

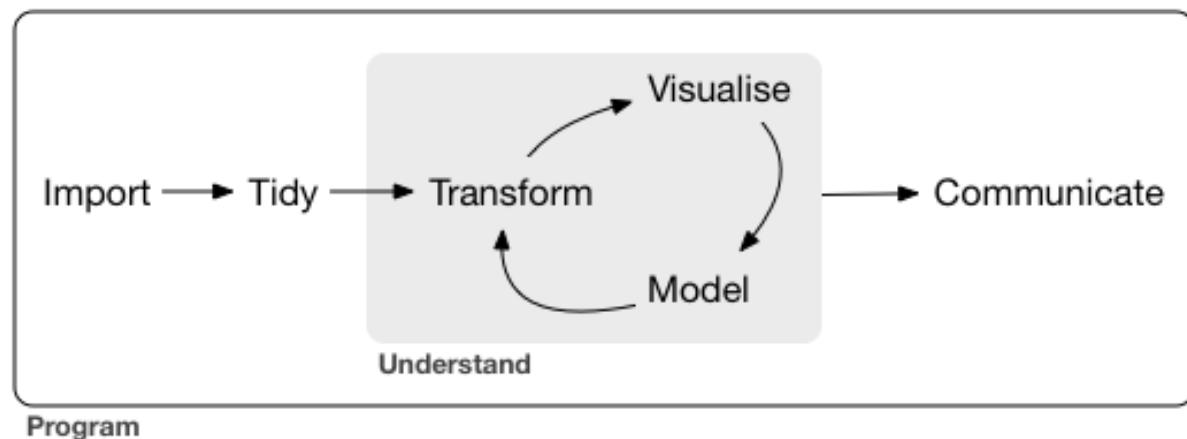
R for Data Science

1 Introduction

Data science is an exciting discipline that allows you to turn raw data into understanding, insight, and knowledge. The goal of “R for Data Science” is to help you learn the most important tools in R that will allow you to do data science. After reading this book, you’ll have the tools to tackle a wide variety of data science challenges, using the best parts of R.

1.1 What you will learn

Data science is a **huge field, and there's no way you can master it by reading a single book**. The goal of this book is to give you a solid foundation in the most important tools. Our model of the tools needed in a typical data science project looks something like this:



First you must import your data into R. This typically means that you take data stored in a file, database, or web API, and load it into a data frame in R. If you can't get your data into R, you can't do data science on it!

Once you've imported your data, it is a good idea to tidy it. Tidying your data means storing it in a consistent form that matches the semantics of the dataset with the way it is stored. In brief, when your data is tidy, each column is a variable, and each row is an observation. Tidy data is important because the consistent structure lets you focus your struggle on questions about the data, not fighting to get the data into the right form for different functions.

Once you have tidy data, a common first step is to transform it. Transformation includes narrowing in on observations of interest (like all people in one city, or all data from the last year), creating new variables that are functions of existing variables (like computing velocity from speed and time), and calculating a set of summary statistics (like counts or means). Together, tidying and transforming are called wrangling, because getting your data in a form that's natural to work with often feels like a fight!

Once you have tidy data with the variables you need, there are two main engines of knowledge generation: visualisation and modelling. These have complementary strengths and weaknesses so any real analysis will iterate between them many times.

Visualisation is a fundamentally human activity. A good visualisation will show you things that you did not expect, or raise new questions about the data. A good visualisation might also hint **that you're asking the wrong question, or you need to collect different data. Visualisations can surprise you, but don't scale particularly well** because they require a human to interpret them.

Models are complementary tools to visualisation. Once you have made your questions sufficiently precise, you can use a model to answer them. Models are a fundamentally mathematical or computational tool, so they **generally scale well. Even when they don't, it's** usually cheaper to buy more computers than it is to buy more brains! But every model makes assumptions, and by its very nature a model cannot question its own assumptions. That means a model cannot fundamentally surprise you.

The last step of data science is communication, an absolutely critical part of any data analysis **project. It doesn't matter how well your models and visualisation have led you to understand the** data unless you can also communicate your results to others.

Surrounding all these tools is programming. Programming is a cross-cutting tool that you use in **every part of the project. You don't need to be an expert programmer to be a data scientist, but** learning more about programming pays off because becoming a better programmer allows you to automate common tasks, and solve new problems with greater ease.

You'll use these tools in every data science project, but for most projects they're not enough. There's a rough 80-20 rule at play: you can tackle about 80% of every project using the tools that you'll learn in this book, but you'll need other tools to tackle the remaining 20%. Throughout this book we'll point you to resources where you can learn more.

1.2 How this book is organised

The previous description of the tools of data science is organised roughly according to the order **in which you use them in an analysis (although of course you'll iterate through them multiple times).** In our experience, however, this is not the best way to learn them:

- Starting with data ingest and tidying is sub-optimal because 80% of the time it's routine and boring, and the other 20% of the time it's weird and frustrating. That's a bad place to start learning a new subject! Instead, we'll start with visualisation and transformation of data that's already been imported and tidied. That way, when you ingest and tidy your own data, your motivation will stay high because you know the pain is worth it.
- Some topics are best explained with other tools. For example, we believe that it's easier to understand how models work if you already know about visualisation, tidy data, and programming.
- Programming tools are not necessarily interesting in their own right, but do allow you to tackle considerably more challenging problems. We'll give you a selection of programming tools in the middle of the book, and then you'll see how they can combine with the data science tools to tackle interesting modelling problems.

Within each chapter, we try and stick to a similar pattern: start with some motivating examples so you can see the bigger picture, and then dive into the details. Each section of the book is paired

with exercises to help you practice what you've learned. While it's tempting to skip the exercises, there's no better way to learn than practicing on real problems.

1.3 What you won't learn

There are some important topics that this book doesn't cover. We believe it's important to stay ruthlessly focused on the essentials so you can get up and running as quickly as possible. That means **this book can't cover every important topic**.

1.3.1 Big data

This book proudly focuses on small, in-memory datasets. This is the right place to start because **you can't tackle big data unless you have experience with small data**. The tools you learn in this book will easily handle hundreds of megabytes of data, and with a little care you can typically use them to work with 1-2 Gb of data. If you're routinely working with larger data (10-100 Gb, say), you should learn more about `data.table`. This book doesn't teach `data.table` because it has a very **concise interface which makes it harder to learn since it offers fewer linguistic cues**. But if you're working with large data, the performance payoff is worth the extra effort required to learn it.

If your data is bigger than this, carefully consider if your big data problem might actually be a small data problem in disguise. While the complete data might be big, often the data needed to answer a specific question is small. You might be able to find a subset, subsample, or summary **that fits in memory and still allows you to answer the question that you're interested in**. The challenge here is finding the right small data, which often requires a lot of iteration.

Another possibility is that your big data problem is actually a large number of small data problems. Each individual problem might fit in memory, but you have millions of them. For example, you might want to fit a model to each person in your dataset. That would be trivial if you had just 10 or 100 people, but instead you have a million. Fortunately each problem is independent of the others (a setup that is sometimes called embarrassingly parallel), so you just need a system (like Hadoop or Spark) that allows you to send different datasets to different **computers for processing**. Once you've figured out how to answer the question for a single subset using the tools described in this book, you learn new tools like `sparklyr`, `rhive`, and `ddr` to solve it for the full dataset.

1.3.2 Python, Julia, and friends

In this book, you won't learn anything about Python, Julia, or any other programming language useful for data science. This isn't because we think these tools are bad. They're not! And in practice, most data science teams use a mix of languages, often at least R and Python.

However, we strongly believe that it's best to master one tool at a time. You will get better faster if you dive deep, rather than spreading yourself thinly over many topics. This doesn't mean you should only know one thing, just that you'll generally learn faster if you stick to one thing at a time. You should strive to learn new things throughout your career, but make sure your understanding is solid before you move on to the next interesting thing.

We think R is a great place to start your data science journey because it is an environment designed from the ground up to support data science. R is not just a programming language, but it is also an interactive environment for doing data science. To support interaction, R is a much

more flexible language than many of its peers. This flexibility comes with its downsides, but the big upside is how easy it is to evolve tailored grammars for specific parts of the data science process. These mini languages help you think about problems as a data scientist, while supporting fluent interaction between your brain and the computer.

1.3.3 Non-rectangular data

This book focuses exclusively on rectangular data: collections of values that are each associated with a variable and an observation. There are lots of datasets that do not naturally fit in this paradigm: including images, sounds, trees, and text. But rectangular data frames are extremely common in science and industry, and we believe that they are a great place to start your data science journey.

1.3.4 Hypothesis confirmation

It's possible to divide data analysis into two camps: hypothesis generation and hypothesis confirmation (sometimes called confirmatory analysis). The focus of this book is unabashedly on **hypothesis generation, or data exploration. Here you'll look deeply at the data and, in combination with your subject knowledge, generate many interesting hypotheses to help explain why the data behaves the way it does.** You evaluate the hypotheses informally, using your scepticism to challenge the data in multiple ways.

The complement of hypothesis generation is hypothesis confirmation. Hypothesis confirmation is hard for two reasons:

1. You need a precise mathematical model in order to generate falsifiable predictions. This often requires considerable statistical sophistication.
2. You can only use an observation once to confirm a hypothesis. As soon as you use it **more than once you're back to doing exploratory analysis. This means to do hypothesis confirmation you need to "preregister" (write out in advance) your analysis plan, and not deviate from it even when you have seen the data. We'll talk a little about some strategies you can use to make this easier in [modelling](#).**

It's common to think about modelling as a tool for hypothesis confirmation, and visualisation as a tool for hypothesis generation. But that's a false dichotomy: models are often used for exploration, and with a little care you can use visualisation for confirmation. The key difference is how often do you look at each observation: if you look only once, it's confirmation; if you look more than once, it's exploration.

1.4 Prerequisites

We've made a few assumptions about what you already know in order to get the most out of this book. You should be generally numerically literate, and it's helpful if you have some programming experience already. If you've never programmed before, you might find [Hands on Programming with R](#) by Garrett to be a useful adjunct to this book.

There are four things you need to run the code in this book: R, RStudio, a collection of R packages called the tidyverse, and a handful of other packages. Packages are the fundamental units of reproducible R code. They include reusable functions, the documentation that describes how to use them, and sample data.

1.4.1 R

To download R, go to CRAN, the comprehensive R archive network. CRAN is composed of a set of mirror servers **distributed around the world and is used to distribute R and R packages. Don't try and pick a mirror that's close to you: instead use the cloud mirror**, <https://cloud.r-project.org>, which automatically figures it out for you.

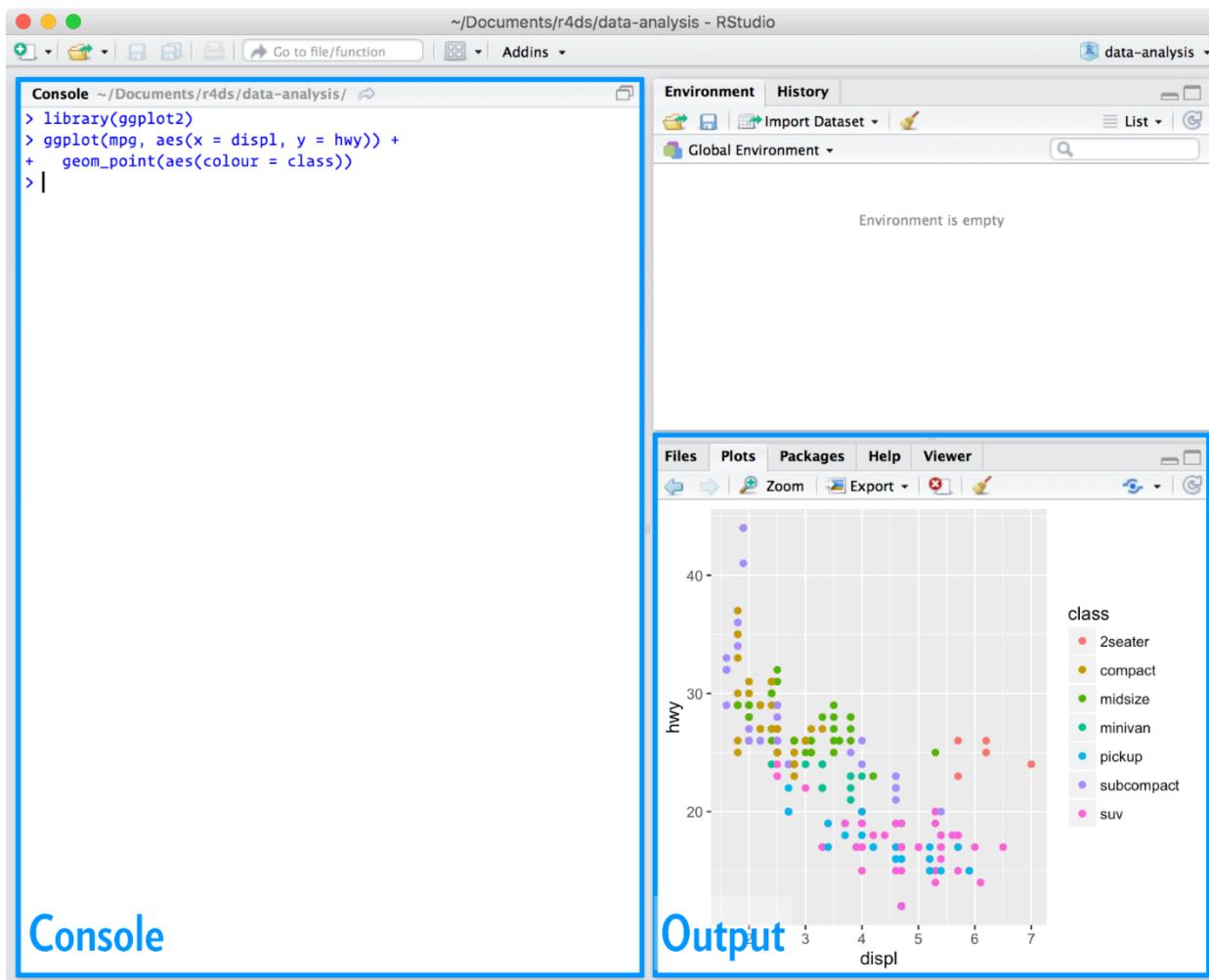
A new major version of R comes out once a year, and there are **2-3 minor releases each year**. It's a good idea to update regularly. Upgrading can be a bit of a hassle, especially for major versions, which require you to reinstall all your packages, but putting it off only makes it worse.

1.4.2 RStudio

RStudio is an integrated development environment, or IDE, for R programming. Download and install it from <http://www.rstudio.com/download>. RStudio is updated a couple of times a year.

When a new version is available, RStudio will let you know. It's a good idea to upgrade regularly so you can take advantage of the latest and greatest features. For this book, make sure you have RStudio 1.0.0.

When you **start RStudio**, you'll see two key regions in the interface:



For now, all you need to know is that you type R code in the console pane, and press enter to run it. You'll learn more as we go along!

1.4.3 The tidyverse

You'll also need to install some R packages. An R package is a collection of functions, data, and documentation that extends the capabilities of base R. Using packages is key to the successful use of R. The majority of the packages that you will learn in this book are part of the so-called tidyverse. The packages in the tidyverse share a common philosophy of data and R programming, and are designed to work together naturally.

You can install the complete tidyverse with a single line of code:

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

On your own computer, type that line of code in the console, and then press enter to run it. R will download the packages from CRAN and install them on to your computer. If you have problems installing, make sure that you are connected to the internet, and that <https://cloud.r-project.org/> isn't blocked by your firewall or proxy.

You will not be able to use the functions, objects, and help files in a package until you load it with `library()`. Once you have installed a package, you can load it with the `library()` function:

```
library(tidyverse)
#> Loading tidyverse: ggplot2
#> Loading tidyverse: tibble
#> Loading tidyverse: tidyverse
#> Loading tidyverse: readr
#> Loading tidyverse: purrr
#> Loading tidyverse: dplyr
#> Conflicts with tidy packages -----
-----
#> filter(): dplyr, stats
#> lag():     dplyr, stats
```

This tells you that tidyverse is loading the ggplot2, tibble, tidyverse, readr, purrr, and dplyr packages. These are considered to be the core of the tidyverse **because you'll use them in almost every analysis.**

Packages in the tidyverse change fairly frequently. You can see if updates are available, and optionally install them, by running `tidyverse_update()`.

1.4.4 Other packages

There are many other excellent packages that are not part of the tidyverse, because they solve problems in a different domain, or are designed with a different set of underlying principles. This **doesn't make them better or worse, just different. In other words, the complement to the tidyverse** is not the messyverse, but many other universes of interrelated packages. As you tackle more data science projects with R, you'll learn new packages and new ways of thinking about data.

In this book we'll use three data packages from outside the tidyverse:

```
install.packages(c("nycflights13", "gapminder", "Lahman"))
```

These packages provide data on airline flights, world development, and baseball that we'll use to illustrate key data science ideas.

1.5 Running R code

The previous section showed you a couple of examples of running R code. Code in the book looks like this:

```
1 + 2
#> [1] 3
#> [1] 3
```

If you run the same code in your local console, it will look like this:

```
> 1 + 2
[1] 3
```

There are two main differences. In your console, you type after the >, called the prompt; **we don't** show the prompt in the book. In the book, output is commented out with #>; in your console it **appears directly after your code**. **These two differences mean that if you're working with an** electronic version of the book, you can easily copy code out of the book and into the console.

Throughout the book we use a consistent set of conventions to refer to code:

- Functions are in a code font and followed by parentheses, like `sum()`, or `mean()`.
- Other R objects (like data or function arguments) are in a code font, without parentheses, like `flights` or `x`.
- **If we want to make it clear what package an object comes from, we'll use the package** name followed by two colons, like `dplyr::mutate()`, or `nycflights13::flights`. This is also valid R code.

1.6 Getting help and learning more

This book is not an island; there is no single resource that will allow you to master R. As you start to apply the techniques described in this book to your own data you will soon find questions that I do not answer. This section describes a few tips on how to get help, and to help you keep learning.

If you get stuck, start with Google. Typically adding “R” to a query is enough to restrict it to relevant results: if the search isn’t useful, it often means that there aren’t any R-specific results available. Google is particularly useful for error messages. If you get an error message and you have no idea what it means, try googling it! Chances are that someone else has been confused by it in the past, and there will be help somewhere on the web. (If the error message isn’t in English, run `Sys.setenv(LANGUAGE = "en")` and re-run the code; you’re more likely to find help for English error messages.)

If Google doesn’t help, try stackoverflow. Start by spending a little time searching for an existing answer, including [R] to restrict your search to questions and answers that use R. If you don’t find anything useful, prepare a minimal reproducible example or reprex. A good reprex makes it easier for other people to help you, and often you’ll figure out the problem yourself in the course of making it.

There are three things you need to include to make your example reproducible: required packages, data, and code.

1. Packages should be loaded at the top of the script, so it’s easy to see which ones the example needs. This is a good time to check that you’re using the latest version of each package; it’s possible you’ve discovered a bug that’s been fixed since you installed the package. For packages in the tidyverse, the easiest way to check is to run `tidyverse_update()`.
2. The easiest way to include data in a question is to use `dput()` to generate the R code to recreate it. For example, to recreate the `mtcars` dataset in R, I’d perform the following steps:
 1. Run `dput(mtcars)` in R
 2. Copy the output
 3. In my reproducible script, type `mtcars <-` then paste.

Try and find the smallest subset of your data that still reveals the problem.

3. Spend a little bit of time ensuring that your code is easy for others to read:
 - o **Make sure you've** used spaces and your variable names are concise, yet informative.
 - o Use comments to indicate where your problem lies.
 - o Do your best to remove everything that is not related to the problem.
The shorter your code is, the easier it is to understand, and the easier it is to fix.

Finish by checking that you have actually made a reproducible example by starting a fresh R session and copying and pasting your script in.

You should also spend some time preparing yourself to solve problems before they occur. Investing a little time in learning R each day will pay off handsomely in the long run. One way is to follow what Hadley, Garrett, and everyone else at RStudio are doing on the [RStudio blog](#). This is where we post announcements about new packages, new IDE features, and in-person courses. You might also want to follow Hadley ([@hadleywickham](#)) or Garrett ([@statgarrett](#)) on Twitter, or follow [@rstudiotips](#) to keep up with new features in the IDE.

To keep up with the R community more broadly, we recommend reading [http://www.r-bloggers.com](#): it aggregates over 500 blogs about R from around the world. If you're an active Twitter user, follow the #rstats hashtag. Twitter is one of the key tools that Hadley uses to keep up with new developments in the community.

1.7 Acknowledgements

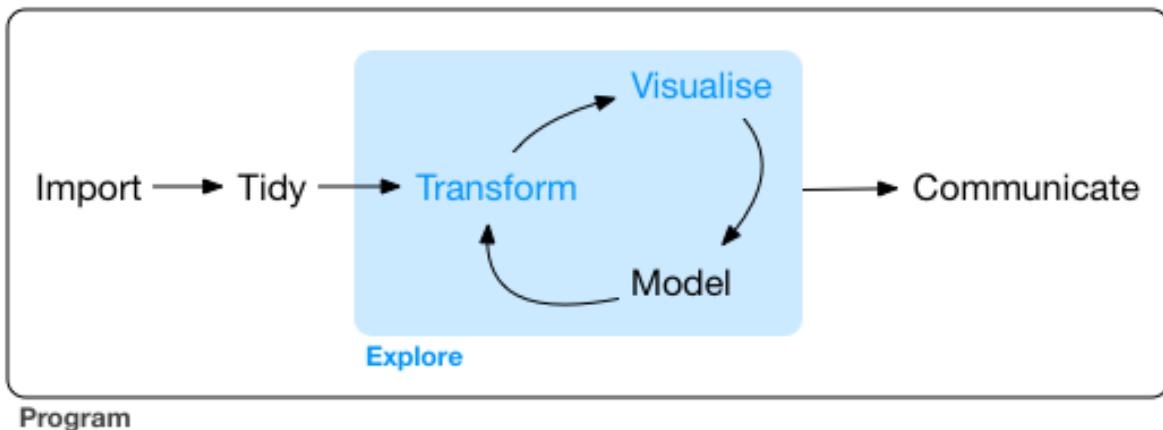
This book isn't just the product of Hadley and Garrett, but is the result of many conversations (in person and online) that we've had with the many people in the R community. There are a few people we'd like to thank in particular, because they have spent many hours answering our dumb questions and helping us to better think about data science:

- Jenny Bryan and Lionel Henry for many helpful discussions around working with lists and list-columns.
- The three chapters on workflow were adapted (with permission), from [http://stat545.com/block002_hello-r-workspace-wd-project.html](#) by Jenny Bryan.
- Genevera Allen for discussions about models, modelling, the statistical learning perspective, and the difference between hypothesis generation and hypothesis confirmation.
- Yihui Xie for his work on the [bookdown](#) package, and for tirelessly responding to my feature requests.
- Bill Behrman for his thoughtful reading of the entire book, and for trying it out with his data science class at Stanford.
- The #rstats twitter community who reviewed all of the draft chapters and provided tons of useful feedback.
- Tal Galili for augmenting his dendextend package to support a section on clustering that did not make it into the final draft.

This book was written in the open, and many people contributed pull requests to fix minor problems. Special thanks goes to everyone who contributed via GitHub:

2 Introduction

The goal of the first part of this book is to get you up to speed with the basic tools of data exploration as quickly as possible. Data exploration is the art of looking at your data, rapidly generating hypotheses, quickly testing them, then repeating again and again and again. The goal of data exploration is to generate many promising leads that you can later explore in more depth.

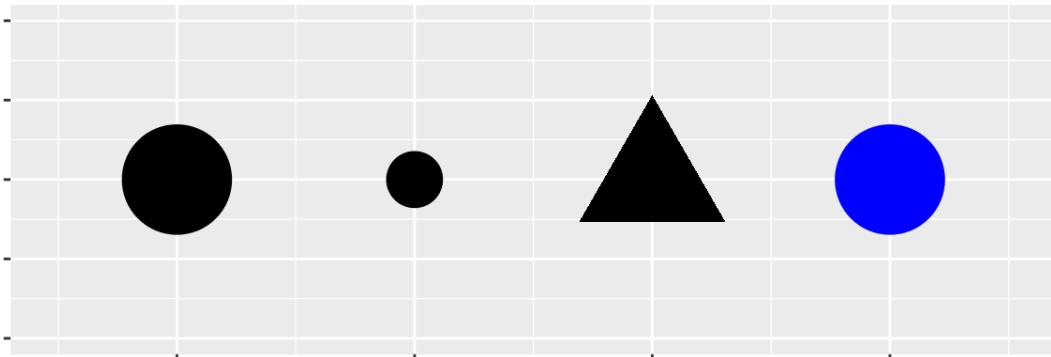


In this part of the book you will learn some useful tools that have an immediate payoff:

- Visualisation is a great place to start with R programming, because the payoff is so clear: you get to make elegant and informative plots that help you understand data. In [data visualisation](#) you'll dive into visualisation, learning the basic structure of a ggplot2 plot, and powerful techniques for turning data into plots.
- Visualisation alone is typically not enough, so in [data transformation](#) you'll learn the key verbs that allow you to select important variables, filter out key observations, create new variables, and compute summaries.
- Finally, in [exploratory data analysis](#), you'll combine visualisation and transformation with your curiosity and scepticism to ask and answer interesting questions about data.

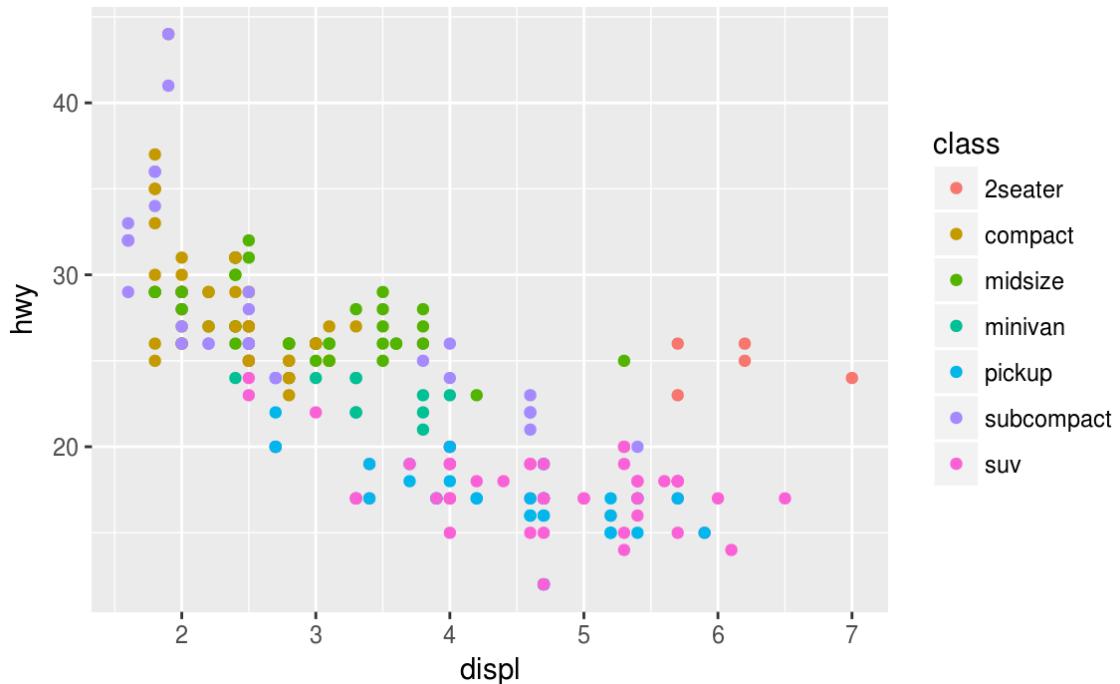
Modelling is an important part of the exploratory process, but you don't have the skills to effectively learn or apply it yet. We'll come back to it in [modelling](#), once you're better equipped with more data wrangling and programming tools.

Nestled among these three chapters that teach you the tools of exploration are three chapters that focus on your R workflow. In [workflow: basics](#), [workflow: scripts](#), and [workflow: projects](#) you'll learn good practices for writing and organising your R code. These will set you up for success in the long run, as they'll give you the tools to stay organised when you tackle real projects.



You can convey information about your data by mapping the aesthetics in your plot to the variables in your dataset. For example, you can map the colors of your points to the `class` variable to reveal the class of each car.

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy, color = class))
```



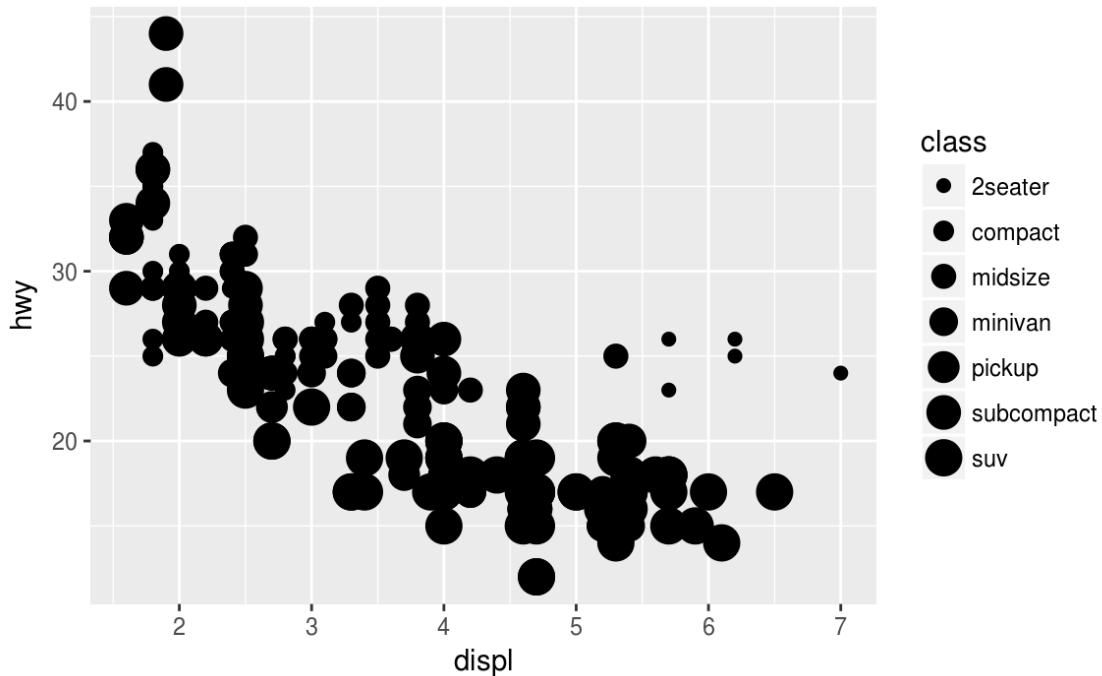
(If you prefer British English, like Hadley, you can use `colour` instead of `color`.)

To map an aesthetic to a variable, associate the name of the aesthetic to the name of the variable inside `aes()`. `ggplot2` will automatically assign a unique level of the aesthetic (here a unique color) to each unique value of the variable, a process known as scaling. `ggplot2` will also add a legend that explains which levels correspond to which values.

The colors reveal that many of the unusual points are two-seater cars. These cars don't seem like hybrids, and are, in fact, sports cars! Sports cars have large engines like SUVs and pickup trucks, but small bodies like midsize and compact cars, which improves their gas mileage. In hindsight, these cars were unlikely to be hybrids since they have large engines.

In the above example, we mapped `class` to the color aesthetic, but we could have mapped `class` to the size aesthetic in the same way. In this case, the exact size of each point would reveal its class affiliation. We get a `warning` here, because mapping an unordered variable (`class`) to an ordered aesthetic (`size`) is not a good idea.

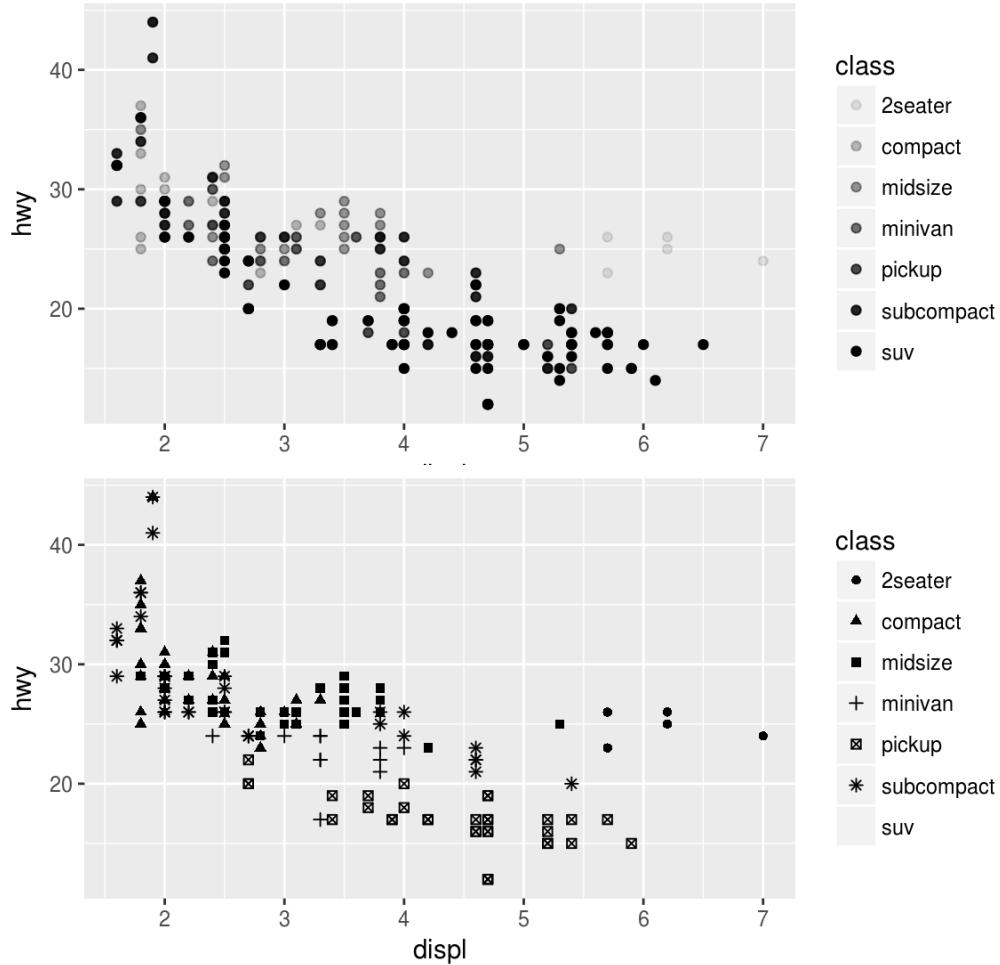
```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy, size = class))
#> Warning: Using size for a discrete variable is not advised.
```



Or we could have mapped `class` to the `alpha` aesthetic, which controls the transparency of the points, or the shape of the points.

```
# Left
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy, alpha = class))

# Right
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy, shape = class))
```



What happened to the SUVs? ggplot2 will only use six shapes at a time. By default, additional groups will go unplotted when you use the shape aesthetic.

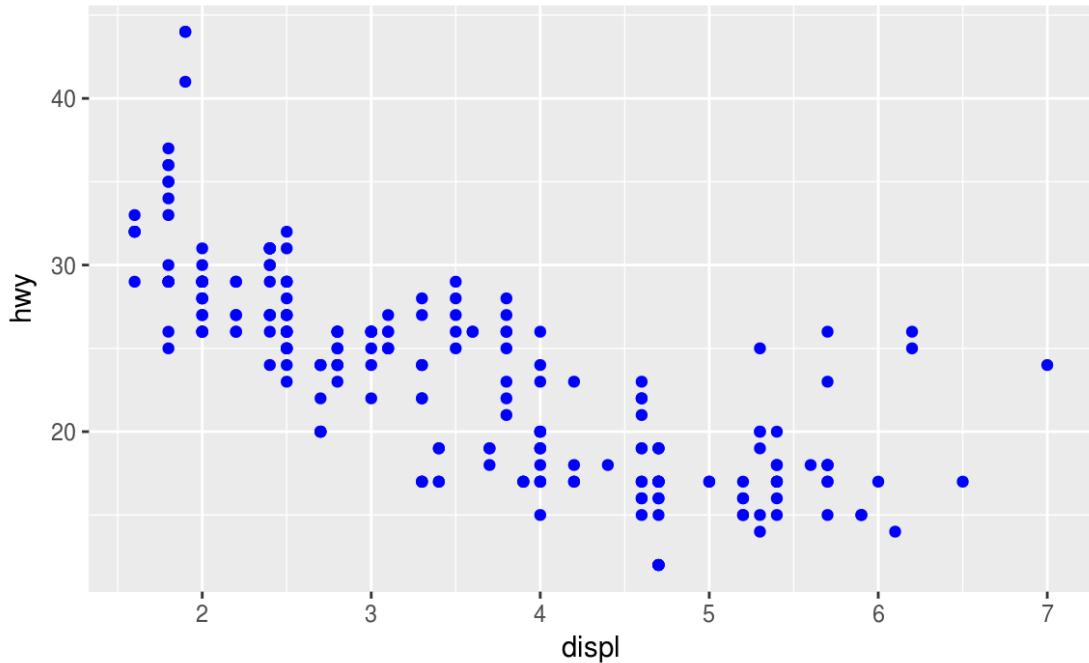
For each aesthetic, you use `aes()` to associate the name of the aesthetic with a variable to display. The `aes()` function gathers together each of the aesthetic mappings used by a layer and **passes them to the layer's mapping argument**. The syntax highlights a useful insight about `x` and `y`: the `x` and `y` locations of a point are themselves aesthetics, visual properties that you can map to variables to display information about the data.

Once you map an aesthetic, ggplot2 takes care of the rest. It selects a reasonable scale to use with the aesthetic, and it constructs a legend that explains the mapping between levels and values. For `x` and `y` aesthetics, ggplot2 does not create a legend, but it creates an axis line with tick marks and a label. The axis line acts as a legend; it explains the mapping between locations and values.

You can also set the aesthetic properties of your geom manually. For example, we can make all of the points in our plot blue:

```
ggplot(data = mpg) +
```

```
geom_point(mapping = aes(x = displ, y = hwy), color = "blue")
```



Here, the color doesn't convey information about a variable, but only changes the appearance of the plot. To set an aesthetic manually, set the aesthetic by name as an argument of your geom function; i.e. it goes `outside` of `aes()`. You'll need to pick a value that makes sense for that aesthetic:

- The name of a color as a character string.
- The size of a point in mm.
- The shape of a point as a number, as shown in Figure 3.1.

□ 0	× 4	⊕ 10	■ 15	■ 22
○ 1	▽ 6	△ 11	● 16	● 21
△ 2	⊗ 7	⊕ 12	▲ 17	▲ 24
◇ 5	* 8	⊗ 13	◆ 18	◆ 23
+ 3	◇ 9	□ 14	● 19	● 20

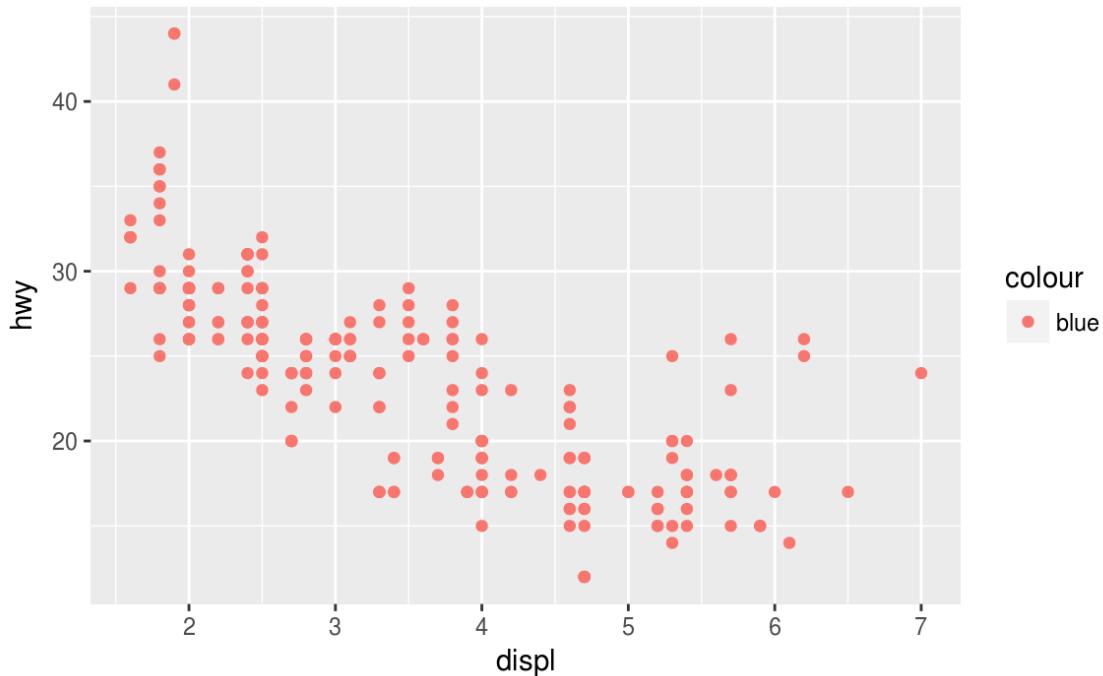
Figure 3.1: R has 25 built in shapes that are identified by numbers. There are some seeming duplicates: for example, 0, 15, and 22 are all squares. The difference comes from the interaction

of the `colour` and `fill` aesthetics. The hollow shapes (0–14) have a border determined by `colour`; the solid shapes (15–18) are filled with `colour`; the filled shapes (21–24) have a border of `colour` and are filled with `fill`.

3.3.1 Exercises

1. What's gone wrong with this code? Why are the points not blue?
2.

```
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy, color = "blue"))
```



3. Which variables in `mpg` are categorical? Which variables are continuous? (Hint: type `?mpg` to read the documentation for the dataset). How can you see this information when you run `mpg`?
4. Map a continuous variable to `color`, `size`, and `shape`. How do these aesthetics behave differently for categorical vs. continuous variables?
5. What happens if you map the same variable to multiple aesthetics?
6. What does the `stroke` aesthetic do? What shapes does it work with? (Hint: use `?geom_point`)
7. What happens if you map an aesthetic to something other than a variable name, like `aes(colour = displ < 5)`?

3.4 Common problems

As you start to run R code, you're likely to run into problems. Don't worry — it happens to everyone. I have been writing R code for years, and every day I still write code that doesn't work!

Start by carefully comparing the code that you're running to the code in the book. R is extremely picky, and a misplaced character can make all the difference. Make sure that every `(` is matched

with a) and every " is paired with another ". Sometimes you'll run the code and nothing happens. Check the left-hand of your console: if it's a +, it means that R doesn't think you've typed a complete expression and it's waiting for you to finish it. In this case, it's usually easy to start from scratch again by pressing ESCAPE to abort processing the current command.

One common problem when creating ggplot2 graphics is to put the + in the wrong place: it has to come at the end of the line, not the start. In other words, make sure you haven't accidentally written code like this:

```
ggplot(data = mpg)  
+ geom_point(mapping = aes(x = displ, y = hwy))
```

If you're still stuck, try the help. You can get help about any R function by running `?function_name` in the console, or selecting the function name and pressing F1 in RStudio. Don't worry if the help doesn't seem that helpful - instead skip down to the examples and look for code that matches what you're trying to do.

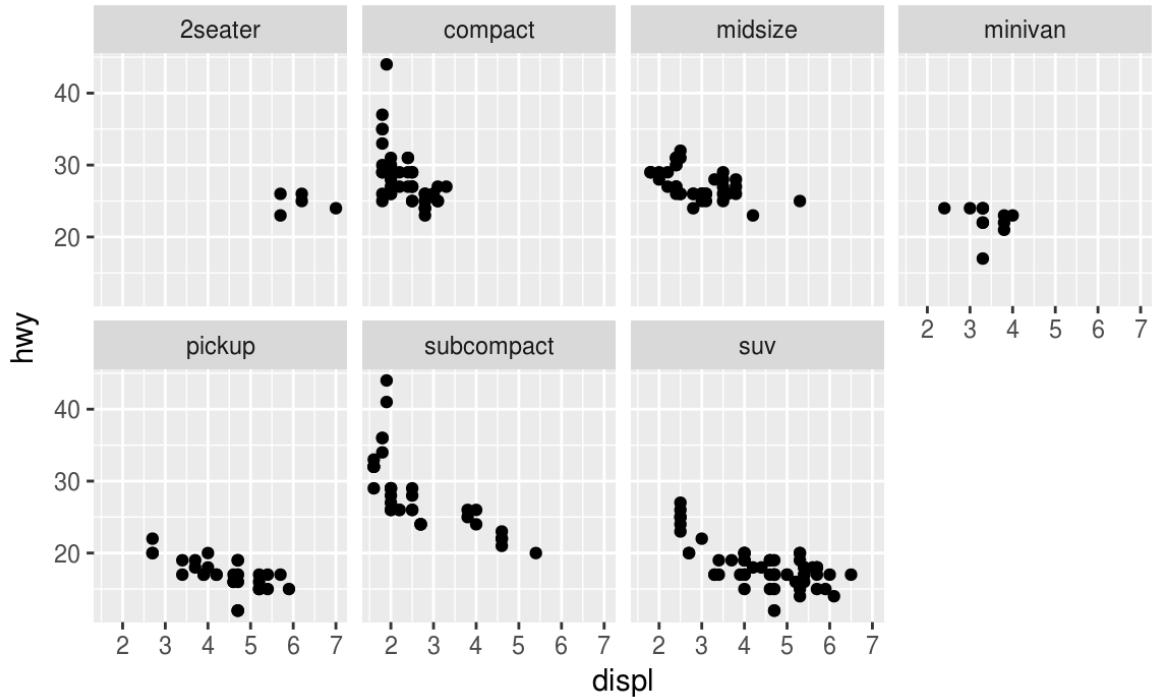
If that doesn't help, carefully read the error message. Sometimes the answer will be buried there! But when you're new to R, the answer might be in the error message but you don't yet know how to understand it. Another great tool is Google: try googling the error message, as it's likely someone else has had the same problem, and has gotten help online.

3.5 Facets

One way to add additional variables is with aesthetics. Another way, particularly useful for categorical variables, is to split your plot into facets, subplots that each display one subset of the data.

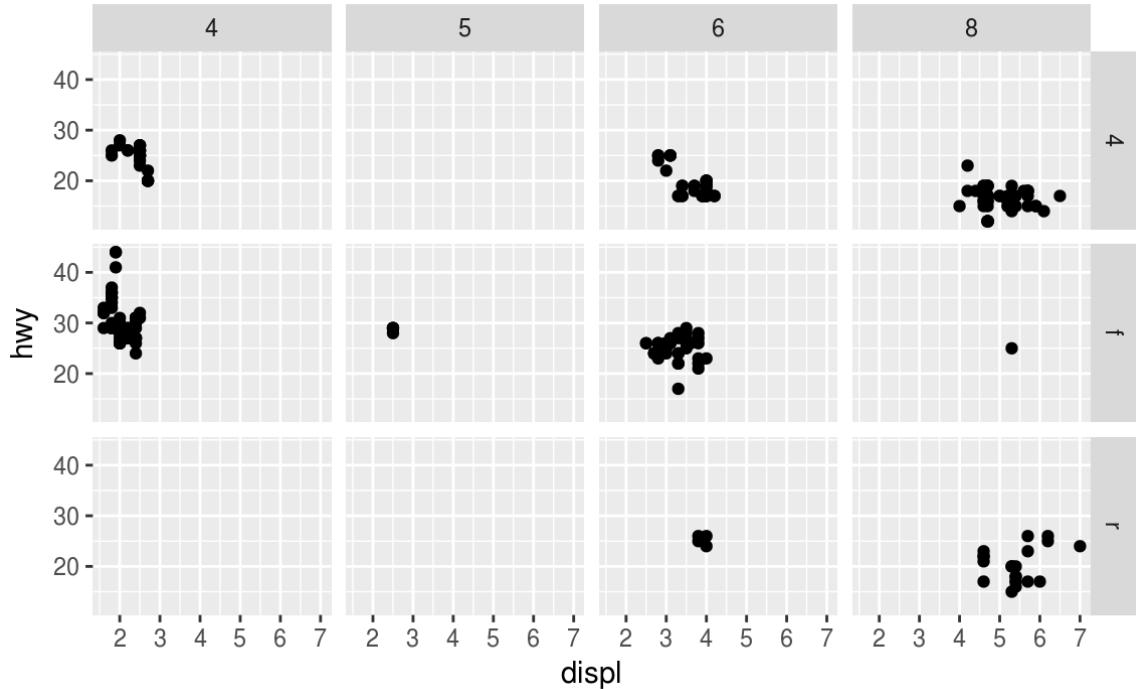
To facet your plot by a single variable, use `facet_wrap()`. The first argument of `facet_wrap()` should be a formula, which you create with ~ followed by a variable name (here "formula" is the name of a data structure in R, not a synonym for "equation"). The variable that you pass to `facet_wrap()` should be discrete.

```
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy)) +  
  facet_wrap(~ class, nrow = 2)
```



To facet your plot on the combination of two variables, add `facet_grid()` to your plot call. The first argument of `facet_grid()` is also a formula. This time the formula should contain two variable names separated by a ~.

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_grid(drv ~ cyl)
```



If you prefer to not facet in the rows or columns dimension, use a `.` instead of a variable name, e.g. `+ facet_grid(. ~ cyl)`.

3.5.1 Exercises

1. What happens if you facet on a continuous variable?
2. What do the empty cells in plot with `facet_grid(drv ~ cyl)` mean? How do they relate to this plot?
3.

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy))
```
4. What plots does the following code make? What does `.` do?
5.

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_grid(drv ~ .)
```
- 6.
7.

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_grid(. ~ cyl)
```
- 8.
9.

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_grid(. ~ cyl)
```
- 10.
11. Take the first faceted plot in this section:
12.

