aesthetic mapping much like `x-position`; thus it must be included within the parentheses of the `aes()` mapping. The following code, where the `fill` aesthetic is specified outside the `aes()` mapping will yield an error. This is a fairly common error that new `ggplot` users make:

```
ggplot(data = flights, mapping = aes(x = carrier), fill = origin) +
  geom_bar()
```

Second, the `fill` aesthetic corresponds to the color used to fill the bars, while the `color` aesthetic corresponds to the color of the outline of the bars. Observe in Figure 2.29 that mapping `origin` to `color` and not `fill` yields grey bars with different colored outlines.

```
ggplot(data = flights, mapping = aes(x = carrier, color = origin)) +
  geom_bar()
```

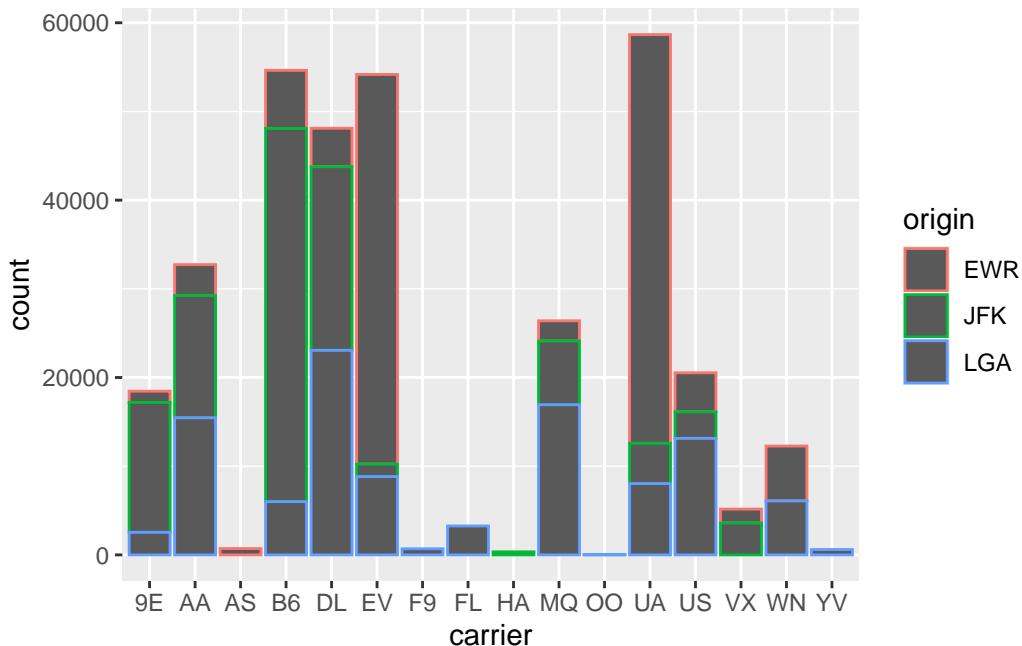


Figure 2.29: Stacked barplot with color aesthetic used instead of fill

#### Learning Check 2.32

What kinds of questions are not easily answered by looking at the above figure?

#### Learning Check 2.33

What can you say, if anything, about the relationship between airline and airport in NYC in 2013 in regards to the number of departing flights?

Another alternative to stacked barplots are *side-by-side barplots*, also known as a *dodged barplot*. The code to created a side-by-side barplot is identical to the code to create a stacked barplot, but with a `position = "dodge"` argument added to `geom_bar()`. In other words, we are overriding the default barplot type, which is a stacked barplot, and specifying it to be a side-by-side barplot.

```
ggplot(data = flights, mapping = aes(x = carrier, fill = origin)) +
  geom_bar(position = "dodge")
```

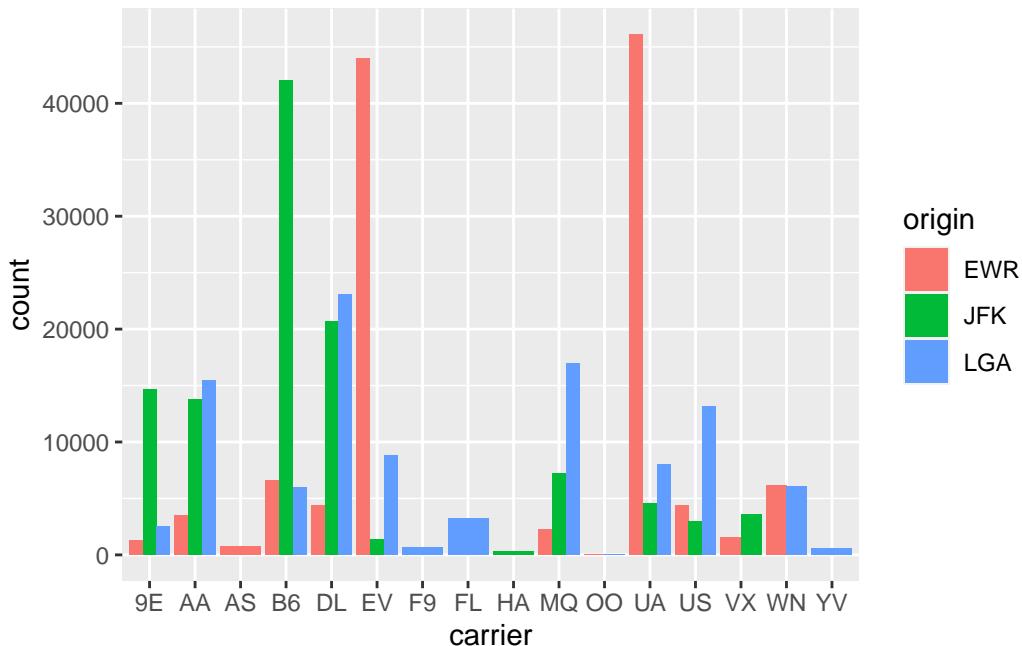


Figure 2.30: Side-by-side AKA dodged barplot comparing the number of flights by carrier and origin

#### Learning Check 2.34

Why might the side-by-side (AKA dodged) barplot be preferable to a stacked barplot in this case?

#### Learning Check 2.35

What are the disadvantages of using a side-by-side (AKA dodged) barplot, in general?

Lastly, another type of barplot is a *faceted barplot*. Recall in Section 2.6 we visualized the distribution of hourly temperatures at the 3 NYC airports *split* by month using facets. We apply the same principle to our barplot visualizing the frequency of `carrier` split by `origin`: instead of mapping `origin`

```
ggplot(data = flights, mapping = aes(x = carrier)) +
  geom_bar() +
  facet_wrap(~ origin, ncol = 1)
```