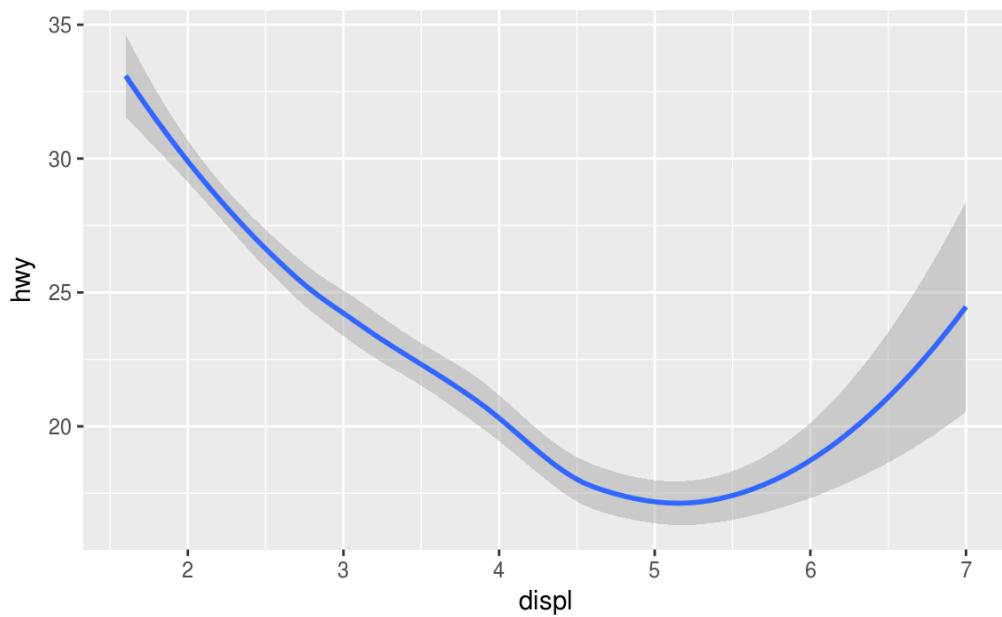
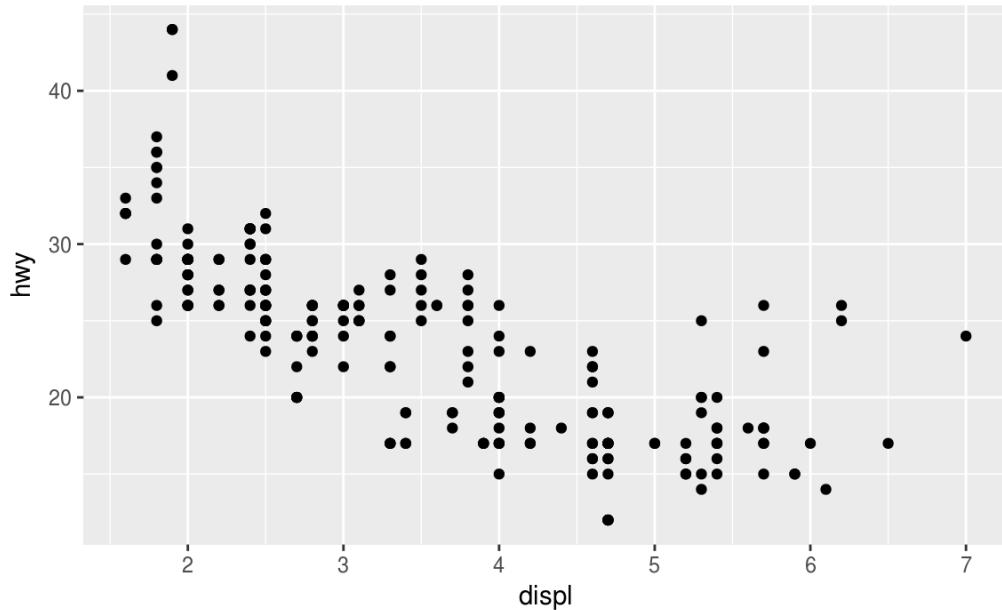
```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_wrap(~ class, nrow = 2)
```
- 13.

What are the advantages to using faceting instead of the colour aesthetic? What are the disadvantages? How might the balance change if you had a larger dataset?

14. Read `?facet_wrap`. What does `nrow` do? What does `ncol` do? What other options control the layout of the individual panels? Why doesn't `facet_grid()` have `nrow` and `ncol` argument?
15. When using `facet_grid()` you should usually put the variable with more unique levels in the columns. Why?

3.6 Geometric objects

How are these two plots similar?



Both plots contain the same x variable, the same y variable, and both describe the same data. But the plots are not identical. Each plot uses a different visual object to represent the data. In ggplot2 syntax, we say that they use different geoms.

A geom is the geometrical object that a plot uses to represent data. People often describe plots by the type of geom that the plot uses. For example, bar charts use bar geoms, line charts use line geoms, boxplots use boxplot geoms, and so on. Scatterplots break the trend; they use the point geom. As we see above, you can use different geoms to plot the same data. The plot on the left uses the point geom, and the plot on the right uses the smooth geom, a smooth line fitted to the data.

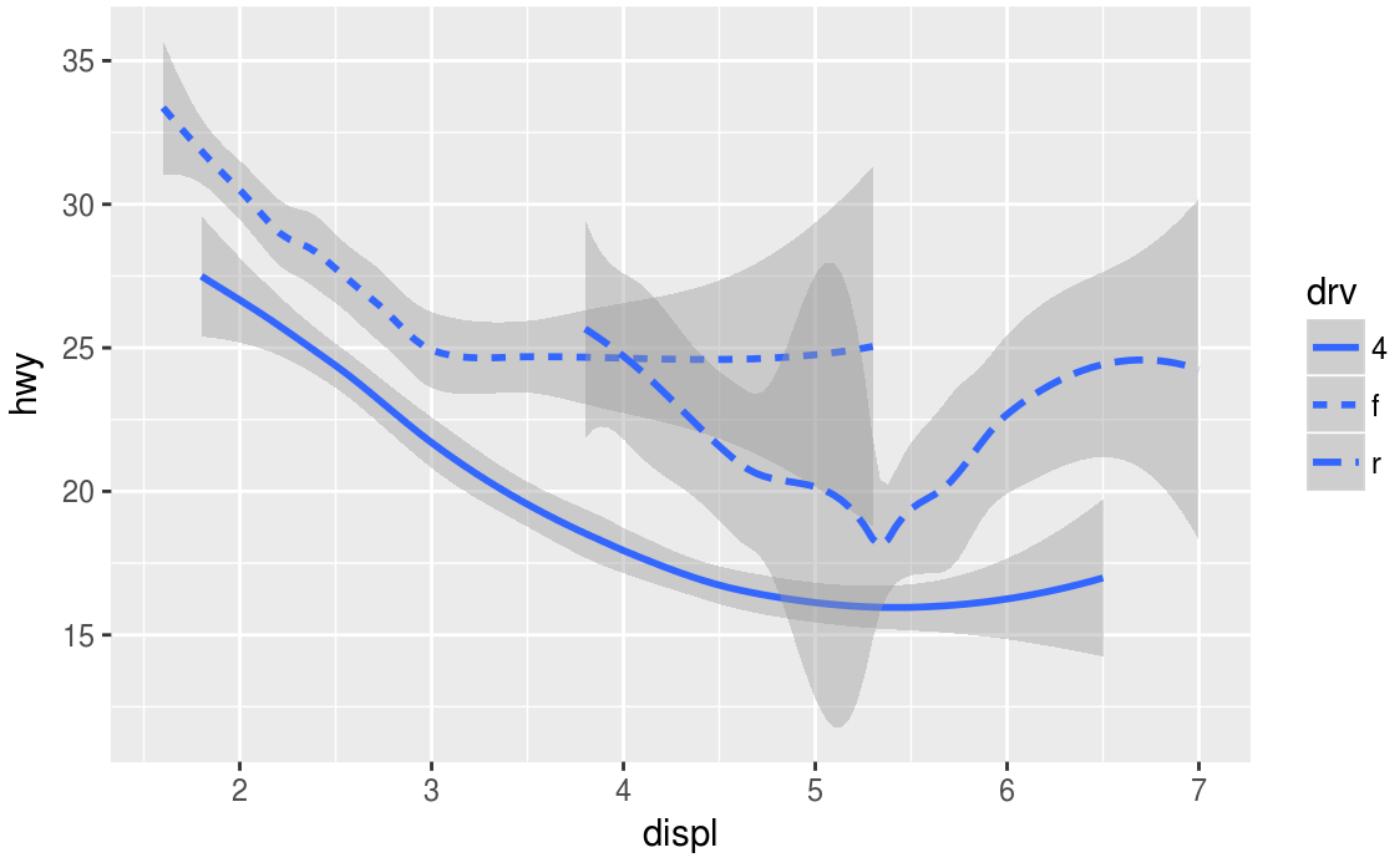
To change the geom in your plot, change the geom function that you add to `ggplot()`. For instance, to make the plots above, you can use this code:

```
# left
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy))

# right
ggplot(data = mpg) +
  geom_smooth(mapping = aes(x = displ, y = hwy))
```

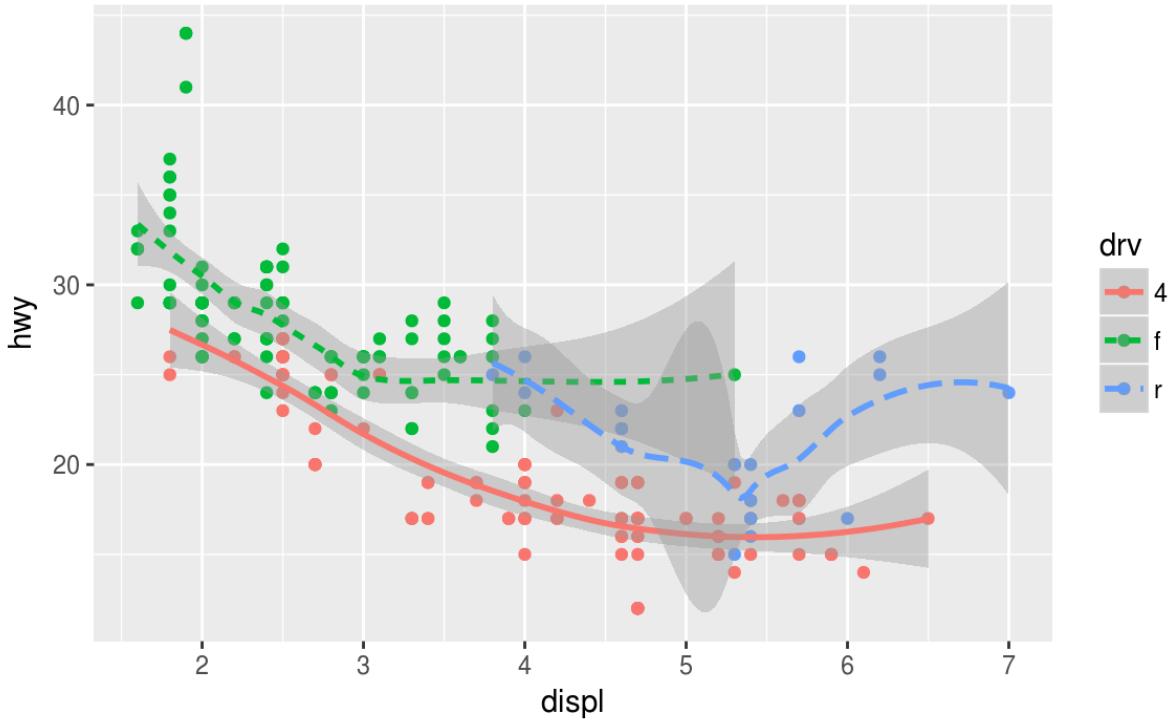
Every geom function in ggplot2 takes a `mapping` argument. However, not every aesthetic works **with every geom**. You could set the shape of a point, but you couldn't set the "shape" of a line. On the other hand, you could set the linetype of a line. `geom_smooth()` will draw a different line, with a different linetype, for each unique value of the variable that you map to linetype.

```
ggplot(data = mpg) +
  geom_smooth(mapping = aes(x = displ, y = hwy, linetype = drv))
```



Here `geom_smooth()` separates the cars into three lines based on their `drv` value, which **describes a car's drivetrain**. One line describes all of the points with a `4` value, one line describes all of the points with an `f` value, and one line describes all of the points with an `r` value. Here, `4` stands for four-wheel drive, `f` for front-wheel drive, and `r` for rear-wheel drive.

If this sounds strange, we can make it more clear by overlaying the lines on top of the raw data and then coloring everything according to `drv`.



Notice that this plot contains two geoms in the same graph! If this makes you excited, buckle up. In the next section, we will learn how to place multiple geoms in the same plot.

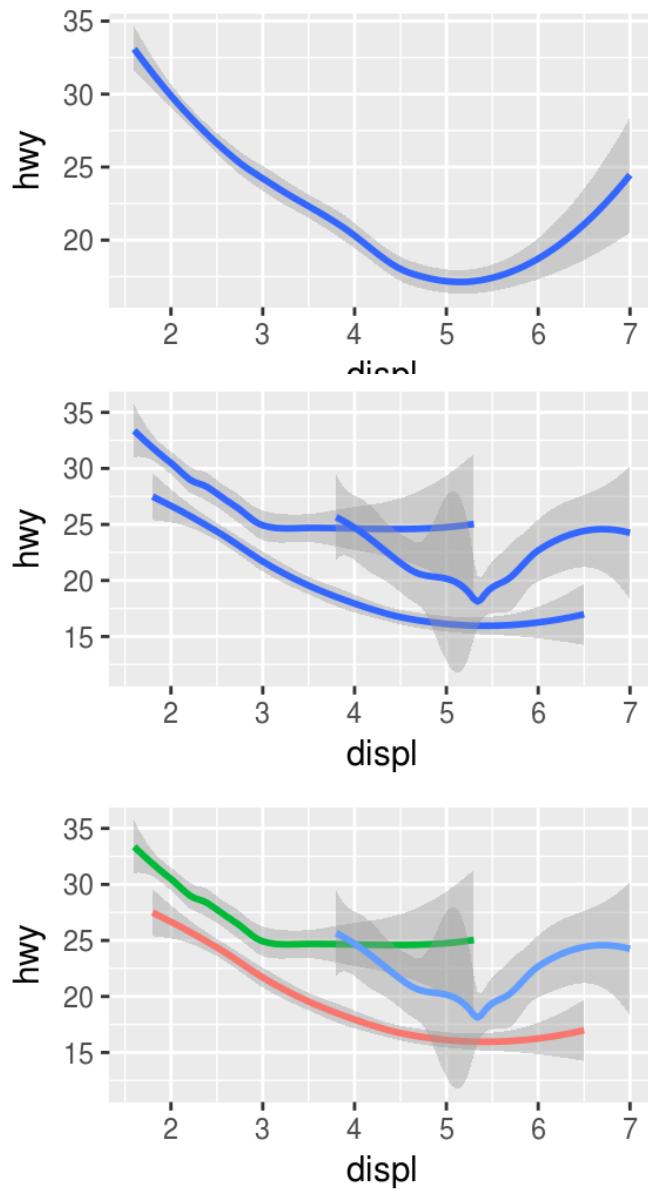
ggplot2 provides over 30 geoms, and extension packages provide even more (see <https://www.ggplot2-exts.org> for a sampling). The best way to get a comprehensive overview is the ggplot2 cheatsheet, which you can find at <http://rstudio.com/cheatsheets>. To learn more about any single geom, use help: `?geom_smooth`.

Many geoms, like `geom_smooth()`, use a single geometric object to display multiple rows of data. For these geoms, you can set the `group` aesthetic to a categorical variable to draw multiple objects. ggplot2 will draw a separate object for each unique value of the grouping variable. In practice, ggplot2 will automatically group the data for these geoms whenever you map an aesthetic to a discrete variable (as in the `linetype` example). It is convenient to rely on this feature because the `group` aesthetic by itself does not add a legend or distinguishing features to the geoms.

```
ggplot(data = mpg) +
  geom_smooth(mapping = aes(x = displ, y = hwy))

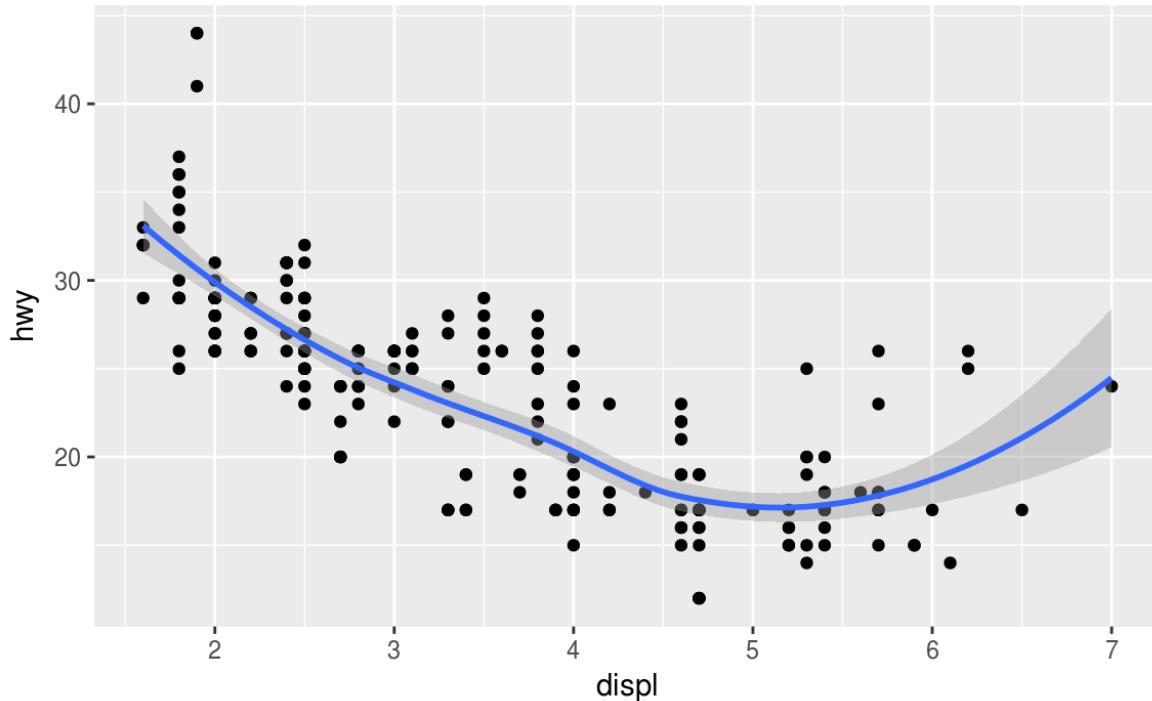
ggplot(data = mpg) +
  geom_smooth(mapping = aes(x = displ, y = hwy, group = drv))

ggplot(data = mpg) +
  geom_smooth(
    mapping = aes(x = displ, y = hwy, color = drv),
    show.legend = FALSE
  )
```



To display multiple geoms in the same plot, add multiple geom functions to `ggplot()`:

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  geom_smooth(mapping = aes(x = displ, y = hwy))
```

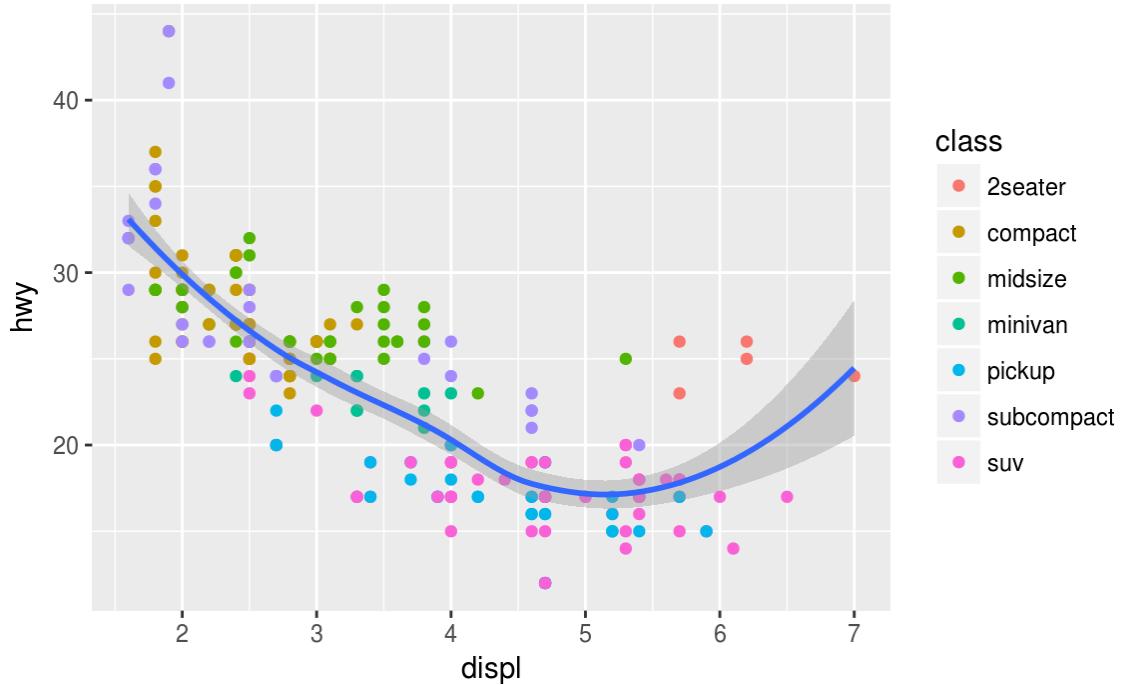


This, however, introduces some duplication in our code. Imagine if you wanted to change the y-axis to display `cty` instead of `hwy`. **You'd need to change the variable in two places, and you** might forget to update one. You can avoid this type of repetition by passing a set of mappings to `ggplot()`. `ggplot2` will treat these mappings as global mappings that apply to each geom in the graph. In other words, this code will produce the same plot as the previous code:

```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +
  geom_point() +
  geom_smooth()
```

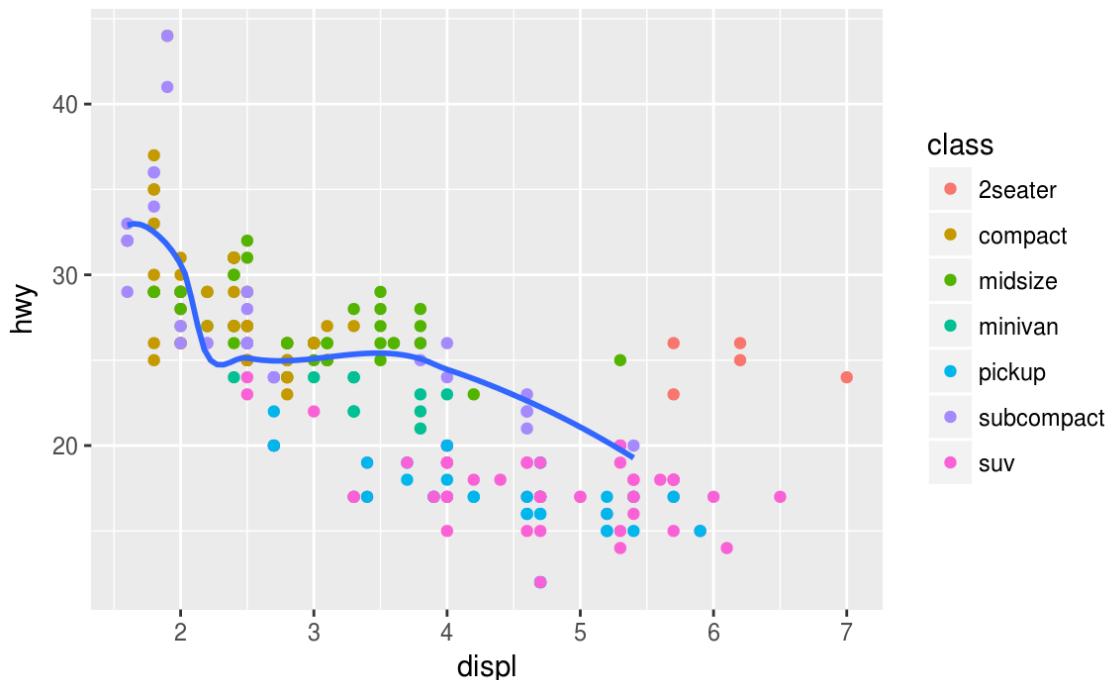
If you place mappings in a geom function, `ggplot2` will treat them as local mappings for the layer. It will use these mappings to extend or overwrite the global mappings **for that layer only**. This makes it possible to display different aesthetics in different layers.

```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +
  geom_point(mapping = aes(color = class)) +
  geom_smooth()
```



You can use the same idea to specify different data for each layer. Here, our smooth line displays just a subset of the mpg dataset, the subcompact cars. The local data argument in `geom_smooth()` overrides the global data argument in `ggplot()` for that layer only.

```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +
  geom_point(mapping = aes(color = class)) +
  geom_smooth(data = filter(mpg, class == "subcompact"), se = FALSE)
```



(You'll learn how `filter()` works in the next chapter: for now, just know that this command selects only the subcompact cars.)

3.6.1 Exercises

1. What geom would you use to draw a line chart? A boxplot? A histogram? An area chart?
2. Run this code in your head and predict what the output will look like. Then, run the code in R and check your predictions.
3.

```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy, color = drv)) +
```
4.

```
geom_point() +  
geom_smooth(se = FALSE)
```
5. What does `show.legend = FALSE` do? What happens if you remove it? Why do you think I used it earlier in the chapter?
6. What does the `se` argument to `geom_smooth()` do?
7. Will these two graphs look different? Why/why not?
8.

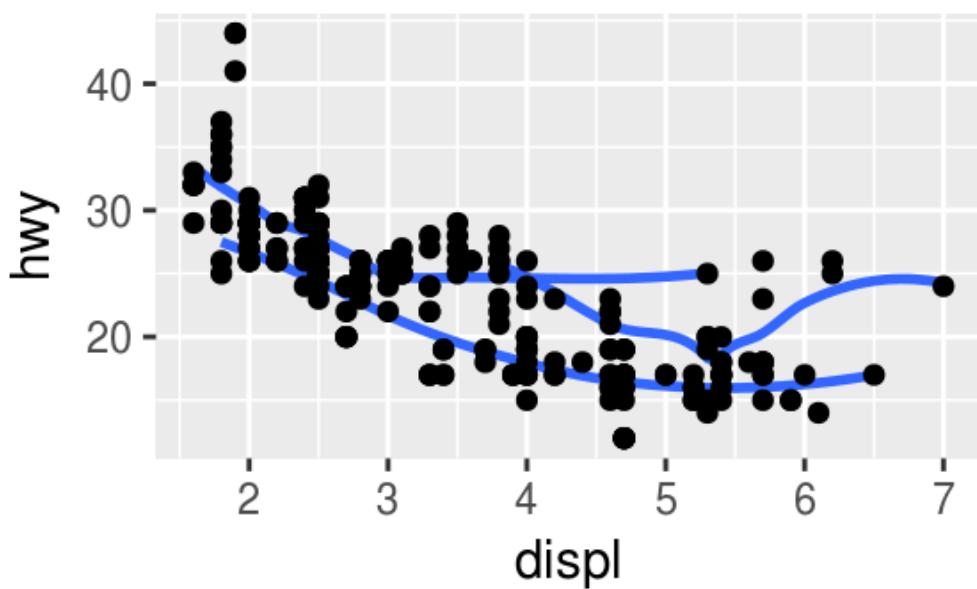
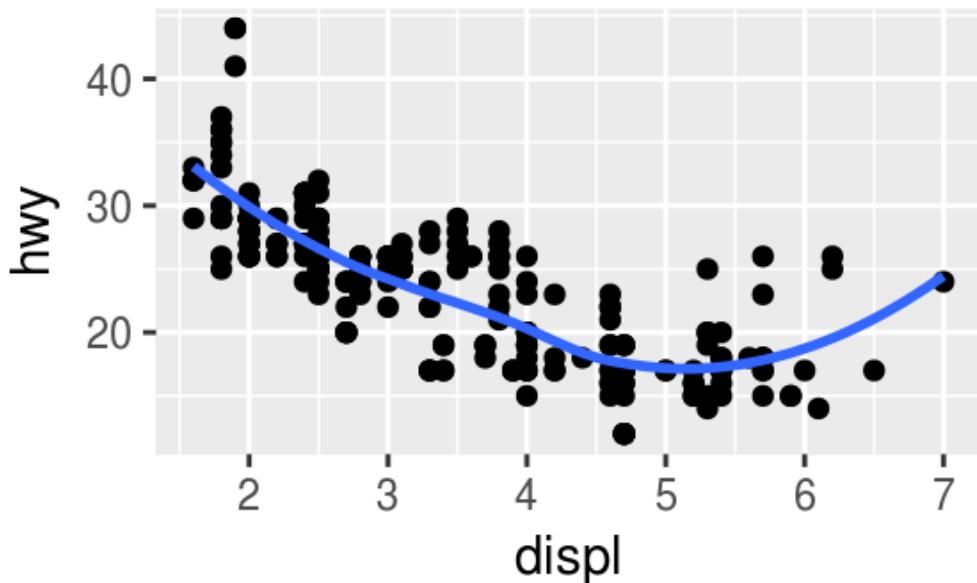
```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +
```
9.

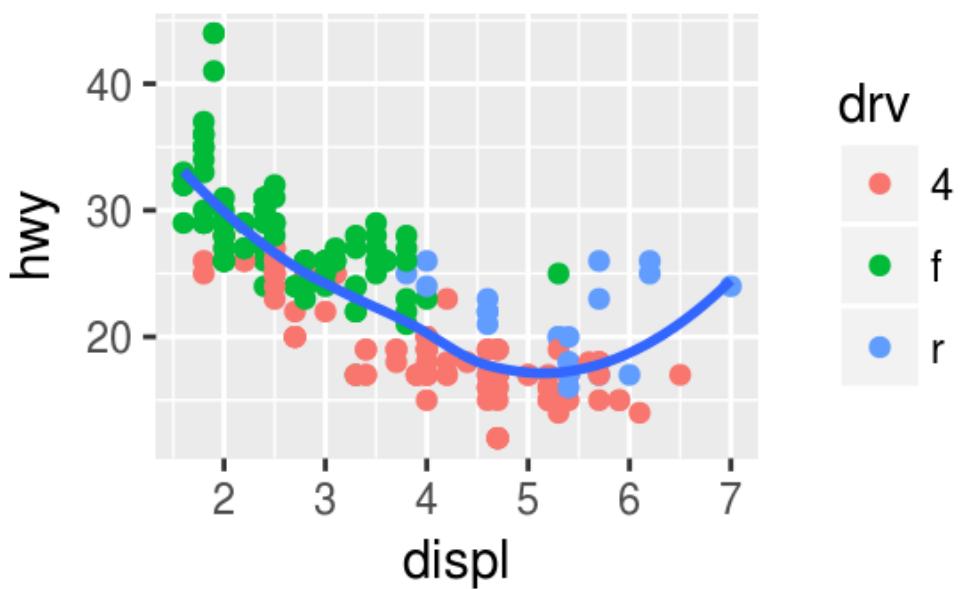
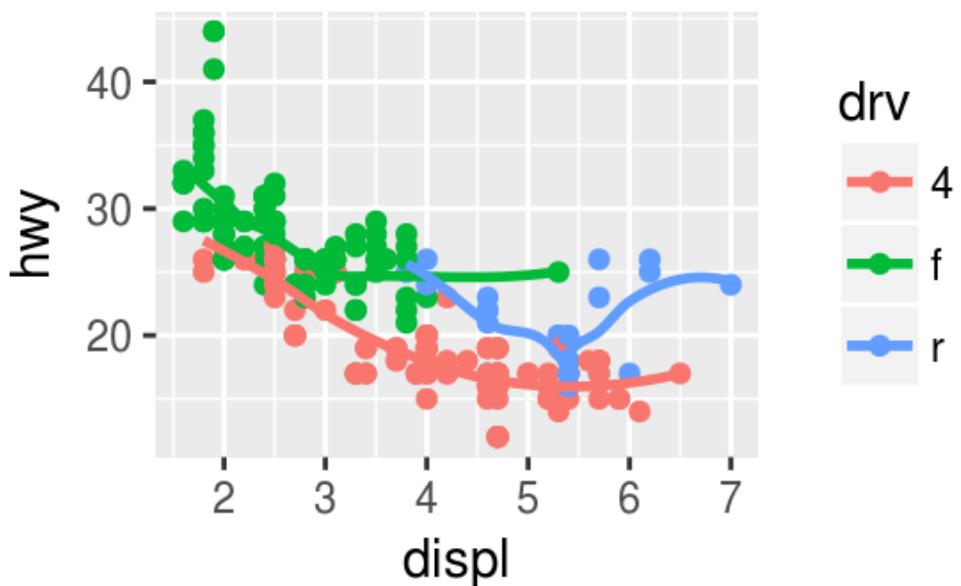
```
geom_point() +
```
10.

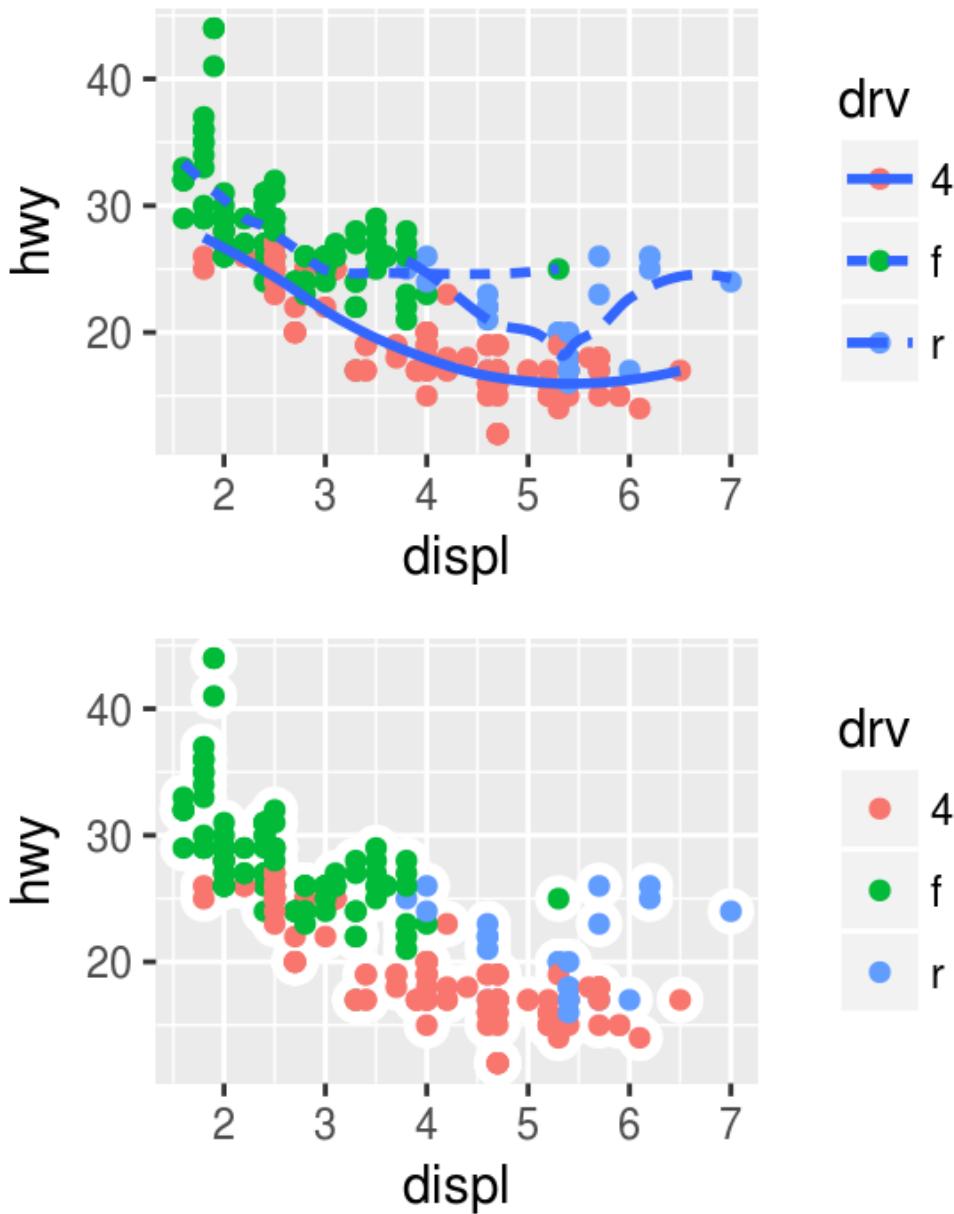
```
geom_smooth()
```
- 11.
12.

```
ggplot() +
```
13.

```
geom_point(data = mpg, mapping = aes(x = displ, y = hwy)) +  
geom_smooth(data = mpg, mapping = aes(x = displ, y = hwy))
```
14. Recreate the R code necessary to generate the following graphs.



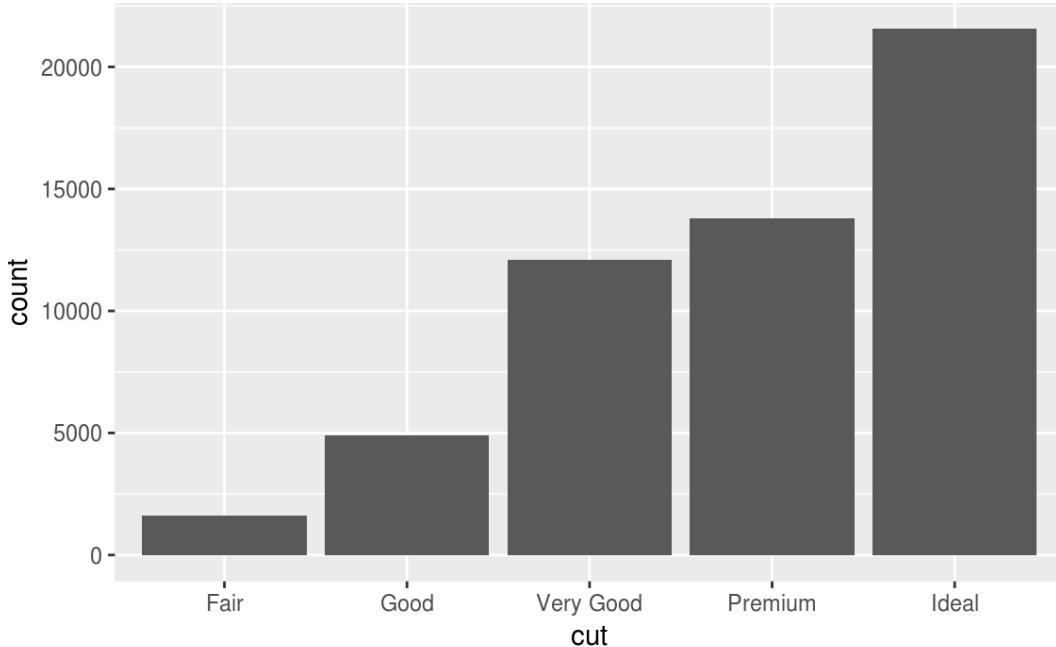




3.7 Statistical transformations

Next, let's take a look at a bar chart. Bar charts seem simple, but they are interesting because they reveal something subtle about plots. Consider a basic bar chart, as drawn with `geom_bar()`. The following chart displays the total number of diamonds in the `diamonds` dataset, grouped by `cut`. The `diamonds` dataset comes in ggplot2 and contains information about ~54,000 diamonds, including the `price`, `carat`, `color`, `clarity`, and `cut` of each diamond. The chart shows that more diamonds are available with high quality cuts than with low quality cuts.

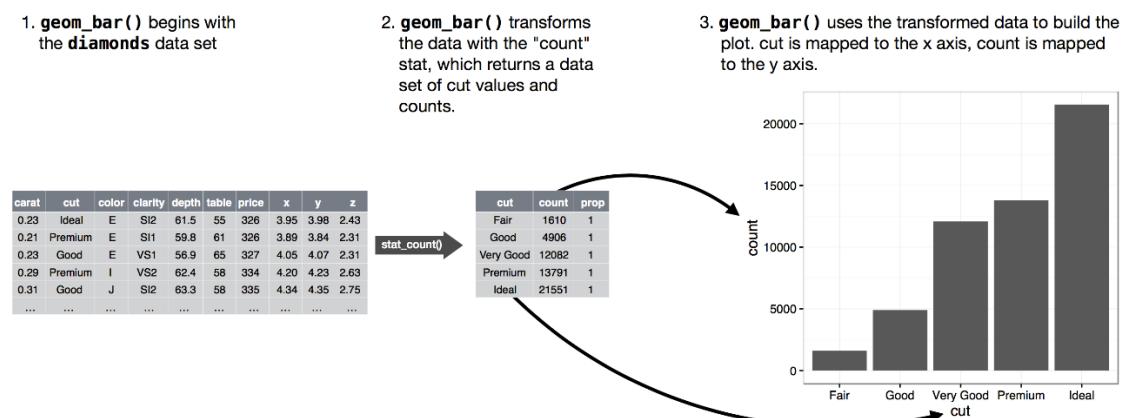
```
ggplot(data = diamonds) +
  geom_bar(mapping = aes(x = cut))
```



On the x-axis, the chart displays `cut`, a variable from `diamonds`. On the y-axis, it displays `count`, but `count` is not a variable in `diamonds`! Where does `count` come from? Many graphs, like scatterplots, plot the raw values of your dataset. Other graphs, like bar charts, calculate new values to plot:

- bar charts, histograms, and frequency polygons bin your data and then plot bin counts, the number of points that fall in each bin.
- smoothers fit a model to your data and then plot predictions from the model.
- boxplots compute a robust summary of the distribution and then display a specially formatted box.

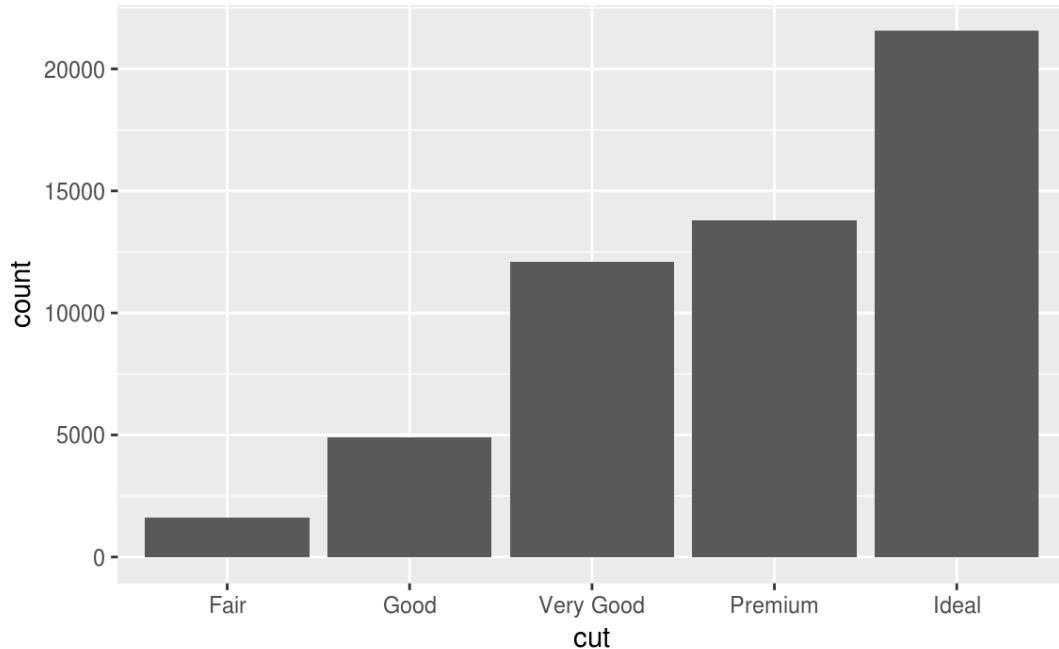
The algorithm used to calculate new values for a graph is called a stat, short for statistical transformation. The figure below describes how this process works with `geom_bar()`.



You can learn which stat a geom uses by inspecting the default value for the `stat` argument. For example, `?geom_bar` shows that the default value for `stat` is “count”, which means that `geom_bar()` uses `stat_count()`. `stat_count()` is documented on the same page as `geom_bar()`, and if you scroll down you can find a section called “Computed variables”. That describes how it computes two new variables: `count` and `prop`.

You can generally use geoms and stats interchangeably. For example, you can recreate the previous plot using `stat_count()` instead of `geom_bar()`:

```
ggplot(data = diamonds) +
  stat_count(mapping = aes(x = cut))
```



This works because every geom has a default stat; and every stat has a default geom. This means that you can typically use geoms without worrying about the underlying statistical transformation. There are three reasons you might need to use a stat explicitly:

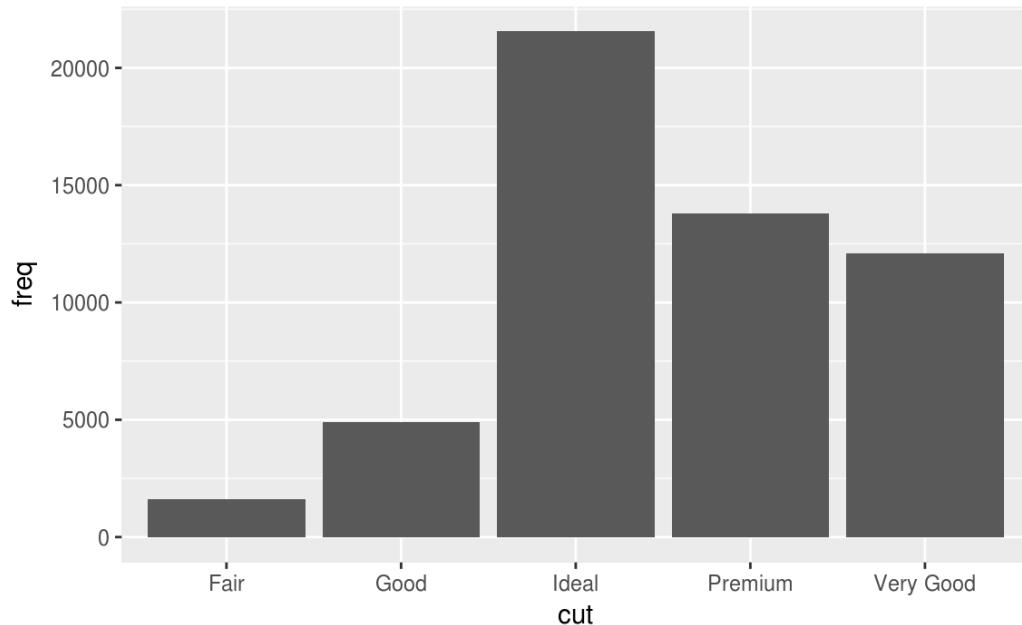
1. You might want to override the default stat. In the code below, I change the stat of `geom_bar()` from `count` (the default) to `identity`. This lets me map the height of the bars to the raw values of a `yy` variable. Unfortunately when people talk about bar charts casually, they might be referring to this type of bar chart, where the height of the bar is already present in the data, or the previous bar chart where the height of the bar is generated by counting rows.