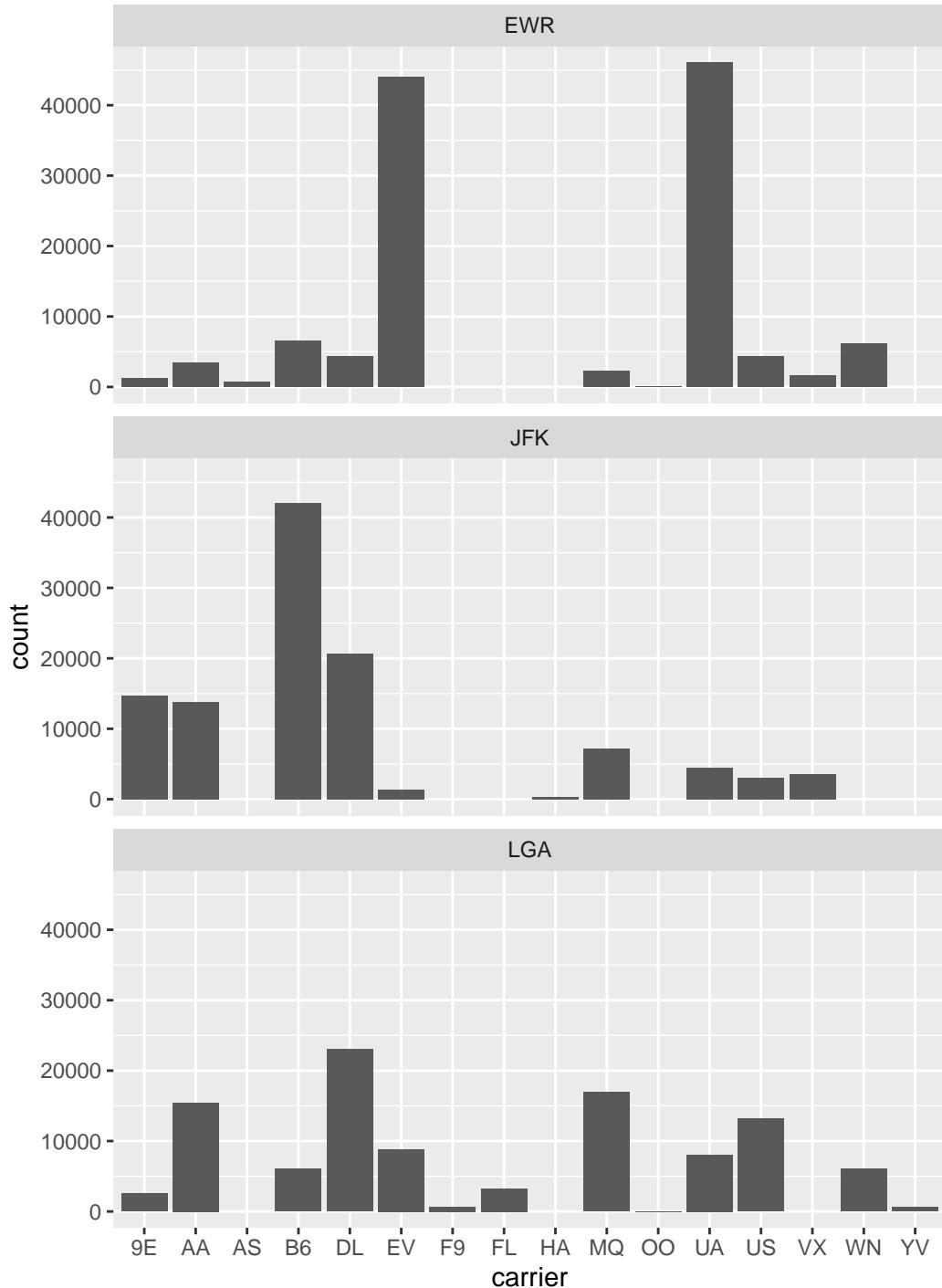


Figure 2.31: Faceted barplot comparing the number of flights by carrier and origin

#### Learning Check 2.36

Why is the faceted barplot preferred to the side-by-side and stacked barplots in this case?

#### Learning Check 2.37

What information about the different carriers at different airports is more easily seen in the faceted barplot?

### 2.8.4 Summary

Barplots are the preferred way of displaying the distribution of a categorical variable, or in other words the frequency with which the different categories called *levels* occur. They are easy to understand and make it easy to make comparisons across levels. When trying to visualize two categorical variables, you have many options: stacked barplots, side-by-side barplots, and faceted barplots. Depending on what aspect of the joint distribution you are trying to emphasize, you will need to make a choice between these three types of barplots.

## 2.9 Conclusion

### 2.9.1 Summary table

Let's recap all five of the Five Named Graphs (5NG) in Table 2.4 summarizing their differences. Using these 5NG, you'll be able to visualize the distributions and relationships of variables contained in a wide array of datasets. This will be even more the case as we start to map more variables to more of each geometric object's aesthetic attribute options, further unlocking the awesome power of the `ggplot2` package.

Table 2.4: Summary of 5NG

Named graph Shows	Geometric object	Notes
Scatter Relationship between 2 numerical variables	<code>geom_point()</code>	
Linegraph Relationship between 2 numerical variables	<code>geom_line()</code>	Used when there is a sequential order to x-variable e.g. time
Histogram Distribution of 1 numerical variable	<code>geom_histogram()</code>	Facetted histograms show the distribution of 1 numerical variable split by the values of another variable

Named graph Shows	Geometric object	Notes
Boxplot Distribution of 1 numerical variable split by the values of another variable	<code>geom_boxplot()</code>	
Barplot Distribution of 1 categorical variable	<code>geom_bar()</code> when counts are not pre-counted, <code>geom_col()</code> when counts are pre-counted	Stacked, side-by-side, and faceted barplots show the joint distribution of 2 categorical variables

## 2.9.2 Argument specification

Run the following two segments of code. First this:

```
ggplot(data = flights, mapping = aes(x = carrier)) +
  geom_bar()
```

then this:

```
ggplot(flights, aes(x = carrier)) +
  geom_bar()
```

You'll notice that both code segments create the same barplot, even though in the second segment we omitted the `data =` and `mapping =` code argument names. This is because the `ggplot()` by default assumes that the `data` argument comes first and the `mapping` argument comes second. So as long as you specify the data frame in question first and the `aes()` mapping second, you can omit the explicit statement of the argument names `data =` and `mapping =`.