```
2. demo <- tribble(
  3.   ~cut,           ~freq,
  4.   "Fair",         1610,
  5.   "Good",         4906,
  6.   "Very Good",   12082,
  7.   "Premium",     13791,
  8.   "Ideal",        21551
  9. )
```

```

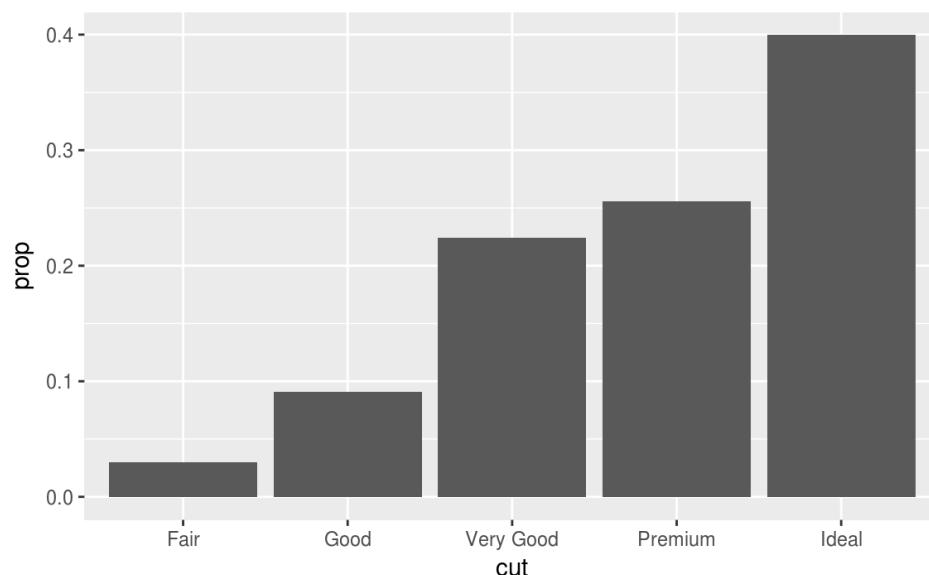
10.
11. ggplot(data = demo) +
    geom_bar(mapping = aes(x = cut, y = freq), stat = "identity")

```



(Don't worry that you haven't seen `<-` or `tribble()` before. You might be able to guess at their meaning from the context, and you'll learn exactly what they do soon!)

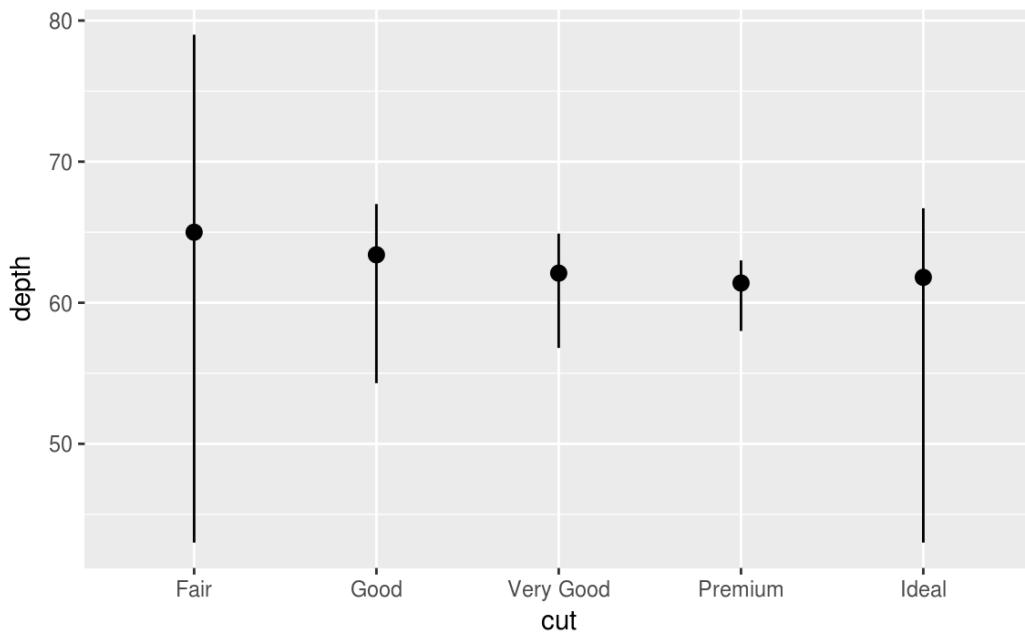
12. You might want to override the default mapping from transformed variables to aesthetics.
 For example, you might want to display a bar chart of proportion, rather than count:
 13. `ggplot(data = diamonds) +
 geom_bar(mapping = aes(x = cut, y = ..prop.., group = 1))`



To find the variables computed by the stat, look for the help section titled “computed variables”.

14. You might want to draw greater attention to the statistical transformation in your code. For example, you might use `stat_summary()`, which summarises the y values for each **unique x value, to draw attention to the summary that you’re computing**:

```
15. ggplot(data = diamonds) +  
16.   stat_summary(  
17.     mapping = aes(x = cut, y = depth),  
18.     fun.ymin = min,  
19.     fun.ymax = max,  
20.     fun.y = median  
)
```



ggplot2 provides over 20 stats for you to use. Each stat is a function, so you can get help in the usual way, e.g. `?stat_bin`. To see a complete list of stats, try the ggplot2 cheatsheet.

3.7.1 Exercises

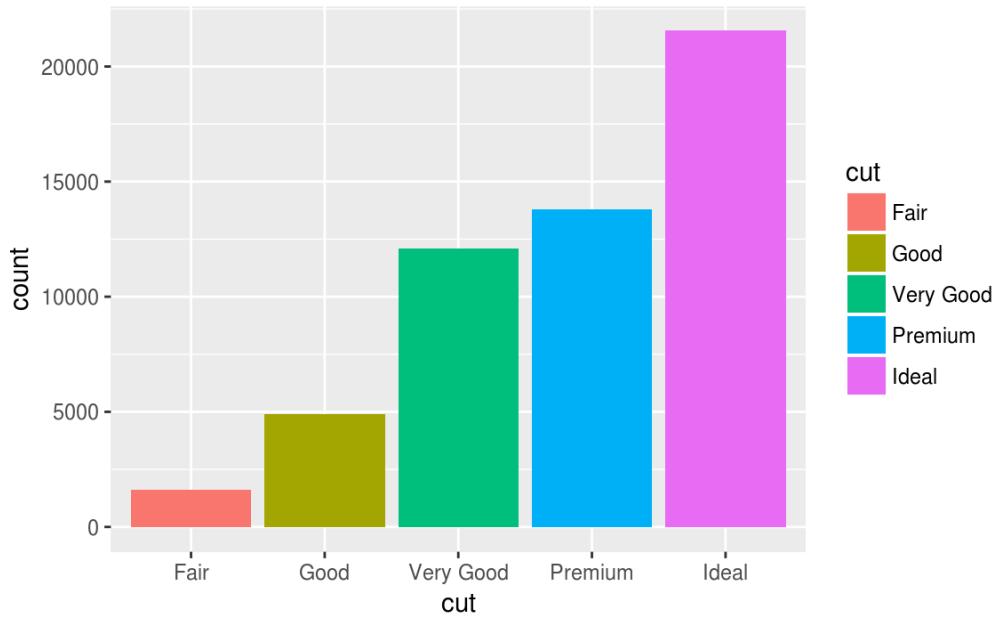
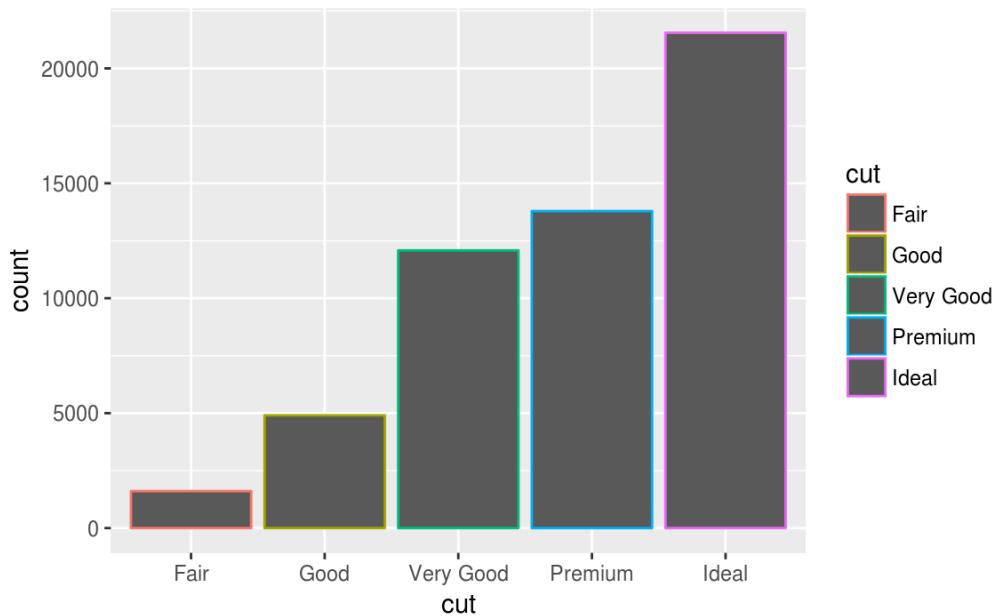
1. What is the default geom associated with `stat_summary()`? How could you rewrite the previous plot to use that geom function instead of the stat function?
2. What does `geom_col()` do? How is it different to `geom_bar()`?
3. Most geoms and stats come in pairs that are almost always used in concert. Read through the documentation and make a list of all the pairs. What do they have in common?
4. What variables does `stat_smooth()` compute? What parameters control its behaviour?
5. In our proportion bar chart, we need to set `group = 1`. Why? In other words what is the problem with these two graphs?
6. `ggplot(data = diamonds) +`
7. `geom_bar(mapping = aes(x = cut, y = ..prop..))`
8. `ggplot(data = diamonds) +`

```
geom_bar(mapping = aes(x = cut, fill = color, y = ..prop..))
```

3.8 Position adjustments

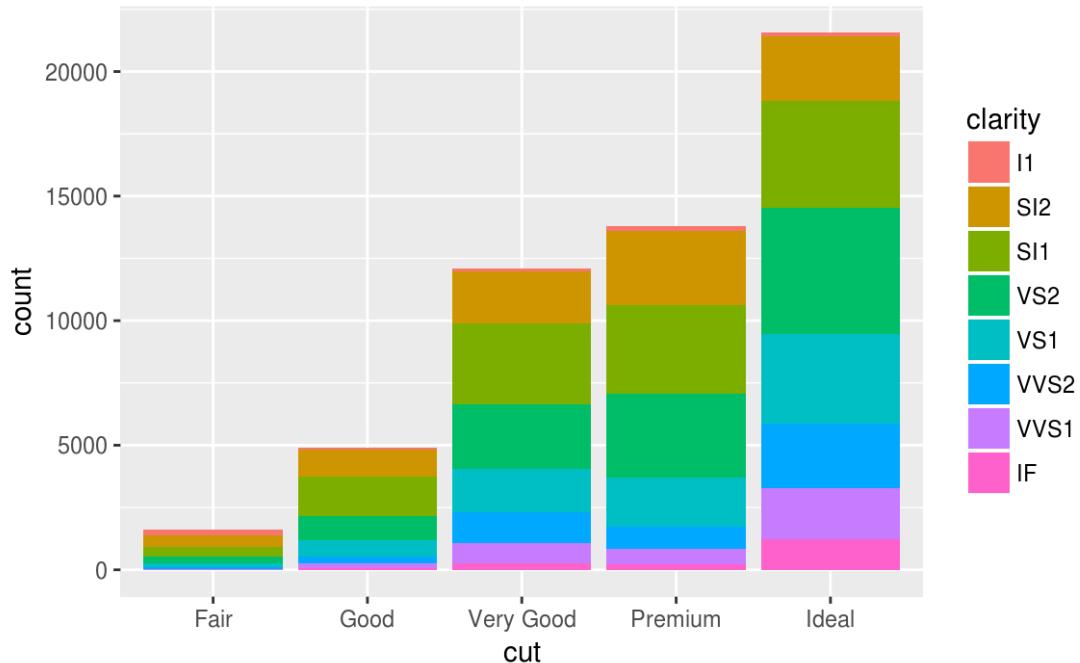
There's one more piece of magic associated with bar charts. You can colour a bar chart using either the `colour` aesthetic, or, more usefully, `fill`:

```
ggplot(data = diamonds) +  
  geom_bar(mapping = aes(x = cut, colour = cut))  
ggplot(data = diamonds) +  
  geom_bar(mapping = aes(x = cut, fill = cut))
```



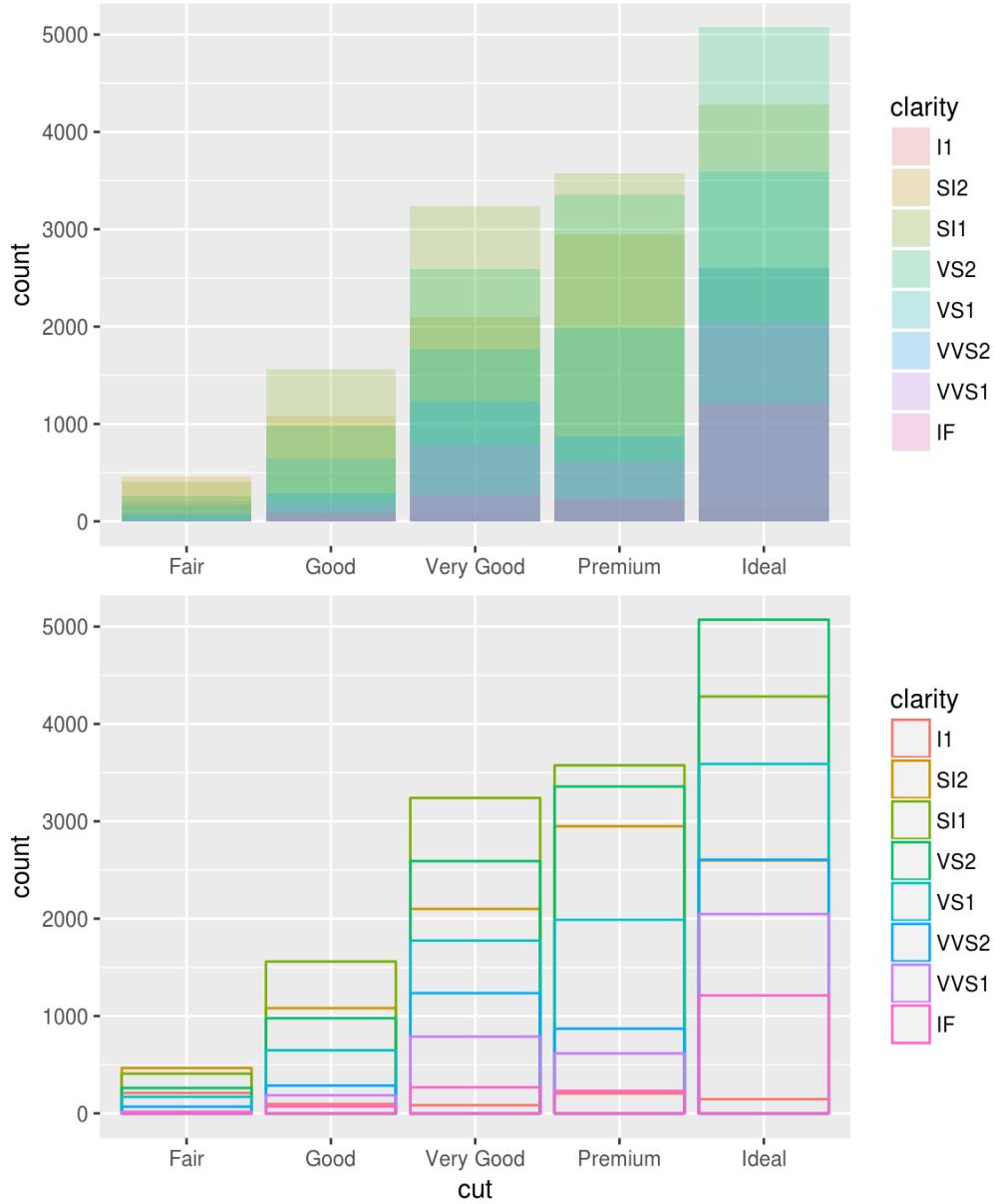
Note what happens if you map the fill aesthetic to another variable, like `clarity`: the bars are automatically stacked. Each colored rectangle represents a combination of `cut` and `clarity`.

```
ggplot(data = diamonds) +
  geom_bar(mapping = aes(x = cut, fill = clarity))
```



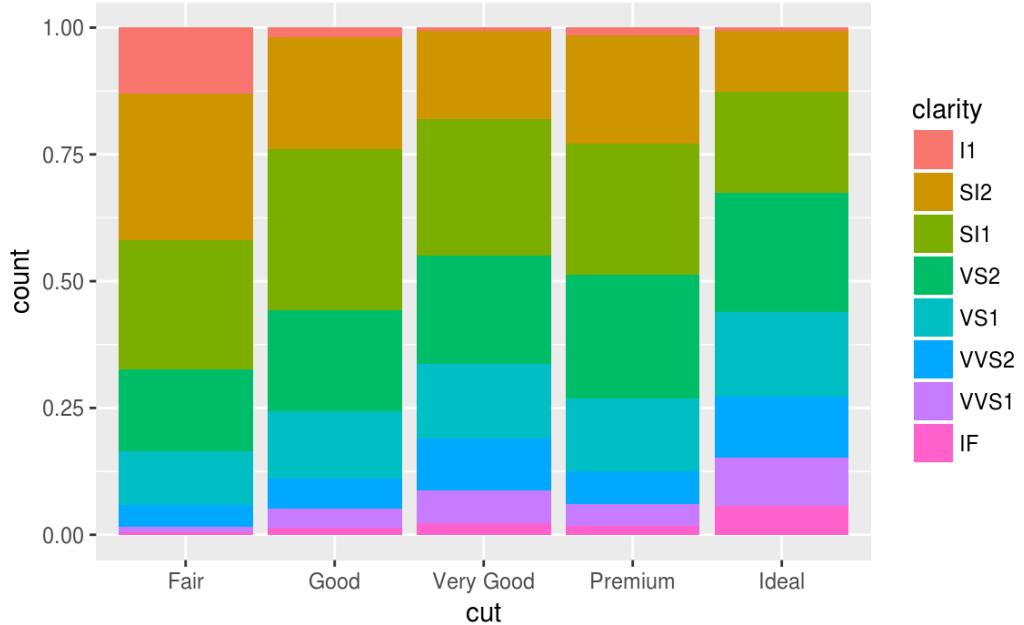
The stacking is performed automatically by the position adjustment specified by the `position` argument. If you don't want a stacked bar chart, you can use one of three other options: "identity", "dodge" or "fill".

- `position = "identity"` will place each object exactly where it falls in the context of the graph. This is not very useful for bars, because it overlaps them. To see that overlapping we either need to make the bars slightly transparent by setting `alpha` to a small value, or completely transparent by setting `fill = NA`.
- `ggplot(data = diamonds, mapping = aes(x = cut, fill = clarity)) + geom_bar(alpha = 1/5, position = "identity")`
- `ggplot(data = diamonds, mapping = aes(x = cut, colour = clarity)) + geom_bar(fill = NA, position = "identity")`

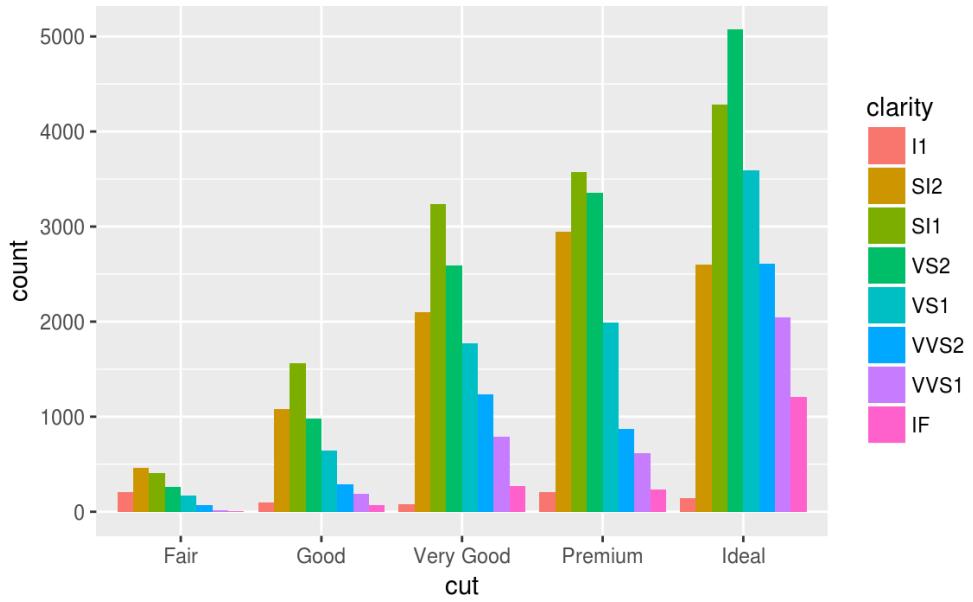


The identity position adjustment is more useful for 2d geoms, like points, where it is the default.

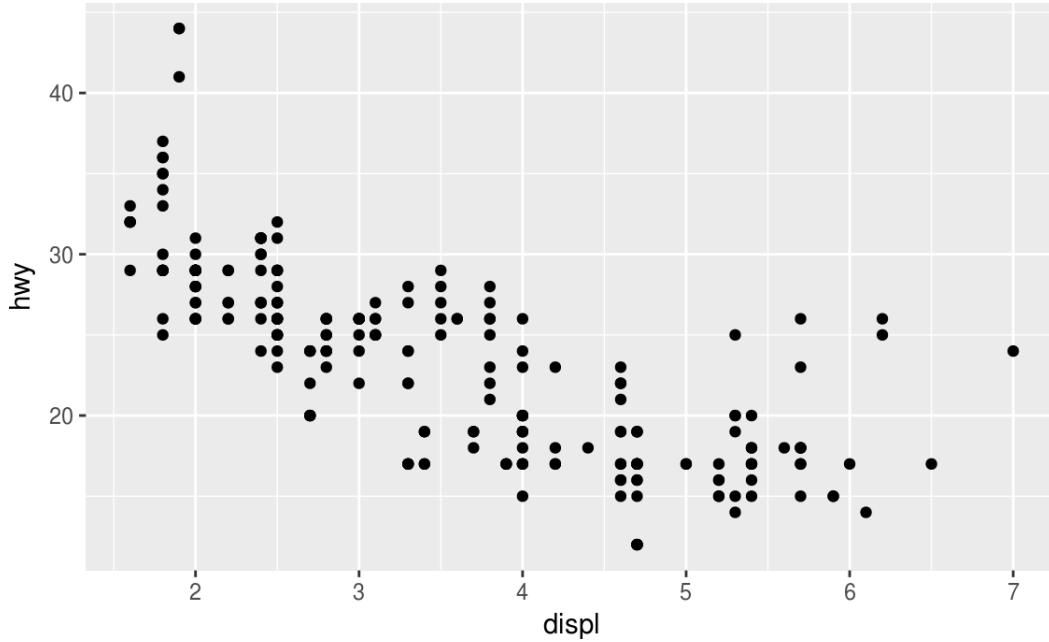
- `position = "fill"` works like stacking, but makes each set of stacked bars the same height. This makes it easier to compare proportions across groups.
- `ggplot(data = diamonds) +
 geom_bar(mapping = aes(x = cut, fill = clarity), position =
 "fill")`



- `position = "dodge"` places overlapping objects directly `beside` one another. This makes it easier to compare individual values.
- ```
ggplot(data = diamonds) +
 geom_bar(mapping = aes(x = cut, fill = clarity), position = "dodge")
```



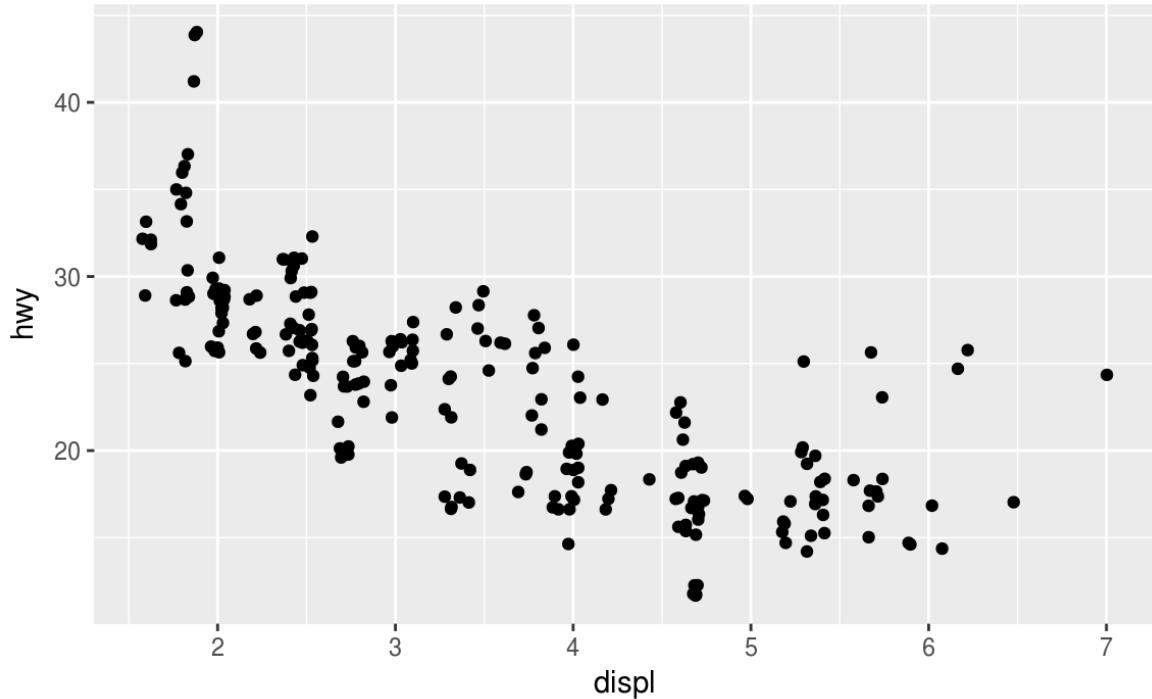
**There's one other type of adjustment that's not useful for bar charts, but it can be very useful for scatterplots.** Recall our first scatterplot. Did you notice that the plot displays only 126 points, even though there are 234 observations in the dataset?



The values of `hwy` and `displ` are rounded so the points appear on a grid and many points overlap each other. This problem is known as overplotting. This arrangement makes it hard to see where the mass of the data is. Are the data points spread equally throughout the graph, or is there one special combination of `hwy` and `displ` that contains 109 values?

You can avoid this gridding by setting the position adjustment to “jitter”. `position = "jitter"` adds a small amount of random noise to each point. This spreads the points out because no two points are likely to receive the same amount of random noise.

```
ggplot(data = mpg) +
 geom_point(mapping = aes(x = displ, y = hwy), position = "jitter")
```

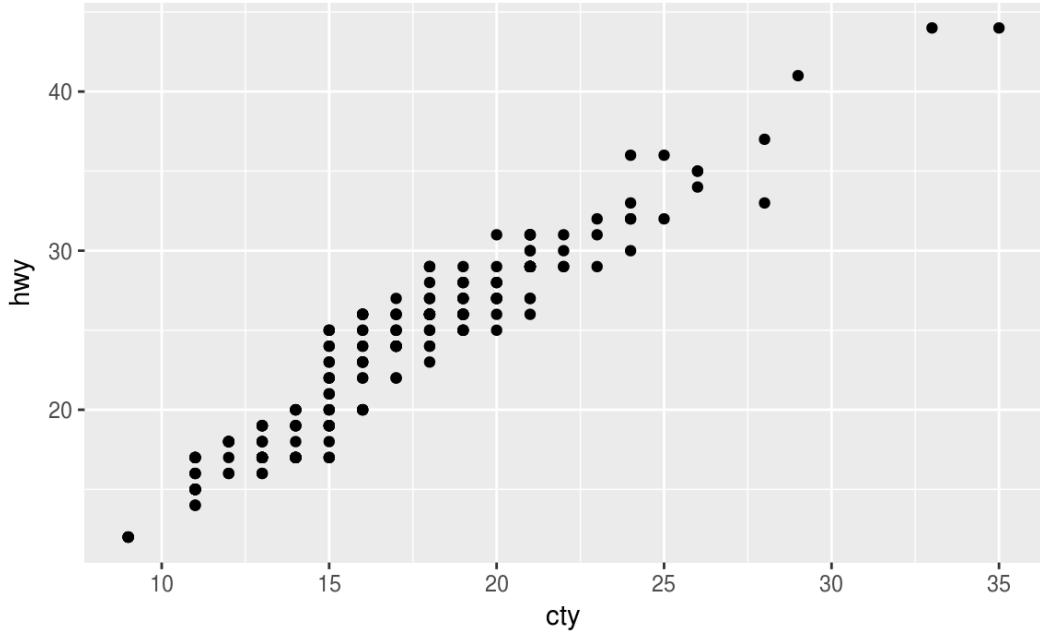


Adding randomness seems like a strange way to improve your plot, but while it makes your graph less accurate at small scales, it makes your graph ~~more~~ revealing at large scales. Because this is such a useful operation, ggplot2 comes with a shorthand for `geom_point(position = "jitter")`: `geom_jitter()`.

To learn more about a position adjustment, look up the help page associated with each adjustment: `?position_dodge`, `?position_fill`, `?position_identity`, `?position_jitter`, and `?position_stack`.

### 3.8.1 Exercises

1. What is the problem with this plot? How could you improve it?
2. `ggplot(data = mpg, mapping = aes(x = cty, y = hwy)) + geom_point()`

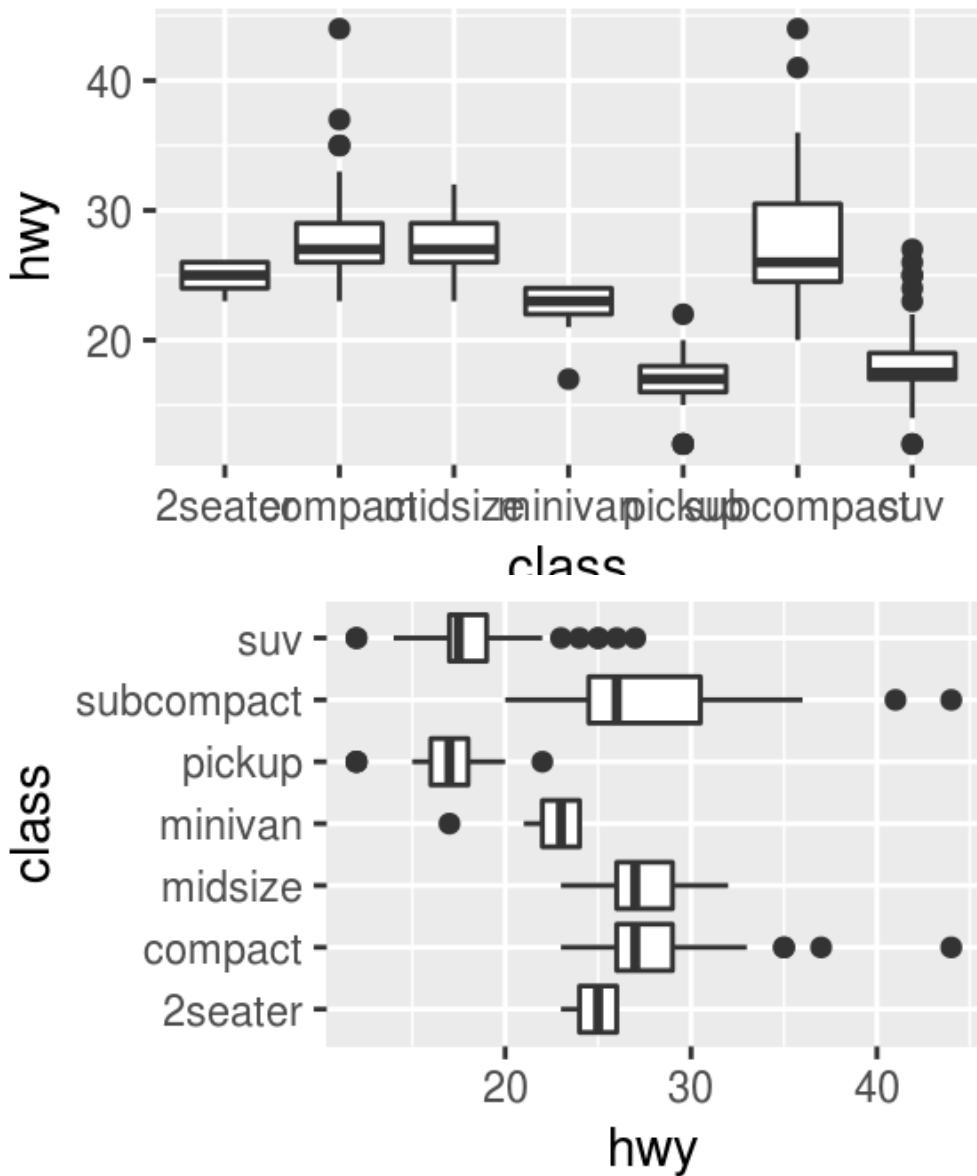


3. What parameters to `geom_jitter()` control the amount of jittering?
4. Compare and contrast `geom_jitter()` with `geom_count()`.
5. **What's the default position adjustment for `geom_boxplot()`?** Create a visualisation of the `mpg` dataset that demonstrates it.

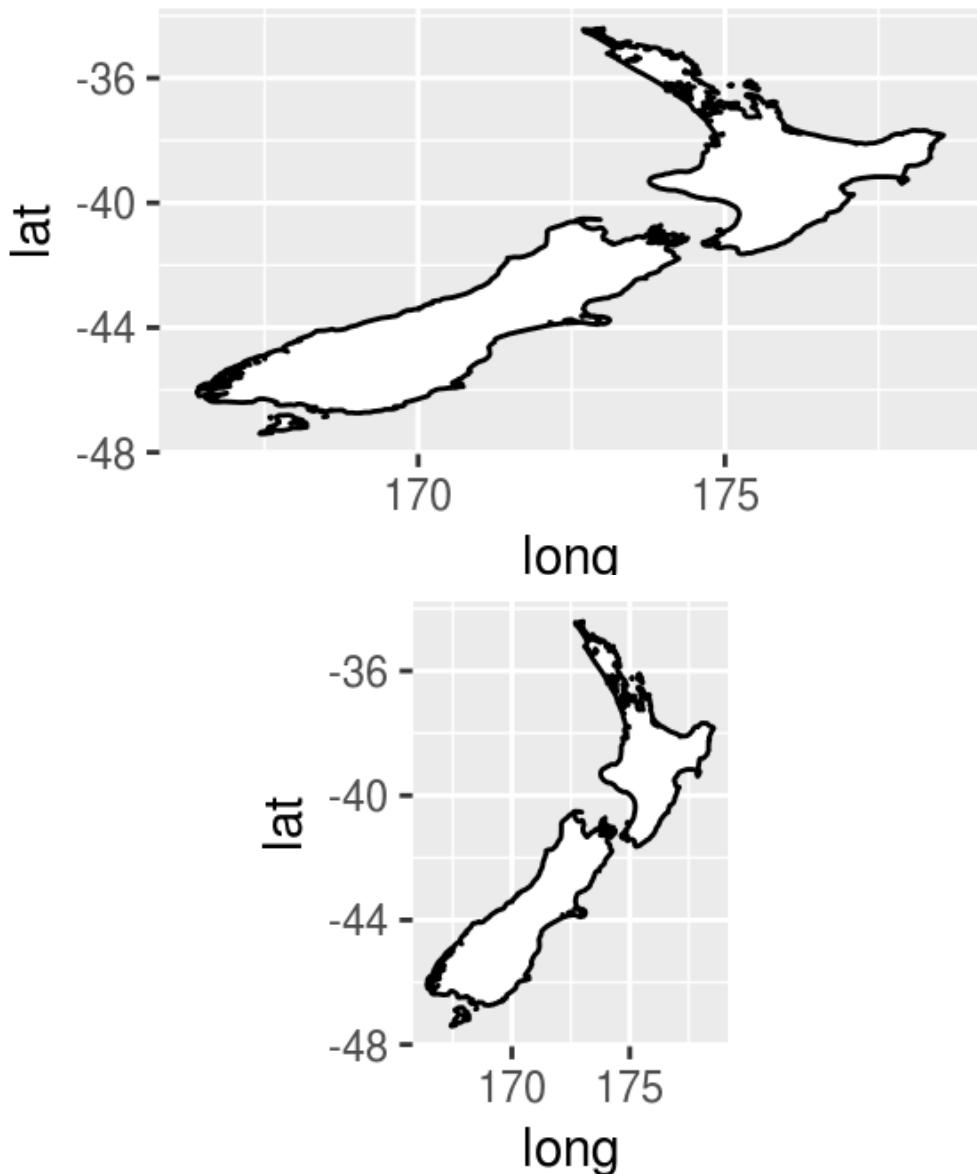
## 3.9 Coordinate systems

Coordinate systems are probably the most complicated part of ggplot2. The default coordinate system is the Cartesian coordinate system where the x and y positions act independently to determine the location of each point. There are a number of other coordinate systems that are occasionally helpful.

- `coord_flip()` switches the x and y axes. This is useful (for example), if you want **horizontal boxplots. It's also useful for long labels: it's hard to get them to fit without overlapping on the x-axis.**
- `ggplot(data = mpg, mapping = aes(x = class, y = hwy)) +  
 geom_boxplot()`
- `ggplot(data = mpg, mapping = aes(x = class, y = hwy)) +  
 geom_boxplot() +  
 coord_flip()`

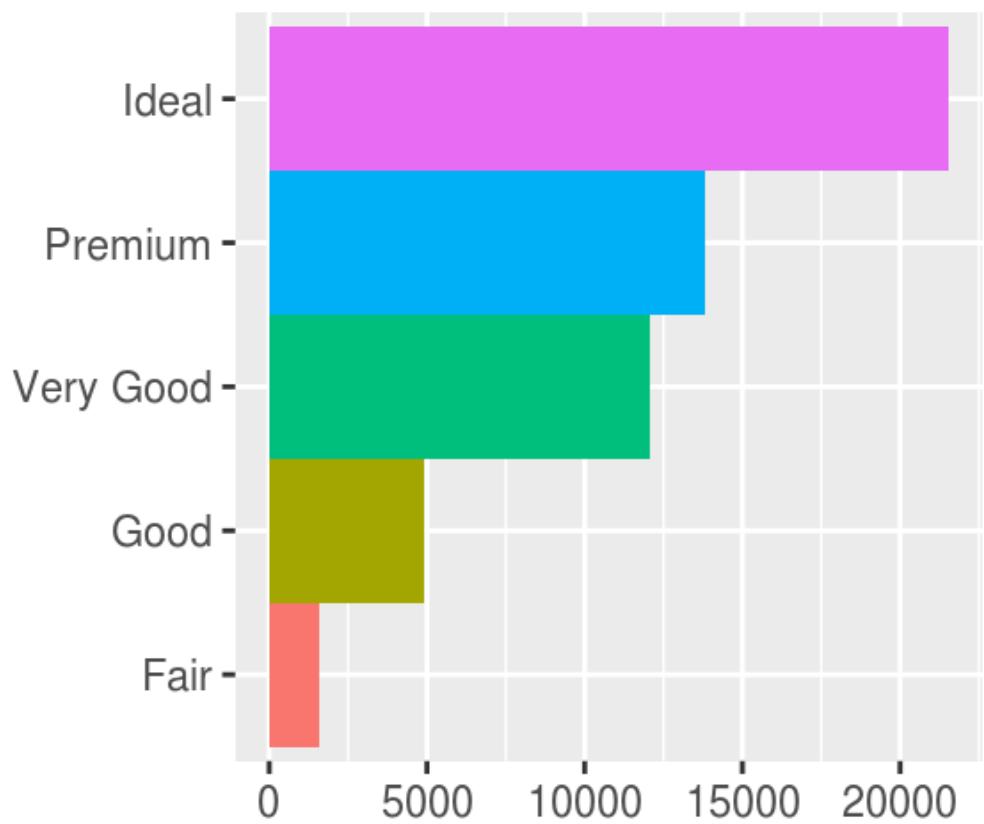


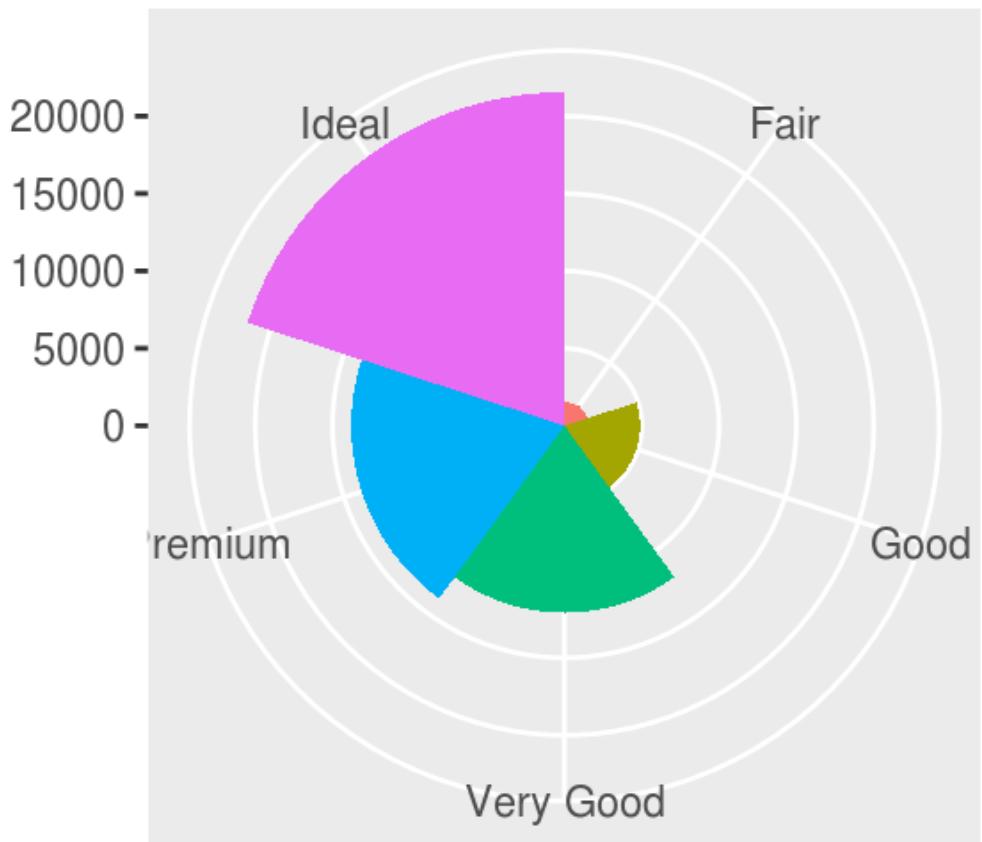
- `coord_quickmap()` sets the aspect ratio correctly for maps. This is very important if you're plotting spatial data with `ggplot2` (which unfortunately we don't have the space to cover in this book).
- `nz <- map_data("nz")`
- 
- `ggplot(nz, aes(long, lat, group = group)) +`
- `geom_polygon(fill = "white", colour = "black")`
- 
- `ggplot(nz, aes(long, lat, group = group)) +`
- `geom_polygon(fill = "white", colour = "black") +`
- `coord_quickmap()`



- `coord_polar()` uses polar coordinates. Polar coordinates reveal an interesting connection between a bar chart and a Coxcomb chart.
- `bar <- ggplot(data = diamonds) +  
 geom_bar(  
 mapping = aes(x = cut, fill = cut),  
 show.legend = FALSE,  
 width = 1  
 ) +  
 theme(aspect.ratio = 1) +  
 labs(x = NULL, y = NULL)  
•  
•`
- `bar + coord_flip()`

```
bar + coord_polar()
```

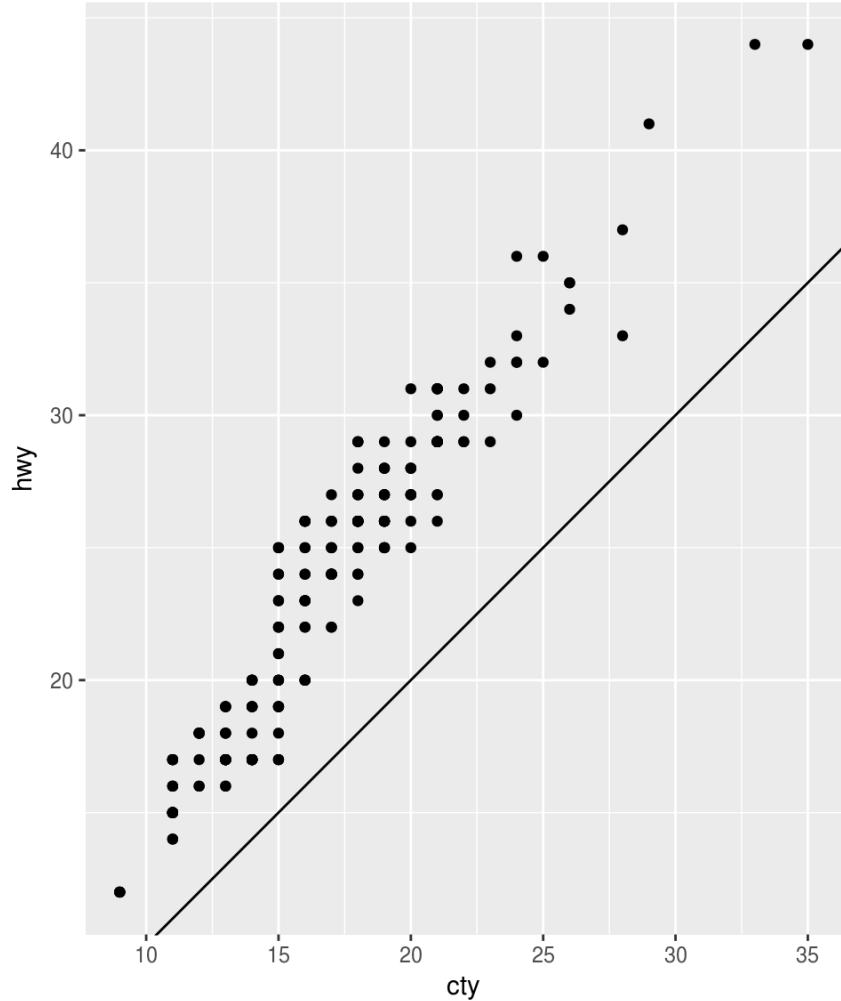




### 3.9.1 Exercises

1. Turn a stacked bar chart into a pie chart using `coord_polar()`.
2. What does `labs()` do? Read the documentation.
3. **What's** the difference between `coord_quickmap()` and `coord_map()`?
4. What does the plot below tell you about the relationship between city and highway mpg? Why is `coord_fixed()` important? What does `geom_abline()` do?

```
5. ggplot(data = mpg, mapping = aes(x = cty, y = hwy)) +
6. geom_point() +
7. geom_abline() +
 coord_fixed()
```



## 3.10 The layered grammar of graphics

In the previous sections, you learned much more than how to make scatterplots, bar charts, and boxplots. You learned a foundation that you can use to make **any** type of plot with ggplot2. To see this, let's add position adjustments, stats, coordinate systems, and faceting to our code template:

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

Our new template takes seven parameters, the bracketed words that appear in the template. In practice, you rarely need to supply all seven parameters to make a graph because ggplot2 will provide useful defaults for everything except the data, the mappings, and the geom function.

The seven parameters in the template compose the grammar of graphics, a formal system for building plots. The grammar of graphics is based on the insight that you can uniquely describe **any** plot as a combination of a dataset, a geom, a set of mappings, a stat, a position adjustment, a coordinate system, and a faceting scheme.

To see how this works, consider how you could build a basic plot from scratch: you could start with a dataset and then transform it into the information that you want to display (with a stat).

**1. Begin with the `diamonds` data set**

| carat | cut     | color | clarity | depth | table | price | x    | y    | z    |
|-------|---------|-------|---------|-------|-------|-------|------|------|------|
| 0.23  | Ideal   | E     | SI2     | 61.5  | 55    | 326   | 3.95 | 3.98 | 2.43 |
| 0.21  | Premium | E     | SI1     | 59.8  | 61    | 326   | 3.89 | 3.84 | 2.31 |
| 0.23  | Good    | E     | VS1     | 56.9  | 65    | 327   | 4.05 | 4.07 | 2.31 |
| 0.29  | Premium | I     | VS2     | 62.4  | 58    | 334   | 4.20 | 4.23 | 2.63 |
| 0.31  | Good    | J     | SI2     | 63.3  | 58    | 335   | 4.34 | 4.35 | 2.75 |
| ...   | ...     | ...   | ...     | ...   | ...   | ...   | ...  | ...  | ...  |

**2. Compute counts for each cut value with `stat_count()`.**

| cut       | count | prop |
|-----------|-------|------|
| Fair      | 1610  | 1    |
| Good      | 4906  | 1    |
| Very Good | 12082 | 1    |
| Premium   | 13791 | 1    |
| Ideal     | 21551 | 1    |

`stat_count()` →

Next, you could choose a geometric object to represent each observation in the transformed data. You could then use the aesthetic properties of the geoms to represent variables in the data. You would map the values of each variable to the levels of an aesthetic.

**3. Represent each observation with a bar.**

**4. Map the `fill` of each bar to the `..count..` variable.**

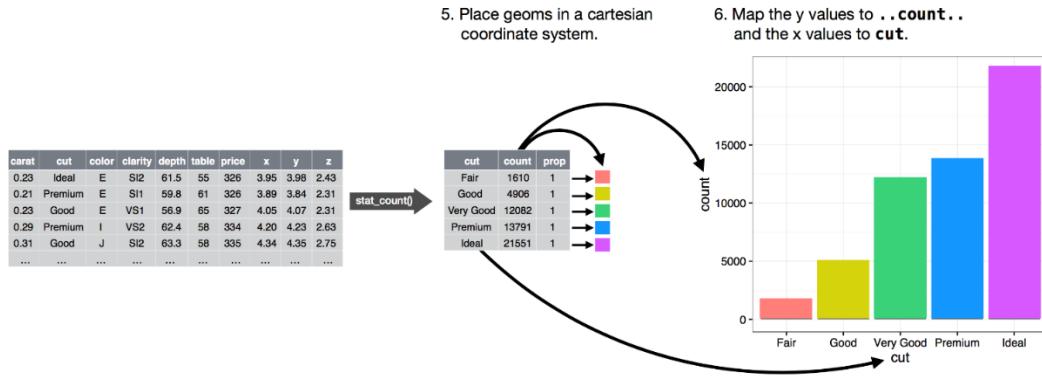
| carat | cut     | color | clarity | depth | table | price | x    | y    | z    |
|-------|---------|-------|---------|-------|-------|-------|------|------|------|
| 0.23  | Ideal   | E     | SI2     | 61.5  | 55    | 326   | 3.95 | 3.98 | 2.43 |
| 0.21  | Premium | E     | SI1     | 59.8  | 61    | 326   | 3.89 | 3.84 | 2.31 |
| 0.23  | Good    | E     | VS1     | 56.9  | 65    | 327   | 4.05 | 4.07 | 2.31 |
| 0.29  | Premium | I     | VS2     | 62.4  | 58    | 334   | 4.20 | 4.23 | 2.63 |
| 0.31  | Good    | J     | SI2     | 63.3  | 58    | 335   | 4.34 | 4.35 | 2.75 |
| ...   | ...     | ...   | ...     | ...   | ...   | ...   | ...  | ...  | ...  |

`stat_count()` →

| cut       | count | prop |
|-----------|-------|------|
| Fair      | 1610  | 1    |
| Good      | 4906  | 1    |
| Very Good | 12082 | 1    |
| Premium   | 13791 | 1    |
| Ideal     | 21551 | 1    |



You'd then select a coordinate system to place the geoms into. You'd use the location of the objects (which is itself an aesthetic property) to display the values of the x and y variables. At that point, you would have a complete graph, but you could further adjust the positions of the geoms within the coordinate system (a position adjustment) or split the graph into subplots (faceting). You could also extend the plot by adding one or more additional layers, where each additional layer uses a dataset, a geom, a set of mappings, a stat, and a position adjustment.



You could use this method to build **any** plot that you imagine. In other words, you can use the code template that you've learned in this chapter to build hundreds of thousands of unique

# 4 Workflow: basics

You now have some experience running R code. I didn't give you many details, but you've obviously figured out the basics, or you would've thrown this book away in frustration! Frustration is natural when you start programming in R, because it is such a stickler for punctuation, and even one character out of place will cause it to complain. But while you should expect to be a little frustrated, take comfort in that it's both typical and temporary: it happens to everyone, and the only way to get over it is to keep trying.

Before we go any further, let's make sure you've got a solid foundation in running R code, and that you know about some of the most helpful RStudio features.

## 4.1 Coding basics

Let's review some basics we've so far omitted in the interests of getting you plotting as quickly as possible. You can use R as a calculator:

```
1 / 200 * 30
#> [1] 0.15
(59 + 73 + 2) / 3
#> [1] 44.7
sin(pi / 2)
#> [1] 1
```

You can create new objects with <-:

```
x <- 3 * 4
```

All R statements where you create objects, assignment statements, have the same form:

```
object_name <- value
```

When reading that code say "object name gets value" in your head.

You will make lots of assignments and <- is a pain to type. Don't be lazy and use =: it will work, but it will cause confusion later. Instead, use RStudio's keyboard shortcut: Alt + - (the minus sign). Notice that RStudio automagically surrounds <- with spaces, which is a good code formatting practice. Code is miserable to read on a good day, so giveyoureyesabreak and use spaces.

## 4.2 What's in a name?

Object names must start with a letter, and can only contain letters, numbers, \_ and .. You want your object names to be descriptive, so you'll need a convention for multiple words. I recommend snake\_case where you separate lowercase words with \_.

```
i_use_snake_case
otherPeopleUseCamelCase
```

```
some.people.use.periods
And_aFew.People_RENOUNCEconvention
```

We'll come back to code style later, in [functions](#).

You can inspect an object by typing its name:

```
x
#> [1] 12
```

Make another assignment:

```
this_is_a_really_long_name <- 2.5
```

To inspect this object, try out RStudio's completion facility: type "this", press TAB, add characters until you have a unique prefix, then press return.

Ooops, you made a mistake! `this_is_a_really_long_name` should have value 3.5 not 2.5. Use another keyboard shortcut to help you fix it. Type "this" then press Cmd/Ctrl + ↑. That will list all the commands you've typed that start those letters. Use the arrow keys to navigate, then press enter to retype the command. Change 2.5 to 3.5 and rerun.

Make yet another assignment:

```
r_rocks <- 2 ^ 3
```

Let's try to inspect it:

```
r_rock
#> Error: object 'r_rock' not found
R_rocks
#> Error: object 'R_rocks' not found
```

There's an implied contract between you and R: it will do the tedious computation for you, but in return, you must be completely precise in your instructions. Typos matter. Case matters.

## 4.3 Calling functions

R has a large collection of built-in functions that are called like this:

```
function_name(arg1 = val1, arg2 = val2, ...)
```

Let's try using `seq()` which makes regular sequences of numbers and, while we're at it, learn more helpful features of RStudio. Type `se` and hit TAB. A popup shows you possible completions. Specify `seq()` by typing more (a "q") to disambiguate, or by using ↑/↓ arrows to select. Notice the floating tooltip that pops up, reminding you of the function's arguments and purpose. If you want more help, press F1 to get all the details in the help tab in the lower right pane.

Press TAB once more when you've selected the function you want. RStudio will add matching opening `( )` and closing `( )` parentheses for you. Type the arguments `1, 10` and hit return.

```
seq(1, 10)
#> [1] 1 2 3 4 5 6 7 8 9 10
```

Type this code and notice you get similar assistance with the paired quotation marks:

```
x <- "hello world"
```

Quotation marks and parentheses must always come in a pair. RStudio does its best to help you, but it's still possible to mess up and end up with a mismatch. If this happens, R will show you the continuation character "+":

```
> x <- "hello
+
```

The + tells you that R is waiting for more input; it doesn't think you're done yet. Usually that means you've forgotten either a " or a ). Either add the missing pair, or press ESCAPE to abort the expression and try again.

If you make an assignment, you don't get to see the value. You're then tempted to immediately double-check the result:

```
y <- seq(1, 10, length.out = 5)
y
#> [1] 1.00 3.25 5.50 7.75 10.00
```

This common action can be shortened by surrounding the assignment with parentheses, which causes assignment and "print to screen" to happen.

```
(y <- seq(1, 10, length.out = 5))
#> [1] 1.00 3.25 5.50 7.75 10.00
```

Now look at your environment in the upper right pane:

The screenshot shows the RStudio interface with the 'Environment' tab selected. The 'Global Environment' dropdown is open, showing a list of objects: r\_rocks, this\_is\_a\_really\_long\_variable, x, and y. The 'Values' section below displays the following data:

| Object                         | Value                        |
|--------------------------------|------------------------------|
| r_rocks                        | 8                            |
| this_is_a_really_long_variable | 2.5                          |
| x                              | "hello world"                |
| y                              | num [1:5] 1 3.25 5.5 7.75 10 |

Here you can see all of the objects that you've created.

## 4.4 Practice

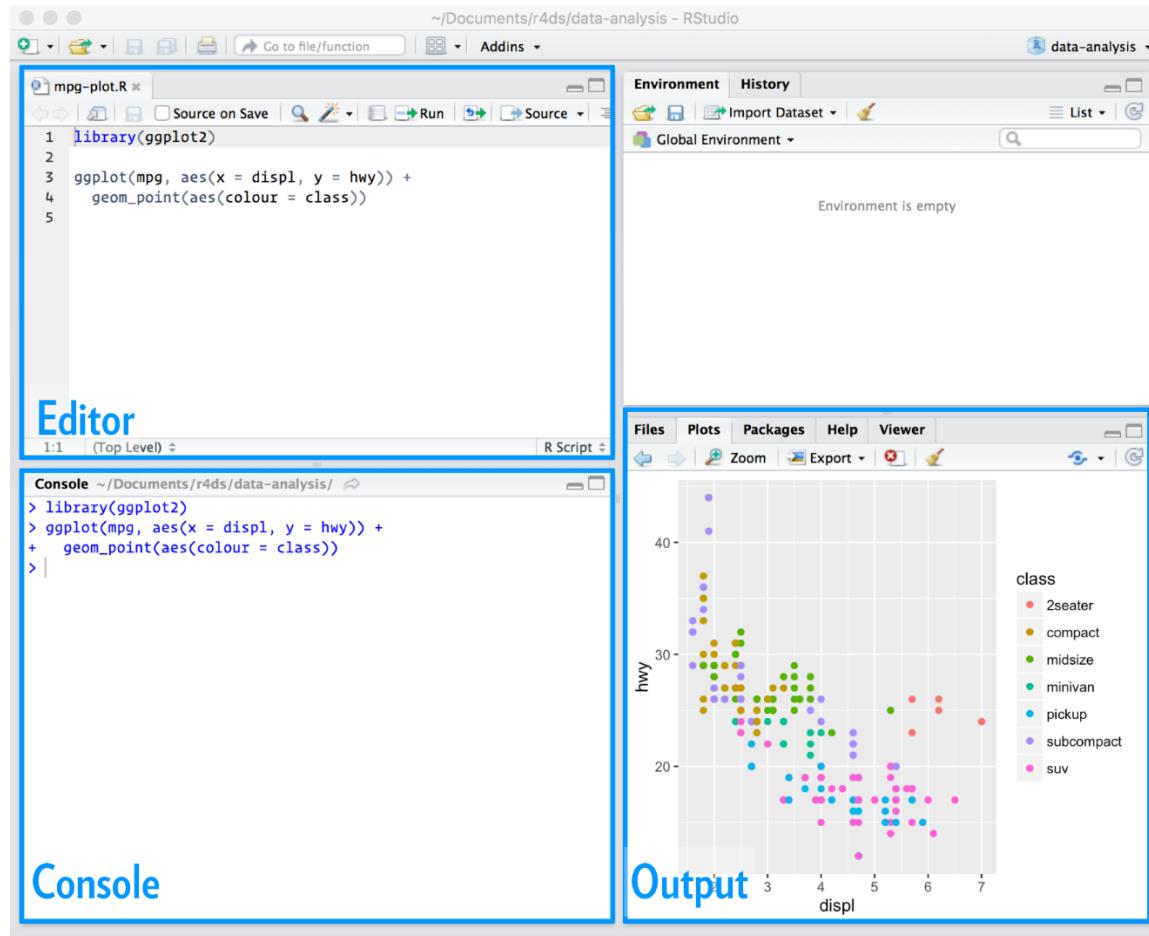
1. Why does this code not work?
2. `my_variable <- 10`
3. `my_variable`  
`#> Error in eval(expr, envir, enclos): object 'my_variable' not found`

Look carefully! (This may seem like an exercise in pointlessness, but training your brain to notice even the tiniest difference will pay off when programming.)

4. Tweak each of the following R commands so that they run correctly:
  5. `library(tidyverse)`
  - 6.
  7. `ggplot(data = mpg) +`
  8.   `geom_point(mapping = aes(x = displ, y = hwy))`
  - 9.
  10. `filter(mpg, cyl = 8)`  
`filter(diamond, carat > 3)`
11. Press Alt + Shift + K. What happens? How can you get to the same place using the menus?

# 6 Workflow: scripts

So far you've been using the console to run code. That's a great place to start, but you'll find it gets cramped pretty quickly as you create more complex ggplot2 graphics and dplyr pipes. To give yourself **more room to work, it's a great idea to use the script editor**. Open it up either by clicking the File menu, and selecting New File, then R script, or using the keyboard shortcut **Cmd/Ctrl + Shift + N**. Now you'll see four panes:



The script editor is a great place to put code you care about. Keep experimenting in the console, but once you have written code that works and does what you want, put it in the script editor. RStudio will automatically save the contents of the editor when you quit RStudio, and will automatically load it when you re-open. **Nevertheless, it's a good idea to save your scripts regularly and to back them up.**

## 6.1 Running code

The script editor is also a great place to build up complex ggplot2 plots or long sequences of dplyr manipulations. The key to using the script editor effectively is to memorise one of the most important keyboard shortcuts: **Cmd/Ctrl + Enter**. This executes the current R expression in the console. For example, take the code below. If your cursor is at █, pressing **Cmd/Ctrl + Enter** will

run the complete command that generates `not_cancelled`. It will also move the cursor to the next statement (beginning with `not_cancelled %>%`). That makes it easy to run your complete script by repeatedly pressing Cmd/Ctrl + Enter.

```
library(dplyr)
library(nycflights13)

not_cancelled <- flights %>%
 filter(!is.na(dep_delay), !is.na(arr_delay))

not_cancelled %>%
 group_by(year, month, day) %>%
 summarise(mean = mean(dep_delay))
```

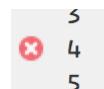
Instead of running expression-by-expression, you can also execute the complete script in one step: Cmd/Ctrl + Shift + S. Doing this **regularly is a great way to check that you've captured all the important parts of your code in the script.**

I recommend that you always start your script with the packages that you need. That way, if you share your code with others, they can easily see what packages they need to install. Note, however, that you should never include `install.packages()` or `setwd()` in a script that you share. **It's very antisocial to change settings on someone else's computer!**

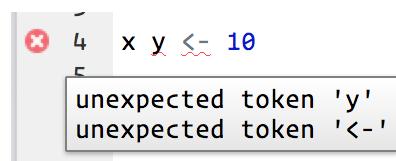
When working through future chapters, I highly recommend starting in the editor and practicing your keyboard shortcuts. Over time, sending code to the console in this way will become so natural that you won't even think about it.

## 6.2 RStudio diagnostics

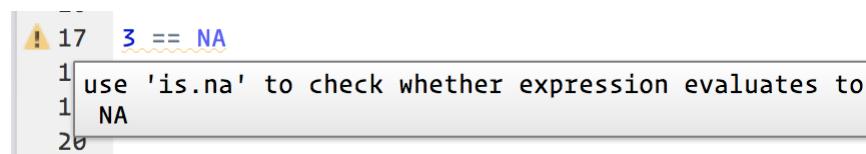
The script editor will also highlight syntax errors with a red squiggly line and a cross in the sidebar:



Hover over the cross to see what the problem is:



RStudio will also let you know about potential problems:



## 6.3 Practice

1. Go to the RStudio Tips twitter account, <https://twitter.com/rstudiotips> and find one tip that looks interesting. Practice using it!
2. What other common mistakes will RStudio diagnostics report? Read <https://support.rstudio.com/hc/en-us/articles/205753617-Code-Diagnostics> to find out.

# 7 Exploratory Data Analysis

## 7.1 Introduction

This chapter will show you how to use visualisation and transformation to explore your data in a systematic way, a task that statisticians call exploratory data analysis, or EDA for short. EDA is an iterative cycle. You:

1. Generate questions about your data.
2. Search for answers by visualising, transforming, and modelling your data.
3. Use what you learn to refine your questions and/or generate new questions.

EDA is not a formal process with a strict set of rules. More than anything, EDA is a state of mind. During the initial phases of EDA you should feel free to investigate every idea that occurs to you. Some of these ideas will pan out, and some will be dead ends. As your exploration continues, you **will home in on a few particularly productive areas that you'll eventually write up and communicate to others.**

EDA is an important part of any data analysis, even if the questions are handed to you on a platter, because you always need to investigate the quality of your data. Data cleaning is just one application of EDA: you ask questions about whether your data meets your expectations or not. **To do data cleaning, you'll need to deploy all the tools of EDA: visualisation, transformation, and modelling.**

### 7.1.1 Prerequisites

In this chapter we'll combine what you've learned about `dplyr` and `ggplot2` to interactively ask questions, answer them with data, and then ask new questions.

```
library(tidyverse)
```

## 7.2 Questions

“There are no routine statistical questions, only questionable statistical routines.” — Sir David Cox

“Far better an approximate answer to the right question, which is often vague, than an exact answer to the wrong question, which can always be made precise.” — John Tukey

Your goal during EDA is to develop an understanding of your data. The easiest way to do this is to use questions as tools to guide your investigation. When you ask a question, the question focuses your attention on a specific part of your dataset and helps you decide which graphs, models, or transformations to make.

EDA is fundamentally a creative process. And like most creative processes, the key to asking **quality** questions is to generate a large **quantity** of questions. It is difficult to ask revealing questions at the start of your analysis because you do not know what insights are contained in your dataset. On the other hand, each new question that you ask will expose you to a new aspect

of your data and increase your chance of making a discovery. You can quickly drill down into the most interesting parts of your data—and develop a set of thought-provoking questions—if you follow up each question with a new question based on what you find.

There is no rule about which questions you should ask to guide your research. However, two types of questions will always be useful for making discoveries within your data. You can loosely word these questions as:

1. What type of variation occurs within my variables?
2. What type of covariation occurs between my variables?

The rest of this chapter will look at these two questions. I'll explain what variation and covariation are, and I'll show you several ways to answer each question. To make the discussion easier, let's define some terms:

- A variable is a quantity, quality, or property that you can measure.
- A value is the state of a variable when you measure it. The value of a variable may change from measurement to measurement.
- An observation is a set of measurements made under similar conditions (you usually make all of the measurements in an observation at the same time and on the same object). An observation will contain several values, each associated with a different variable. I'll sometimes refer to an observation as a data point.
- Tabular data is a set of values, each associated with a variable and an observation. Tabular data is **tidy** if each value is placed in its own “cell”, each variable in its own column, and each observation in its own row.

So far, all of the data that you've seen has been tidy. In real-life, most data isn't tidy, so we'll come back to these ideas again in [tidy data](#).

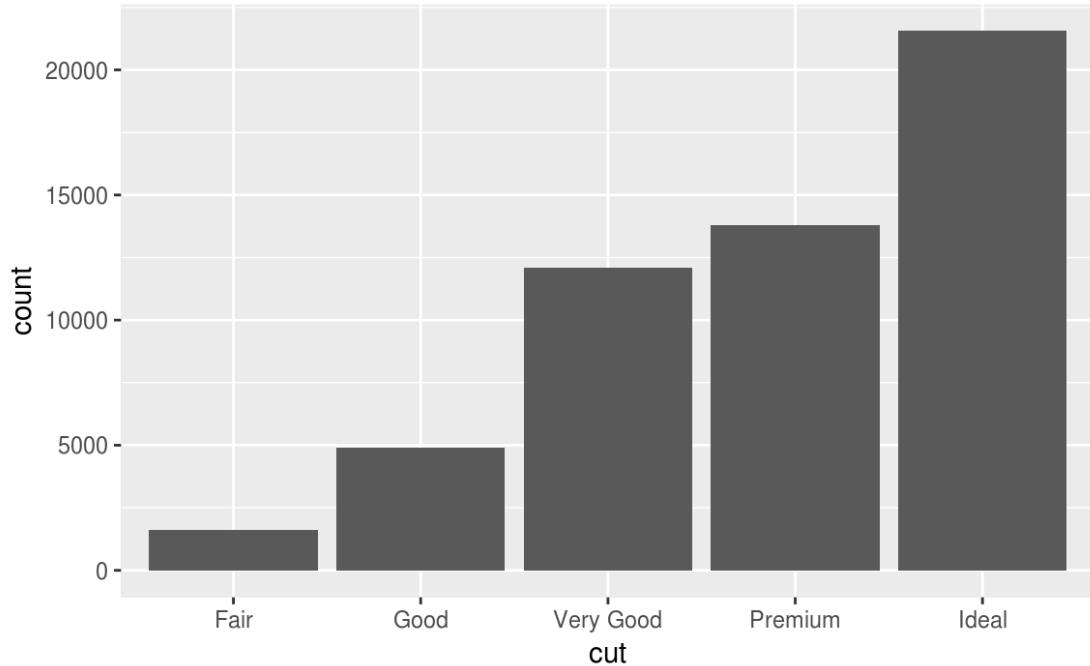
## 7.3 Variation

Variation is the tendency of the values of a variable to change from measurement to measurement. You can see variation easily in real life; if you measure any continuous variable twice, you will get two different results. This is true even if you measure quantities that are constant, like the speed of light. Each of your measurements will include a small amount of error that varies from measurement to measurement. Categorical variables can also vary if you measure across different subjects (e.g. the eye colors of different people), or different times (e.g. the energy levels of an electron at different moments). Every variable has its own pattern of variation, which can reveal interesting information. The best way to understand that pattern is to **visualise the distribution of the variable's values**.

### 7.3.1 Visualising distributions

How you visualise the distribution of a variable will depend on whether the variable is categorical or continuous. A variable is categorical if it can only take one of a small set of values. In R, categorical variables are usually saved as factors or character vectors. To examine the distribution of a categorical variable, use a bar chart:

```
ggplot(data = diamonds) +
 geom_bar(mapping = aes(x = cut))
```

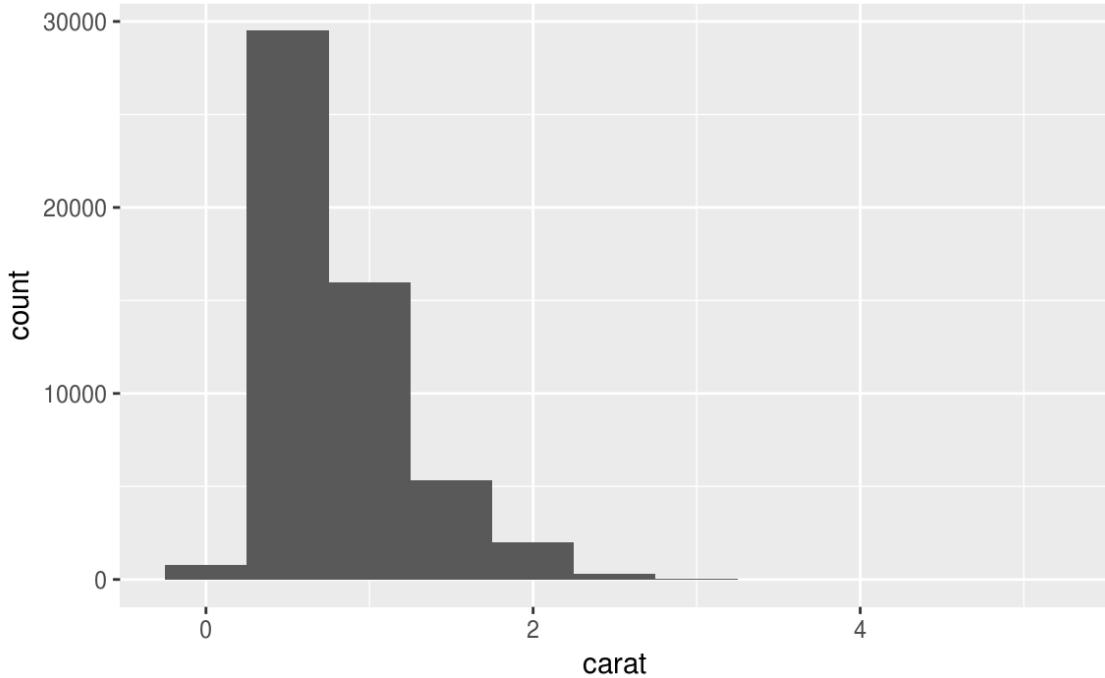


The height of the bars displays how many observations occurred with each x value. You can compute these values manually with `dplyr::count()`:

```
diamonds %>%
 count(cut)
#> # A tibble: 5 × 2
#> cut n
#> <ord> <int>
#> 1 Fair 1610
#> 2 Good 4906
#> 3 Very Good 12082
#> 4 Premium 13791
#> 5 Ideal 21551
```

A variable is continuous if it can take any of an infinite set of ordered values. Numbers and date-times are two examples of continuous variables. To examine the distribution of a continuous variable, use a histogram:

```
ggplot(data = diamonds) +
 geom_histogram(mapping = aes(x = carat), binwidth = 0.5)
```



You can compute this by hand by combining `dplyr::count()` and `ggplot2::cut_width()`:

```
diamonds %>%
 count(cut_width(carat, 0.5))
#> # A tibble: 11 × 2
#> `cut_width(carat, 0.5)` n
#> <fctr> <int>
#> 1 [-0.25,0.25] 785
#> 2 (0.25,0.75] 29498
#> 3 (0.75,1.25] 15977
#> 4 (1.25,1.75] 5313
#> 5 (1.75,2.25] 2002
#> 6 (2.25,2.75] 322
#> # ... with 5 more rows
```

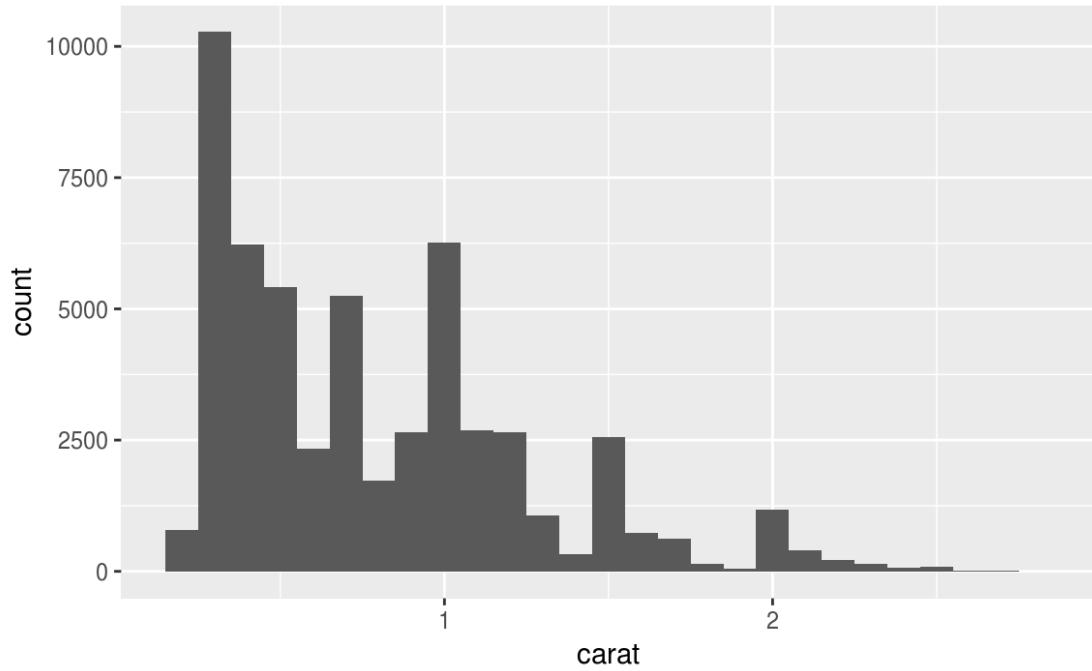
A histogram divides the x-axis into equally spaced bins and then uses the height of a bar to display the number of observations that fall in each bin. In the graph above, the tallest bar shows that almost 30,000 observations have a `carat` value between 0.25 and 0.75, which are the left and right edges of the bar.

You can set the width of the intervals in a histogram with the `binwidth` argument, which is measured in the units of the `x` variable. You should always explore a variety of binwidths when working with histograms, as different binwidths can reveal different patterns. For example, here is how the graph above looks when we zoom into just the diamonds with a size of less than three carats and choose a smaller binwidth.

```
smaller <- diamonds %>%
 filter(carat < 3)

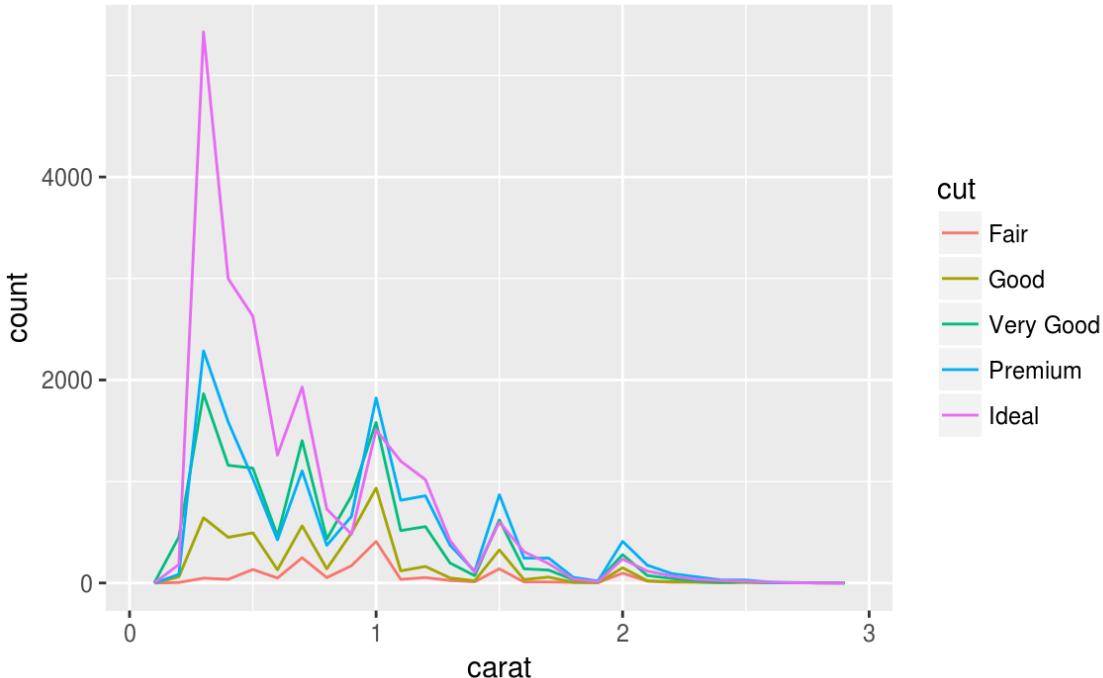
ggplot(data = smaller, mapping = aes(x = carat)) +
```

```
geom_histogram(binwidth = 0.1)
```



If you wish to overlay multiple histograms in the same plot, I recommend using `geom_freqpoly()` instead of `geom_histogram()`. `geom_freqpoly()` performs the same calculation as `geom_histogram()`, but instead of displaying the counts with bars, uses lines instead. **It's much easier to understand overlapping lines than bars.**

```
ggplot(data = smaller, mapping = aes(x = carat, colour = cut)) +
 geom_freqpoly(binwidth = 0.1)
```



There are a few challenges with this type of plot, which we will come back to in [visualising a categorical and a continuous variable](#).

Now that you can visualise variation, what should you look for in your plots? And what type of **follow-up questions should you ask?** I've put together a list below of the most useful types of information that you will find in your graphs, along with some follow-up questions for each type of information. The key to asking good follow-up questions will be to rely on your curiosity (What do you want to learn more about?) as well as your skepticism (How could this be misleading?).

### 7.3.2 Typical values

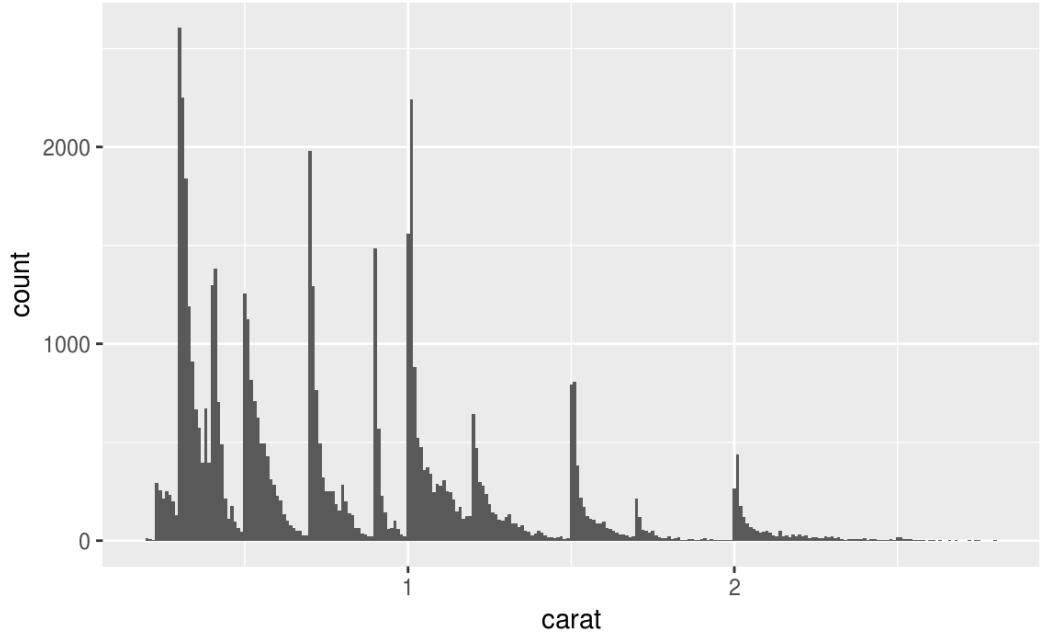
In both bar charts and histograms, tall bars show the common values of a variable, and shorter bars show less-common values. Places that do not have bars reveal values that were not seen in your data. To turn this information into useful questions, look for anything unexpected:

- Which values are the most common? Why?
- Which values are rare? Why? Does that match your expectations?
- Can you see any unusual patterns? What might explain them?

As an example, the histogram below suggests several interesting questions:

- Why are there more diamonds at whole carats and common fractions of carats?
- Why are there more diamonds slightly to the right of each peak than there are slightly to the left of each peak?
- Why are there no diamonds bigger than 3 carats?

```
ggplot(data = smaller, mapping = aes(x = carat)) +
 geom_histogram(binwidth = 0.01)
```

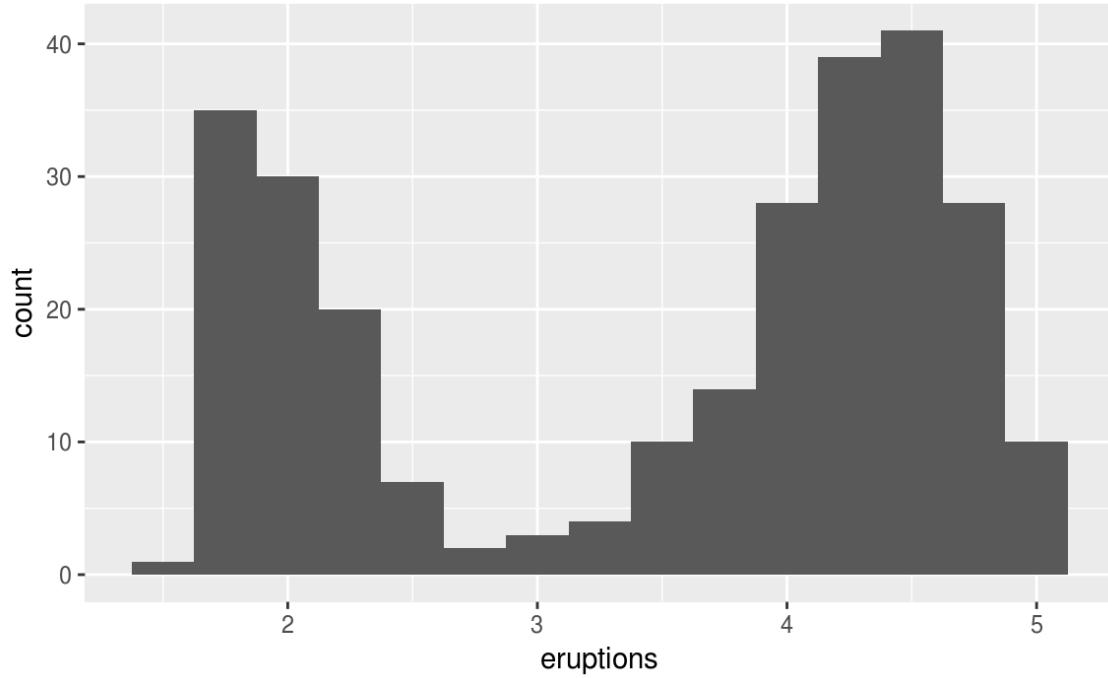


Clusters of similar values suggest that subgroups exist in your data. To understand the subgroups, ask:

- How are the observations within each cluster similar to each other?
- How are the observations in separate clusters different from each other?
- How can you explain or describe the clusters?
- Why might the appearance of clusters be misleading?

The histogram below shows the length (in minutes) of 272 eruptions of the Old Faithful Geyser in Yellowstone National Park. Eruption times appear to be clustered into two groups: there are short eruptions (of around 2 minutes) and long eruptions (4-5 minutes), but little in between.

```
ggplot(data = faithful, mapping = aes(x = eruptions)) +
 geom_histogram(binwidth = 0.25)
```



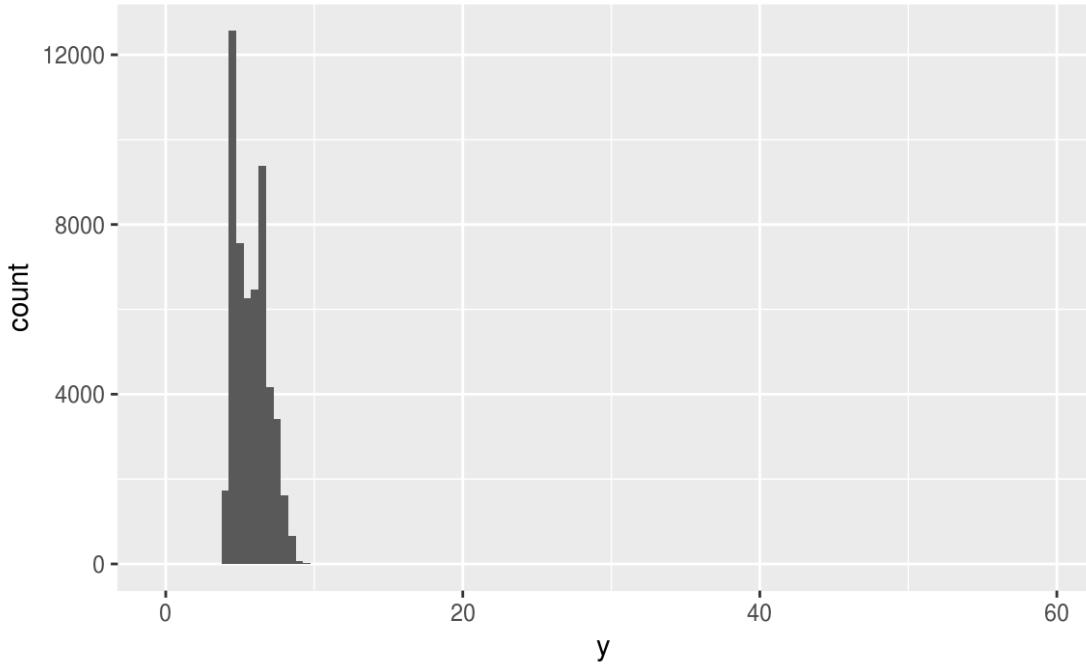
Many of the questions above will prompt you to explore a relationship **between** variables, for **example**, to see if the values of one variable can explain the behavior of another variable. We'll get to that shortly.

### 7.3.3 Unusual values

**Outliers** are observations that are unusual; data points that don't seem to fit the pattern.

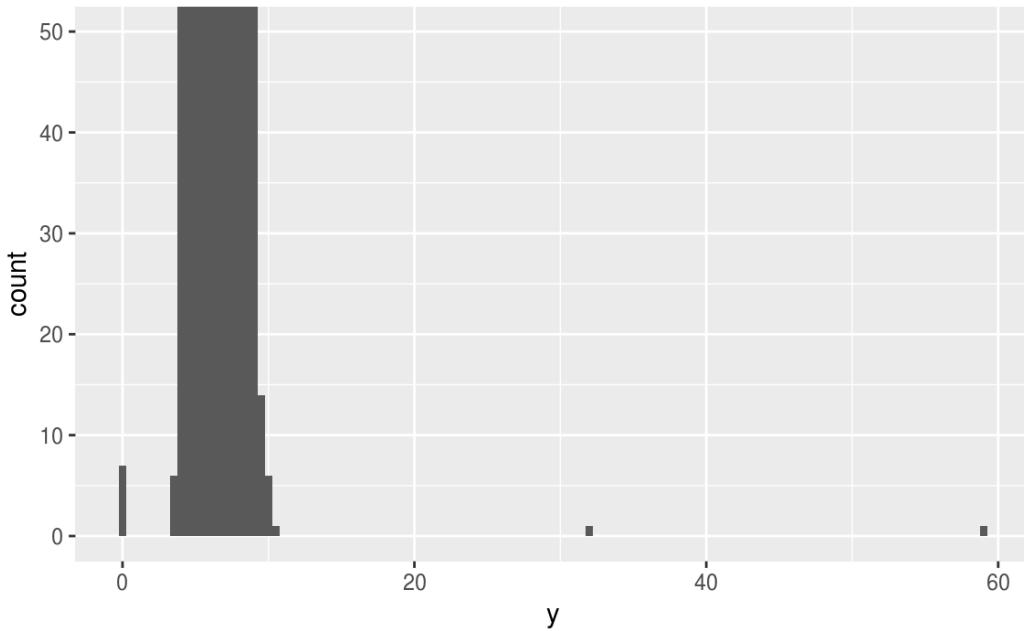
Sometimes outliers are data entry errors; other times outliers suggest important new science. When you have a lot of data, outliers are sometimes difficult to see in a histogram. For example, take the distribution of the `y` variable from the diamonds dataset. The only evidence of outliers is the unusually wide limits on the x-axis.

```
ggplot(diamonds) +
 geom_histogram(mapping = aes(x = y), binwidth = 0.5)
```



There are so many observations in the common bins that the rare bins are so short that you can't see them (although maybe if you stare intently at 0 you'll spot something). To make it easy to see the unusual values, we need to zoom to small values of the y-axis with `coord_cartesian()`:

```
ggplot(diamonds) +
 geom_histogram(mapping = aes(x = y), binwidth = 0.5) +
 coord_cartesian(ylim = c(0, 50))
```



(`coord_cartesian()` also has an `xlim()` argument for when you need to zoom into the x-axis. `ggplot2` also has `xlim()` and `ylim()` functions that work slightly differently: they throw away the data outside the limits.)

This allows us to see that there are three unusual values: 0, ~30, and ~60. We pluck them out with `dplyr`:

```
unusual <- diamonds %>%
 filter(y < 3 | y > 20) %>%
 select(price, x, y, z) %>%
 arrange(y)
unusual
#> # A tibble: 9 × 4
#> price x y z
#> <int> <dbl> <dbl> <dbl>
#> 1 5139 0.00 0.0 0.00
#> 2 6381 0.00 0.0 0.00
#> 3 12800 0.00 0.0 0.00
#> 4 15686 0.00 0.0 0.00
#> 5 18034 0.00 0.0 0.00
#> 6 2130 0.00 0.0 0.00
#> 7 2130 0.00 0.0 0.00
#> 8 2075 5.15 31.8 5.12
#> 9 12210 8.09 58.9 8.06
```

The `y` variable measures one of the three dimensions of these diamonds, in mm. We know that **diamonds can't have a width of 0mm, so these values must be incorrect. We might also suspect** that measurements of 32mm and 59mm are implausible: those diamonds are over an inch long, **but don't cost hundreds of thousands of dollars!**

**It's good practice to repeat your analysis with and without the outliers. If they have minimal effect on the results, and you can't figure out why they're there, it's reasonable to replace them with missing values, and move on. However, if they have a substantial effect on your results, you shouldn't drop them without justification. You'll need to figure out what caused them (e.g. a data entry error) and disclose that you removed them in your write-up.**

### 7.3.4 Exercises

1. Explore the distribution of each of the `x`, `y`, and `z` variables in `diamonds`. What do you learn? Think about a diamond and how you might decide which dimension is the length, width, and depth.
2. Explore the distribution of `price`. Do you discover anything unusual or surprising? (Hint: Carefully think about the `binwidth` and make sure you try a wide range of values.)
3. How many diamonds are 0.99 carat? How many are 1 carat? What do you think is the cause of the difference?
4. Compare and contrast `coord_cartesian()` vs `xlim()` or `ylim()` when zooming in on a histogram. What happens if you leave `binwidth` unset? What happens if you try and zoom so only half a bar shows?

## 7.4 Missing values

If you've encountered unusual values in your dataset, and simply want to move on to the rest of your analysis, you have two options.

1. Drop the entire row with the strange values:
2. 

```
diamonds2 <- diamonds %>%
 filter(between(y, 3, 20))
```

I don't recommend this option because just because one measurement is invalid, doesn't mean all the measurements are. Additionally, if you have low quality data, by time that you've applied this approach to every variable you might find that you don't have any data left!

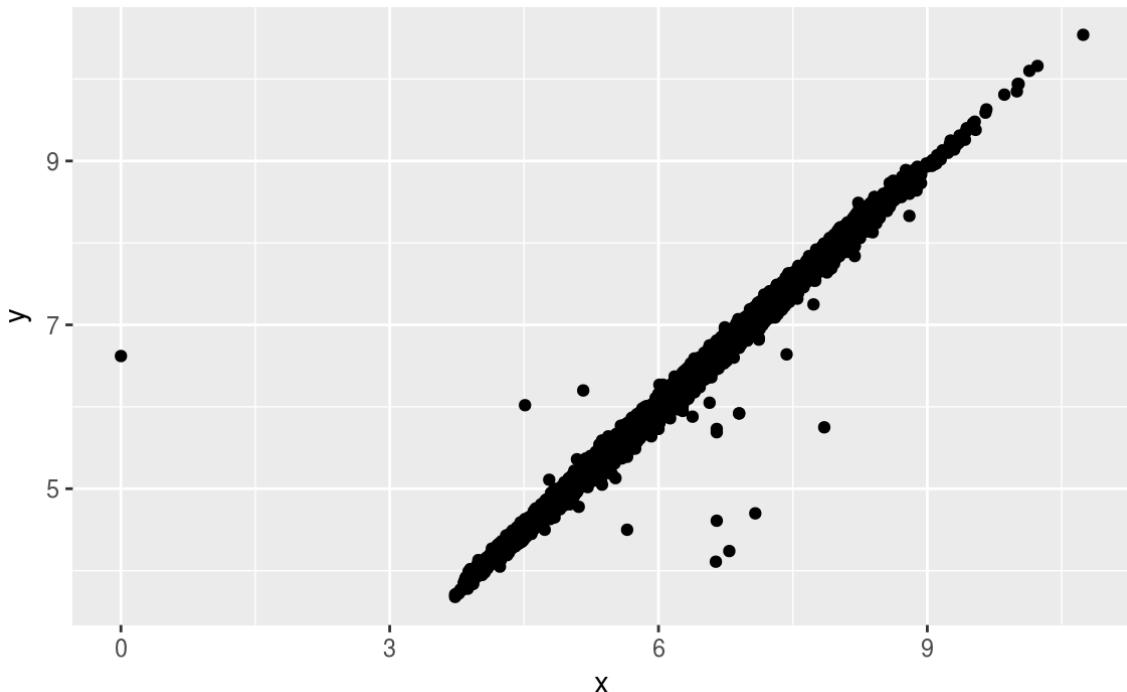
3. Instead, I recommend replacing the unusual values with missing values. The easiest way to do this is to use `mutate()` to replace the variable with a modified copy. You can use the `ifelse()` function to replace unusual values with `NA`:
4. 

```
diamonds2 <- diamonds %>%
 mutate(y = ifelse(y < 3 | y > 20, NA, y))
```

`ifelse()` has three arguments. The first argument `test` should be a logical vector. The result will contain the value of the second argument, `yes`, when `test` is `TRUE`, and the value of the third argument, `no`, when it is false.

Like R, ggplot2 subscribes to the philosophy that missing values should never silently go missing. It's not obvious where you should plot missing values, so ggplot2 doesn't include them in the plot, but it does warn that they've been removed:

```
ggplot(data = diamonds2, mapping = aes(x = x, y = y)) +
 geom_point()
#> Warning: Removed 9 rows containing missing values (geom_point).
```

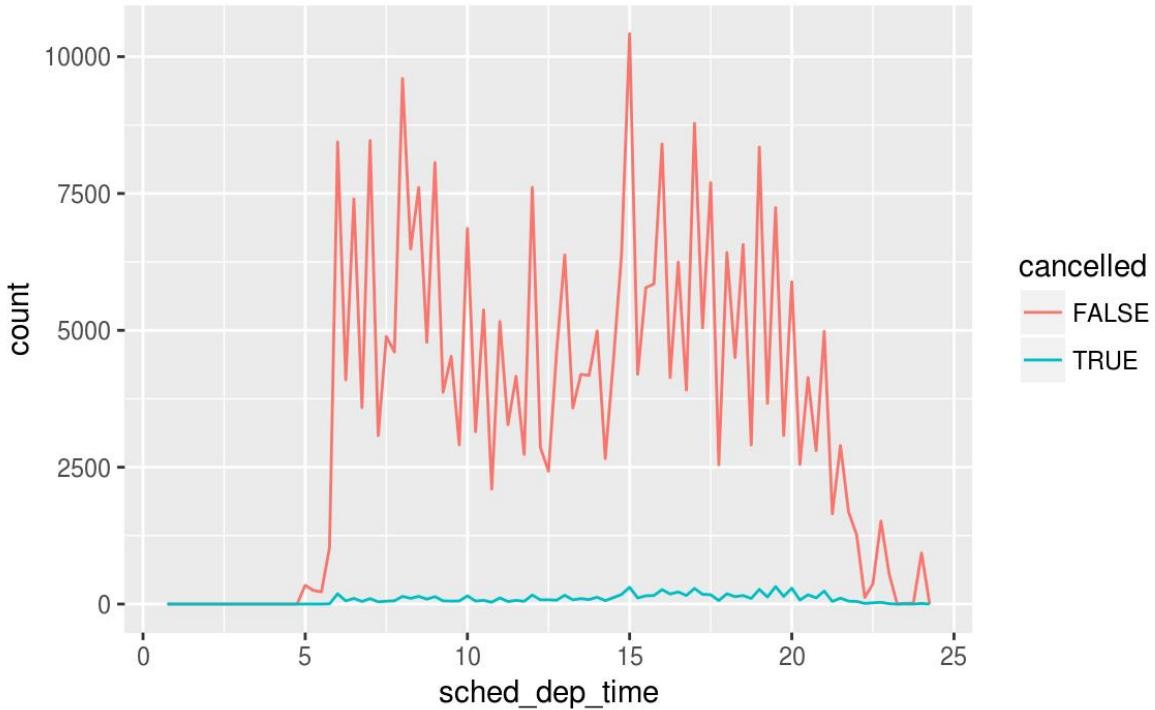


To suppress that warning, set `na.rm = TRUE`:

```
ggplot(data = diamonds2, mapping = aes(x = x, y = y)) +
 geom_point(na.rm = TRUE)
```

Other times you want to understand what makes observations with missing values different to observations with recorded values. For example, in `nycflights13::flights`, missing values in the `dep_time` variable indicate that the flight was cancelled. So you might want to compare the scheduled departure times for cancelled and non-cancelled times. You can do this by making a new variable with `is.na()`.