Going forward for the rest of this book, all `ggplot()` will be like the second segment above: with the `data =` and `mapping =` explicit naming of the argument omitted and the default ordering of arguments respected.

## 2.9.3 Additional resources

If you want to further unlock the power of the `ggplot2` package for data visualization, we suggest you that you check out RStudio's "Data Visualization with `ggplot2`" cheatsheet. This cheatsheet summarizes much more than what we've discussed in this chapter, in particular the many more than the 5 `geom` geometric objects we covered in this Chapter, while providing quick and easy to read visual descriptions.

You can access this cheatsheet by going to the RStudio Menu Bar -> Help -> Cheatsheets -> “Data Visualization with ggplot2”:

## Data Visualization with ggplot2 :: CHEAT SHEET

### Basics

`ggplot2` is based on the grammar of graphics, the idea that you can build every graph from the same components: a **data set**, a **coordinate system**, and **geoms**—visual marks that represent data points.



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like `size`, `color`, and `x` and `y` locations.



Complete the template below to build a graph.

```
ggplot(data = DATA) +
  <GEOM_FUNCTION>(mapping = aes(MAPPINGS),
  stat = <STAT>, position = <POSITION>) +
  <COORDINATE_FUNCTION>+
  <FACET_FUNCTION>+
  <SCALE_FUNCTION>+
  <THEME_FUNCTION>
```

`ggplot(data = mpg, aes(x = cyl, y = hwy))` Begins a plot that you finish by adding layers to. Add one geom function per layer.

`aesthetic mappings` `data` `geom`  
`plot(x = cyl, y = hwy, data = mpg, geom = "point")` Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

`last_plot()` Returns the last plot

`ggsave("plot.png", width = 5, height = 5)` Saves last plot as 5" x 5" file named "plot.png" in working directory. Matches file type to file extension.



### Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

#### GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))

a + geom_blank()
# (Useful for expanding limits)

b + geom_curve(aes(yend = lat + 1,
  xend = long + 1, curvature = 2))
  , x, yend, y, yend,
  alpha, angle, color, curvature, linetype, size
  , lineheight, size, vjust

a + geom_path(linewidth = "butt", linejoin = "round",
  linemiter = 1)
  , x, y, alpha, color, group, linetype, size

a + geom_polygon(aes(group = group))
  , x, y, alpha, color, fill, group, linetype, size

b + geom_rect(aes(min = long, ymin = lat, xmax =
  long + 1, ymax = lat + 1))
  , xmax, xmin, ymax,
  ymin, alpha, color, fill, linetype, size

a + geom_ribbon(aes(ymin = unemploy - 900,
  alpha, color, fill, group, linetype, size
  , alpha = unemploy + 900))
  , x, y, alpha, color, fill, group, linetype, size
```

#### LINE SEGMENTS

common aesthetics: `x`, `y`, `alpha`, `color`, `linetype`, `size`

```
b + geom_abline(aes(intercept = 0, slope = 1))
  , b + geom_hline(aes(intercept = lat))
  , b + geom_vline(aes(intercept = long))

b + geom_segment(aes(yend = lat + 1, xend = long + 1))
  , b + geom_spoke(aes(angle = 1155, radius = 1))
```

#### ONE VARIABLE continuous

```
c <- ggplot(mpg, aes(hwy)); c <- ggplot(mpg)

c + geom_area(stat = "bin")
  , x, y, alpha, color, fill, linetype, size
  , c + geom_density(kernel = "gaussian")
  , x, y, alpha, color, fill, group, linetype, size, weight

c + geom_dotplot()
  , x, y, alpha, color, fill, group, linetype, size
  , c + geom_freqpoly()
  , x, y, alpha, color, group,
  linetype, size
  , c + geom_histogram(binwidth = 5)
  , x, y, alpha, color, fill, linetype, size, weight
  , c2 + geom_aq(aes(sample = hwy))
```

#### discrete