```
nycflights13::flights %>%
 mutate(
 cancelled = is.na(dep_time),
 sched_hour = sched_dep_time %/% 100,
 sched_min = sched_dep_time %% 100,
 sched_dep_time = sched_hour + sched_min / 60
) %>%
 ggplot(mapping = aes(sched_dep_time)) +
 geom_freqpoly(mapping = aes(colour = cancelled), binwidth = 1/4)
```



However this plot isn't great because there are many more non-cancelled flights than cancelled flights. In the next section we'll explore some techniques for improving this comparison.

#### 7.4.1 Exercises

1. What happens to missing values in a histogram? What happens to missing values in a bar chart? Why is there a difference?

2. What does `na.rm = TRUE` do in `mean()` and `sum()`?

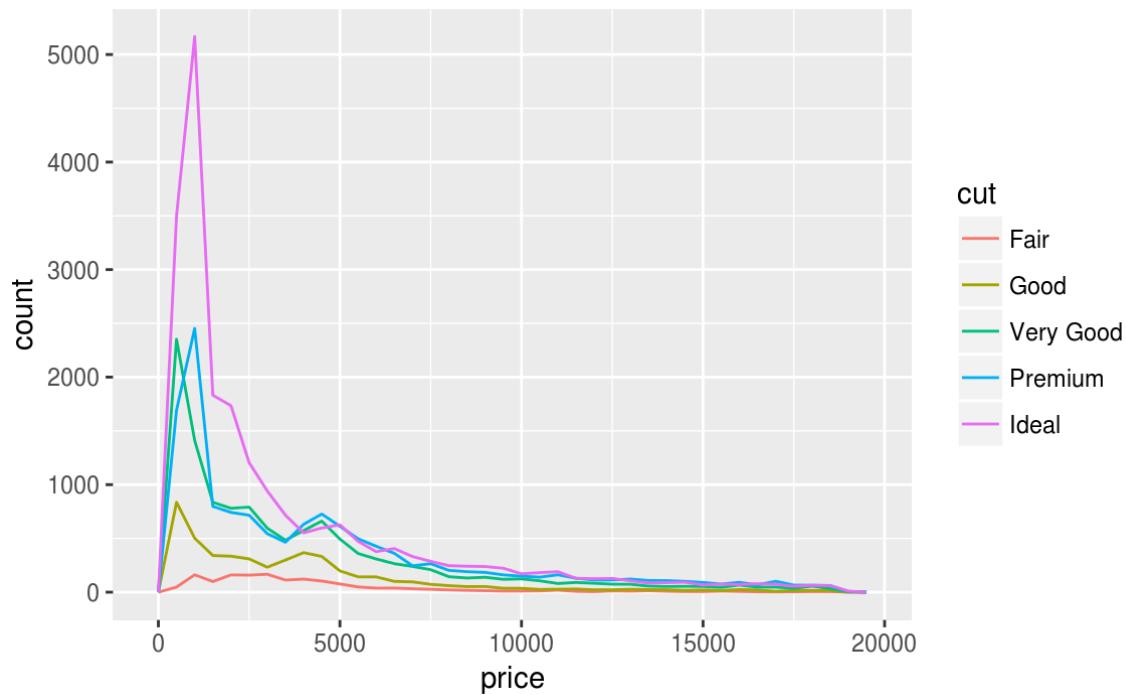
## 7.5 Covariation

If variation describes the behavior **within** a variable, covariation describes the behavior **between** variables. Covariation is the tendency for the values of two or more variables to vary together in a related way. The best way to spot covariation is to visualise the relationship between two or more variables. How you do that should again depend on the type of variables involved.

### 7.5.1 A categorical and continuous variable

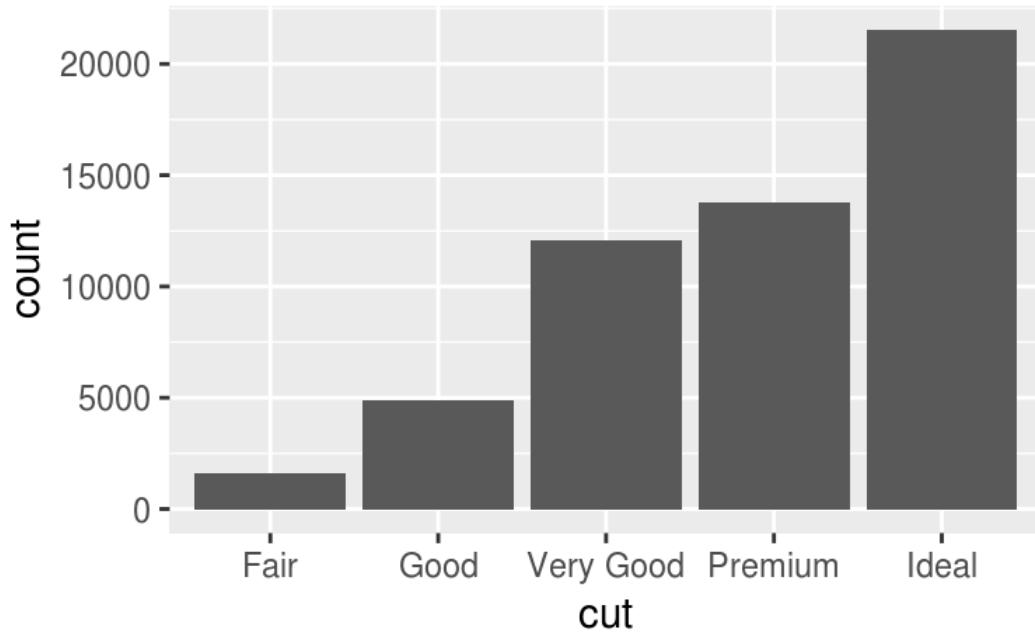
**It's common to** want to explore the distribution of a continuous variable broken down by a categorical variable, as in the previous frequency polygon. The default appearance of `geom_freqpoly()` is not that useful for that sort of comparison because the height is given by the count. That means if one of the groups is much smaller than the others, it's hard to see the differences in shape. For example, let's explore how the price of a diamond varies with its quality:

```
ggplot(data = diamonds, mapping = aes(x = price)) +
 geom_freqpoly(mapping = aes(colour = cut), binwidth = 500)
```



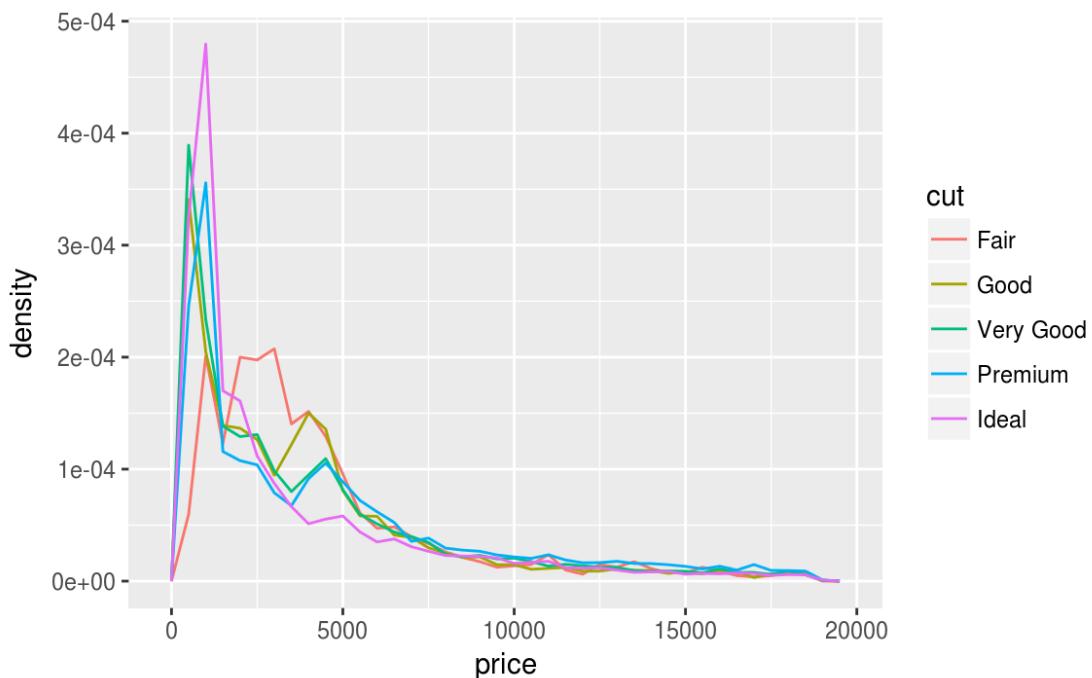
It's hard to see the difference in distribution because the overall counts differ so much:

```
ggplot(diamonds) +
 geom_bar(mapping = aes(x = cut))
```



To make the comparison easier we need to swap what is displayed on the y-axis. Instead of **displaying count**, we'll **display density**, which is the count standardised so that the area under each frequency polygon is one.

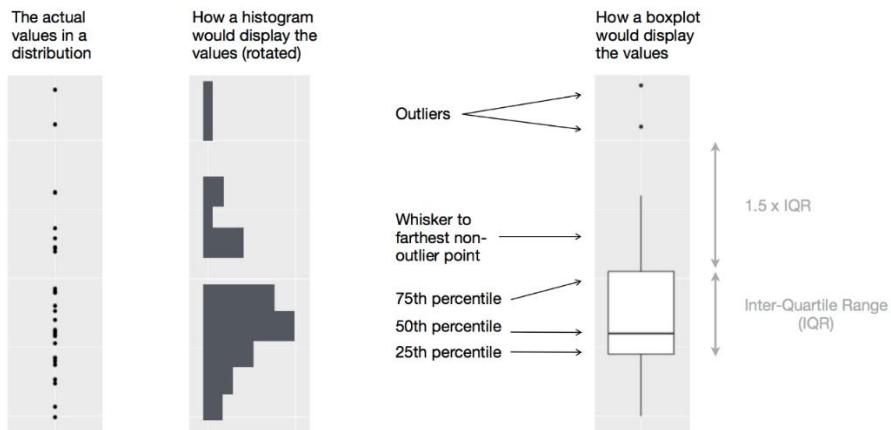
```
ggplot(data = diamonds, mapping = aes(x = price, y = ..density..)) +
 geom_freqpoly(mapping = aes(colour = cut), binwidth = 500)
```



There's something rather surprising about this plot - it appears that fair diamonds (the lowest quality) have the highest average price! But maybe that's because frequency polygons are a little hard to interpret - there's a lot going on in this plot.

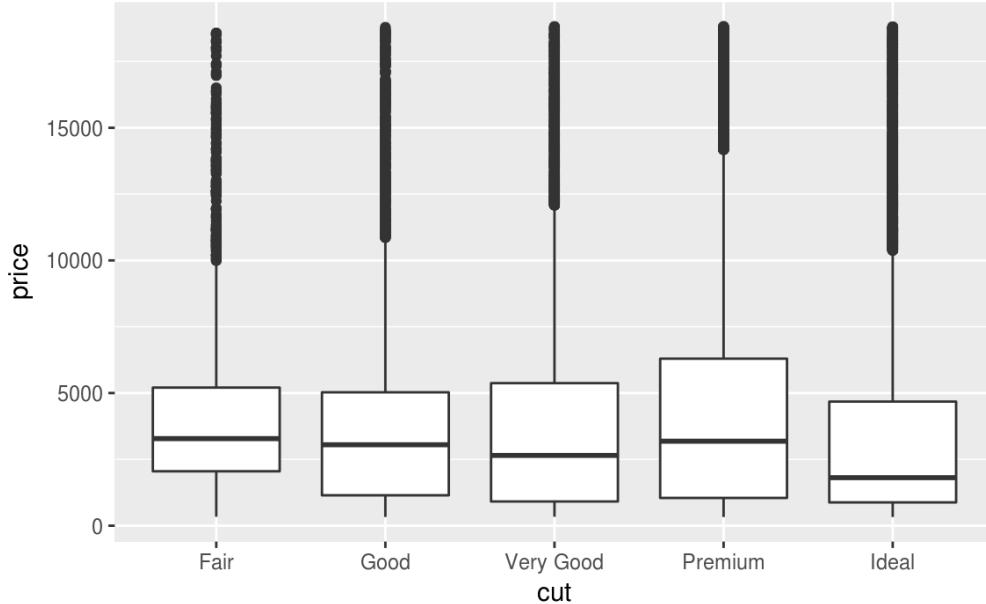
Another alternative to display the distribution of a continuous variable broken down by a categorical variable is the boxplot. A boxplot is a type of visual shorthand for a distribution of values that is popular among statisticians. Each boxplot consists of:

- A box that stretches from the 25th percentile of the distribution to the 75th percentile, a distance known as the interquartile range (IQR). In the middle of the box is a line that displays the median, i.e. 50th percentile, of the distribution. These three lines give you a sense of the spread of the distribution and whether or not the distribution is symmetric about the median or skewed to one side.
- Visual points that display observations that fall more than 1.5 times the IQR from either edge of the box. These outlying points are unusual so are plotted individually.
- A line (or whisker) that extends from each end of the box and goes to the farthest non-outlier point in the distribution.



Let's take a look at the distribution of price by cut using `geom_boxplot()`:

```
ggplot(data = diamonds, mapping = aes(x = cut, y = price)) +
 geom_boxplot()
```

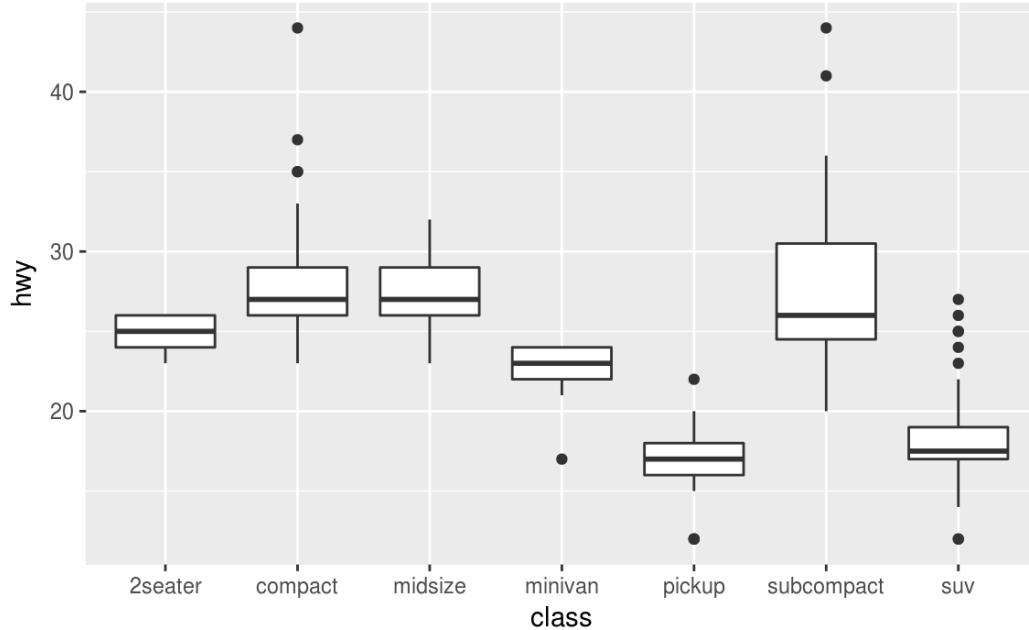


We see much less information about the distribution, but the boxplots are much more compact so we can more easily compare them (and fit more on one plot). It supports the counterintuitive **finding that better quality diamonds are cheaper on average! In the exercises, you'll be challenged to figure out why.**

`cut` is an ordered factor: fair is worse than good, which is worse than very good and so on. Many **categorical variables don't have such an intrinsic order, so you might want to reorder them to** make a more informative display. One way to do that is with the `reorder()` function.

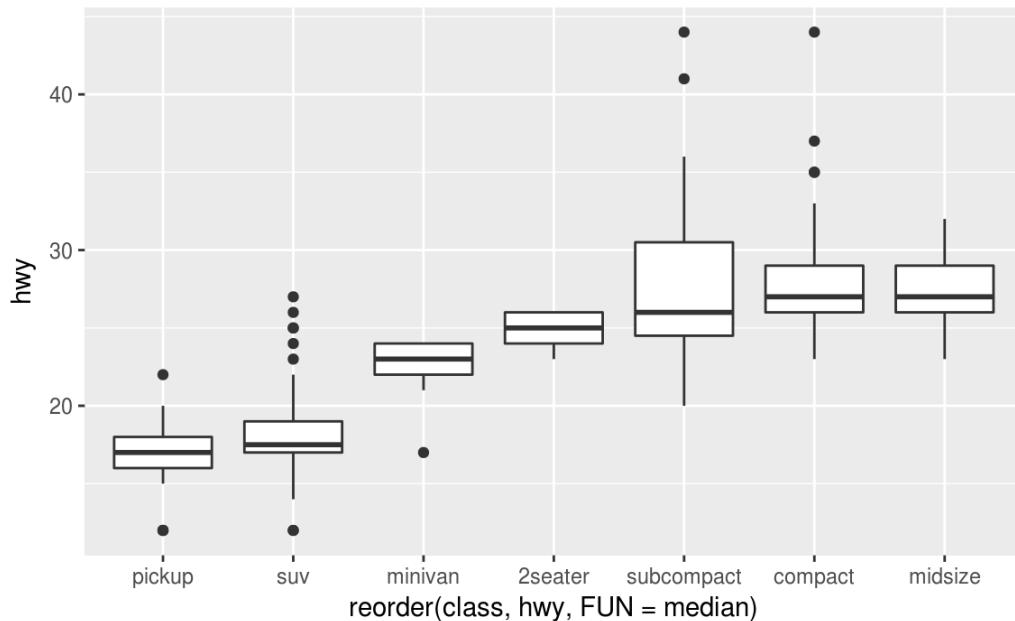
For example, take the `class` variable in the `mpg` dataset. You might be interested to know how highway mileage varies across classes:

```
ggplot(data = mpg, mapping = aes(x = class, y = hwy)) +
 geom_boxplot()
```



To make the trend easier to see, we can reorder `class` based on the median value of `hwy`:

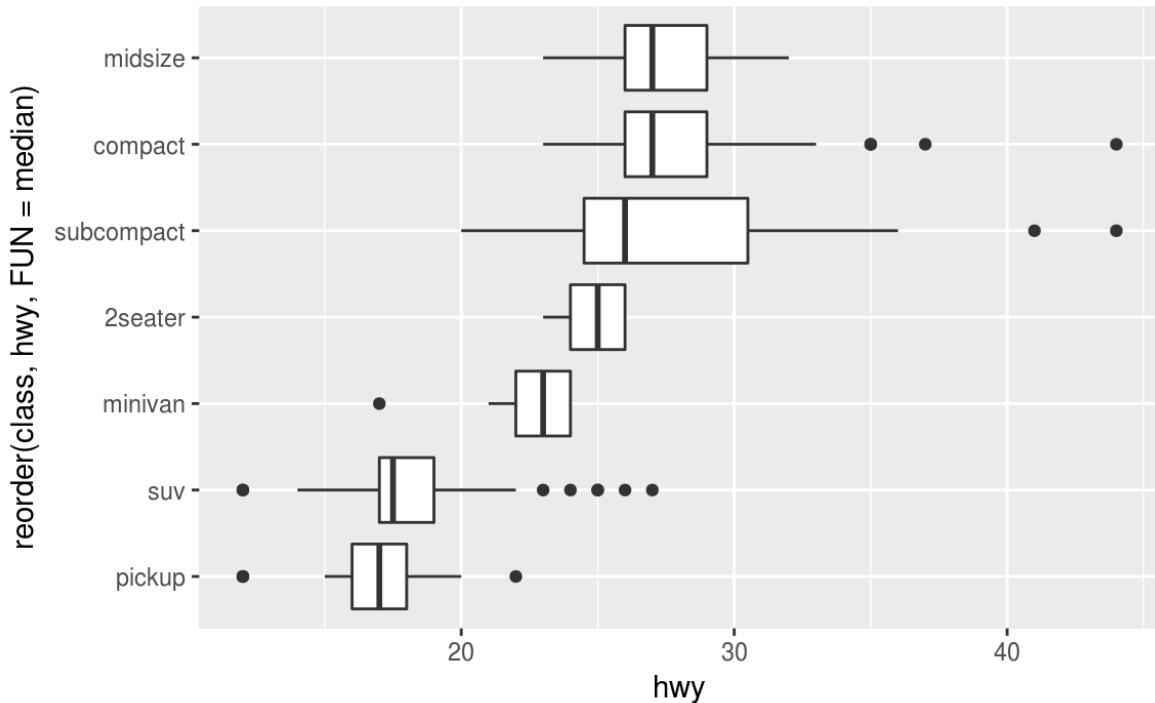
```
ggplot(data = mpg) +
 geom_boxplot(mapping = aes(x = reorder(class, hwy, FUN = median), y = hwy))
```



If you have long variable names, `geom_boxplot()` will work better if you flip it 90°. You can do that with `coord_flip()`.

```
ggplot(data = mpg) +
```

```
geom_boxplot(mapping = aes(x = reorder(class, hwy, FUN = median), y =
hwy)) +
coord_flip()
```



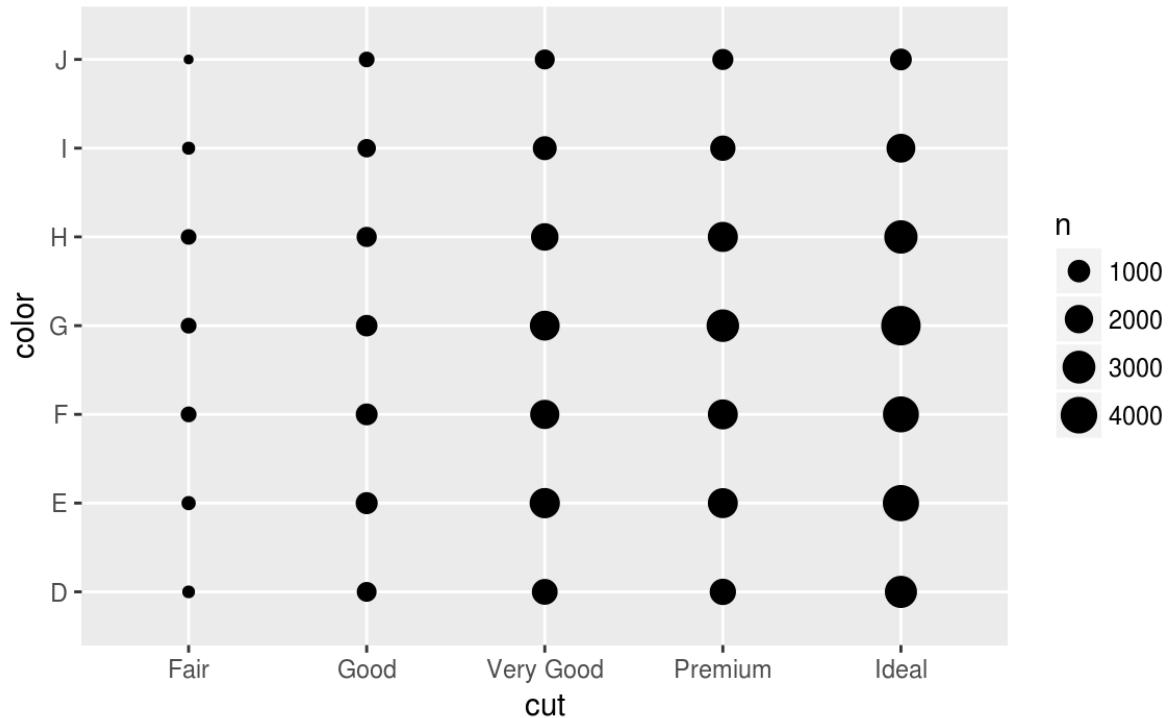
#### 7.5.1.1 Exercises

1. Use what you've learned to improve the visualisation of the departure times of cancelled vs. non-cancelled flights.
2. What variable in the diamonds dataset is most important for predicting the price of a diamond? How is that variable correlated with cut? Why does the combination of those two relationships lead to lower quality diamonds being more expensive?
3. Install the ggstance package, and create a horizontal boxplot. How does this compare to using `coord_flip()`?
4. One problem with boxplots is that they were developed in an era of much smaller datasets **and tend to display a prohibitively large number of “outlying values”**. One approach to remedy this problem is the letter value plot. Install the lvplot package, and try using `geom_lv()` to display the distribution of price vs cut. What do you learn? How do you interpret the plots?
5. Compare and contrast `geom_violin()` with a faceted `geom_histogram()`, or a coloured `geom_freqpoly()`. What are the pros and cons of each method?
6. If you have a small dataset, it's sometimes useful to use `geom_jitter()` to see the relationship between a continuous and categorical variable. The ggbeeswarm package provides a number of methods similar to `geom_jitter()`. List them and briefly describe what each one does.

#### 7.5.2 Two categorical variables

To visualise the covariation between categorical variables, you'll need to count the number of observations for each combination. One way to do that is to rely on the built-in `geom_count()`:

```
ggplot(data = diamonds) +
 geom_count(mapping = aes(x = cut, y = color))
```



The size of each circle in the plot displays how many observations occurred at each combination of values. Covariation will appear as a strong correlation between specific x values and specific y values.

Another approach is to compute the count with dplyr:

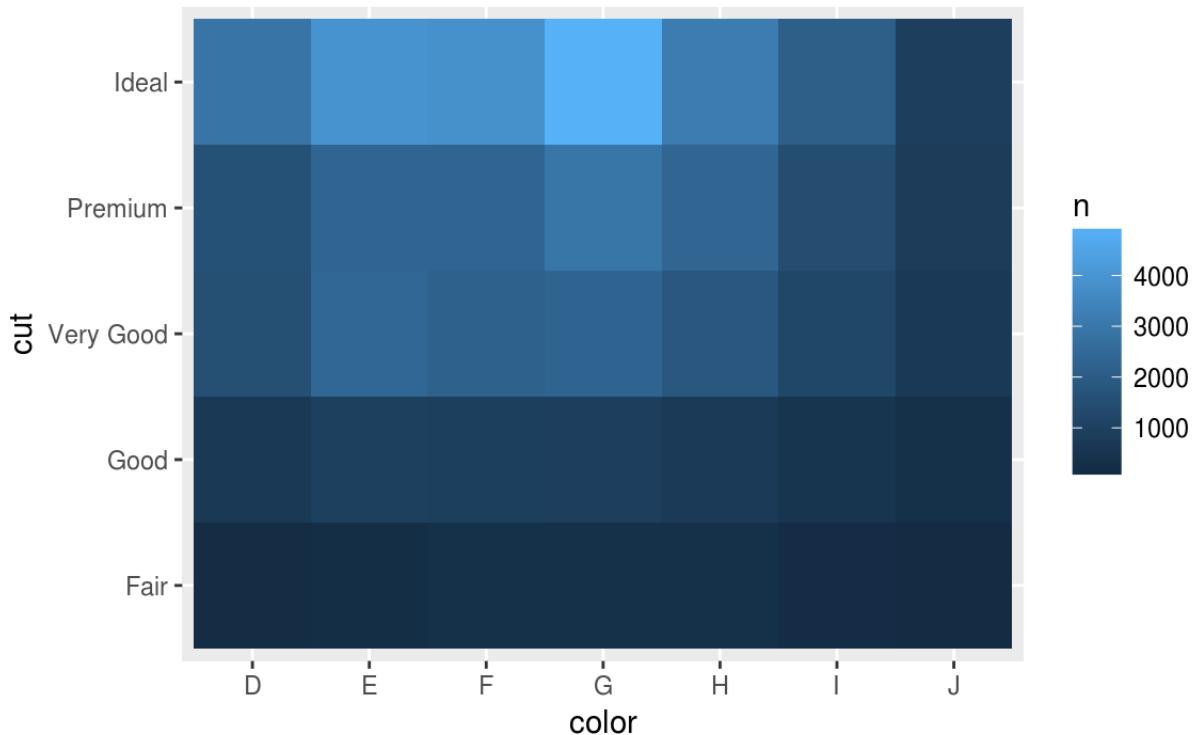
```
diamonds %>%
 count(color, cut)
#> Source: local data frame [35 x 3]
#> Groups: color [?]
#>
#> color cut n
#> <ord> <ord> <int>
#> 1 D Fair 163
#> 2 D Good 662
#> 3 D Very Good 1513
#> 4 D Premium 1603
#> 5 D Ideal 2834
#> 6 E Fair 224
#> # ... with 29 more rows
```

Then visualise with `geom_tile()` and the fill aesthetic:

```

diamonds %>%
 count(color, cut) %>%
 ggplot(mapping = aes(x = color, y = cut)) +
 geom_tile(mapping = aes(fill = n))

```



If the categorical variables are unordered, you might want to use the `seriation` package to simultaneously reorder the rows and columns in order to more clearly reveal interesting patterns. For larger plots, you might want to try the `d3heatmap` or `heatmaply` packages, which create interactive plots.

### 7.5.2.1 Exercises

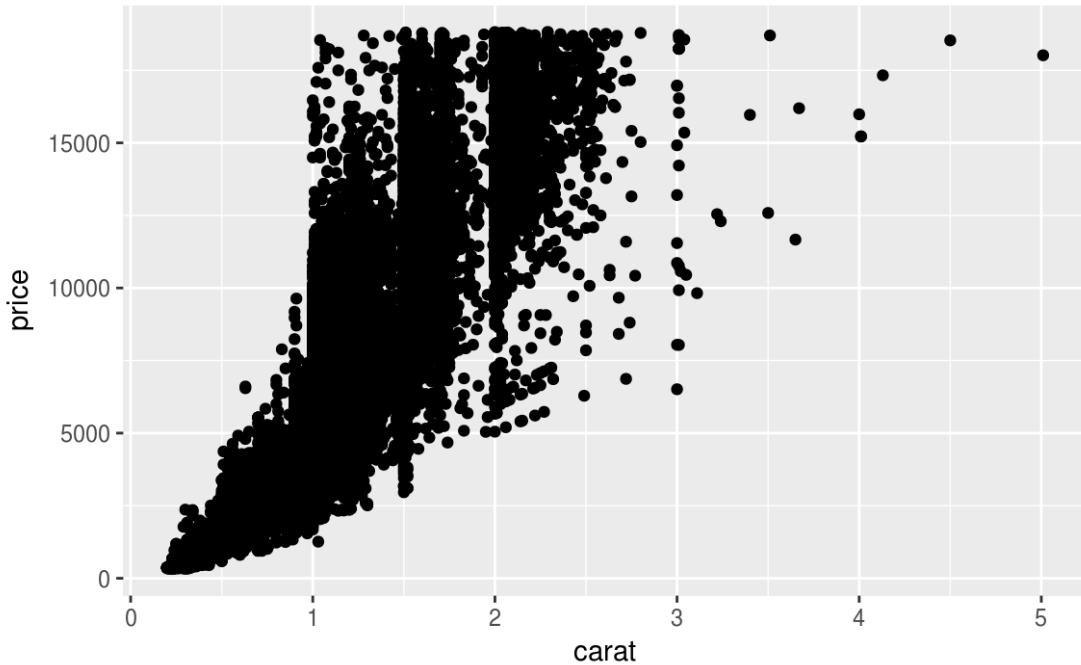
1. How could you rescale the count dataset above to more clearly show the distribution of `cut` within `color`, or `color` within `cut`?
2. Use `geom_tile()` together with `dplyr` to explore how average flight delays vary by destination and month of year. What makes the plot difficult to read? How could you improve it?
3. Why is it slightly better to use `aes(x = color, y = cut)` rather than `aes(x = cut, y = color)` in the example above?

### 7.5.3 Two continuous variables

**You've already seen one great way to visualise the covariation between two continuous variables:** draw a scatterplot with `geom_point()`. You can see covariation as a pattern in the points. For example, you can see an exponential relationship between the carat size and price of a diamond.

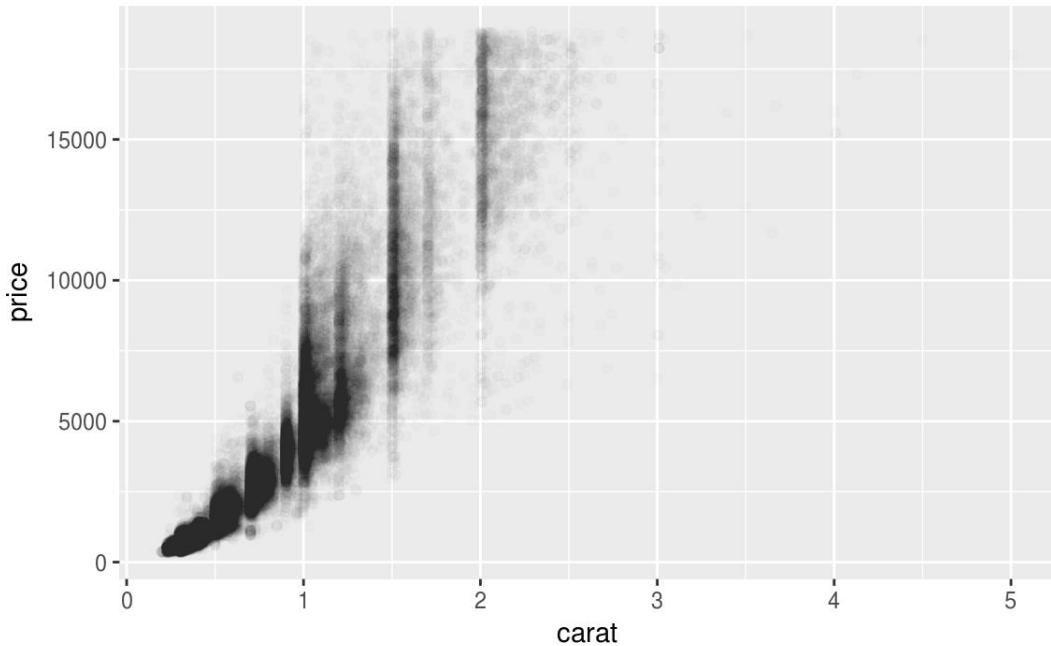
```
ggplot(data = diamonds) +
```

```
geom_point(mapping = aes(x = carat, y = price))
```



Scatterplots become less useful as the size of your dataset grows, because points begin to **overplot**, and pile up into **areas of uniform black (as above)**. You've already seen one way to fix the problem: using the `alpha` aesthetic to add transparency.

```
ggplot(data = diamonds) +
 geom_point(mapping = aes(x = carat, y = price), alpha = 1 / 100)
```

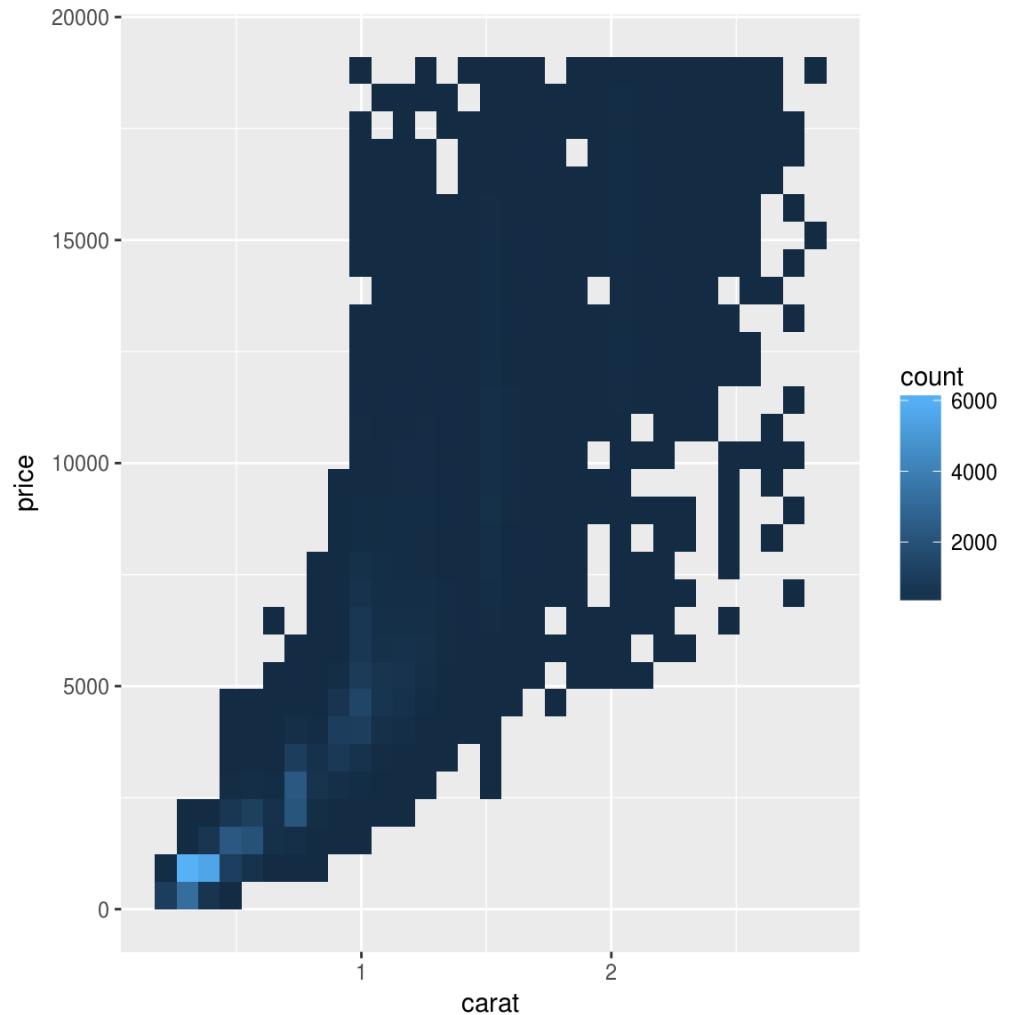


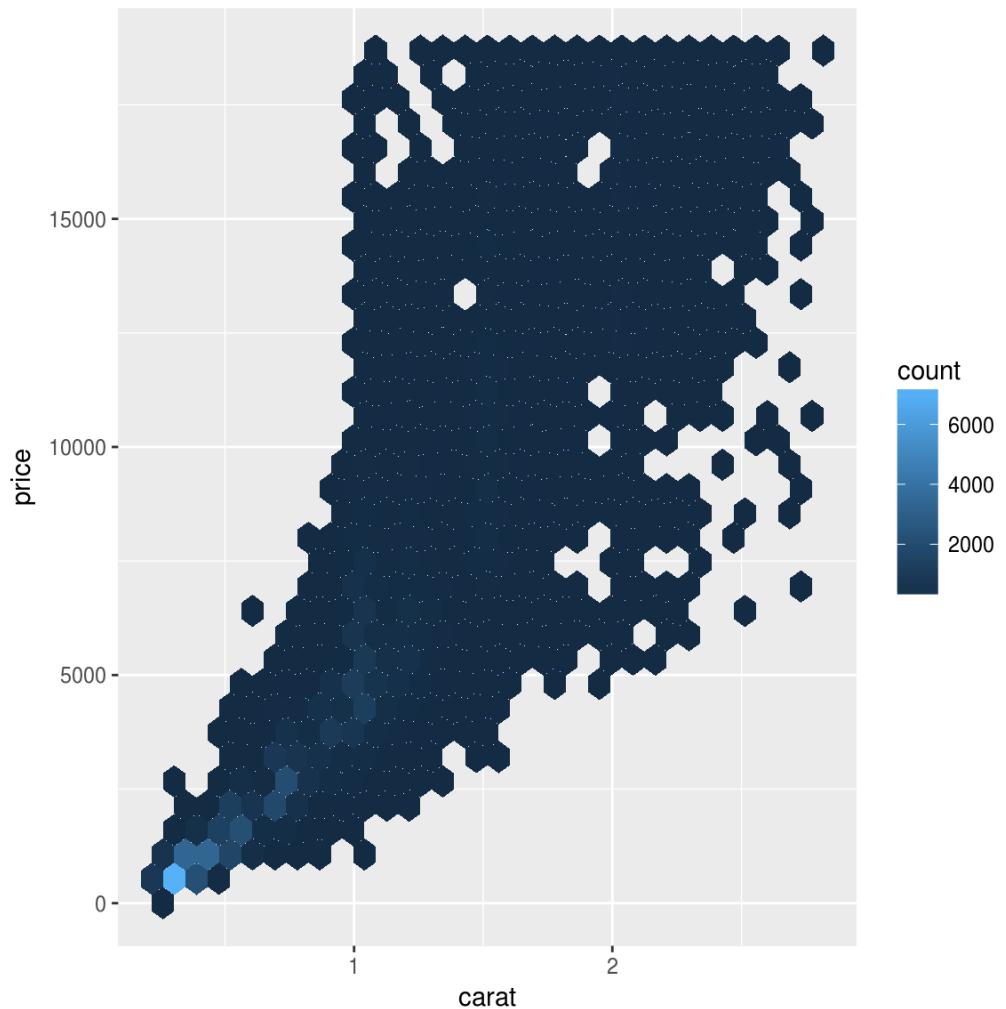
But using transparency can be challenging for very large datasets. Another solution is to use bin. Previously you used `geom_histogram()` and `geom_freqpoly()` to bin in one dimension. Now **you'll learn how to use** `geom_bin2d()` and `geom_hex()` to bin in two dimensions.

`geom_bin2d()` and `geom_hex()` divide the coordinate plane into 2d bins and then use a fill color to display how many points fall into each bin. `geom_bin2d()` creates rectangular bins. `geom_hex()` creates hexagonal bins. You will need to install the `hexbin` package to use `geom_hex()`.

```
ggplot(data = smaller) +
 geom_bin2d(mapping = aes(x = carat, y = price))

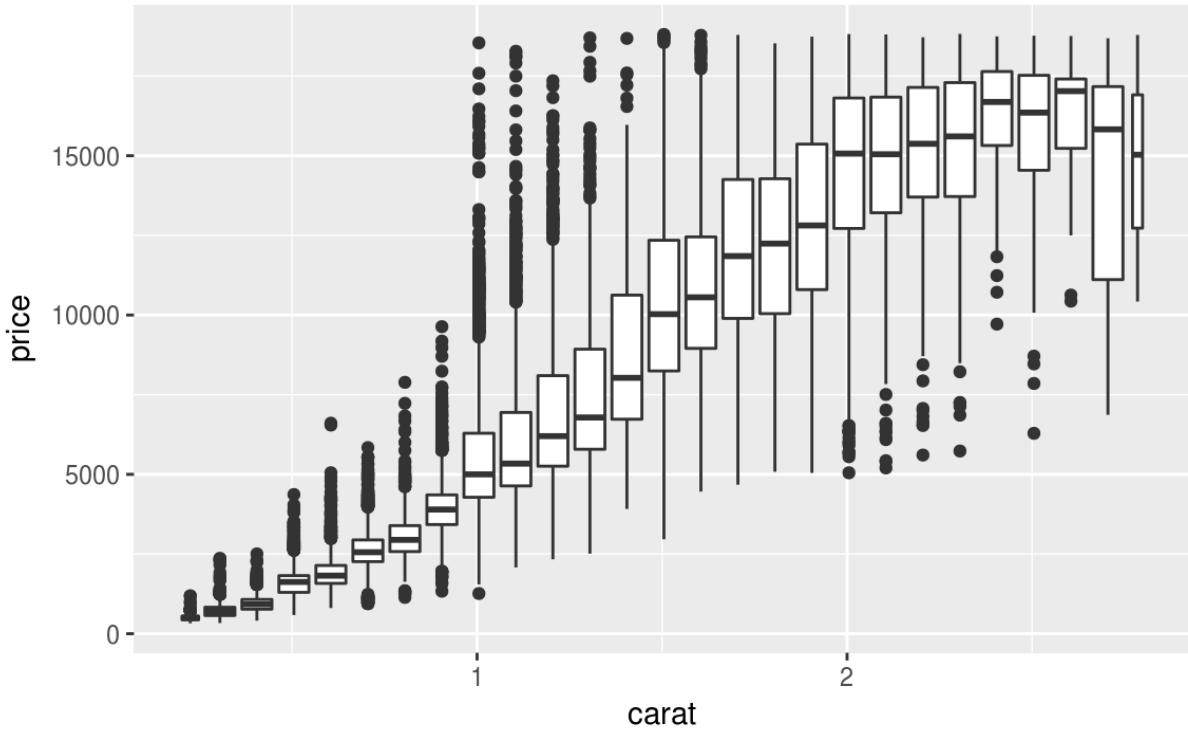
install.packages("hexbin")
ggplot(data = smaller) +
 geom_hex(mapping = aes(x = carat, y = price))
```





Another option is to bin one continuous variable so it acts like a categorical variable. Then you can use one of the techniques for visualising the combination of a categorical and a continuous variable that you learned about. For example, you could bin `carat` and then for each group, display a boxplot:

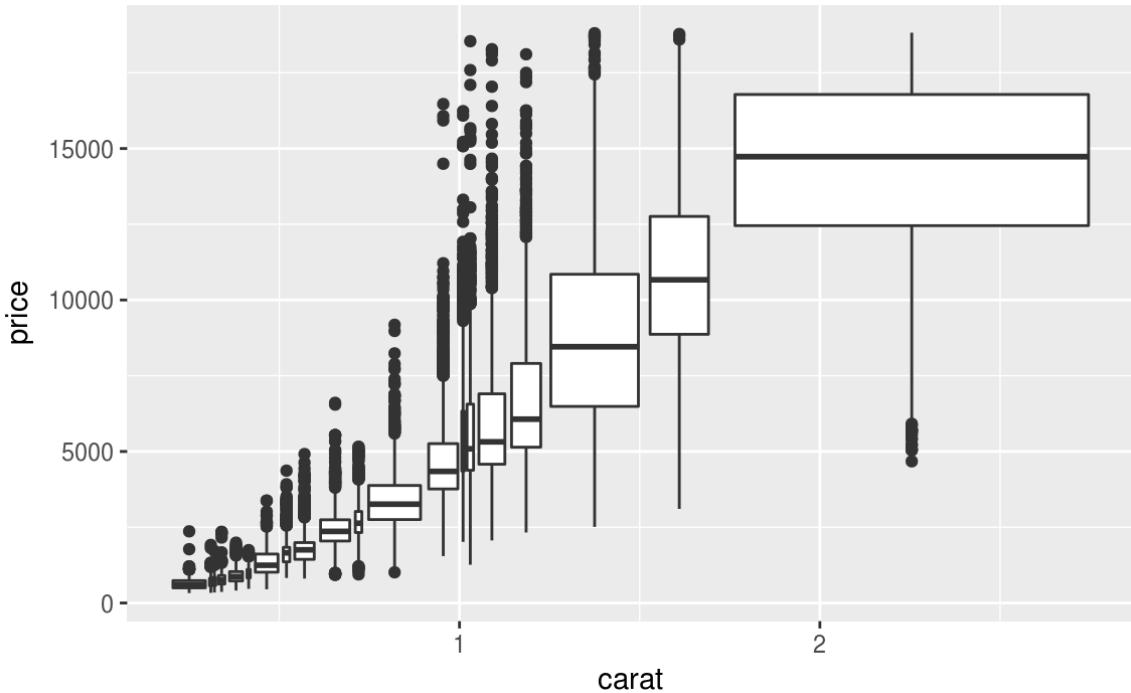
```
ggplot(data = smaller, mapping = aes(x = carat, y = price)) +
 geom_boxplot(mapping = aes(group = cut_width(carat, 0.1)))
```



`cut_width(x, width)`, as used above, divides `x` into bins of width `width`. By default, boxplots look roughly the same (apart from number of outliers) regardless of how many **observations there are, so it's difficult to tell that each boxplot summarises a different number of points**. One way to show that is to make the width of the boxplot proportional to the number of points with `varwidth = TRUE`.

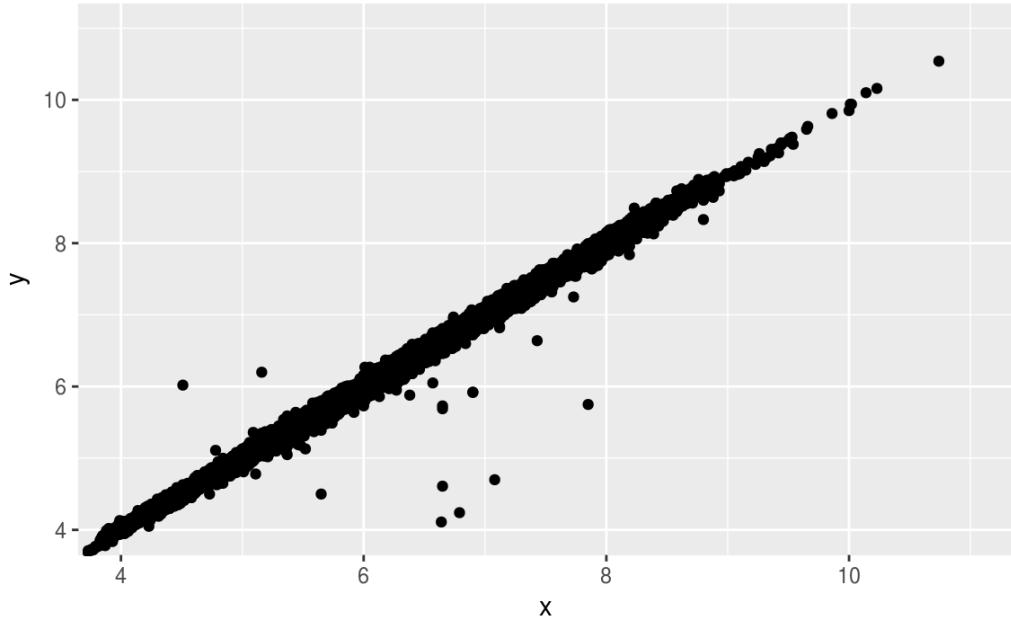
Another approach is to display approximately the same number of points in each bin. That's the job of `cut_number()`:

```
ggplot(data = smaller, mapping = aes(x = carat, y = price)) +
 geom_boxplot(mapping = aes(group = cut_number(carat, 20)))
```



#### 7.5.3.1 Exercises

1. Instead of summarising the conditional distribution with a boxplot, you could use a frequency polygon. What do you need to consider when using `cut_width()` vs `cut_number()`? How does that impact a visualisation of the 2d distribution of `carat` and `price`?
2. Visualise the distribution of `carat`, partitioned by `price`.
3. How does the price distribution of very large diamonds compare to small diamonds. Is it as you expect, or does it surprise you?
4. Combine two **of the techniques you've learned to visualise the combined distribution of cut, carat, and price.**
5. Two dimensional plots reveal outliers that are not visible in one dimensional plots. For example, some points in the plot below have an unusual combination of `x` and `y` values, which makes the points outliers even though their `x` and `y` values appear normal when examined separately.
6. `ggplot(data = diamonds) +`
7.   `geom_point(mapping = aes(x = x, y = y)) +`  
`coord_cartesian(xlim = c(4, 11), ylim = c(4, 11))`



Why is a scatterplot a better display than a binned plot for this case?

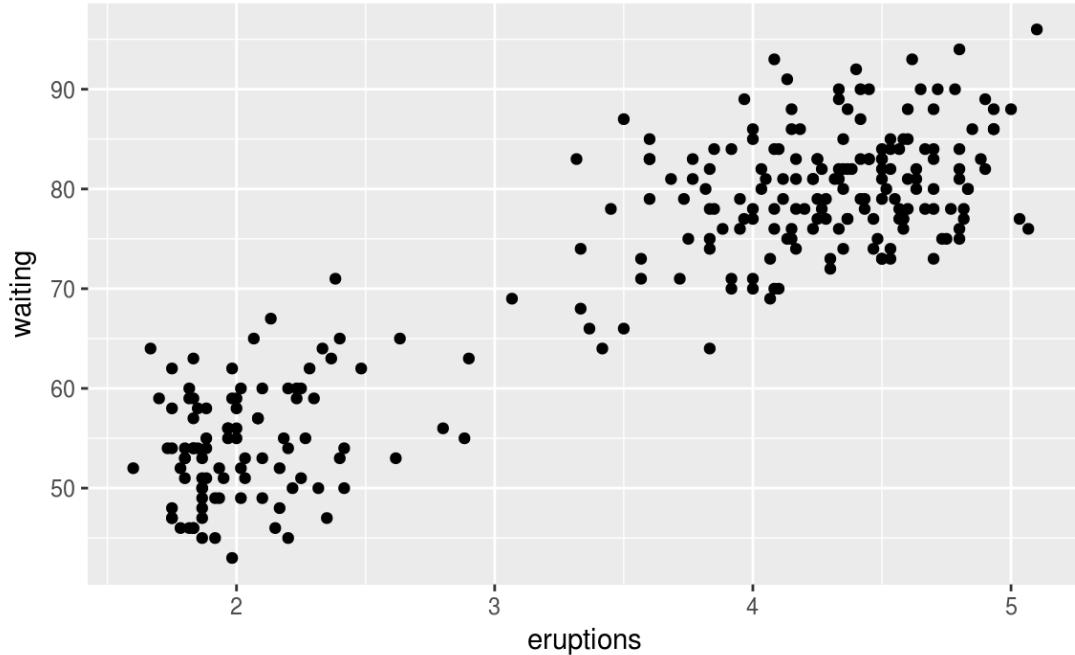
## 7.6 Patterns and models

Patterns in your data provide clues about relationships. If a systematic relationship exists between two variables it will appear as a pattern in the data. If you spot a pattern, ask yourself:

- Could this pattern be due to coincidence (i.e. random chance)?
- How can you describe the relationship implied by the pattern?
- How strong is the relationship implied by the pattern?
- What other variables might affect the relationship?
- Does the relationship change if you look at individual subgroups of the data?

A scatterplot of Old Faithful eruption lengths versus the wait time between eruptions shows a pattern: longer wait times are associated with longer eruptions. The scatterplot also displays the two clusters that we noticed above.

```
ggplot(data = faithful) +
 geom_point(mapping = aes(x = eruptions, y = waiting))
```



Patterns provide one of the most useful tools for data scientists because they reveal covariation. If you think of variation as a phenomenon that creates uncertainty, covariation is a phenomenon that reduces it. If two variables covary, you can use the values of one variable to make better predictions about the values of the second. If the covariation is due to a causal relationship (a special case), then you can use the value of one variable to control the value of the second.

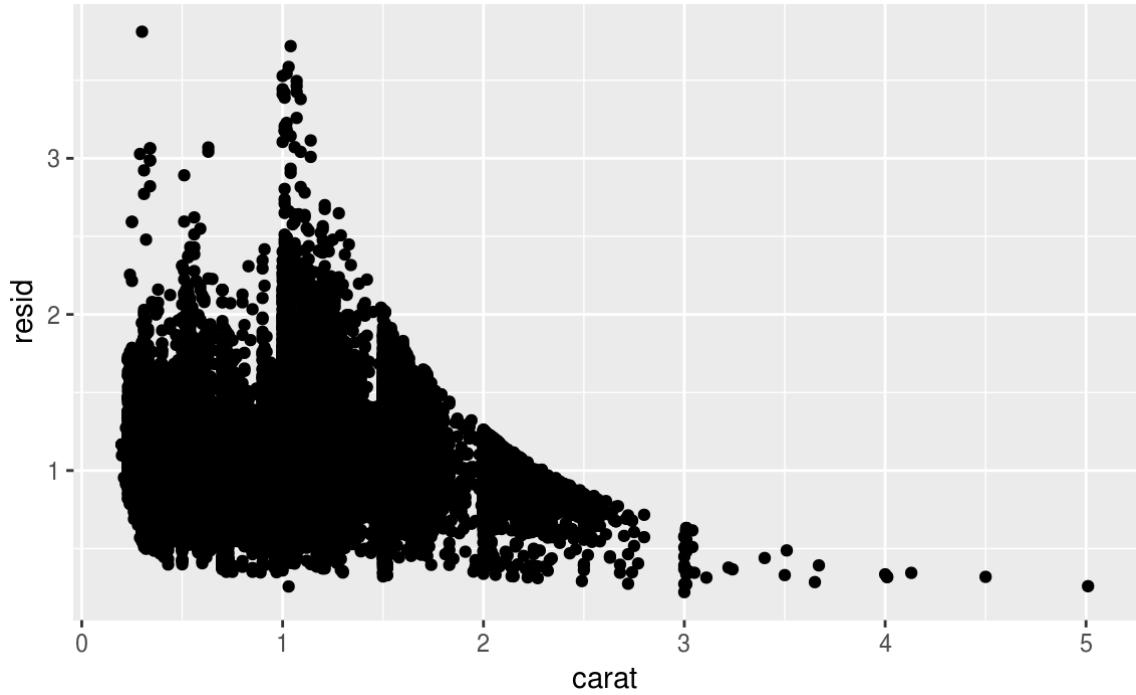
Models are a tool for extracting patterns out of data. For example, consider the diamonds data. It's hard to understand the relationship between `cut` and `price`, because `cut` and `carat`, and `carat` and `price` are tightly related. It's possible to use a model to remove the very strong relationship between `price` and `carat` so we can explore the subtleties that remain. The following code fits a model that predicts `price` from `carat` and then computes the residuals (the difference between the predicted value and the actual value). The residuals give us a view of the price of the diamond, once the effect of `carat` has been removed.

```
library(modelr)

mod <- lm(log(price) ~ log(carat), data = diamonds)

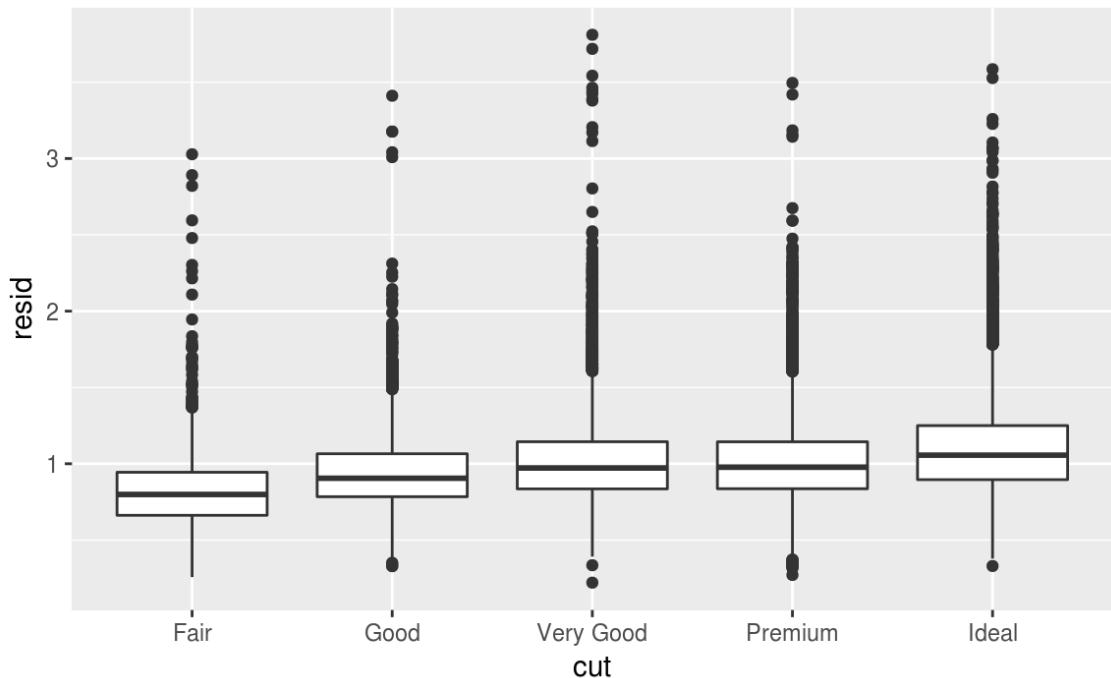
diamonds2 <- diamonds %>%
 add_residuals(mod) %>%
 mutate(resid = exp(resid))

ggplot(data = diamonds2) +
 geom_point(mapping = aes(x = carat, y = resid))
```



Once you've removed the strong relationship between carat and price, you can see what you expect in the relationship between cut and price: relative to their size, better quality diamonds are more expensive.

```
ggplot(data = diamonds2) +
 geom_boxplot(mapping = aes(x = cut, y = resid))
```



You'll learn how models, and the `modelr` package, work in the final part of the book, [model](#). We're saving modelling for later because understanding what models are and how they work is easiest once you have tools of data wrangling and programming in hand.

## 7.7 ggplot2 calls

As we move on from these introductory chapters, we'll transition to a more concise expression of `ggplot2` code. So far we've been very explicit, which is helpful when you are learning:

```
ggplot(data = faithful, mapping = aes(x = eruptions)) +
 geom_freqpoly(binwidth = 0.25)
```

Typically, the first one or two arguments to a function are so important that you should know them by heart. The first two arguments to `ggplot()` are `data` and `mapping`, and the first two arguments to `aes()` are `x` and `y`. In the remainder of the book, we won't supply those names. That saves typing, and, by reducing the amount of boilerplate, makes it easier to see what's different between plots. That's a really important programming concern that we'll come back in [functions](#).

Rewriting the previous plot more concisely yields:

```
ggplot(faithful, aes(eruptions)) +
 geom_freqpoly(binwidth = 0.25)
```

Sometimes we'll turn the end of a pipeline of data transformation into a plot. Watch for the transition from `%>%` to `+`. I wish this transition wasn't necessary but unfortunately `ggplot2` was created before the pipe was discovered.

```
diamonds %>%
 count(cut, clarity) %>%
 ggplot(aes(clarity, cut, fill = n)) +
 geom_tile()
```

## 7.8 Learning more

If you want to learn more about the mechanics of `ggplot2`, I'd highly recommend grabbing a copy of the `ggplot2` book: <https://amzn.com/331924275X>. It's been recently updated, so it includes `dplyr` and `tidyverse` code, and has much more space to explore all the facets of visualisation. Unfortunately the book isn't generally available for free, but if you have a connection to a university you can probably get an electronic version for free through SpringerLink.

Another useful resource is the `R Graphics Cookbook` by Winston Chang. Much of the contents are available online at <http://www.cookbook-r.com/Graphs/>.

I also recommend `Graphical Data Analysis with R`, by Antony Unwin. This is a book-length treatment similar to the material covered in this chapter, but has the space to go into much greater depth.

# 8 Workflow: projects

One day you will need to quit R, go do something else and return to your analysis the next day. One day you will be working on multiple analyses simultaneously that all use R and you want to keep them separate. One day you will need to bring data from the outside world into R and send numerical results and figures from R back out into the world. To handle these real life situations, you need to make two decisions:

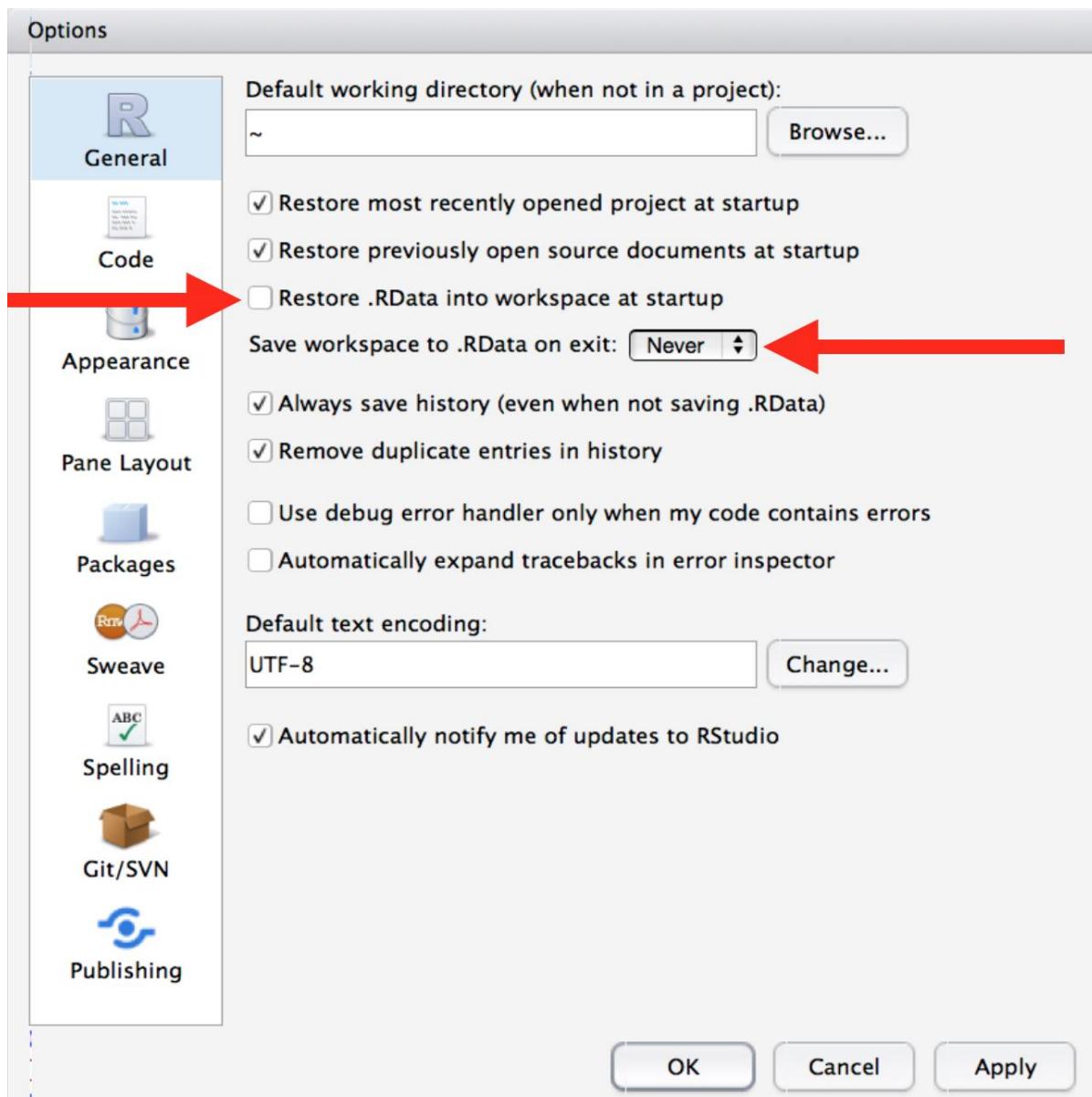
1. **What about your analysis is “real”, i.e.** what will you save as your lasting record of what happened?
2. **Where does your analysis “live”?**

## 8.1 What is real?

As a beginning R user, it’s OK to consider your environment (i.e. the objects listed in the environment pane) “real”. However, in the long run, you’ll be much better off if you consider your R scripts as “real”.

With your R scripts (and your data files), you can recreate the environment. It’s much harder to recreate your R scripts from your environment! You’ll either have to retype a lot of code from memory (making mistakes all the way) or you’ll have to carefully mine your R history.

To foster this behaviour, I highly recommend that you instruct RStudio not to preserve your workspace between sessions:



This will cause you some short-term pain, because now when you restart RStudio it will not remember the results of the code that you ran last time. But this short-term pain will save you long-term agony because it forces you to capture all important interactions in your code. There's nothing worse than discovering three months after the fact that you've only stored the results of an important calculation in your workspace, not the calculation itself in your code.

There is a great pair of keyboard shortcuts that will work together to make sure you've captured the important parts of your code in the editor:

1. Press Cmd/Ctrl + Shift + F10 to restart RStudio.
2. Press Cmd/Ctrl + Shift + S to rerun the current script.

I use this pattern hundreds of times a week.

## 8.2 Where does your analysis live?

R has a powerful notion of the working directory. This is where R looks for files that you ask it to load, and where it will put any files that you ask it to save. RStudio shows your current working directory at the top of the console:



And you can print this out in R code by running `getwd()`:

```
getwd()
#> [1] "/Users/hadley/Documents/r4ds/r4ds"
```

As a beginning R user, it's OK to let your home directory, documents directory, or any other weird directory on your computer be R's working directory. But you're six chapters into this book, and you're no longer a rank beginner. Very soon now you should evolve to organising your analytical projects into directories and, when working on a project, setting R's working directory to the associated directory.

I do not recommend it, but you can also set the working directory from within R:

```
setwd("/path/to/my/CoolProject")
```

But you should never do this because there's a better way; a way that also puts you on the path to managing your R work like an expert.

## 8.3 Paths and directories

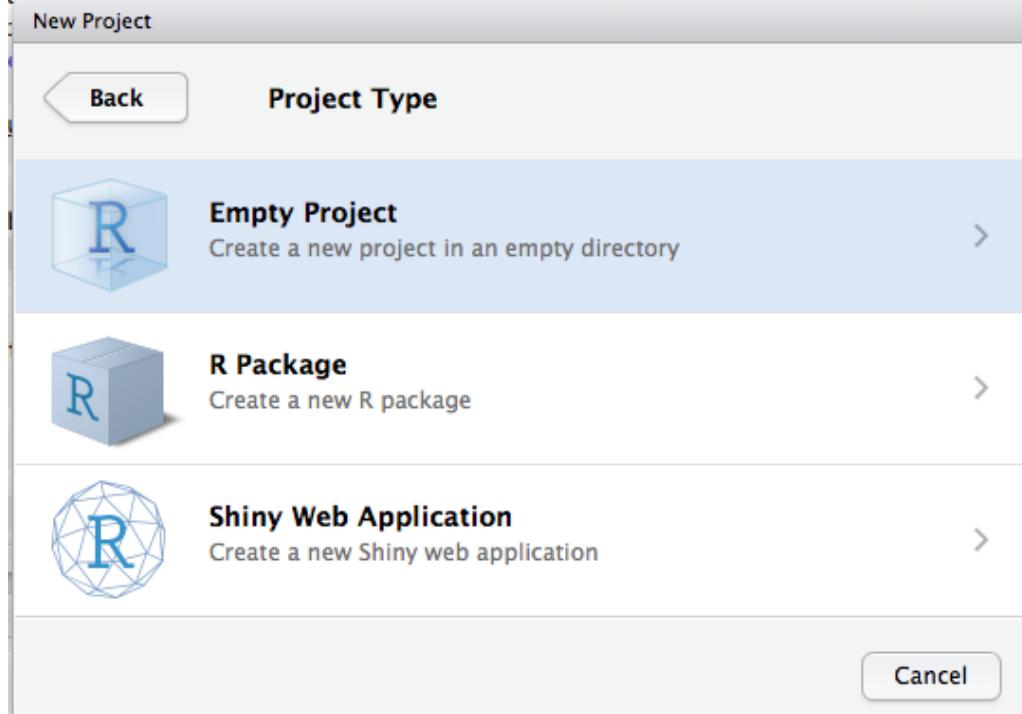
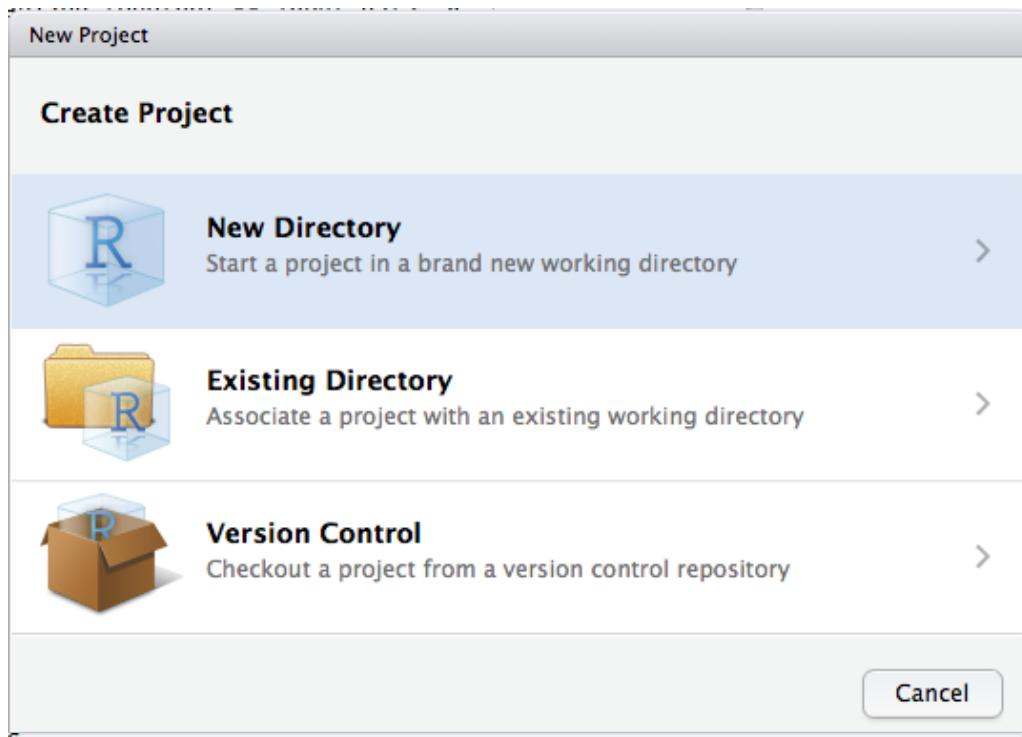
Paths and directories are a little complicated because there are two basic styles of paths: Mac/Linux and Windows. There are three chief ways in which they differ:

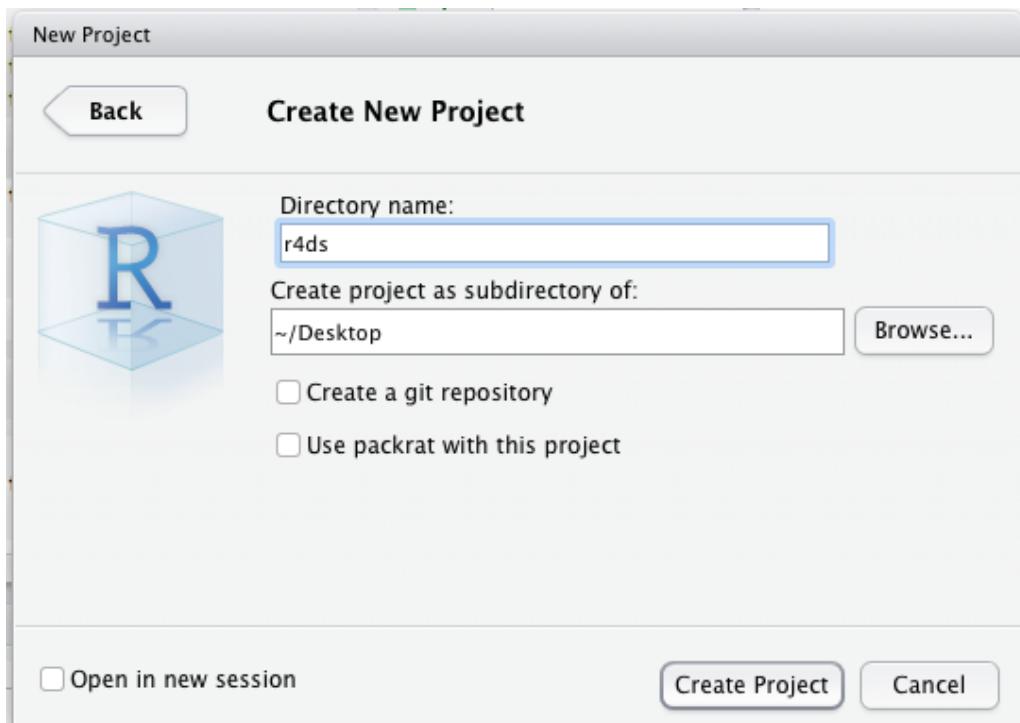
1. The most important difference is how you separate the components of the path. Mac and Linux uses slashes (e.g. `plots/diamonds.pdf`) and Windows uses backslashes (e.g. `plots\diamonds.pdf`). **R can work with either type (no matter what platform you're currently using)**, but unfortunately, backslashes mean something special to R, and to get a single backslash in the path, you need to type two backslashes! That makes life frustrating, so I recommend always using the Linux/Mac style with forward slashes.
2. Absolute paths (i.e. paths that point to the same place regardless of your working directory) look different. In Windows they start with a drive letter (e.g. `C:`) or two backslashes (e.g. `\servername`) **and in Mac/Linux they start with a slash “/” (e.g. /users/hadley)**. You should never use absolute paths in your scripts, because they hinder sharing: no one else will have exactly the same directory configuration as you.
3. The last minor difference is the place that `~` points to. `~` is a convenient shortcut to your **home directory**. **Windows doesn't really have the notion of a home directory, so it instead points to your documents directory**.

## 8.4 RStudio projects

R experts keep all the files associated with a project together — input data, R scripts, analytical results, figures. This is such a wise and common practice that RStudio has built-in support for this via projects.

**Let's make a project for you to use while you're working through the rest of this book. Click File > New Project, then:**





Call your project `r4ds` and think carefully about which `subdirectory` you put the project in. If you **don't store it somewhere sensible, it will be** hard to find it in the future!

Once this process is complete, you'll get a new RStudio project just for this book. Check that the "home" directory of your project is the current working directory:

```
getwd()
#> [1] /Users/hadley/Documents/r4ds/r4ds
```

Whenever you refer to a file with a relative path it will look for it here.

**Now enter the following commands in the script editor, and save the file, calling it "diamonds.R". Next, run the complete script which will save a PDF and CSV file into your project directory. Don't worry about the details, you'll learn them later in the book.**

```
library(tidyverse)

ggplot(diamonds, aes(carat, price)) +
 geom_hex()
ggsave("diamonds.pdf")

write_csv(diamonds, "diamonds.csv")
```

Quit RStudio. Inspect the folder associated with your project — notice the `.Rproj` file. Double-click that file to **re-open the project. Notice you get back to where you left off: it's the same** working directory and command history, and all the files you were working on are still open. Because you followed my instructions above, you will, however, have a completely fresh **environment, guaranteeing that you're starting with a clean slate.**

In your favorite OS-specific way, search your computer for `diamonds.pdf` and you will find the PDF (no surprise) but also the script that created it (`diamonds.R`). This is huge win! One day you will want to remake a figure or just understand where it came from. If you rigorously save figures to files with R code and never with the mouse or the clipboard, you will be able to reproduce old work with ease!

## 8.5 Summary

In summary, RStudio projects give you a solid workflow that will serve you well in the future:

- Create an RStudio project for each data analysis project.
- **Keep data files there; we'll talk about loading them into R in [data import](#).**
- Keep scripts there; edit them, run them in bits or as a whole.
- Save your outputs (plots and cleaned data) there.
- Only ever use relative paths, not absolute paths.

Everything you need is in one place, and cleanly separated from all the other projects that you are working on.

# 10 Tibbles

## 10.1 Introduction

Throughout this book we work with “tibbles” instead of R’s traditional `data.frame`. Tibbles are data frames, but they tweak some older behaviours to make life a little easier. R is an old language, and some things that were useful 10 or 20 years ago now get in your way. It’s difficult to change base R without breaking existing code, so most innovation occurs in packages. Here we will describe the `tibble` package, which provides opinionated data frames that make working in the tidyverse a little easier. In most places, I’ll use the term `tibble` and `data frame` interchangeably; when I want to draw particular attention to R’s built-in data frame, I’ll call them `data.frames`.

If this chapter leaves you wanting to learn more about tibbles, you might enjoy `vignette("tibble")`.

### 10.1.1 Prerequisites

In this chapter we’ll explore the `tibble` package, part of the core tidyverse.

```
library(tidyverse)
```

## 10.2 Creating tibbles

Almost all of the functions that you’ll use in this book produce tibbles, as tibbles are one of the unifying features of the tidyverse. Most other R packages use regular data frames, so you might want to coerce a data frame to a tibble. You can do that with `as_tibble()`:

```
as_tibble(iris)
#> # A tibble: 150 × 5
#> Sepal.Length Sepal.Width Petal.Length Petal.Width Species
#> <dbl> <dbl> <dbl> <dbl> <fctr>
#> 1 5.1 3.5 1.4 0.2 setosa
#> 2 4.9 3.0 1.4 0.2 setosa
#> 3 4.7 3.2 1.3 0.2 setosa
#> 4 4.6 3.1 1.5 0.2 setosa
#> 5 5.0 3.6 1.4 0.2 setosa
#> 6 5.4 3.9 1.7 0.4 setosa
#> # ... with 144 more rows
```

You can create a new tibble from individual vectors with `tibble()`. `tibble()` will automatically recycle inputs of length 1, and allows you to refer to variables that you just created, as shown below.

```
tibble(
 x = 1:5,
 y = 1,
 z = x ^ 2 + y
)
```

```
#> # A tibble: 5 × 3
#> x y z
#> <int> <dbl> <dbl>
#> 1 1 1 2
#> 2 2 1 5
#> 3 3 1 10
#> 4 4 1 17
#> 5 5 1 26
```

If you're already familiar with `data.frame()`, note that `tibble()` does much less: it never changes the type of the inputs (e.g. it never converts strings to factors!), it never changes the names of variables, and it never creates row names.

It's possible for a tibble to have column names that are not valid R variable names, aka non-syntactic names. For example, they might not start with a letter, or they might contain unusual characters like a space. To refer to these variables, you need to surround them with backticks, `:

```
tb <- tibble(
 `:`) ` = "smile",
 ` ` = "space",
 `2000` = "number"
)
tb
#> # A tibble: 1 × 3
#> `:` ` ` `2000`
#> <chr> <chr> <chr>
#> 1 smile space number
```

You'll also need the backticks when working with these variables in other packages, like `ggplot2`, `dplyr`, and `tidyverse`.

Another way to create a tibble is with `tribble()`, short for transposed tibble. `tribble()` is customised for data entry in code: column headings are defined by formulas (i.e. they start with ~), and entries are separated by commas. This makes it possible to lay out small amounts of data in easy to read form.

```
tribble(
 ~x, ~y, ~z,
 #---|---|---
 "a", 2, 3.6,
 "b", 1, 8.5
)
#> # A tibble: 2 × 3
#> x y z
#> <chr> <dbl> <dbl>
#> 1 a 2 3.6
#> 2 b 1 8.5
```

I often add a comment (the line starting with #), to make it really clear where the header is.

## 10.3 Tibbles vs. `data.frame`

There are two main differences in the usage of a tibble vs. a classic `data.frame`: printing and subsetting.

### 10.3.1 Printing

Tibbles have a refined print method that shows only the first 10 rows, and all the columns that fit on screen. This makes it much easier to work with large data. In addition to its name, each column reports its type, a nice feature borrowed from `str()`:

```
tibble(
 a = lubridate::now() + runif(1e3) * 86400,
 b = lubridate::today() + runif(1e3) * 30,
 c = 1:1e3,
 d = runif(1e3),
 e = sample(letters, 1e3, replace = TRUE)
)
#> # A tibble: 1,000 × 5
#> a b c d e
#> <dttm> <date> <int> <dbl> <chr>
#> 1 2017-05-04 14:08:35 2017-05-11 1 0.368 h
#> 2 2017-05-05 08:13:45 2017-05-16 2 0.612 n
#> 3 2017-05-05 02:37:24 2017-05-26 3 0.415 l
#> 4 2017-05-04 15:58:41 2017-05-25 4 0.212 x
#> 5 2017-05-04 12:22:58 2017-05-22 5 0.733 a
#> 6 2017-05-04 23:23:55 2017-05-18 6 0.460 v
#> # ... with 994 more rows
```

**Tibbles are designed so that you don't accidentally overwhelm your console when you print large data frames.** But sometimes you need more output than the default display. There are a few options that can help.

First, you can explicitly `print()` the data frame and control the number of rows (`n`) and the width of the display. `width = Inf` will display all columns:

```
nycflights13::flights %>%
 print(n = 10, width = Inf)
```

You can also control the default print behaviour by setting options:

- `options(tibble.print_max = n, tibble.print_min = m)`: if more than `m` rows, print only `n` rows. Use `options(dplyr.print_min = Inf)` to always show all rows.
- Use `options(tibble.width = Inf)` to always print all columns, regardless of the width of the screen.

You can see a complete list of options by looking at the package help with `package?tibble`.

A final option is to use RStudio's built-in data viewer to get a scrollable view of the complete dataset. This is also often useful at the end of a long chain of manipulations.

```
nycflights13::flights %>%
 View()
```

## 10.3.2 Subsetting

So far all the tools you've learned have worked with complete data frames. If you want to pull out a single variable, you need some new tools, `$` and `[`. `[` can extract by name or position; `$` only extracts by name but is a little less typing.

```
df <- tibble(
 x = runif(5),
 y = rnorm(5)
)

Extract by name
df$x
#> [1] 0.434 0.395 0.548 0.762 0.254
df[["x"]]
#> [1] 0.434 0.395 0.548 0.762 0.254

Extract by position
df[[1]]
#> [1] 0.434 0.395 0.548 0.762 0.254
```

To use these in a pipe, you'll need to use the special placeholder `..`:

```
df %>% .$x
#> [1] 0.434 0.395 0.548 0.762 0.254
df %>% .[["x"]]
#> [1] 0.434 0.395 0.548 0.762 0.254
```

Compared to a `data.frame`, tibbles are more strict: they never do partial matching, and they will generate a warning if the column you are trying to access does not exist.

## 10.4 Interacting with older code

Some older functions don't work with tibbles. If you encounter one of these functions, use `as.data.frame()` to turn a tibble back to a `data.frame`:

```
class(as.data.frame(tb))
#> [1] "data.frame"
```

The main reason that some older functions don't work with tibble is the `[` function. We don't use `[` much in this book because `dplyr::filter()` and `dplyr::select()` allow you to solve the same problems with clearer code (but you will learn a little about it in [vector subsetting](#)). With base R data frames, `[` sometimes returns a data frame, and sometimes returns a vector. With tibbles, `[` always returns another tibble.

## 10.5 Exercises

1. How can you tell if an object is a tibble? (Hint: try printing `mtcars`, which is a regular data frame).

2. Compare and contrast the following operations on a `data.frame` and equivalent tibble.  
What is different? Why might the default data frame behaviours cause you frustration?
3. `df <- data.frame(abc = 1, xyz = "a")`
4. `df$x`
5. `df[, "xyz"]`  
`df[, c("abc", "xyz")]`
  
6. If you have the name of a variable stored in an object, e.g. `var <- "mpg"`, how can you extract the reference variable from a tibble?
7. Practice referring to non-syntactic names in the following data frame by:
  1. Extracting the variable called 1.
  2. Plotting a scatterplot of 1 vs 2.
  3. Creating a new column called 3 which is 2 divided by 1.
  4. Renaming the columns to one, two and three.
8. `annoying <- tibble(`
9.   ``1` = 1:10,`
10.   ``2` = `1` * 2 + rnorm(length(`1`))`  
    `)`
  
11. What does `tibble::enframe()` do? When might you use it?
12. What option controls how many additional column names are printed at the footer of a tibble?

# 11 Data import

## 11.1 Introduction

Working with data provided by R packages is a great way to learn the tools of data science, but at **some point you want to stop learning and start working with your own data. In this chapter, you'll learn how to read plain-text rectangular files into R.** Here, we'll only scratch the surface of data import, but many of the principles will translate to **other forms of data.** We'll finish with a few pointers to packages that are useful for other types of data.

### 11.1.1 Prerequisites

In this chapter, you'll learn how to load flat files in R with the `readr` package, which is part of the core tidyverse.

```
library(tidyverse)
```

## 11.2 Getting started

Most of `readr`'s functions are concerned with turning flat files into data frames:

- `read_csv()` reads comma delimited files, `read_csv2()` reads semicolon separated files (common in countries where `,` is used as the decimal place), `read_tsv()` reads tab delimited files, and `read_delim()` reads in files with any delimiter.
- `read_fwf()` reads fixed width files. You can specify fields either by their widths with `fwf_widths()` or their position with `fwf_positions()`. `read_table()` reads a common variation of fixed width files where columns are separated by white space.
- `read_log()` reads Apache style log files. (But also check out `webreadr` which is built on top of `read_log()` and provides many more helpful tools.)

These functions all have similar syntax: once you've mastered one, you can use the others with ease. For the rest of this chapter we'll focus on `read_csv()`. Not only are csv files one of the most common forms of data storage, but once you understand `read_csv()`, you can easily apply your knowledge to all the other functions in `readr`.

The first argument to `read_csv()` is the most important: it's the path to the file to read.

```
heights <- read_csv("data/heights.csv")
#> Parsed with column specification:
#> cols(
#> earn = col_double(),
#> height = col_double(),
#> sex = col_character(),
#> ed = col_integer(),
#> age = col_integer(),
#> race = col_character()
#>)
```

When you run `read_csv()` it prints out a column specification that gives the name and type of each column. That's an important part of `readr`, which we'll come back to in [parsing a file](#).

You can also supply an inline csv file. This is useful for experimenting with `readr` and for creating reproducible examples to share with others:

```
read_csv("a,b,c
1,2,3
4,5,6")
#> # A tibble: 2 × 3
#> a b c
#> <int> <int> <int>
#> 1 1 2 3
#> 2 4 5 6
```

In both cases `read_csv()` uses the first line of the data for the column names, which is a very common convention. There are two cases where you might want to tweak this behaviour:

1. Sometimes there are a few lines of metadata at the top of the file. You can use `skip = n` to skip the first `n` lines; or use `comment = "#"` to drop all lines that start with (e.g.) `#`.  
1. `read_csv("The first line of metadata`  
2. `The second line of metadata`  
3. `x,y,z`  
4. `1,2,3", skip = 2)`  
5. `#> # A tibble: 1 × 3`  
6. `#> x y z`  
7. `#> <int> <int> <int>`  
8. `#> 1 1 2 3`  
9. `#> 10.`  
10. `11. read_csv("# A comment I want to skip`  
11. `12. x,y,z`  
12. `13. 1,2,3", comment = "#")`  
13. `#> # A tibble: 1 × 3`  
14. `#> x y z`  
15. `#> <int> <int> <int>`  
16. `#> 1 1 2 3`  
17. `#> 17.`  
17. The data might not have column names. You can use `col_names = FALSE` to tell `read_csv()` not to treat the first row as headings, and instead label them sequentially from X1 to Xn:  
18. `read_csv("1,2,3\n4,5,6", col_names = FALSE)`  
19. `#> # A tibble: 2 × 3`  
20. `#> X1 X2 X3`  
21. `#> <int> <int> <int>`  
22. `#> 1 1 2 3`  
23. `#> 2 4 5 6`

("\n" is a convenient shortcut for adding a new line. You'll learn more about it and other types of string escape in [string basics](#).)

Alternatively you can pass `col_names` a character vector which will be used as the column names:

```

read_csv("1,2,3\n4,5,6", col_names = c("x", "y", "z"))
#> # A tibble: 2 × 3
#> x y z
#> <int> <int> <int>
#> 1 1 2 3
#> 2 4 5 6

```

Another option that commonly needs tweaking is `na`: this specifies the value (or values) that are used to represent missing values in your file:

```

read_csv("a,b,c\n1,2,.", na = ".")
#> # A tibble: 1 × 3
#> a b c
#> <int> <int> <chr>
#> 1 1 2 <NA>

```

This is all you need to know to read ~75% of CSV files that you'll encounter in practice. You can also easily adapt what you've learned to read tab separated files with `read_tsv()` and fixed width files with `read_fwf()`. To read in more challenging files, you'll need to learn more about how `readr` parses each column, turning them into R vectors.

### 11.2.1 Compared to base R

If you've used R before, you might wonder why we're not using `read.csv()`. There are a few good reasons to favour `readr` functions over the base equivalents:

- They are typically much faster (~10x) than their base equivalents. Long running jobs have a **progress bar**, so you can see what's happening. If you're looking for raw speed, try `data.table::fread()`. It doesn't fit quite so well into the tidyverse, but it can be quite a bit faster.
- They produce **tibbles**, they don't convert character vectors to factors, use row names, or munge the column names. These are common sources of frustration with the base R functions.
- They are more reproducible. Base R functions inherit some behaviour from your operating system and environment variables, so import code that works on your computer might not work on someone else's.

### 11.2.2 Exercises

1. What function would you use to read a file where fields were separated with "`|`"?
2. Apart from `file`, `skip`, and `comment`, what other arguments do `read_csv()` and `read_tsv()` have in common?
3. What are the most important arguments to `read_fwf()`?
4. Sometimes strings in a CSV file contain commas. To prevent them from causing problems they need to be surrounded by a quoting character, like `"` or `'`. By convention, `read_csv()` assumes that the quoting character will be `"`, and if you want to change it you'll need to use `read_delim()` instead. What arguments do you need to specify to read the following text into a data frame?

```

"x,y\n1,'a,b'"
5. Identify what is wrong with each of the following inline CSV files. What happens when you
run the code?
6. read_csv("a,b\n1,2,3\n4,5,6")
7. read_csv("a,b,c\n1,2\n1,2,3,4")
8. read_csv("a,b\nn\"1")
9. read_csv("a,b\n1,2\nna,b")
read_csv("a;b\n1;3")

```

## 11.3 Parsing a vector

Before we get into the details of how `readr` reads files from disk, we need to take a little detour to talk about the `parse_*` functions. These functions take a character vector and return a more specialised vector like a logical, integer, or date:

```

str(parse_logical(c("TRUE", "FALSE", "NA")))
#> logi [1:3] TRUE FALSE NA
str(parse_integer(c("1", "2", "3")))
#> int [1:3] 1 2 3
str(parse_date(c("2010-01-01", "1979-10-14")))
#> Date[1:2], format: "2010-01-01" "1979-10-14"

```

These functions are useful in their own right, but are also an important building block for `readr`. **Once you've learned how the individual parsers work in this section, we'll circle back and see** how they fit together to parse a complete file in the next section.

Like all functions in the tidyverse, the `parse_*` functions are uniform: the first argument is a character vector to parse, and the `na` argument specifies which strings should be treated as missing:

```

parse_integer(c("1", "231", ".", "456"), na = ".")
#> [1] 1 231 NA 456

```

If parsing fails, you'll get a warning:

```

x <- parse_integer(c("123", "345", "abc", "123.45"))
#> Warning: 2 parsing failures.
#> row col expected actual
#> 3 -- an integer abc
#> 4 -- no trailing characters .45

```

And the failures will be missing in the output:

```

x
#> [1] 123 345 NA NA
#> attr(",problems")
#> # A tibble: 2 × 4
#> row col expected actual
#> <int> <int> <chr> <chr>
#> 1 3 NA an integer abc

```

```
#> 2 4 NA no trailing characters .45
```

If there are many parsing failures, you'll need to use `problems()` to get the complete set. This returns a tibble, which you can then manipulate with `dplyr`.

```
problems(x)
#> # A tibble: 2 × 4
#> row col expected actual
#> <int> <int> <chr> <chr>
#> 1 3 NA an integer abc
#> 2 4 NA no trailing characters .45
```

Using parsers is mostly a matter of understanding what's available and how they deal with different types of input. There are eight particularly important parsers:

1. `parse_logical()` and `parse_integer()` parse logicals and integers respectively. **There's basically nothing that can go wrong with these parsers so I won't describe them here further.**
2. `parse_double()` is a strict numeric parser, and `parse_number()` is a flexible numeric parser. These are more complicated than you might expect because different parts of the world write numbers in different ways.
3. `parse_character()` **seems so simple that it shouldn't be necessary. But one** complication makes it quite important: character encodings.
4. `parse_factor()` create factors, the data structure that R uses to represent categorical variables with fixed and known values.
5. `parse_datetime()`, `parse_date()`, and `parse_time()` allow you to parse various date & time specifications. These are the most complicated because there are so many different ways of writing dates.

The following sections describe these parsers in more detail.

### 11.3.1 Numbers

It seems like it should be straightforward to parse a number, but three problems make it tricky:

1. People write numbers differently in different parts of the world. For example, some countries use `.` in between the integer and fractional parts of a real number, while others use `,`.
2. Numbers are often surrounded by other characters that provide some context, like `"$1000"` or `"10%"`.
3. **Numbers often contain “grouping” characters to make them easier to read, like** `"1,000,000"`, **and these grouping characters vary around the world.**

To address the first problem, `readr` has the notion of a “locale”, an object that specifies parsing options that differ from place to place. When parsing numbers, the most important option is the character you use for the decimal mark. You can override the default value of `.` by creating a new locale and setting the `decimal_mark` argument:

```
parse_double("1.23")
#> [1] 1.23
parse_double("1,23", locale = locale(decimal_mark = ","))
```

```
#> [1] 1.23
```

**readr's default locale is US**-centric, because generally R is US-centric (i.e. the documentation of base R is written in American English). An alternative approach would be to try and guess the defaults from your operating system. This is hard to do well, and, more importantly, makes your code fragile: even if it works on your computer, it might fail when you email it to a colleague in another country.

`parse_number()` addresses the second problem: it ignores non-numeric characters before and after the number. This is particularly useful for currencies and percentages, but also works to extract numbers embedded in text.

```
parse_number("$100")
#> [1] 100
parse_number("20%")
#> [1] 20
parse_number("It cost $123.45")
#> [1] 123
```

The final problem is addressed by the combination of `parse_number()` and the locale as `parse_number()` will ignore the “grouping mark”:

```
Used in America
parse_number("$123,456,789")
#> [1] 1.23e+08

Used in many parts of Europe
parse_number("123.456.789", locale = locale(grouping_mark = "."))
#> [1] 1.23e+08

Used in Switzerland
parse_number("123'456'789", locale = locale(grouping_mark = "'"))
#> [1] 1.23e+08
```

### 11.3.2 Strings

It seems like `parse_character()` should be really simple — it could just return its input. **Unfortunately life isn't so simple, as there are multiple ways to represent the same string. To understand what's going on, we need to dive into** the details of how computers represent strings. In R, we can get at the underlying representation of a string using `charToRaw()`:

```
charToRaw("Hadley")
#> [1] 48 61 64 6c 65 79
```

Each hexadecimal number represents a byte of information: 48 is H, 61 is a, and so on. The mapping from hexadecimal number to character is called the encoding, and in this case the **encoding is called ASCII. ASCII does a great job of representing English characters, because it's the American Standard Code for Information Interchange.**

Things get more complicated for languages other than English. In the early days of computing there were many competing standards for encoding non-English characters, and to correctly

interpret a string you needed to know both the values and the encoding. For example, two common encodings are Latin1 (aka ISO-8859-1, used for Western European languages) and Latin2 (aka ISO-8859-2, used for Eastern European languages). In Latin1, the byte b1 is “±”, but in Latin2, it’s “ä”! **Fortunately, today there is one standard that is supported almost everywhere:** UTF-8. UTF-8 can encode just about every character used by humans today, as well as many extra symbols (like emoji!).

readr uses UTF-8 everywhere: it assumes your data is UTF-8 encoded when you read it, and always uses it when writing. This is a good default, but will fail for data produced by older **systems that don’t understand UTF-8**. If this happens to you, your strings will look weird when you print them. Sometimes just **one or two characters might be messed up; other times you’ll get complete gibberish**. For example:

```
x1 <- "El Ni\xf1o was particularly bad this year"
x2 <- "\x82\xb1\x82\xf1\x82\xc9\x82\xbf\x82\xcd"

x1
#> [1] "El Ni\xf1o was particularly bad this year"
x2
#> [1] "\x82\xb1\x82\xf1\x82\xbf\x82\xcd"
```

To fix the problem you need to specify the encoding in `parse_character()`:

```
parse_character(x1, locale = locale(encoding = "Latin1"))
#> [1] "El Niño was particularly bad this year"
parse_character(x2, locale = locale(encoding = "Shift-JIS"))
#> [1] "こんにちは"
```

**How do you find the correct encoding?** If you’re lucky, it’ll be included somewhere in the data documentation. Unfortunately, that’s rarely the case, so readr provides `guess_encoding()` to help you figure it out. It’s not foolproof, and it works better when you have lots of text (unlike here), but it’s a reasonable place to start. Expect to try a few different encodings before you find the right one.

```
guess_encoding(charToRaw(x1))
#> # A tibble: 2 × 2
#> encoding confidence
#> <chr> <dbl>
#> 1 ISO-8859-1 0.46
#> 2 ISO-8859-9 0.23
guess_encoding(charToRaw(x2))
#> # A tibble: 1 × 2
#> encoding confidence
#> <chr> <dbl>
#> 1 KOI8-R 0.42
```

The first argument to `guess_encoding()` can either be a path to a file, or, as in this case, a raw vector (useful if the strings are already in R).

Encodings are a rich and complex topic, and I’ve only scratched the surface here. If you’d like to learn more I’d recommend reading the detailed explanation at <http://kunststube.net/encoding/>.

### 11.3.3 Factors

R uses factors to represent categorical variables that have a known set of possible values. Give `parse_factor()` a vector of known levels to generate a warning whenever an unexpected value is present:

```
fruit <- c("apple", "banana")
parse_factor(c("apple", "banana", "bananana"), levels = fruit)
#> Warning: 1 parsing failure.
#> row col expected actual
#> 3 -- value in level set banana
#> [1] apple banana <NA>
#> attr(,"problems")
#> # A tibble: 1 × 4
#> row col expected actual
#> <int> <int> <chr> <chr>
#> 1 3 NA value in level set banana
#> Levels: apple banana
```

But if you have many problematic entries, it's often easier to leave as character vectors and then use the tools you'll learn about in [strings](#) and [factors](#) to clean them up.

### 11.3.4 Dates, date-times, and times

You pick between three parsers depending on whether you want a date (the number of days since 1970-01-01), a date-time (the number of seconds since midnight 1970-01-01), or a time (the number of seconds since midnight). When called without any additional arguments:

- `parse_datetime()` expects an ISO8601 date-time. ISO8601 is an international standard in which the components of a date are organised from biggest to smallest: year, month, day, hour, minute, second.
- `parse_datetime("2010-10-01T2010")`
- `#> [1] "2010-10-01 20:10:00 UTC"`
- `# If time is omitted, it will be set to midnight`
- `parse_datetime("20101010")`  
`#> [1] "2010-10-10 UTC"`

This is the most important date/time standard, and if you work with dates and times frequently, I recommend reading [https://en.wikipedia.org/wiki/ISO\\_8601](https://en.wikipedia.org/wiki/ISO_8601)

- `parse_date()` expects a four digit year, a - or /, the month, a - or /, then the day:
- `parse_date("2010-10-01")`  
`#> [1] "2010-10-01"`
- `parse_time()` expects the hour, :, minutes, optionally : and seconds, and an optional am/pm specifier:
- `library(hms)`
- `parse_time("01:10 am")`
- `#> 01:10:00`
- `parse_time("20:10:01")`

```
#> 20:10:01
```

Base R doesn't have a great built in class for time data, so we use the one provided in the `hms` package.

If these defaults don't work for your data you can supply your own date-time format, built up of the following pieces:

Year

`%Y` (4 digits).  
`%y` (2 digits); 00-69 -> 2000-2069, 70-99 -> 1970-1999.

Month

`%m` (2 digits).  
`%b` (abbreviated name, like "Jan").  
`%B` (full name, "January").

Day

`%d` (2 digits).  
`%e` (optional leading space).

Time

`%H` 0-23 hour.  
`%I` 0-12, must be used with `%p`.  
`%p` AM/PM indicator.  
`%M` minutes.  
`%S` integer seconds.  
`%OS` real seconds.  
`%z` Time zone (as name, e.g. America/Chicago). Beware of abbreviations: if you're American, note that "EST" is a Canadian time zone that does not have daylight savings time. It is not Eastern Standard Time! We'll come back to this [time zones](#).  
`%Z` (as offset from UTC, e.g. +0800).

Non-digits

`%.` skips one non-digit character.  
`%*` skips any number of non-digits.

The best way to figure out the correct format is to create a few examples in a character vector, and test with one of the parsing functions. For example:

```
parse_date("01/02/15", "%m/%d/%y")
#> [1] "2015-01-02"
parse_date("01/02/15", "%d/%m/%y")
#> [1] "2015-02-01"
parse_date("01/02/15", "%y/%m/%d")
#> [1] "2001-02-15"
```

If you're using `%b` or `%B` with non-English month names, you'll need to set the `lang` argument to `locale()`. See the list of built-in languages in `date_names_langs()`, or if your language is not already included, create your own with `date_names()`.