```
d <- ggplot(mpg, aes(f1))
  , d + geom_bar()
  , x, alpha, color, fill, linetype, size, weight
```

#### TWO VARIABLES

```
continuous x, continuous y
e <- ggplot(mpg, aes(cty, hwy))

e + geom_label(aes(label = cty), nudge_x = 1,
  nudge_y = 1, check_overlap = TRUE)
  , x, y, label, alpha, angle, color, family, fontface, hjust,
  lineheight, size, vjust

e + geom_jitter(height = 2, width = 2)
  , x, y, alpha, color, fill, shape, size, stroke

e + geom_point(), x, y, alpha, color, fill, shape,
  linetype, size, weight

e + geom_quantile(), x, y, alpha, color, group,
  linetype, size, weight

e + geom_rug(sides = "bl")
  , x, y, alpha, color, group, linetype, size, weight

e + geom_smooth(method = lm)
  , x, y, alpha, color, fill, group, linetype, size, weight

e + geom_text(aes(label = cyl), nudge_x = 1,
  nudge_y = 1, check_overlap = TRUE)
  , x, y, label, alpha, angle, color, family, fontface, hjust,
  lineheight, size, vjust
```

#### discrete x, continuous y

```
f <- ggplot(mpg, aes(class, hwy))

f + geom_col()
  , x, y, alpha, color, fill, group, linetype, size, weight
  , f + geom_boxplot()
  , x, y, lower, middle, upper,
  ymax, ymin, alpha, color, fill, group, linetype, size, weight
  , f + geom_dotplot(binaxis = "y", stackdir =
  "center")
  , x, y, alpha, color, fill, group
  , f + geom_violin(scale = "area")
  , x, y, alpha, color, fill, group, linetype, size, weight
```

#### discrete x, discrete y

```
g <- ggplot(diamonds, aes(cut, color))
  , g + geom_count()
  , x, y, alpha, color, fill, shape, size, stroke
```

#### THREE VARIABLES

```
seals$z <- with(seals, sqrt(delta_long^2 + delta_lat^2))
l <- ggplot(seals, aes(long, lat))

l + geom_contour(aes(z = z))
  , x, y, z, alpha, color, group, linetype, size, weight
  , l + geom_raster(aes(fill = z))
  , x, y, label, alpha, fill, interpolates = FALSE
  , k + geom_map(aes(map_id = state), map = map)
  , expand_limits(x = -120, y = 30)
  , map_id, alpha, color, fill, group, linetype, size
```

#### continuous bivariate distribution

```
h <- ggplot(diamonds, aes(carat, price))

h + geom_bin2d(binwidth = c(0.25, 500))
  , x, y, alpha, color, fill, linetype, size, weight
  , h + geom_density2d()
  , x, y, alpha, color, group, linetype, size
  , h + geom_hex()
  , x, y, alpha, colour, fill, size
```

#### continuous function

```
i <- ggplot(economics, aes(date, unemploy))
```

```
i + geom_ribbon()
  , x, y, alpha, color, fill, linetype, size
```

```
i + geom_line()
  , x, y, alpha, color, group, linetype, size
```

```
i + geom_step(direction = "hv")
  , x, y, alpha, color, group, linetype, size
```

**visualizing error**

```
j <- data.frame(grp = c("A", "B"), fit = 4.5, se = 1.2)
```

```
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
  , j + geom_crossbar(fatten = 2)
  , x, y, max, ymn, alpha, color, fill, group, linetype, size
```

```
j + geom_errorbar()
  , x, y, max, ymn, alpha, color, fill, group, linetype, size
```

```
j + geom_linerange()
  , x, ymn, ymax, alpha, color, group, linetype, size
```

```
j + geom_pointrange()
  , x, ymn, ymax, alpha, color, fill, group, linetype, size
```

**maps**

```
data <- data.frame(murder = USArests$Murder,
  state = tolower(rownames(USArests)))
```