```
parse_date("1 janvier 2015", "%d %B %Y", locale = locale("fr"))
#> [1] "2015-01-01"
```

### 11.3.5 Exercises

1. What are the most important arguments to `locale()`?
2. What happens if you try and set `decimal_mark` and `grouping_mark` to the same character? What happens to the default value of `grouping_mark` when you set `decimal_mark` to ","? **What happens to the default value of `decimal_mark` when you set the `grouping_mark` to "."?**
3. I didn't discuss the `date_format` and `time_format` options to `locale()`. What do they do? Construct an example that shows when they might be useful.
4. If you live outside the US, create a new locale object that encapsulates the settings for the types of file you read most commonly.
5. **What's the difference between `read_csv()` and `read_csv2()`?**
6. What are the most common encodings used in Europe? What are the most common encodings used in Asia? Do some googling to find out.
7. Generate the correct format string to parse each of the following dates and times:
  8. `d1 <- "January 1, 2010"`
  9. `d2 <- "2015-Mar-07"`
  10. `d3 <- "06-Jun-2017"`
  11. `d4 <- c("August 19 (2015)", "July 1 (2015)")`
  12. `d5 <- "12/30/14" # Dec 30, 2014`
  13. `t1 <- "1705"`  
`t2 <- "11:15:10.12 PM"`

## 11.4 Parsing a file

Now that you've learned how to parse an individual vector, it's time to return to the beginning and explore how `readr` parses a file. There are two new things that you'll learn about in this section:

1. How `readr` automatically guesses the type of each column.
2. How to override the default specification.

### 11.4.1 Strategy

`readr` uses a heuristic to figure out the type of each column: it reads the first 1000 rows and uses some (moderately conservative) heuristics to figure out the type of each column. You can emulate this process with a character vector using `guess_parser()`, which returns `readr`'s best guess, and `parse_guess()` which uses that guess to parse the column:

```
guess_parser("2010-10-01")
#> [1] "date"
guess_parser("15:01")
#> [1] "time"
guess_parser(c("TRUE", "FALSE"))
#> [1] "logical"
guess_parser(c("1", "5", "9"))
#> [1] "integer"
guess_parser(c("12,352,561"))
#> [1] "number"

str(parse_guess("2010-10-10"))
```

```
#> Date[1:1], format: "2010-10-10"
```

The heuristic tries each of the following types, stopping when it finds a match:

- logical: contains only “F”, “T”, “FALSE”, or “TRUE”.
- integer: contains only numeric characters (and –).
- double: contains only valid doubles (including numbers like  $4.5e-5$ ).
- number: contains valid doubles with the grouping mark inside.
- time: matches the default `time_format`.
- date: matches the default `date_format`.
- date-time: any ISO8601 date.

If none of these rules apply, then the column will stay as a vector of strings.

## 11.4.2 Problems

These defaults don't always work for larger files. There are two basic problems:

1. The first thousand rows might be a special case, and `readr` guesses a type that is not sufficiently general. For example, you might have a column of doubles that only contains integers in the first 1000 rows.
2. The column might contain a lot of missing values. If the first 1000 rows contain only `NAs`, **readr will guess that it's a character vector, whereas you probably want to parse it as something more specific.**

`readr` contains a challenging CSV that illustrates both of these problems:

```
challenge <- read_csv(readr_example("challenge.csv"))
#> Parsed with column specification:
#> cols(
#> x = col_integer(),
#> y = col_character()
#>)
#> Warning: 1000 parsing failures.
#> row col expected actual
file
#> 1001 x no trailing characters .23837975086644292
'./home/travis/R/Library/readr/extdata/challenge.csv'
#> 1002 x no trailing characters .41167997173033655
'./home/travis/R/Library/readr/extdata/challenge.csv'
#> 1003 x no trailing characters .7460716762579978
'./home/travis/R/Library/readr/extdata/challenge.csv'
#> 1004 x no trailing characters .723450553836301
'./home/travis/R/Library/readr/extdata/challenge.csv'
#> 1005 x no trailing characters .614524137461558
'./home/travis/R/Library/readr/extdata/challenge.csv'
#>
.....
#> See problems(...) for more details.
```

(Note the use of `readr_example()` which finds the path to one of the files included with the package)

There are two printed outputs: the column specification generated by looking at the first 1000 rows, and the first five **parsing failures**. It's always a good idea to explicitly pull out the `problems()`, so you can explore them in more depth:

```
problems(challenge)
#> # A tibble: 1,000 × 5
#> row col expected actual
#> <int> <chr> <chr> <chr>
#> 1 1001 x no trailing characters .23837975086644292
#> 2 1002 x no trailing characters .41167997173033655
#> 3 1003 x no trailing characters .7460716762579978
#> 4 1004 x no trailing characters .723450553836301
#> 5 1005 x no trailing characters .614524137461558
#> 6 1006 x no trailing characters .473980569280684
#> # ... with 994 more rows, and 1 more variables: file <chr>
```

A good strategy is to work column by column until there are no problems remaining. Here we can see that there are a lot of parsing problems with the `x` column - there are trailing characters after the integer value. That suggests we need to use a double parser instead.

To fix the call, start by copying and pasting the column specification into your original call:

```
challenge <- read_csv(
 readr_example("challenge.csv"),
 col_types = cols(
 x = col_integer(),
 y = col_character()
)
)
```

Then you can tweak the type of the `x` column:

```
challenge <- read_csv(
 readr_example("challenge.csv"),
 col_types = cols(
 x = col_double(),
 y = col_character()
)
)
```

That fixes the first problem, but if we look at the last few rows, you'll see that they're dates stored in a character vector:

```
tail(challenge)
#> # A tibble: 6 × 2
#> x y
#> <dbl> <chr>
#> 1 0.805 2019-11-21
#> 2 0.164 2018-03-29
```

```
#> 3 0.472 2014-08-04
#> 4 0.718 2015-08-16
#> 5 0.270 2020-02-04
#> 6 0.608 2019-01-06
```

You can fix that by specifying that `y` is a date column:

```
challenge <- read_csv(
 readr_example("challenge.csv"),
 col_types = cols(
 x = col_double(),
 y = col_date()
)
)
tail(challenge)
#> # A tibble: 6 × 2
#> x y
#> <dbl> <date>
#> 1 0.805 2019-11-21
#> 2 0.164 2018-03-29
#> 3 0.472 2014-08-04
#> 4 0.718 2015-08-16
#> 5 0.270 2020-02-04
#> 6 0.608 2019-01-06
```

Every `parse_xyz()` function has a corresponding `col_xyz()` function. You use `parse_xyz()` when the data is in a character vector in R already; you use `col_xyz()` when you want to tell `readr` how to load the data.

I highly recommend always supplying `col_types`, building up from the print-out provided by `readr`. This ensures that you have a consistent and reproducible data import script. If you rely on the default guesses and your data changes, `readr` will continue to read it in. If you want to be really strict, use `stop_for_problems()`: that will throw an error and stop your script if there are any parsing problems.

### 11.4.3 Other strategies

There are a few other general strategies to help you parse files:

- In the previous example, we just got unlucky: if we look at just one more row than the default, we can correctly parse in one shot:
- `challenge2 <- read_csv(readr_example("challenge.csv"), guess_max = 1001)`
- `#> Parsed with column specification:`
- `#> cols(`
- `#> x = col_double(),`
- `#> y = col_date(format = "")`
- `#> )`
- `challenge2`
- `#> # A tibble: 2,000 × 2`
- `#> x y`

- #> <dbl> <date>
- #> 1 404 <NA>
- #> 2 4172 <NA>
- #> 3 3004 <NA>
- #> 4 787 <NA>
- #> 5 37 <NA>
- #> 6 2332 <NA>
- #> # ... with 1,994 more rows
- **Sometimes it's easier to diagnose problems if you just read in all the columns as character vectors:**
- challenge2 <- read\_csv(readr\_example("challenge.csv"),
- col\_types = cols(.default = col\_character())
- )

This is particularly useful in conjunction with `type_convert()`, which applies the parsing heuristics to the character columns in a data frame.

```
df <- tribble(
 ~x, ~y,
 "1", "1.21",
 "2", "2.32",
 "3", "4.56"
)
df
#> # A tibble: 3 × 2
#> x y
#> <chr> <chr>
#> 1 1 1.21
#> 2 2 2.32
#> 3 3 4.56

Note the column types
type_convert(df)
#> Parsed with column specification:
#> cols(
#> x = col_integer(),
#> y = col_double()
#>)
#> # A tibble: 3 × 2
#> x y
#> <int> <dbl>
#> 1 1 1.21
#> 2 2 2.32
#> 3 3 4.56
```

- **If you're reading a very large file, you might want to set `n_max` to a smallish number like 10,000 or 100,000.** That will accelerate your iterations while you eliminate common problems.
- **If you're having major parsing problems, sometimes it's easier** to just read into a character vector of lines with `read_lines()`, or even a character vector of length 1 with

`read_file()`. Then you can use the string parsing skills you'll learn later to parse more exotic formats.

## 11.5 Writing to a file

readr also comes with two useful functions for writing data back to disk: `write_csv()` and `write_tsv()`. Both functions increase the chances of the output file being read back in correctly by:

- Always encoding strings in UTF-8.
- Saving dates and date-times in ISO8601 format so they are easily parsed elsewhere.

If you want to export a csv file to Excel, use `write_excel_csv()` — this writes a special character (a “byte order mark”) at the start of the file which tells Excel that you’re using the UTF-8 encoding.

The most important arguments are `x` (the data frame to save), and `path` (the location to save it). You can also specify how missing values are written with `na`, and if you want to append to an existing file.

```
write_csv(challenge, "challenge.csv")
```

Note that the type information is lost when you save to csv:

```
challenge
#> # A tibble: 2,000 × 2
#> x y
#> <dbl> <date>
#> 1 404 <NA>
#> 2 4172 <NA>
#> 3 3004 <NA>
#> 4 787 <NA>
#> 5 37 <NA>
#> 6 2332 <NA>
#> # ... with 1,994 more rows
write_csv(challenge, "challenge-2.csv")
read_csv("challenge-2.csv")
#> Parsed with column specification:
#> cols(
#> x = col_integer(),
#> y = col_character()
#>)
#> # A tibble: 2,000 × 2
#> x y
#> <int> <chr>
#> 1 404 <NA>
#> 2 4172 <NA>
#> 3 3004 <NA>
#> 4 787 <NA>
#> 5 37 <NA>
#> 6 2332 <NA>
```

```
#> # ... with 1,994 more rows
```

This makes CSVs a little unreliable for caching interim results—you need to recreate the column specification every time you load in. There are two alternatives:

1. `write_rds()` and `read_rds()` are uniform wrappers around the base functions `readRDS()` and `saveRDS()`. These store data in R's custom binary format called **RDS**:
2. `write_rds(challenge, "challenge.rds")`
3. `read_rds("challenge.rds")`
4. #> # A tibble: 2,000 × 2
5. #> x y
6. #> <dbl> <date>
7. #> 1 404 <NA>
8. #> 2 4172 <NA>
9. #> 3 3004 <NA>
10. #> 4 787 <NA>
11. #> 5 37 <NA>
12. #> 6 2332 <NA>
- #> # ... with 1,994 more rows

13. The feather package implements a fast binary file format that can be shared across programming languages:

14. `library(feather)`
15. `write_feather(challenge, "challenge.feather")`
16. `read_feather("challenge.feather")`
17. #> # A tibble: 2,000 × 2
18. #> x y
19. #> <dbl> <date>
20. #> 1 404 <NA>
21. #> 2 4172 <NA>
22. #> 3 3004 <NA>
23. #> 4 787 <NA>
24. #> 5 37 <NA>
25. #> 6 2332 <NA>
- #> # ... with 1,994 more rows

Feather tends to be faster than RDS and is usable outside of R. RDS supports list-columns (**which you'll learn about in many models**); feather currently does not.

## 11.6 Other types of data

To get other types of data into R, we recommend starting with the tidyverse packages listed below. They're certainly not perfect, but they are a good place to start. For rectangular data:

- `haven` reads SPSS, Stata, and SAS files.
- `readxl` reads excel files (both `.xls` and `.xlsx`).
- `DBI`, along with a database specific backend (e.g. `RMySQL`, `RSQLite`, `RPostgreSQL` etc) allows you to run SQL queries against a database and return a data frame.

For hierarchical data: use `jsonlite` (by Jeroen Ooms) for json, and `xml2` for XML. Jenny Bryan has some excellent worked examples at <https://jennybc.github.io/purrr-tutorial/>.

For other file types, try the [R data import/export manual](#) and the `rio` package.

# 12 Tidy data

## 12.1 Introduction

“Happy families are all alike; every unhappy family is unhappy in its own way.” — Leo Tolstoy

“Tidy datasets are all alike, but every messy dataset is messy in its own way.” — Hadley Wickham

In this chapter, you will learn a consistent way to organise your data in R, an organisation called tidy data. Getting your data into this format requires some upfront work, but that work pays off in the long term. Once you have tidy data and the tidy tools provided by packages in the tidyverse, you will spend much less time munging data from one representation to another, allowing you to spend more time on the analytic questions at hand.

This chapter will give you a practical introduction to tidy data and the accompanying tools in the **tidyr package**. If you’d like to learn more about the underlying theory, you might enjoy the **Tidy Data** paper published in the Journal of Statistical Software, <http://www.jstatsoft.org/v59/i10/paper>.

### 12.1.1 Prerequisites

In this chapter we’ll focus on **tidyverse**, a package that provides a bunch of tools to help tidy up your messy datasets. **tidyverse** is a member of the core tidyverse.

```
library(tidyverse)
```

## 12.2 Tidy data

You can represent the same underlying data in multiple ways. The example below shows the same data organised in four different ways. Each dataset shows the same values of four variables **country**, **year**, **population**, and **cases**, but each dataset organises the values in a different way.

```
table1
#> # A tibble: 6 × 4
#> country year cases population
#> <chr> <int> <int> <int>
#> 1 Afghanistan 1999 745 19987071
#> 2 Afghanistan 2000 2666 20595360
#> 3 Brazil 1999 37737 172006362
#> 4 Brazil 2000 80488 174504898
#> 5 China 1999 212258 1272915272
#> 6 China 2000 213766 1280428583
table2
#> # A tibble: 12 × 4
#> country year type count
#> <chr> <int> <chr> <int>
#> 1 Afghanistan 1999 cases 745
#> 2 Afghanistan 1999 population 19987071
```

```

#> 3 Afghanistan 2000 cases 2666
#> 4 Afghanistan 2000 population 20595360
#> 5 Brazil 1999 cases 37737
#> 6 Brazil 1999 population 172006362
#> # ... with 6 more rows
table3
#> # A tibble: 6 × 3
#> country year rate
#> <chr> <int> <chr>
#> 1 Afghanistan 1999 745/19987071
#> 2 Afghanistan 2000 2666/20595360
#> 3 Brazil 1999 37737/172006362
#> 4 Brazil 2000 80488/174504898
#> 5 China 1999 212258/1272915272
#> 6 China 2000 213766/1280428583

Spread across two tibbles
table4a # cases
#> # A tibble: 3 × 3
#> country `1999` `2000`
#> <chr> <int> <int>
#> 1 Afghanistan 745 2666
#> 2 Brazil 37737 80488
#> 3 China 212258 213766
table4b # population
#> # A tibble: 3 × 3
#> country `1999` `2000`
#> <chr> <int> <int>
#> 1 Afghanistan 19987071 20595360
#> 2 Brazil 172006362 174504898
#> 3 China 1272915272 1280428583

```

These are all representations of the same underlying data, but they are not equally easy to use. One dataset, the tidy dataset, will be much easier to work with inside the tidyverse.

There are three interrelated rules which make a dataset tidy:

1. Each variable must have its own column.
2. Each observation must have its own row.
3. Each value must have its own cell.

Figure 12.1 shows the rules visually.

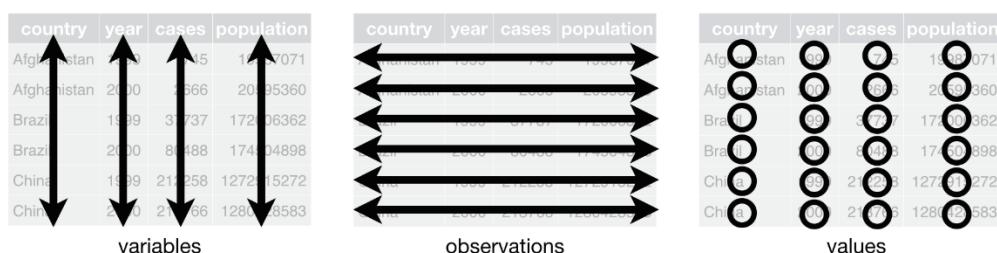


Figure 12.1: Following three rules makes a dataset tidy: variables are in columns, observations are in rows, and values are in cells.

These three rules are interrelated because it's impossible to only satisfy two of the three. That interrelationship leads to an even simpler set of practical instructions:

1. Put each dataset in a tibble.
2. Put each variable in a column.

In this example, only `table1` is tidy. It's the only representation where each column is a variable.

Why ensure that your data is tidy? There are two main advantages:

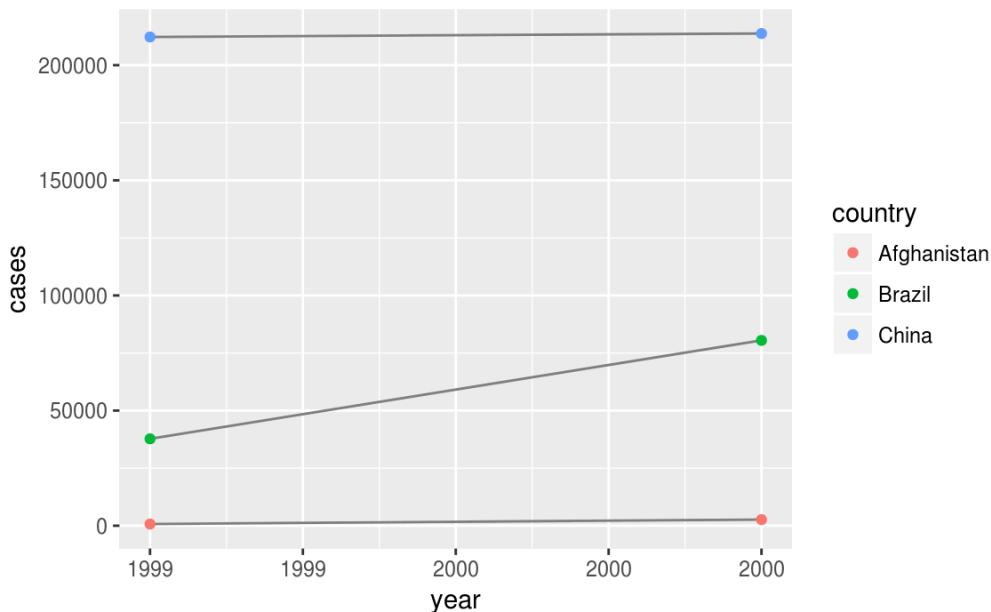
1. There's a general advantage to picking one consistent way of storing data. If you have a consistent data structure, it's easier to learn the tools that work with it because they have an underlying uniformity.
2. There's a specific advantage to placing variables in columns because it allows R's vectorised nature to shine. As you learned in `mutate` and `summary functions`, most built-in R functions work with vectors of values. That makes transforming tidy data feel particularly natural.

`dplyr`, `ggplot2`, and all the other packages in the `tidyverse` are designed to work with tidy data. Here are a couple of small examples showing how you might work with `table1`.

```
Compute rate per 10,000
table1 %>%
 mutate(rate = cases / population * 10000)
#> # A tibble: 6 × 5
#> country year cases population rate
#> <chr> <int> <int> <int> <dbl>
#> 1 Afghanistan 1999 745 19987071 0.373
#> 2 Afghanistan 2000 2666 20595360 1.294
#> 3 Brazil 1999 37737 172006362 2.194
#> 4 Brazil 2000 80488 174504898 4.612
#> 5 China 1999 212258 1272915272 1.667
#> 6 China 2000 213766 1280428583 1.669

Compute cases per year
table1 %>%
 count(year, wt = cases)
#> # A tibble: 2 × 2
#> year n
#> <int> <int>
#> 1 1999 250740
#> 2 2000 296920

Visualise changes over time
library(ggplot2)
ggplot(table1, aes(year, cases)) +
 geom_line(aes(group = country), colour = "grey50") +
 geom_point(aes(colour = country))
```



### 12.2.1 Exercises

1. Using prose, describe how the variables and observations are organised in each of the sample tables.
2. Compute the `rate` for `table2`, and `table4a + table4b`. You will need to perform four operations:
  1. Extract the number of TB cases per country per year.
  2. Extract the matching population per country per year.
  3. Divide cases by population, and multiply by 10000.
  4. Store back in the appropriate place.

Which representation is easiest to work with? Which is hardest? Why?

3. Recreate the plot showing change in cases over time using `table2` instead of `table1`. What do you need to do first?

## 12.3 Spreading and gathering

The principles of tidy data seem so obvious that you might wonder if you'll ever encounter a dataset that isn't tidy. Unfortunately, however, most data that you will encounter will be untidy. There are two main reasons:

1. **Most people aren't familiar with the principles of tidy data, and it's hard to derive them yourself unless you spend a lot of time working with data.**
2. Data is often organised to facilitate some use other than analysis. For example, data is often organised to make entry as easy as possible.

This means for most real analyses, you'll need to do some tidying. The first step is always to figure out what the variables and observations are. Sometimes this is easy; other times you'll

need to consult with the people who originally generated the data. The second step is to resolve one of two common problems:

1. One variable might be spread across multiple columns.
2. One observation might be scattered across multiple rows.

Typically a dataset will only suffer from one of these problems; it'll only suffer from both if you're really unlucky! To fix these problems, you'll need the two most important functions in `tidyR`: `gather()` and `spread()`.

### 12.3.1 Gathering

A common problem is a dataset where some of the column names are not names of variables, but **values** of a variable. Take `table4a`: the column names 1999 and 2000 represent values of the `year` variable, and each row represents two observations, not one.

```
table4a
#> # A tibble: 3 × 3
#> country `1999` `2000`
#> <chr> <int> <int>
#> 1 Afghanistan 745 2666
#> 2 Brazil 37737 80488
#> 3 China 212258 213766
```

To tidy a dataset like this, we need to gather those columns into a new pair of variables. To describe that operation we need three parameters:

- The set of columns that represent values, not variables. In this example, those are the columns 1999 and 2000.
- The name of the variable whose values form the column names. I call that the `key`, and here it is `year`.
- The name of the variable whose values are spread over the cells. I call that `value`, and here it's the number of cases.

Together those parameters generate the call to `gather()`:

```
table4a %>%
 gather(`1999`, `2000`, key = "year", value = "cases")
#> # A tibble: 6 × 3
#> country year cases
#> <chr> <chr> <int>
#> 1 Afghanistan 1999 745
#> 2 Brazil 1999 37737
#> 3 China 1999 212258
#> 4 Afghanistan 2000 2666
#> 5 Brazil 2000 80488
#> 6 China 2000 213766
```

The columns to gather are specified with `dplyr::select()` style notation. Here there are only two columns, so we list them individually. Note that "1999" and "2000" are non-syntactic names

(because they don't start with a letter) so we have to surround them in backticks. To refresh your memory of the other ways to select columns, see [select](#).

| country     | year | cases  | country     | 1999   | 2000   |
|-------------|------|--------|-------------|--------|--------|
| Afghanistan | 1999 | 745    | Afghanistan | 745    | 2666   |
| Afghanistan | 2000 | 2666   | Brazil      | 37737  | 80488  |
| Brazil      | 1999 | 37737  | China       | 212258 | 213766 |
| Brazil      | 2000 | 80488  |             |        |        |
| China       | 1999 | 212258 |             |        |        |
| China       | 2000 | 213766 |             |        |        |

Figure 12.2: Gathering `table4` into a tidy form.

In the final result, the gathered columns are dropped, and we get new `key` and `value` columns. Otherwise, the relationships between the original variables are preserved. Visually, this is shown in Figure 12.2. We can use `gather()` to tidy `table4b` in a similar fashion. The only difference is the variable stored in the cell values:

```
table4b %>%
 gather(`1999`, `2000`, key = "year", value = "population")
#> # A tibble: 6 × 3
#> country year population
#> <chr> <chr> <int>
#> 1 Afghanistan 1999 19987071
#> 2 Brazil 1999 172006362
#> 3 China 1999 1272915272
#> 4 Afghanistan 2000 20595360
#> 5 Brazil 2000 174504898
#> 6 China 2000 1280428583
```

To combine the tidied versions of `table4a` and `table4b` into a single tibble, we need to use `dplyr::left_join()`, which you'll learn about in [relational data](#).

```
tidy4a <- table4a %>%
 gather(`1999`, `2000`, key = "year", value = "cases")
tidy4b <- table4b %>%
 gather(`1999`, `2000`, key = "year", value = "population")
left_join(tidy4a, tidy4b)
#> Joining, by = c("country", "year")
#> # A tibble: 6 × 4
#> country year cases population
#> <chr> <chr> <int> <int>
#> 1 Afghanistan 1999 745 19987071
#> 2 Brazil 1999 37737 172006362
#> 3 China 1999 212258 1272915272
#> 4 Afghanistan 2000 2666 20595360
#> 5 Brazil 2000 80488 174504898
#> 6 China 2000 213766 1280428583
```

### 12.3.2 Spreading

Spreading is the opposite of gathering. You use it when an observation is scattered across multiple rows. For example, take `table2`: an observation is a country in a year, but each observation is spread across two rows.

```
table2
#> # A tibble: 12 × 4
#> country year type count
#> <chr> <int> <chr> <int>
#> 1 Afghanistan 1999 cases 745
#> 2 Afghanistan 1999 population 19987071
#> 3 Afghanistan 2000 cases 2666
#> 4 Afghanistan 2000 population 20595360
#> 5 Brazil 1999 cases 37737
#> 6 Brazil 1999 population 172006362
#> # ... with 6 more rows
```

To tidy this up, we first analyse the representation in similar way to `gather()`. This time, however, we only need two parameters:

- The column that contains variable names, the key **column**. Here, it's `type`.
- The column that contains values forms multiple variables, the value **column**. Here it's `count`.

Once we've figured that out, we can use `spread()`, as shown programmatically below, and visually in Figure 12.3.

```
spread(table2, key = type, value = count)
#> # A tibble: 6 × 4
#> country year cases population
#> <chr> <int> <int> <int>
#> 1 Afghanistan 1999 745 19987071
#> 2 Afghanistan 2000 2666 20595360
#> 3 Brazil 1999 37737 172006362
#> 4 Brazil 2000 80488 174504898
#> 5 China 1999 212258 1272915272
#> 6 China 2000 213766 1280428583
```

The diagram illustrates the transformation of a wide table (table2) into a tidy table. On the left, the wide table has four columns: country, year, key, and value. On the right, the tidy table has four columns: country, year, cases, and population. Arrows point from the country, year, and value columns of the wide table to their respective columns in the tidy table. The value column is transformed into two new columns: cases and population.

| country     | year | key        | value      |
|-------------|------|------------|------------|
| Afghanistan | 1999 | cases      | 745        |
| Afghanistan | 1999 | population | 19987071   |
| Afghanistan | 2000 | cases      | 2666       |
| Afghanistan | 2000 | population | 20595360   |
| Brazil      | 1999 | cases      | 37737      |
| Brazil      | 1999 | population | 172006362  |
| Brazil      | 2000 | cases      | 80488      |
| Brazil      | 2000 | population | 174504898  |
| China       | 1999 | cases      | 212258     |
| China       | 1999 | population | 1272915272 |
| China       | 2000 | cases      | 213766     |
| China       | 2000 | population | 1280428583 |

| country     | year | cases  | population |
|-------------|------|--------|------------|
| Afghanistan | 1999 | 745    | 19987071   |
| Afghanistan | 2000 | 2666   | 20595360   |
| Brazil      | 1999 | 37737  | 172006362  |
| Brazil      | 2000 | 80488  | 174504898  |
| China       | 1999 | 212258 | 1272915272 |
| China       | 2000 | 213766 | 1280428583 |

table2

Figure 12.3: Spreading table2 makes it tidy

As you might have guessed from the common `key` and `value` arguments, `spread()` and `gather()` are complements. `gather()` makes wide tables narrower and longer; `spread()` makes long tables shorter and wider.

### 12.3.3 Exercises

1. Why are `gather()` and `spread()` not perfectly symmetrical?

Carefully consider the following example:

```
2. stocks <- tibble(
 3. year = c(2015, 2015, 2016, 2016),
 4. half = c(1, 2, 1, 2),
 5. return = c(1.88, 0.59, 0.92, 0.17)
 6.)
 7. stocks %>%
 8. spread(year, return) %>%
 gather("year", "return", `2015`:`2016`)
```

(Hint: look at the variable types and think about column names.)

Both `spread()` and `gather()` have a `convert` argument. What does it do?

9. Why does this code fail?

```
10. table4a %>%
 11. gather(1999, 2000, key = "year", value = "cases")
 #> Error in combine_vars(vars, ind_list): Position must be between
 0 and n
```

12. Why does spreading this tibble fail? How could you add a new column to fix the problem?

```
13. people <- tribble(
 14. ~name, ~key, ~value,
 15. #-----|-----|-----
 16. "Phillip Woods", "age", 45,
 17. "Phillip Woods", "height", 186,
 18. "Phillip Woods", "age", 50,
```

```

19. "Jessica Cordero", "age", 37,
20. "Jessica Cordero", "height", 156
)

21. Tidy the simple tibble below. Do you need to spread or gather it? What are the variables?
22. preg <- tribble(
23. ~pregnant, ~male, ~female,
24. "yes", NA, 10,
25. "no", 20, 12
)

```

## 12.4 Separating and uniting

So far you've learned how to `tidy` `table2` and `table4`, but not `table3`. `table3` has a different problem: we have one column (`rate`) that contains two variables (`cases` and `population`). To fix this problem, we'll need the `separate()` function. You'll also learn about the complement of `separate(): unite()`, which you use if a single variable is spread across multiple columns.

### 12.4.1 Separate

`separate()` pulls apart one column into multiple columns, by splitting wherever a separator character appears. Take `table3`:

```
table3
#> # A tibble: 6 × 3
#> country year rate
#> <chr> <int> <chr>
#> 1 Afghanistan 1999 745/19987071
#> 2 Afghanistan 2000 2666/20595360
#> 3 Brazil 1999 37737/172006362
#> 4 Brazil 2000 80488/174504898
#> 5 China 1999 212258/1272915272
#> 6 China 2000 213766/1280428583
```

The `rate` column contains both `cases` and `population` variables, and we need to split it into two variables. `separate()` takes the name of the column to separate, and the names of the columns to separate into, as shown in Figure 12.4 and the code below.

```
table3 %>%
 separate(rate, into = c("cases", "population"))
#> # A tibble: 6 × 4
#> country year cases population
#> <chr> <int> <chr> <chr>
#> 1 Afghanistan 1999 745 19987071
#> 2 Afghanistan 2000 2666 20595360
#> 3 Brazil 1999 37737 172006362
#> 4 Brazil 2000 80488 174504898
#> 5 China 1999 212258 1272915272
#> 6 China 2000 213766 1280428583
```

The diagram shows two tables side-by-side. The left table, labeled 'table3', has columns 'country', 'year', and 'rate'. The right table is the result of separating the 'rate' column into 'cases' and 'population'. An arrow points from the right side of the first table to the left side of the second table.

| country     | year | rate                |
|-------------|------|---------------------|
| Afghanistan | 1999 | 745 / 19987071      |
| Afghanistan | 2000 | 2666 / 20595360     |
| Brazil      | 1999 | 37737 / 172006362   |
| Brazil      | 2000 | 80488 / 174504898   |
| China       | 1999 | 212258 / 1272915272 |
| China       | 2000 | 213766 / 1280428583 |

| country     | year | cases  | population |
|-------------|------|--------|------------|
| Afghanistan | 1999 | 745    | 19987071   |
| Afghanistan | 2000 | 2666   | 20595360   |
| Brazil      | 1999 | 37737  | 172006362  |
| Brazil      | 2000 | 80488  | 174504898  |
| China       | 1999 | 212258 | 1272915272 |
| China       | 2000 | 213766 | 1280428583 |

table3

Figure 12.4: Separating table3 makes it tidy

By default, `separate()` will split values wherever it sees a non-alphanumeric character (i.e. a **character that isn't a number or letter**). For example, in the code above, `separate()` split the values of `rate` at the forward slash characters. If you wish to use a specific character to separate a column, you can pass the character to the `sep` argument of `separate()`. For example, we could rewrite the code above as:

```
table3 %>%
 separate(rate, into = c("cases", "population"), sep = "/")
```

(Formally, `sep` is a **regular expression**, which you'll learn more about in [strings](#).)

**Look carefully at the column types:** you'll notice that `case` and `population` are character columns. This is the default behaviour in `separate()`: it leaves the type of the column as is. **Here, however, it's not very useful as those really are numbers. We can ask `separate()` to try and convert to better types using `convert = TRUE`:**

```
table3 %>%
 separate(rate, into = c("cases", "population"), convert = TRUE)
#> # A tibble: 6 × 4
#> country year cases population
#> <chr> <int> <int> <int>
#> 1 Afghanistan 1999 745 19987071
#> 2 Afghanistan 2000 2666 20595360
#> 3 Brazil 1999 37737 172006362
#> 4 Brazil 2000 80488 174504898
#> 5 China 1999 212258 1272915272
#> 6 China 2000 213766 1280428583
```

You can also pass a vector of integers to `sep`. `separate()` will interpret the integers as positions to split at. Positive values start at 1 on the far-left of the strings; negative value start at -1 on the far-right of the strings. When using integers to separate strings, the length of `sep` should be one less than the number of names in `into`.

You can use this arrangement to separate the last two digits of each year. This make this data less tidy, but is useful in other cases, as you'll see in a little bit.

```
table3 %>%
 separate(year, into = c("century", "year"), sep = 2)
#> # A tibble: 6 × 4
#> country century year rate
#> <chr> <chr> <chr> <chr>
#> 1 Afghanistan 19 99 745/19987071
#> 2 Afghanistan 20 00 2666/20595360
#> 3 Brazil 19 99 37737/172006362
#> 4 Brazil 20 00 80488/174504898
#> 5 China 19 99 212258/1272915272
#> 6 China 20 00 213766/1280428583
```

## 12.4.2 Unite

`unite()` is the inverse of `separate()`: it combines multiple columns into a single column. You'll need it much less frequently than `separate()`, but it's still a useful tool to have in your back pocket.



| country     | year | rate                |
|-------------|------|---------------------|
| Afghanistan | 1999 | 745 / 19987071      |
| Afghanistan | 2000 | 2666 / 20595360     |
| Brazil      | 1999 | 37737 / 172006362   |
| Brazil      | 2000 | 80488 / 174504898   |
| China       | 1999 | 212258 / 1272915272 |
| China       | 2000 | 213766 / 1280428583 |

| country     | century | year | rate                |
|-------------|---------|------|---------------------|
| Afghanistan | 19      | 99   | 745 / 19987071      |
| Afghanistan | 20      | 0    | 2666 / 20595360     |
| Brazil      | 19      | 99   | 37737 / 172006362   |
| Brazil      | 20      | 0    | 80488 / 174504898   |
| China       | 19      | 99   | 212258 / 1272915272 |
| China       | 20      | 0    | 213766 / 1280428583 |

table6

Figure 12.5: Uniting table5 makes it tidy

We can use `unite()` to rejoin the `century` and `year` columns that we created in the last example. That data is saved as `tidyR::table5`. `unite()` takes a data frame, the name of the new variable to create, and a set of columns to combine, again specified in `dplyr::select()` style:

```
table5 %>%
 unite(new, century, year)
#> # A tibble: 6 × 3
#> country new rate
#> <chr> <chr> <chr>
#> 1 Afghanistan 19_99 745/19987071
#> 2 Afghanistan 20_00 2666/20595360
#> 3 Brazil 19_99 37737/172006362
#> 4 Brazil 20_00 80488/174504898
#> 5 China 19_99 212258/1272915272
```

```
#> 6 China 20_00 213766/1280428583
```

In this case we also need to use the `sep` argument. The default will place an underscore (`_`) between the values from different columns. Here we don't want any separator so we use `""`:

```
table5 %>%
 unite(new, century, year, sep = "")
#> # A tibble: 6 × 3
#> country new rate
#> <chr> <chr> <chr>
#> 1 Afghanistan 1999 745/19987071
#> 2 Afghanistan 2000 2666/20595360
#> 3 Brazil 1999 37737/172006362
#> 4 Brazil 2000 80488/174504898
#> 5 China 1999 212258/1272915272
#> 6 China 2000 213766/1280428583
```

### 12.4.3 Exercises

1. What do the `extra` and `fill` arguments do in `separate()`? Experiment with the various options for the following two toy datasets.
2. `tibble(x = c("a,b,c", "d,e,f,g", "h,i,j")) %>% separate(x, c("one", "two", "three"))`
- 3.
- 4.
5. `tibble(x = c("a,b,c", "d,e", "f,g,i")) %>% separate(x, c("one", "two", "three"))`
6. Both `unite()` and `separate()` have a `remove` argument. What does it do? Why would you set it to `FALSE`?
7. Compare and contrast `separate()` and `extract()`. Why are there three variations of separation (by position, by separator, and with groups), but only one `unite`?

## 12.5 Missing values

Changing the representation of a dataset brings up an important subtlety of missing values. Surprisingly, a value can be missing in one of two possible ways:

- Explicitly, i.e. flagged with `NA`.
- Implicitly, i.e. simply not present in the data.

Let's illustrate this idea with a very simple data set:

```
stocks <- tibble(
 year = c(2015, 2015, 2015, 2015, 2016, 2016, 2016),
 qtr = c(1, 2, 3, 4, 2, 3, 4),
 return = c(1.88, 0.59, 0.35, NA, 0.92, 0.17, 2.66)
)
```

There are two missing values in this dataset:

- The return for the fourth quarter of 2015 is explicitly missing, because the cell where its value should be instead contains `NA`.
- The return for the first quarter of 2016 is implicitly missing, because it simply does not appear in the dataset.

One way to think about the difference is with this Zen-like koan: An explicit missing value is the presence of an absence; an implicit missing value is the absence of a presence.

The way that a dataset is represented can make implicit values explicit. For example, we can make the implicit missing value explicit by putting years in the columns:

```
stocks %>%
 spread(year, return)
#> # A tibble: 4 × 3
#> qtr `2015` `2016`
#> * <dbl> <dbl> <dbl>
#> 1 1 1.88 NA
#> 2 2 0.59 0.92
#> 3 3 0.35 0.17
#> 4 4 NA 2.66
```

Because these explicit missing values may not be important in other representations of the data, you can set `na.rm = TRUE` in `gather()` to turn explicit missing values implicit:

```
stocks %>%
 spread(year, return) %>%
 gather(year, return, `2015`:`2016`, na.rm = TRUE)
#> # A tibble: 6 × 3
#> qtr year return
#> * <dbl> <chr> <dbl>
#> 1 1 2015 1.88
#> 2 2 2015 0.59
#> 3 3 2015 0.35
#> 4 2 2016 0.92
#> 5 3 2016 0.17
#> 6 4 2016 2.66
```

Another important tool for making missing values explicit in tidy data is `complete()`:

```
stocks %>%
 complete(year, qtr)
#> # A tibble: 8 × 3
#> year qtr return
#> <dbl> <dbl> <dbl>
#> 1 2015 1 1.88
#> 2 2015 2 0.59
#> 3 2015 3 0.35
#> 4 2015 4 NA
#> 5 2016 1 NA
#> 6 2016 2 0.92
#> # ... with 2 more rows
```

`complete()` takes a set of columns, and finds all unique combinations. It then ensures the original dataset contains all those values, filling in explicit `NAs` where necessary.

**There's one other important tool that you should know for working with missing values.**

Sometimes when a data source has primarily been used for data entry, missing values indicate that the previous value should be carried forward:

```
treatment <- tribble(
 ~ person, ~ treatment, ~response,
 "Derrick Whitmore", 1, 7,
 NA, 2, 10,
 NA, 3, 9,
 "Katherine Burke", 1, 4
)
```

You can fill in these missing values with `fill()`. It takes a set of columns where you want missing values to be replaced by the most recent non-missing value (sometimes called last observation carried forward).

```
treatment %>%
 fill(person)
#> # A tibble: 4 × 3
#> person treatment response
#> <chr> <dbl> <dbl>
#> 1 Derrick Whitmore 1 7
#> 2 Derrick Whitmore 2 10
#> 3 Derrick Whitmore 3 9
#> 4 Katherine Burke 1 4
```

### 12.5.1 Exercises

1. Compare and contrast the `fill` arguments to `spread()` and `complete()`.
2. What does the `direction` argument to `fill()` do?

## 12.6 Case Study

To finish off the chapter, let's pull together everything you've learned to tackle a realistic data tidying problem. The `tidyverse::who` dataset contains tuberculosis (TB) cases broken down by year, country, age, gender, and diagnosis method. The data comes from the `2014 World Health Organization Global Tuberculosis Report`, available at <http://www.who.int/tb/country/data/download/en/>.

There's a wealth of epidemiological information in this dataset, but it's challenging to work with the data in the form that it's provided:

```
who
#> # A tibble: 7,240 × 60
#> country iso2 iso3 year new_sp_m014 new_sp_m1524 new_sp_m2534
#> <chr> <chr> <chr> <int> <int> <int> <int>
#> 1 Afghanistan AF AFG 1980 NA NA NA
```

```

#> 2 Afghanistan AF AFG 1981 NA NA NA
#> 3 Afghanistan AF AFG 1982 NA NA NA
#> 4 Afghanistan AF AFG 1983 NA NA NA
#> 5 Afghanistan AF AFG 1984 NA NA NA
#> 6 Afghanistan AF AFG 1985 NA NA NA
#> # ... with 7,234 more rows, and 53 more variables: new_sp_m3544
<int>,
#> # new_sp_m4554 <int>, new_sp_m5564 <int>, new_sp_m65 <int>,
#> # new_sp_f014 <int>, new_sp_f1524 <int>, new_sp_f2534 <int>,
#> # new_sp_f3544 <int>, new_sp_f4554 <int>, new_sp_f5564 <int>,
#> # new_sp_f65 <int>, new_sn_m014 <int>, new_sn_m1524 <int>,
#> # new_sn_m2534 <int>, new_sn_m3544 <int>, new_sn_m4554 <int>,
#> # new_sn_m5564 <int>, new_sn_m65 <int>, new_sn_f014 <int>,
#> # new_sn_f1524 <int>, new_sn_f2534 <int>, new_sn_f3544 <int>,
#> # new_sn_f4554 <int>, new_sn_f5564 <int>, new_sn_f65 <int>,
#> # new_ep_m014 <int>, new_ep_m1524 <int>, new_ep_m2534 <int>,
#> # new_ep_m3544 <int>, new_ep_m4554 <int>, new_ep_m5564 <int>,
#> # new_ep_m65 <int>, new_ep_f014 <int>, new_ep_f1524 <int>,
#> # new_ep_f2534 <int>, new_ep_f3544 <int>, new_ep_f4554 <int>,
#> # new_ep_f5564 <int>, new_ep_f65 <int>, newrel_m014 <int>,
#> # newrel_m1524 <int>, newrel_m2534 <int>, newrel_m3544 <int>,
#> # newrel_m4554 <int>, newrel_m5564 <int>, newrel_m65 <int>,
#> # newrel_f014 <int>, newrel_f1524 <int>, newrel_f2534 <int>,
#> # newrel_f3544 <int>, newrel_f4554 <int>, newrel_f5564 <int>,
#> # newrel_f65 <int>

```

This is a very typical real-life example dataset. It contains redundant columns, odd variable codes, and many missing values. In short, `who` is messy, and **we'll need multiple steps to tidy it**. Like `dplyr`, `tidyverse` is designed so that each function does one thing well. That means in real-life situations you'll usually need to string together multiple verbs into a pipeline.

The best place to start is almost always to gather together the columns that are not variables. Let's have a look at what we've got:

- It looks like `country`, `iso2`, and `iso3` are three variables that redundantly specify the country.
- `year` is clearly also a variable.
- **We don't know what all the other columns are yet**, but given the structure in the variable names (e.g. `new_sp_m014`, `new_ep_m014`, `new_ep_f014`) these are likely to be values, not variables.

So we need to gather together all the columns from `new_sp_m014` to `newrel_f65`. **We don't know what those values represent yet, so we'll give them the generic name "key"**. We know the **cells represent the count of cases**, so we'll use the variable `cases`. There are a lot of missing **values in the current representation**, so for now we'll use `na.rm` just so we can focus on the values that are present.

```

who1 <- who %>%
 gather(new_sp_m014:newrel_f65, key = "key", value = "cases", na.rm =
TRUE)
who1
#> # A tibble: 76,046 × 6

```

```

#> country iso2 iso3 year key cases
#> * <chr> <chr> <chr> <int> <chr> <int>
#> 1 Afghanistan AF AFG 1997 new_sp_m014 0
#> 2 Afghanistan AF AFG 1998 new_sp_m014 30
#> 3 Afghanistan AF AFG 1999 new_sp_m014 8
#> 4 Afghanistan AF AFG 2000 new_sp_m014 52
#> 5 Afghanistan AF AFG 2001 new_sp_m014 129
#> 6 Afghanistan AF AFG 2002 new_sp_m014 90
#> # ... with 7.604e+04 more rows

```

We can get some hint of the structure of the values in the new `key` column by counting them:

```

who1 %>%
 count(key)
#> # A tibble: 56 × 2
#> key n
#> <chr> <int>
#> 1 new_ep_f014 1032
#> 2 new_ep_f1524 1021
#> 3 new_ep_f2534 1021
#> 4 new_ep_f3544 1021
#> 5 new_ep_f4554 1017
#> 6 new_ep_f5564 1017
#> # ... with 50 more rows

```

You might be able to parse this out by yourself with a little thought and some experimentation, but luckily we have the data dictionary handy. It tells us:

1. The first three letters of each column denote whether the column contains new or old cases of TB. In this dataset, each column contains new cases.
2. The next two letters describe the type of TB:
  - o `rel` stands for cases of relapse
  - o `ep` stands for cases of extrapulmonary TB
  - o `sn` stands for cases of pulmonary TB that could not be diagnosed by a pulmonary smear (smear negative)
  - o `sp` stands for cases of pulmonary TB that could be diagnosed by a pulmonary smear (smear positive)
3. The sixth letter gives the sex of TB patients. The dataset groups cases by males (`m`) and females (`f`).
4. The remaining numbers gives the age group. The dataset groups cases into seven age groups:
  - o `014` = 0 – 14 years old
  - o `1524` = 15 – 24 years old
  - o `2534` = 25 – 34 years old
  - o `3544` = 35 – 44 years old
  - o `4554` = 45 – 54 years old
  - o `5564` = 55 – 64 years old
  - o `65` = 65 or older

We need to make a minor fix to the format of the column names: unfortunately the names are slightly inconsistent because instead of `new_rel` we have `newrel` (**it's hard to spot this here but**

if you don't fix it we'll get errors in subsequent steps). You'll learn about `str_replace()` in [strings](#), but the basic idea is pretty simple: replace the characters "newrel" with "new\_rel". This makes all variable names consistent.

```
who2 <- whol %>%
 mutate(key = stringr::str_replace(key, "newrel", "new_rel"))
who2
#> # A tibble: 76,046 × 6
#> country iso2 iso3 year key cases
#> <chr> <chr> <chr> <int> <chr> <int>
#> 1 Afghanistan AF AFG 1997 new_sp_m014 0
#> 2 Afghanistan AF AFG 1998 new_sp_m014 30
#> 3 Afghanistan AF AFG 1999 new_sp_m014 8
#> 4 Afghanistan AF AFG 2000 new_sp_m014 52
#> 5 Afghanistan AF AFG 2001 new_sp_m014 129
#> 6 Afghanistan AF AFG 2002 new_sp_m014 90
#> # ... with 7.604e+04 more rows
```

We can separate the values in each code with two passes of `separate()`. The first pass will split the codes at each underscore.

```
who3 <- who2 %>%
 separate(key, c("new", "type", "sexage"), sep = "_")
who3
#> # A tibble: 76,046 × 8
#> country iso2 iso3 year new type sexage cases
#> <chr> <chr> <chr> <int> <chr> <chr> <chr> <int>
#> 1 Afghanistan AF AFG 1997 new sp m014 0
#> 2 Afghanistan AF AFG 1998 new sp m014 30
#> 3 Afghanistan AF AFG 1999 new sp m014 8
#> 4 Afghanistan AF AFG 2000 new sp m014 52
#> 5 Afghanistan AF AFG 2001 new sp m014 129
#> 6 Afghanistan AF AFG 2002 new sp m014 90
#> # ... with 7.604e+04 more rows
```

Then we might as well drop the new column because it's constant in this dataset. While we're dropping columns, let's also drop `iso2` and `iso3` since they're redundant.

```
who3 %>%
 count(new)
#> # A tibble: 1 × 2
#> new n
#> <chr> <int>
#> 1 new 76046
who4 <- who3 %>%
 select(-new, -iso2, -iso3)
```

Next we'll separate `sexage` into `sex` and `age` by splitting after the first character:

```
who5 <- who4 %>%
 separate(sexage, c("sex", "age"), sep = 1)
who5
```

```
#> # A tibble: 76,046 × 6
#> country year type sex age cases
#> <chr> <int> <chr> <chr> <chr> <int>
#> 1 Afghanistan 1997 sp m 014 0
#> 2 Afghanistan 1998 sp m 014 30
#> 3 Afghanistan 1999 sp m 014 8
#> 4 Afghanistan 2000 sp m 014 52
#> 5 Afghanistan 2001 sp m 014 129
#> 6 Afghanistan 2002 sp m 014 90
#> # ... with 7.604e+04 more rows
```

The `who` dataset is now tidy!

I've shown you the code a piece at a time, assigning each interim result to a new variable. This typically isn't how you'd work interactively. Instead, you'd gradually build up a complex pipe:

```
who %>%
 gather(code, value, new_sp_m014:newrel_f65, na.rm = TRUE) %>%
 mutate(code = stringr::str_replace(code, "newrel", "new_rel")) %>%
 separate(code, c("new", "var", "sexage")) %>%
 select(-new, -iso2, -iso3) %>%
 separate(sexage, c("sex", "age"), sep = 1)
```

## 12.6.1 Exercises

1. In this case study I set `na.rm = TRUE` just to make it easier to check that we had the correct values. Is this reasonable? Think about how missing values are represented in this dataset. Are there implicit missing values? What's the difference between an NA and zero?
2. What happens if you neglect the `mutate()` step? (`mutate(key = stringr::str_replace(key, "newrel", "new_rel"))`)
3. I claimed that `iso2` and `iso3` were redundant with `country`. Confirm this claim.
4. For each country, year, and sex compute the total number of cases of TB. Make an informative visualisation of the data.

## 12.7 Non-tidy data

Before we continue on to other topics, it's worth talking briefly about non-tidy data. Earlier in the chapter, I used the pejorative term “messy” to refer to non-tidy data. That's an oversimplification: there are lots of useful and well-founded data structures that are not tidy data. There are two main reasons to use other data structures:

- Alternative representations may have substantial performance or space advantages.
- Specialised fields have evolved their own conventions for storing data that may be quite different to the conventions of tidy data.

Either of these reasons means you'll need something other than a tibble (or data frame). If your data does fit naturally into a rectangular structure composed of observations and variables, I think tidy data should be your default choice. But there are good reasons to use other structures; tidy data is not the only way.

If you'd like to learn more about non-tidy data, I'd highly recommend this thoughtful blog post by Jeff Leek: <http://simplystatistics.org/2016/02/17/non-tidy-data/>

# 13 Relational data

## 13.1 Introduction

It's rare that a data analysis involves only a single table of data. Typically you have many tables of data, and you must combine them to answer the questions that you're interested in. Collectively, multiple tables of data are called relational data because it is the relations, not just the individual datasets, that are important.

Relations are always defined between a pair of tables. All other relations are built up from this simple idea: the relations of three or more tables are always a property of the relations between each pair. Sometimes both elements of a pair can be the same table! This is needed if, for example, you have a table of people, and each person has a reference to their parents.

To work with relational data you need verbs that work with pairs of tables. There are three families of verbs designed to work with relational data:

- Mutating joins, which add new variables to one data frame from matching observations in another.
- Filtering joins, which filter observations from one data frame based on whether or not they match an observation in the other table.
- Set operations, which treat observations as if they were set elements.

The most common place to find relational data is in a relational database management system (or RDBMS), a term that encompasses almost all modern databases. If you've used a database before, you've almost certainly used SQL. If so, you should find the concepts in this chapter familiar, although their expression in dplyr is a little different. Generally, dplyr is a little easier to use than SQL because dplyr is specialised to do data analysis: it makes common data analysis operations easier, at the expense of making it more difficult to do other things that aren't commonly needed for data analysis.

### 13.1.1 Prerequisites

We will explore relational data from `nycflights13` using the two-table verbs from dplyr.