```
map <- map_data("state")
k <- ggplot(data, aes(fill = murder))
```

```
k + geom_map(aes(map_id = state), map = map)
  , expand_limits(x = -120, y = 30)
  , map_id, alpha, color, fill, group, linetype, size
```

```
l + geom_raster(aes(fill = z))
  , x, y, label, alpha, fill, interpolates = FALSE
  , l + geom_tile(aes(fill = z))
  , x, y, alpha, color, fill, linetype, size, width
```

### 2.9.4 What's to come

Recall in Figure 2.2 in Section 2.3 we visualized the relationship between departure delay and arrival delay for Alaska Airlines flights. This necessitated paring or filtering down the `flights` data frame to a new data frame `alaska_flights` consisting of only `carrier == AS` flights first:

```
alaska_flights <- flights %>%
  filter(carrier == "AS")

ggplot(data = alaska_flights, mapping = aes(x = dep_delay, y = arr_delay)) +
  geom_point()
```

Furthermore recall in Figure 2.8 in Section 2.4 we visualized hourly temperature recordings at Newark airport only for the first 15 days of January 2013. This necessitated paring or filtering down the `weather` data frame to a new data frame `early_january_weather` consisting of hourly temperature recordings only for `origin == "EWR"`, `month == 1`, and day less than or equal to 15 first:

```
early_january_weather <- weather %>%
  filter(origin == "EWR" & month == 1 & day <= 15)

ggplot(data = early_january_weather, mapping = aes(x = time_hour, y = temp)) +
  geom_line()
```

These two code segments were a preview of Chapter 3 on data wrangling where we'll delve further into the `dplyr` package. Data wrangling is the process of transforming and modifying existing data with the intent of making it more appropriate for analysis purposes. For example, the two code segments used the `filter()` function to create new data frames (`alaska_flights` and `early_january_weather`) by choosing only a subset of rows of existing data frames (`flights` and `weather`). In this next chapter, we'll formally introduce the `filter()` and other data wrangling functions as well as the *pipe operator* `%>%` which allows you to combine multiple data wrangling actions into a single sequential *chain* of actions. On to Chapter 3 on data wrangling!

## 3 Data Wrangling



Currently working on content transfer from previous version of the book.

## 4 Data Importing & “Tidy Data”



Currently working on content transfer from previous version of the book.

# **Part III**

# **Data Modeling**

# 5 Basic Regression



Currently working on content transfer from previous version of the book.

# 6 Multiple Regression



Currently working on content transfer from previous version of the book.

# **Part IV**

# **Statistical Theory**

# 7 Randomization and Causality



Currently working on content transfer from previous version of the book.

## 8 Populations and Generalizability



Currently working on content transfer from previous version of the book.

# 9 Sampling Distributions



Currently working on content transfer from previous version of the book.

## **Part V**

# **Statistical Inference**

# 10 Confidence Intervals



Currently working on content transfer from previous version of the book.

# 11 P-values



Under Construction

Currently working on content transfer from previous version of the book.

# 12 Hypothesis tests



Under Construction

Currently working on content transfer from previous version of the book.

## 13 Putting it all together



Currently working on content transfer from previous version of the book.

# References

- Grolemund, Garrett, and Hadley Wickham. 2016. *R for Data Science*. <http://r4ds.had.co.nz/>.
- Ismay, Chester. 2016. *Getting Used to r, RStudio, and r Markdown*. <http://ismayc.github.io/rbasics-book>.
- Robbins, Naomi. 2013. *Creating More Effective Graphs*. Chart House.
- Wickham, Hadley. 2021. *Nycflights13: Flights That Departed NYC in 2013*. <https://github.com/hadley/nycflights13>.
- Wickham, Hadley, Winston Chang, Lionel Henry, Thomas Lin Pedersen, Kohske Takahashi, Claus Wilke, Kara Woo, Hiroaki Yutani, and Dewey Dunnington. 2022. *Ggplot2: Create Elegant Data Visualisations Using the Grammar of Graphics*. <https://CRAN.R-project.org/package=ggplot2>.
- Wilkinson, Leland. 2005. *The Grammar of Graphics (Statistics and Computing)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc.

# A Statistical Background



Currently working on content transfer from previous version of the book.