```
library(tidyverse)
library(nycflights13)
```

## 13.2 nycflights13

We will use the `nycflights13` package to learn about relational data. `nycflights13` contains four tibbles that are related to the `flights` table that you used in [data transformation](#):

- `airlines` lets you look up the full carrier name from its abbreviated code:
- `airlines`
- #> # A tibble: 16 × 2
- #> carrier name
 <chr>
- #> <chr>

- #> 1 9E Endeavor Air Inc.
- #> 2 AA American Airlines Inc.
- #> 3 AS Alaska Airlines Inc.
- #> 4 B6 JetBlue Airways
- #> 5 DL Delta Air Lines Inc.
- #> 6 EV ExpressJet Airlines Inc.
- #> # ... with 10 more rows
- airports gives information about each airport, identified by the faa airport code:
- airports
- #> # A tibble: 1,458 × 8

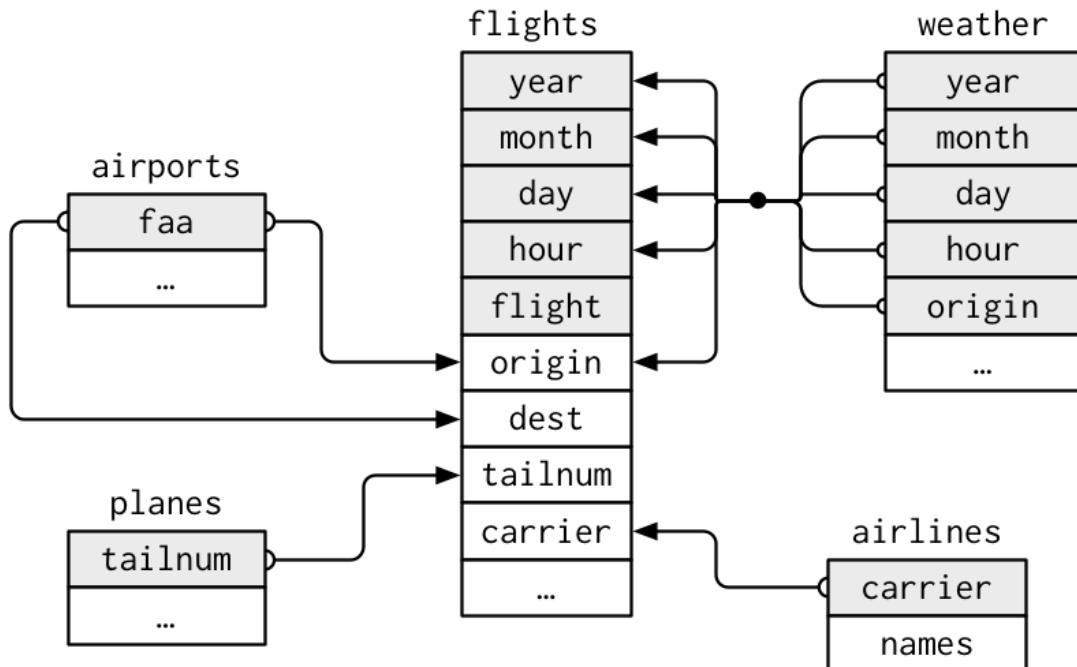
|   |       |                           | name                               | lat   | lon   | alt   | tz    |
|---|-------|---------------------------|------------------------------------|-------|-------|-------|-------|
| • | #>    | faa                       |                                    |       |       |       |       |
|   | dst   |                           |                                    |       |       |       |       |
| • | #>    | <chr>                     |                                    | <chr> | <dbl> | <dbl> | <int> |
|   | <chr> |                           |                                    |       |       |       | <dbl> |
| • | #> 1  | 04G                       | Lansdowne Airport                  | 41.1  | -80.6 | 1044  | -5    |
|   | A     |                           |                                    |       |       |       |       |
| • | #> 2  | 06A                       | Moton Field Municipal Airport      | 32.5  | -85.7 | 264   | -6    |
|   | A     |                           |                                    |       |       |       |       |
| • | #> 3  | 06C                       | Schaumburg Regional                | 42.0  | -88.1 | 801   | -6    |
|   | A     |                           |                                    |       |       |       |       |
| • | #> 4  | 06N                       | Randall Airport                    | 41.4  | -74.4 | 523   | -5    |
|   | A     |                           |                                    |       |       |       |       |
| • | #> 5  | 09J                       | Jekyll Island Airport              | 31.1  | -81.4 | 11    | -5    |
|   | A     |                           |                                    |       |       |       |       |
| • | #> 6  | 0A9                       | Elizabethton Municipal Airport     | 36.4  | -82.2 | 1593  | -5    |
|   | A     |                           |                                    |       |       |       |       |
| • | #> #  | ... with 1,452 more rows, | and 1 more variables: tzzone <chr> |       |       |       |       |

- planes gives information about each plane, identified by its tailnum:
- planes
- #> # A tibble: 3,322 × 9

|   |       |                           | tailnum                            | year                    | type             | manufacturer |         |      |
|---|-------|---------------------------|------------------------------------|-------------------------|------------------|--------------|---------|------|
| • | #>    | tailnum                   | year                               |                         | type             | manufacturer |         |      |
|   | model | engines                   |                                    |                         |                  |              |         |      |
| • | #>    | <chr>                     | <int>                              |                         | <chr>            |              | <chr>   |      |
|   | <chr> | <int>                     |                                    |                         |                  |              |         |      |
| • | #> 1  | N10156                    | 2004                               | Fixed wing multi engine |                  |              | EMBRAER | EMB- |
|   | 145XR | 2                         |                                    |                         |                  |              |         |      |
| • | #> 2  | N102UW                    | 1998                               | Fixed wing multi engine | AIRBUS INDUSTRIE |              | A320-   |      |
|   | 214   | 2                         |                                    |                         |                  |              |         |      |
| • | #> 3  | N103US                    | 1999                               | Fixed wing multi engine | AIRBUS INDUSTRIE |              | A320-   |      |
|   | 214   | 2                         |                                    |                         |                  |              |         |      |
| • | #> 4  | N104UW                    | 1999                               | Fixed wing multi engine | AIRBUS INDUSTRIE |              | A320-   |      |
|   | 214   | 2                         |                                    |                         |                  |              |         |      |
| • | #> 5  | N10575                    | 2002                               | Fixed wing multi engine |                  |              | EMBRAER | EMB- |
|   | 145LR | 2                         |                                    |                         |                  |              |         |      |
| • | #> 6  | N105UW                    | 1999                               | Fixed wing multi engine | AIRBUS INDUSTRIE |              | A320-   |      |
|   | 214   | 2                         |                                    |                         |                  |              |         |      |
| • | #> #  | ... with 3,316 more rows, | and 3 more variables: seats <int>, |                         |                  |              |         |      |
|   | #> #  | speed <int>,              | engine <chr>                       |                         |                  |              |         |      |

- `weather` gives the weather at each NYC airport for each hour:
- `weather`
- #> # A tibble: 26,130 × 15
- #> origin year month day hour temp dewp humid wind\_dir  
wind\_speed
- #> <chr> <dbl> <dbl> <int> <int> <dbl> <dbl> <dbl> <dbl> <dbl>  
<dbl>
- #> 1 EWR 2013 1 1 0 37.0 21.9 54.0 230  
10.4
- #> 2 EWR 2013 1 1 1 37.0 21.9 54.0 230  
13.8
- #> 3 EWR 2013 1 1 2 37.9 21.9 52.1 230  
12.7
- #> 4 EWR 2013 1 1 3 37.9 23.0 54.5 230  
13.8
- #> 5 EWR 2013 1 1 4 37.9 24.1 57.0 240  
15.0
- #> 6 EWR 2013 1 1 6 39.0 26.1 59.4 270  
10.4
- #> # ... with 2.612e+04 more rows, and 5 more variables: wind\_gust  
<dbl>,  
#> # precip <dbl>, pressure <dbl>, visib <dbl>, time\_hour <dttm>

One way to show the relationships between the different tables is with a drawing:



This diagram is a little overwhelming, but it's simple compared to some you'll see in the wild! The key to understanding diagrams like this is to remember each relation always concerns a pair of tables. You don't need to understand the whole thing; you just need to understand the chain of relations between the tables that you are interested in.

For nycflights13:

- `flights` connects to `planes` via a single variable, `tailnum`.
- `flights` connects to `airlines` through the `carrier` variable.
- `flights` connects to `airports` in two ways: via the `origin` and `dest` variables.
- `flights` connects to `weather` via `origin` (the location), and `year`, `month`, `day` and `hour` (the time).

### 13.2.1 Exercises

1. Imagine you wanted to draw (approximately) the route each plane flies from its origin to its destination. What variables would you need? What tables would you need to combine?
2. I forgot to draw the relationship between `weather` and `airports`. What is the relationship and how should it appear in the diagram?
3. `weather` only contains information for the origin (NYC) airports. If it contained weather records for all airports in the USA, what additional relation would it define with `flights`?
4. **We know that some days of the year are “special”, and fewer people than usual fly on them.** How might you represent that data as a data frame? What would be the primary keys of that table? How would it connect to the existing tables?

## 13.3 Keys

The variables used to connect each pair of tables are called **keys**. A key is a variable (or set of variables) that uniquely identifies an observation. In simple cases, a single variable is sufficient to identify an observation. For example, each plane is uniquely identified by its `tailnum`. In other cases, multiple variables may be needed. For example, to identify an observation in `weather` you need five variables: `year`, `month`, `day`, `hour`, and `origin`.

There are two types of keys:

- A **primary key** uniquely identifies an observation in its own table. For example, `planes$tailnum` is a primary key because it uniquely identifies each plane in the `planes` table.
- A **foreign key** uniquely identifies an observation in another table. For example, the `flights$tailnum` is a foreign key because it appears in the `flights` table where it matches each flight to a unique plane.

A variable can be both a primary key **and** a foreign key. For example, `origin` is part of the `weather` primary key, and is also a foreign key for the `airport` table.

**Once you’ve identified the primary keys in your tables, it’s good practice to verify that they do indeed uniquely identify each observation.** One way to do that is to `count()` the primary keys and look for entries where `n` is greater than one:

```
planes %>%
 count(tailnum) %>%
 filter(n > 1)
#> # A tibble: 0 × 2
#> # ... with 2 variables: tailnum <chr>, n <int>
```

```

weather %>%
 count(year, month, day, hour, origin) %>%
 filter(n > 1)
#> Source: local data frame [0 x 6]
#> Groups: year, month, day, hour [0]
#>
#> # ... with 6 variables: year <dbl>, month <dbl>, day <int>, hour
#> # origin <chr>, n <int>

```

Sometimes a table doesn't have an explicit primary key: each row is an observation, but no combination of variables reliably identifies it. For example, what's the primary key in the `flights` table? You might think it would be the date plus the flight or tail number, but neither of those are unique:

```

flights %>%
 count(year, month, day, flight) %>%
 filter(n > 1)
#> Source: local data frame [29,768 x 5]
#> Groups: year, month, day [365]
#>
#> year month day flight n
#> <int> <int> <int> <int> <int>
#> 1 2013 1 1 1 2
#> 2 2013 1 1 3 2
#> 3 2013 1 1 4 2
#> 4 2013 1 1 11 3
#> 5 2013 1 1 15 2
#> 6 2013 1 1 21 2
#> # ... with 2.976e+04 more rows

flights %>%
 count(year, month, day, tailnum) %>%
 filter(n > 1)
#> Source: local data frame [64,928 x 5]
#> Groups: year, month, day [365]
#>
#> year month day tailnum n
#> <int> <int> <int> <chr> <int>
#> 1 2013 1 1 N0EGMQ 2
#> 2 2013 1 1 N11189 2
#> 3 2013 1 1 N11536 2
#> 4 2013 1 1 N11544 3
#> 5 2013 1 1 N11551 2
#> 6 2013 1 1 N12540 2
#> # ... with 6.492e+04 more rows

```

When starting to work with this data, I had naively assumed that each flight number would be only used once per day: that would make it much easier to communicate problems with a specific flight. Unfortunately that is not the case! If a table lacks a primary key, it's sometimes useful to add one with `mutate()` and `row_number()`. That makes it easier to match observations if

you've done some filtering and want to check back in with the original data. This is called a surrogate key.

A primary key and the corresponding foreign key in another table form a relation. Relations are typically one-to-many. For example, each flight has one plane, but each plane has many flights. In other data, you'll occasionally see a 1-to-1 relationship. You can think of this as a special case of 1-to-many. You can model many-to-many relations with a many-to-1 relation plus a 1-to-many relation. For example, in this data there's a many-to-many relationship between airlines and airports: each airline flies to many airports; each airport hosts many airlines.

### 13.3.1 Exercises

1. Add a surrogate key to `flights`.
2. Identify the keys in the following datasets
  1. `Lahman::Batting`,
  2. `babynames::babynames`
  3. `nasaweather::atmos`
  4. `fueleconomy::vehicles`
  5. `ggplot2::diamonds`

(You might need to install some packages and read some documentation.)

3. Draw a diagram illustrating the connections between the `Batting`, `Master`, and `Salaries` tables in the `Lahman` package. Draw another diagram that shows the relationship between `Master`, `Managers`, `AwardsManagers`.

How would you characterise the relationship between the `Batting`, `Pitching`, and `Fielding` tables?

## 13.4 Mutating joins

The first tool we'll look at for combining a pair of tables is the mutating join. A mutating join allows you to combine variables from two tables. It first matches observations by their keys, then copies across variables from one table to the other.

Like `mutate()`, the join functions add variables to the right, so if you have a lot of variables already, the new variables won't get printed out. For these examples, we'll make it easier to see what's going on in the examples by creating a narrower dataset:

```
flights2 <- flights %>%
 select(year:day, hour, origin, dest, tailnum, carrier)
flights2
#> # A tibble: 336,776 × 8
#> year month day hour origin dest tailnum carrier
#> <int> <int> <int> <dbl> <chr> <chr> <chr> <chr>
#> 1 2013 1 1 5 EWR IAH N14228 UA
#> 2 2013 1 1 5 LGA IAH N24211 UA
#> 3 2013 1 1 5 JFK MIA N619AA AA
#> 4 2013 1 1 5 JFK BQN N804JB B6
#> 5 2013 1 1 6 LGA ATL N668DN DL
```

```
#> 6 2013 1 1 5 EWR ORD N39463 UA
#> # ... with 3.368e+05 more rows
```

(Remember, when you're in RStudio, you can also use `View()` to avoid this problem.)

Imagine you want to add the full airline name to the `flights2` data. You can combine the `airlines` and `flights2` data frames with `left_join()`:

```
flights2 %>%
 select(-origin, -dest) %>%
 left_join(airlines, by = "carrier")
#> # A tibble: 336,776 × 7
#> year month day hour tailnum carrier name
#> <int> <int> <int> <dbl> <chr> <chr>
#> 1 2013 1 1 5 N14228 UA United Air Lines Inc.
#> 2 2013 1 1 5 N24211 UA United Air Lines Inc.
#> 3 2013 1 1 5 N619AA AA American Airlines Inc.
#> 4 2013 1 1 5 N804JB B6 JetBlue Airways
#> 5 2013 1 1 6 N668DN DL Delta Air Lines Inc.
#> 6 2013 1 1 5 N39463 UA United Air Lines Inc.
#> # ... with 3.368e+05 more rows
```

The result of joining `airlines` to `flights2` is an additional variable: `name`. This is why I call this type of join a mutating join. In this case, you could have got to the same place using `mutate()` and R's base subsetting:

```
flights2 %>%
 select(-origin, -dest) %>%
 mutate(name = airlines$name[match(carrier, airlines$carrier)])
#> # A tibble: 336,776 × 7
#> year month day hour tailnum carrier name
#> <int> <int> <int> <dbl> <chr> <chr>
#> 1 2013 1 1 5 N14228 UA United Air Lines Inc.
#> 2 2013 1 1 5 N24211 UA United Air Lines Inc.
#> 3 2013 1 1 5 N619AA AA American Airlines Inc.
#> 4 2013 1 1 5 N804JB B6 JetBlue Airways
#> 5 2013 1 1 6 N668DN DL Delta Air Lines Inc.
#> 6 2013 1 1 5 N39463 UA United Air Lines Inc.
#> # ... with 3.368e+05 more rows
```

But this is hard to generalise when you need to match multiple variables, and takes close reading to figure out the overall intent.

The following sections explain, in detail, how mutating joins work. You'll start by learning a useful visual representation of joins. We'll then use that to explain the four mutating join functions: the inner join, and the three outer joins. When working with real data, keys don't always uniquely identify observations, so next we'll talk about what happens when there isn't a unique match. Finally, you'll learn how to tell dplyr which variables are the keys for a given join.

### 13.4.1 Understanding joins

To help you learn how joins work, I'm going to use a visual representation:

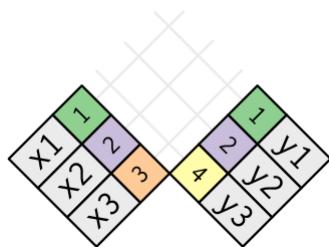
| x | y  |
|---|----|
| 1 | x1 |
| 2 | x2 |
| 3 | x3 |

| y |    |
|---|----|
| 1 | y1 |
| 2 | y2 |
| 4 | y3 |

```
x <- tribble(
 ~key, ~val_x,
 1, "x1",
 2, "x2",
 3, "x3"
)
y <- tribble(
 ~key, ~val_y,
 1, "y1",
 2, "y2",
 4, "y3"
)
```

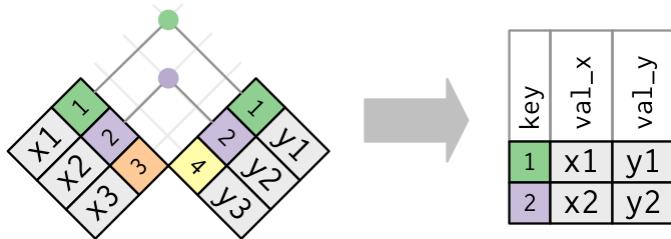
The coloured column represents the “key” variable: these are used to match the rows between the tables. The grey column represents the “value” column that is carried along for the ride. In these examples I’ll show a single key variable, but the idea generalises in a straightforward way to multiple keys and multiple values.

A join is a way of connecting each row in `x` to zero, one, or more rows in `y`. The following diagram shows each potential match as an intersection of a pair of lines.



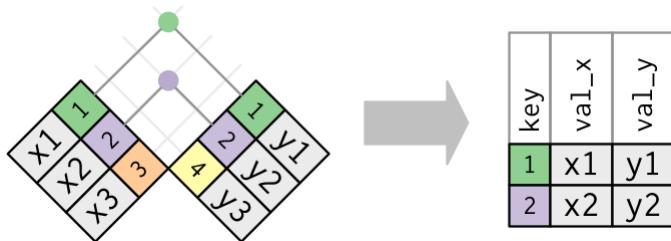
(If you look closely, you might notice that we've switched the order of the key and value columns in `x`. This is to emphasise that joins match based on the key; the value is just carried along for the ride.)

In an actual join, matches will be indicated with dots. The number of dots = the number of matches = the number of rows in the output.



### 13.4.2 Inner join

The simplest type of join is the inner join. An inner join matches pairs of observations whenever their keys are equal:



(To be precise, this is an inner equijoin because the keys are matched using the equality operator. Since most joins are equijoins we usually drop that specification.)

The output of an inner join is a new data frame that contains the key, the x values, and the y values. We use `by` to tell dplyr which variable is the key:

```
x %>%
 inner_join(y, by = "key")
#> # A tibble: 2 × 3
#> key val_x val_y
#> <dbl> <chr> <chr>
#> 1 1 x1 y1
#> 2 2 x2 y2
```

The most important property of an inner join is that unmatched rows are not included in the result. **This means that generally inner joins are usually not appropriate for use in analysis because it's too easy to lose observations.**

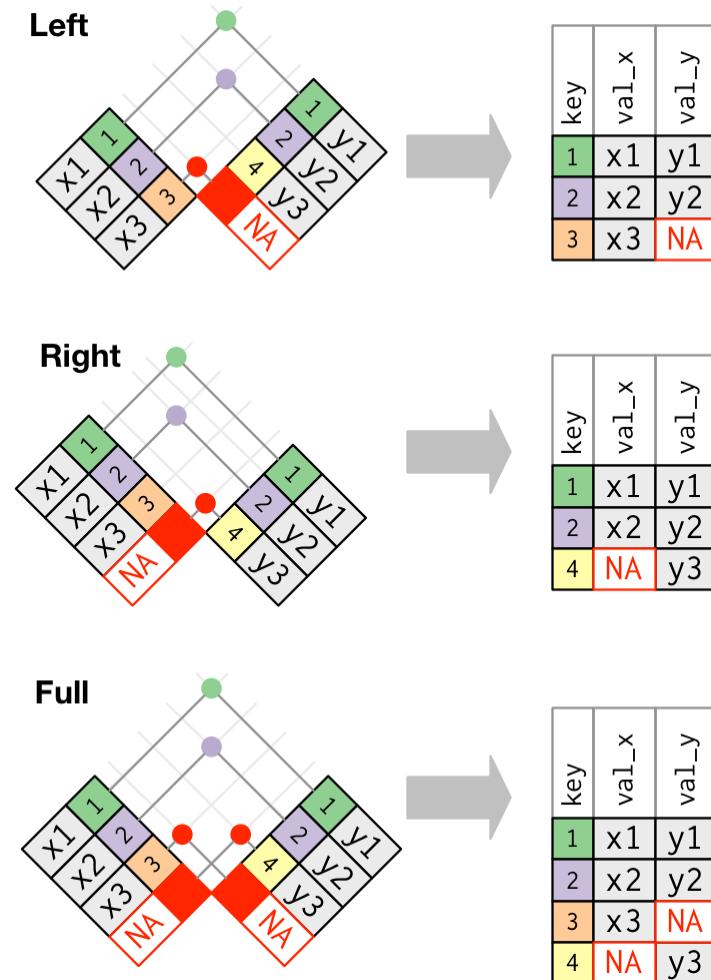
### 13.4.3 Outer joins

An inner join keeps observations that appear in both tables. An outer join keeps observations that appear in at least one of the tables. There are three types of outer joins:

- A left join keeps all observations in `x`.
- A right join keeps all observations in `y`.
- A full join keeps all observations in `x` and `y`.

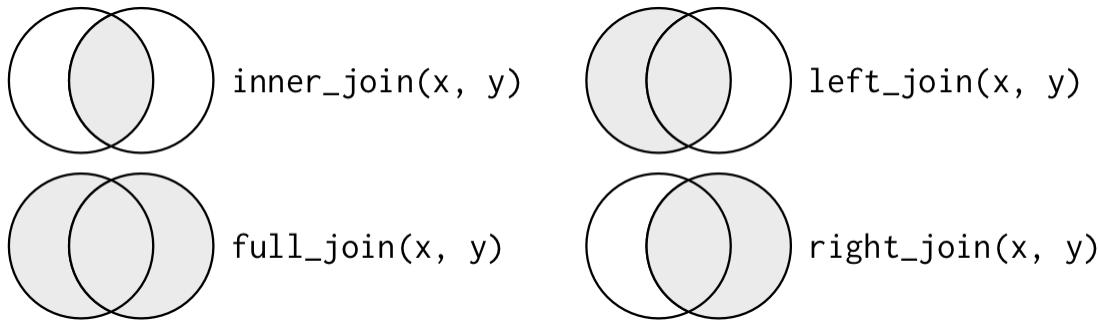
These joins work by adding an additional “virtual” observation to each table. This observation has a key that always matches (if no other key matches), and a value filled with NA.

Graphically, that looks like:



The most commonly used join is the left join: you use this whenever you look up additional data from another table, because it preserves the original observations even when there isn't a match. The left join should be your default join: use it unless you have a strong reason to prefer one of the others.

Another way to depict the different types of joins is with a Venn diagram:

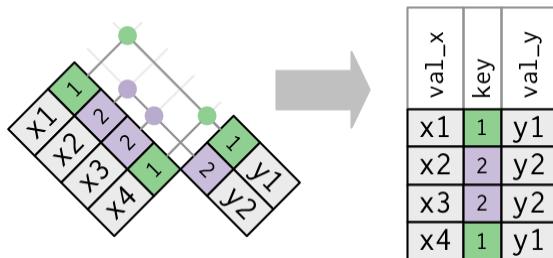


However, this is not a great representation. It might jog your memory about which join preserves the observations in which table, but it suffers from a major limitation: a Venn diagram can't show what happens when keys don't uniquely identify an observation.

### 13.4.4 Duplicate keys

So far all the diagrams have assumed that the keys are unique. But that's not always the case. This section explains what happens when the keys are not unique. There are two possibilities:

1. One table has duplicate keys. This is useful when you want to add in additional information as there is typically a one-to-many relationship.

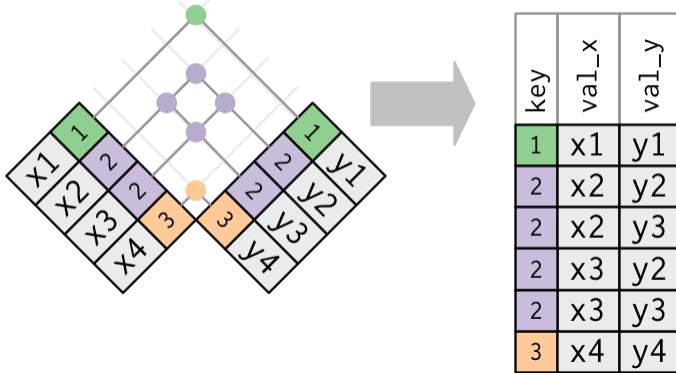


Note that I've put the key column in a slightly different position in the output. This reflects that the key is a primary key in y and a foreign key in x.

```
x <- tribble(
 ~key, ~val_x,
 1, "x1",
 2, "x2",
 2, "x3",
 1, "x4"
)
y <- tribble(
 ~key, ~val_y,
 1, "y1",
 2, "y2"
)
left_join(x, y, by = "key")
#> # A tibble: 4 × 3
#> key val_x val_y
#> <dbl> <chr> <chr>
#> 1 1 x1 y1
#> 2 2 x2 y2
#> 3 2 x3 y2
#> 4 1 x4 y1
```

```
#> 2 2 x2 y2
#> 3 2 x3 y2
#> 4 1 x4 y1
```

2. Both tables have duplicate keys. This is usually an error because in neither table do the keys uniquely identify an observation. When you join duplicated keys, you get all possible combinations, the Cartesian product:



```
x <- tribble(
 ~key, ~val_x,
 1, "x1",
 2, "x2",
 2, "x3",
 3, "x4"
)
y <- tribble(
 ~key, ~val_y,
 1, "y1",
 2, "y2",
 2, "y3",
 3, "y4"
)
left_join(x, y, by = "key")
#> # A tibble: 6 × 3
#> key val_x val_y
#> <dbl> <chr> <chr>
#> 1 1 x1 y1
#> 2 2 x2 y2
#> 3 2 x2 y3
#> 4 2 x3 y2
#> 5 2 x3 y3
#> 6 3 x4 y4
```

### 13.4.5 Defining the key columns

So far, the pairs of tables have always been joined by a single variable, and that variable has the same name in both tables. That constraint was encoded by `by = "key"`. You can use other values for `by` to connect the tables in other ways:

- The default, `by = NULL`, uses all variables that appear in both tables, the so called natural join. For example, the flights and weather tables match on their common variables: `year`, `month`, `day`, `hour` and `origin`.
- ```
flights2 %>%
  left_join(weather)
#> Joining, by = c("year", "month", "day", "hour", "origin")
#> # A tibble: 336,776 × 18
#>   year month   day hour origin dest tailnum carrier temp
#>   <dbl> <dbl> <int> <dbl> <chr> <chr> <chr> <chr> <dbl>
#>   <dbl> <dbl>
#> 1 2013     1     1     5   EWR   IAH   N14228   UA    NA
#> NA   NA
#> 2 2013     1     1     5   LGA   IAH   N24211   UA    NA
#> NA   NA
#> 3 2013     1     1     5   JFK   MIA   N619AA   AA    NA
#> NA   NA
#> 4 2013     1     1     5   JFK   BQN   N804JB   B6    NA
#> NA   NA
#> 5 2013     1     1     6   LGA   ATL   N668DN   DL    39.9
#> 26.1 57.3
#> 6 2013     1     1     5   EWR   ORD   N39463   UA    NA
#> NA   NA
#> # ... with 3.368e+05 more rows, and 7 more variables: wind_dir
#> <dbl>,
#> #   wind_speed <dbl>, wind_gust <dbl>, precip <dbl>, pressure
#> <dbl>,
#> #   visib <dbl>, time_hour <dttm>
```
- A character vector, `by = "x"`. This is like a natural join, but uses only some of the common variables. For example, `flights` and `planes` have `year` variables, but they mean different things so we only want to join by `tailnum`.
- ```
flights2 %>%
 left_join(planes, by = "tailnum")
#> # A tibble: 336,776 × 16
#> year.x month day hour origin dest tailnum carrier year.y
#> <int> <int> <int> <dbl> <chr> <chr> <chr> <chr> <int>
#> 1 2013 1 1 5 EWR IAH N14228 UA 1999
#> 2 2013 1 1 5 LGA IAH N24211 UA 1998
#> 3 2013 1 1 5 JFK MIA N619AA AA 1990
#> 4 2013 1 1 5 JFK BQN N804JB B6 2012
#> 5 2013 1 1 6 LGA ATL N668DN DL 1991
#> 6 2013 1 1 5 EWR ORD N39463 UA 2012
#> # ... with 3.368e+05 more rows, and 7 more variables: type
#> <chr>,
#> # manufacturer <chr>, model <chr>, engines <int>, seats
#> <int>,
#> # speed <int>, engine <chr>
```

Note that the `year` variables (which appear in both input data frames, but are not constrained to be equal) are disambiguated in the output with a suffix.

- A named character vector: `by = c("a" = "b")`. This will match variable `a` in table `x` to variable `b` in table `y`. The variables from `x` will be used in the output.

For example, if we want to draw a map we need to combine the flights data with the airports data which contains the location (`lat` and `long`) of each airport. Each flight has an origin and destination `airport`, so we need to specify which one we want to join to:

```
flights2 %>%
 left_join(airports, c("dest" = "faa"))
#> # A tibble: 336,776 × 15
#> year month day hour origin dest tailnum carrier
#> <int> <int> <int> <dbl> <chr> <chr> <chr> <chr>
#> 1 2013 1 1 5 EWR IAH N14228 UA
#> 2 2013 1 1 5 LGA IAH N24211 UA
#> 3 2013 1 1 5 JFK MIA N619AA AA
#> 4 2013 1 1 5 JFK BQN N804JB B6
#> 5 2013 1 1 6 LGA ATL N668DN DL
#> 6 2013 1 1 5 EWR ORD N39463 UA
#> # ... with 3.368e+05 more rows, and 7 more variables: name
<chr>,
#> # lat <dbl>, lon <dbl>, alt <int>, tz <dbl>, dst <chr>, tzone
<chr>

flights2 %>%
 left_join(airports, c("origin" = "faa"))
#> # A tibble: 336,776 × 15
#> year month day hour origin dest tailnum carrier
name
#> <int> <int> <int> <dbl> <chr> <chr> <chr> <chr>
<chr>
#> 1 2013 1 1 5 EWR IAH N14228 UA Newark
Liberty Intl
#> 2 2013 1 1 5 LGA IAH N24211 UA
La Guardia
#> 3 2013 1 1 5 JFK MIA N619AA AA John F
Kennedy Intl
#> 4 2013 1 1 5 JFK BQN N804JB B6 John F
Kennedy Intl
#> 5 2013 1 1 6 LGA ATL N668DN DL
La Guardia
#> 6 2013 1 1 5 EWR ORD N39463 UA Newark
Liberty Intl
#> # ... with 3.368e+05 more rows, and 6 more variables: lat
<dbl>,
#> # lon <dbl>, alt <int>, tz <dbl>, dst <chr>, tzone <chr>
```

### 13.4.6 Exercises

1. Compute the average delay by destination, then join on the `airports` data frame so you can show the spatial distribution of delays. Here's an easy way to draw a map of the United States:
2. 

```
airports %>%
 semi_join(flights, c("faa" = "dest")) %>%
 ggplot(aes(lon, lat)) +
 borders("state") +
 geom_point() +
 coord_quickmap()
```

(Don't worry if you don't understand what `semi_join()` does — you'll learn about it next.)

You might want to use the `size` or `colour` of the points to display the average delay for each airport.

7. Add the location of the origin and destination (i.e. the `lat` and `lon`) to `flights`.
8. Is there a relationship between the age of a plane and its delays?
9. What weather conditions make it more likely to see a delay?
10. What happened on June 13 2013? Display the spatial pattern of delays, and then use Google to cross-reference with the weather.

### 13.4.7 Other implementations

`base::merge()` can perform all four types of mutating join:

```
dplyr merge
inner_join(x, y) merge(x, y)
left_join(x, y) merge(x, y, all.x = TRUE)
right_join(x, y) merge(x, y, all.y = TRUE),
full_join(x, y) merge(x, y, all.x = TRUE, all.y = TRUE)
```

The advantages of the specific dplyr verbs is that they more clearly convey the intent of your code: the difference between the joins is really important but concealed in the arguments of `merge()`. dplyr's joins are considerably faster and don't mess with the order of the rows.

SQL is the inspiration for dplyr's conventions, so the translation is straightforward:

| dplyr                                   | SQL                                                       |
|-----------------------------------------|-----------------------------------------------------------|
| <code>inner_join(x, y, by = "z")</code> | <code>SELECT * FROM x INNER JOIN y USING (z)</code>       |
| <code>left_join(x, y, by = "z")</code>  | <code>SELECT * FROM x LEFT OUTER JOIN y USING (z)</code>  |
| <code>right_join(x, y, by = "z")</code> | <code>SELECT * FROM x RIGHT OUTER JOIN y USING (z)</code> |
| <code>full_join(x, y, by = "z")</code>  | <code>SELECT * FROM x FULL OUTER JOIN y USING (z)</code>  |

Note that "INNER" and "OUTER" are optional, and often omitted.

Joining different variables between the tables, e.g. `inner_join(x, y, by = c("a" = "b"))` uses a slightly different syntax in SQL: `SELECT * FROM x INNER JOIN y ON x.a =`

y.b. As this syntax suggests, SQL supports a wider range of join types than dplyr because you can connect the tables using constraints other than equality (sometimes called non-equijoins).

## 13.5 Filtering joins

Filtering joins match observations in the same way as mutating joins, but affect the observations, not the variables. There are two types:

- `semi_join(x, y)` keeps all observations in `x` that have a match in `y`.
- `anti_join(x, y)` drops all observations in `x` that have a match in `y`.

Semi-joins are useful for matching filtered summary tables back to the original rows. For example, imagine you've found the top ten most popular destinations:

```
top_dest <- flights %>%
 count(dest, sort = TRUE) %>%
 head(10)
top_dest
#> # A tibble: 10 × 2
#> dest n
#> <chr> <int>
#> 1 ORD 17283
#> 2 ATL 17215
#> 3 LAX 16174
#> 4 BOS 15508
#> 5 MCO 14082
#> 6 CLT 14064
#> # ... with 4 more rows
```

Now you want to find each flight that went to one of those destinations. You could construct a filter yourself:

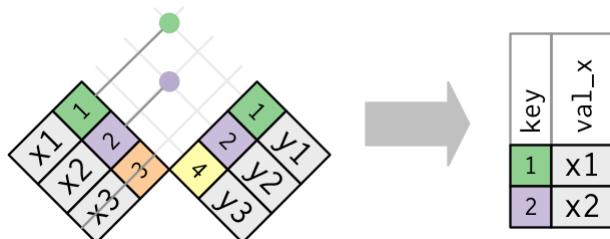
```
flights %>%
 filter(dest %in% top_dest$dest)
#> # A tibble: 141,145 × 19
#> year month day dep_time sched_dep_time dep_delay arr_time
#> <int> <int> <int> <int> <int> <dbl> <int>
#> 1 2013 1 1 542 540 2 923
#> 2 2013 1 1 554 600 -6 812
#> 3 2013 1 1 554 558 -4 740
#> 4 2013 1 1 555 600 -5 913
#> 5 2013 1 1 557 600 -3 838
#> 6 2013 1 1 558 600 -2 753
#> # ... with 1.411e+05 more rows, and 12 more variables:
#> # sched_arr_time <int>, arr_delay <dbl>, carrier <chr>, flight
#> # <int>,
#> # tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
#> # distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
```

But it's difficult to extend that approach to multiple variables. For example, imagine that you'd found the 10 days with highest average delays. How would you construct the filter statement that used `year`, `month`, and `day` to match it back to `flights`?

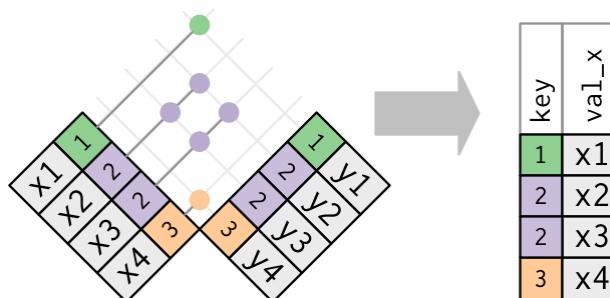
Instead you can use a semi-join, which connects the two tables like a mutating join, but instead of adding new columns, only keeps the rows in `x` that have a match in `y`:

```
flights %>%
 semi_join(top_dest)
#> Joining, by = "dest"
#> # A tibble: 141,145 × 19
#> year month day dep_time sched_dep_time dep_delay arr_time
#> <int> <int> <int> <int> <int> <dbl> <int>
#> 1 2013 1 1 554 558 -4 740
#> 2 2013 1 1 558 600 -2 753
#> 3 2013 1 1 608 600 8 807
#> 4 2013 1 1 629 630 -1 824
#> 5 2013 1 1 656 700 -4 854
#> 6 2013 1 1 709 700 9 852
#> # ... with 1.411e+05 more rows, and 12 more variables:
#> # sched_arr_time <int>, arr_delay <dbl>, carrier <chr>, flight
#> # tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
#> # distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
```

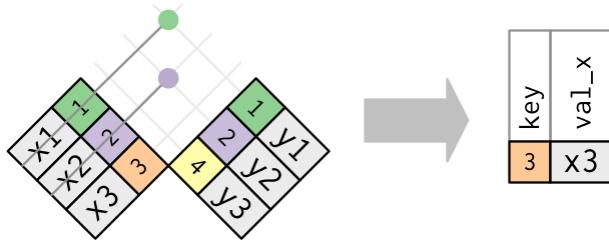
Graphically, a semi-join looks like this:



Only the existence of a match is important; it doesn't matter which observation is matched. This means that filtering joins never duplicate rows like mutating joins do:



The inverse of a semi-join is an anti-join. An anti-join keeps the rows that ***don't*** have a match:



Anti-joins are useful for diagnosing join mismatches. For example, when connecting `flights` and `planes`, you might be interested to know that there are many `flights` **that don't have a match in `planes`**:

```
flights %>%
 anti_join(planes, by = "tailnum") %>%
 count(tailnum, sort = TRUE)
#> # A tibble: 722 × 2
#> tailnum n
#> <chr> <int>
#> 1 <NA> 2512
#> 2 N725MQ 575
#> 3 N722MQ 513
#> 4 N723MQ 507
#> 5 N713MQ 483
#> 6 N735MQ 396
#> # ... with 716 more rows
```

### 13.5.1 Exercises

1. What does it mean for a flight to have a missing `tailnum`? What do the tail numbers that **don't have a matching record in `planes`** have in common? (Hint: one variable explains ~90% of the problems.)
2. Filter flights to only show flights with planes that have flown at least 100 flights.
3. Combine `fueleconomy::vehicles` and `fueleconomy::common` to find only the records for the most common models.
4. Find the 48 hours (over the course of the whole year) that have the worst delays. Cross-reference it with the `weather` data. Can you see any patterns?
5. What does `anti_join(flights, airports, by = c("dest" = "faa"))` tell you? What does `anti_join(airports, flights, by = c("faa" = "dest"))` tell you?
6. **You might expect that there's an implicit relationship between plane and airline**, because each plane is flown by a single airline. Confirm or reject this hypothesis using the tools you've learned above.

## 13.6 Join problems

The data you've been working with in this chapter has been cleaned up so that you'll have as few problems as possible. Your own data is unlikely to be so nice, so there are a few things that you should do with your own data to make your joins go smoothly.

- Start by identifying the variables that form the primary key in each table. You should usually do this based on your understanding of the data, not empirically by looking for a combination of variables that give a unique identifier. If you just look for variables without **thinking about what they mean, you might get (un)lucky and find a combination that's unique** in your current data but the relationship might not be true in general.

For example, the altitude and longitude uniquely identify each airport, but they are not good identifiers!

```
airports %>% count(alt, lon) %>% filter(n > 1)
#> Source: local data frame [0 x 3]
#> Groups: alt [0]
#>
#> # ... with 3 variables: alt <int>, lon <dbl>, n <int>
```

- Check that none of the variables in the primary key are missing. If a value is missing then **it can't identify an observation!**
- Check that your foreign keys match primary keys in another table. The best way to do this is with an `anti_join()`. **It's common for keys not to match because of data entry errors.** Fixing these is often a lot of work.

**If you do have missing keys, you'll need to be thoughtful about your use of inner vs. outer joins, carefully considering whether or not you want to drop rows that don't have a match.**

Be aware that simply checking the number of rows before and after the join is not sufficient to ensure that your join has gone smoothly. If you have an inner join with duplicate keys in both tables, you might get unlucky as the number of dropped rows might exactly equal the number of duplicated rows!

## 13.7 Set operations

The final type of two-table verb are the set operations. Generally, I use these the least frequently, but they are occasionally useful when you want to break a single complex filter into simpler pieces. All these operations work with a complete row, comparing the values of every variable. These expect the `x` and `y` inputs to have the same variables, and treat the observations like sets:

- `intersect(x, y)`: return only observations in both `x` and `y`.
- `union(x, y)`: return unique observations in `x` and `y`.
- `setdiff(x, y)`: return observations in `x`, but not in `y`.

Given this simple data:

```
df1 <- tribble(
 ~x, ~y,
 1, 1,
 2, 1
)
df2 <- tribble(
 ~x, ~y,
 1, 1,
 1, 2
```

)

The four possibilities are:

```
intersect(df1, df2)
#> # A tibble: 1 × 2
#> x y
#> <dbl> <dbl>
#> 1 1 1

Note that we get 3 rows, not 4
union(df1, df2)
#> # A tibble: 3 × 2
#> x y
#> <dbl> <dbl>
#> 1 1 2
#> 2 2 1
#> 3 1 1

setdiff(df1, df2)
#> # A tibble: 1 × 2
#> x y
#> <dbl> <dbl>
#> 1 2 1

setdiff(df2, df1)
#> # A tibble: 1 × 2
#> x y
#> <dbl> <dbl>
#> 1 1 2
```

# 15 Factors

## 15.1 Introduction

In R, factors are used to work with categorical variables, variables that have a fixed and known set of possible values. They are also useful when you want to display character vectors in a non-alphabetical order.

Historically, factors were much easier to work with than characters. As a result, many of the functions in base R automatically convert characters to factors. This means that factors often **crop up in places where they're not actually helpful. Fortunately, you don't need to worry about** that in the tidyverse, and can focus on situations where factors are genuinely useful.

For more historical context on factors, I recommend [stringsAsFactors: An unauthorized biography](#) by Roger Peng, and [stringsAsFactors = <sigh>](#) by Thomas Lumley.

### 15.1.1 Prerequisites

To work with factors, **we'll use the**forcats package, which provides tools for dealing with **categorical variables (and it's an anagram of factors!). It provides a wide range of helpers for** working with factors.forcats is not part of the core tidyverse, so we need to load it explicitly.

```
library(tidyverse)
library(forcats)
```

## 15.2 Creating factors

Imagine that you have a variable that records month:

```
x1 <- c("Dec", "Apr", "Jan", "Mar")
```

Using a string to record this variable has two problems:

1. There are only twelve possible months, and **there's nothing saving you from typos:**

```
x2 <- c("Dec", "Apr", "Jam", "Mar")
```

2. **It doesn't sort in a useful way:**

```
3. sort(x1)
#> [1] "Apr" "Dec" "Jan" "Mar"
```

You can fix both of these problems with a factor. To create a factor you must start by creating a list of the valid levels:

```
month_levels <- c(
 "Jan", "Feb", "Mar", "Apr", "May", "Jun",
 "Jul", "Aug", "Sep", "Oct", "Nov", "Dec"
)
```

Now you can create a factor:

```
y1 <- factor(x1, levels = month_levels)
y1
#> [1] Dec Apr Jan Mar
#> Levels: Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec
sort(y1)
#> [1] Jan Mar Apr Dec
#> Levels: Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec
```

And any values not in the set will be silently converted to NA:

```
y2 <- factor(x2, levels = month_levels)
y2
#> [1] Dec Apr <NA> Mar
#> Levels: Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec
```

If you want a warning, you can use `readr::parse_factor()`:

```
y2 <- parse_factor(x2, levels = month_levels)
#> Warning: 1 parsing failure.
#> row col expected actual
#> 3 -- value in level set Jam
```

If you omit the `levels`, they'll be taken from the data in alphabetical order:

```
factor(x1)
#> [1] Dec Apr Jan Mar
#> Levels: Apr Dec Jan Mar
```

Sometimes you'd prefer that the order of the levels match the order of the first appearance in the data. You can do that when creating the factor by setting `levels` to `unique(x)`, or after the fact, with `fct_inorder()`:

```
f1 <- factor(x1, levels = unique(x1))
f1
#> [1] Dec Apr Jan Mar
#> Levels: Dec Apr Jan Mar

f2 <- x1 %>% factor() %>% fct_inorder()
f2
#> [1] Dec Apr Jan Mar
#> Levels: Dec Apr Jan Mar
```

If you ever need to access the set of valid levels directly, you can do so with `levels()`:

```
levels(f2)
#> [1] "Dec" "Apr" "Jan" "Mar"
```

## 15.3 General Social Survey

For the rest of this chapter, we're going to focus on `forcats::gss_cat`. It's a sample of data from the [General Social Survey](#), which is a long-running US survey conducted by the independent research organization NORC at the University of Chicago. The survey has thousands of questions, so in `gss_cat` I've selected a handful that will illustrate some common challenges you'll encounter when working with factors.

```
gss_cat
#> # A tibble: 21,483 × 9
#> year marital age race rincome partyid
#> <int> <fctr> <int> <fctr> <fctr> <fctr>
#> 1 2000 Never married 26 White $8000 to 9999 Ind,near rep
#> 2 2000 Divorced 48 White $8000 to 9999 Not str republican
#> 3 2000 Widowed 67 White Not applicable Independent
#> 4 2000 Never married 39 White Not applicable Ind,near rep
#> 5 2000 Divorced 25 White Not applicable Not str democrat
#> 6 2000 Married 25 White $20000 - 24999 Strong democrat
#> # ... with 2.148e+04 more rows, and 3 more variables: relig <fctr>,
#> # denom <fctr>, tvhours <int>
```

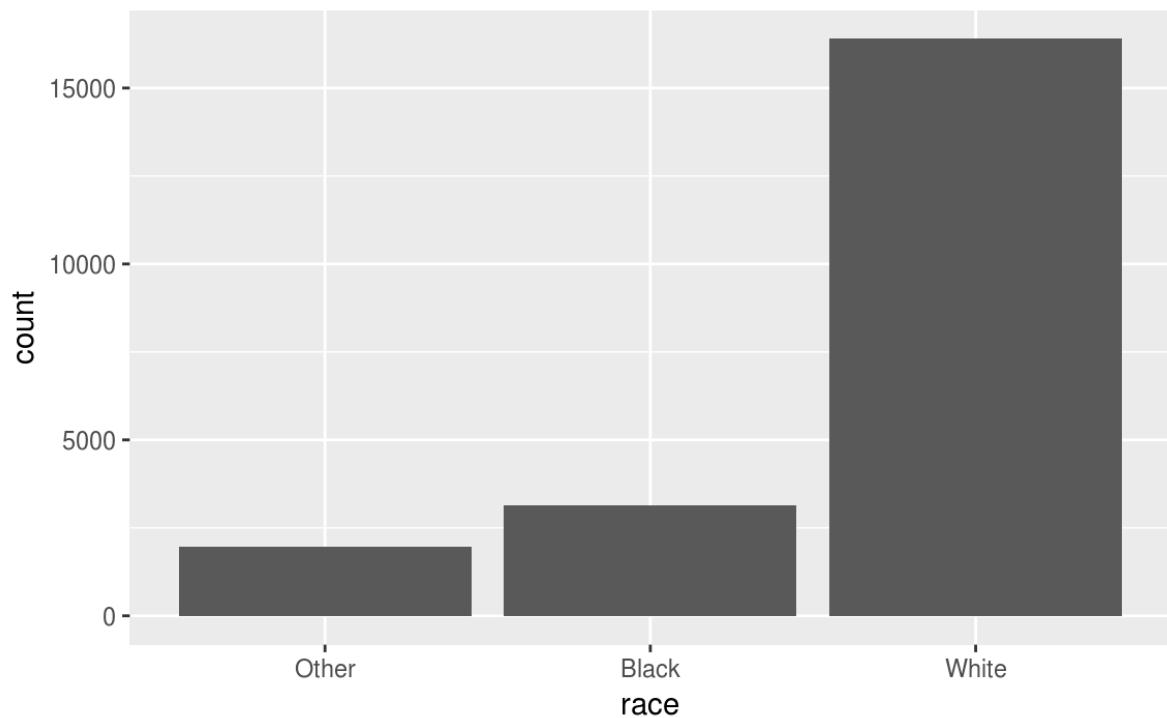
(Remember, since this dataset is provided by a package, you can get more information about the variables with `?gss_cat`.)

When factors are stored in a tibble, you can't see their levels so easily. One way to see them is with `count()`:

```
gss_cat %>%
 count(race)
#> # A tibble: 3 × 2
#> race n
#> <fctr> <int>
#> 1 Other 1959
#> 2 Black 3129
#> 3 White 16395
```

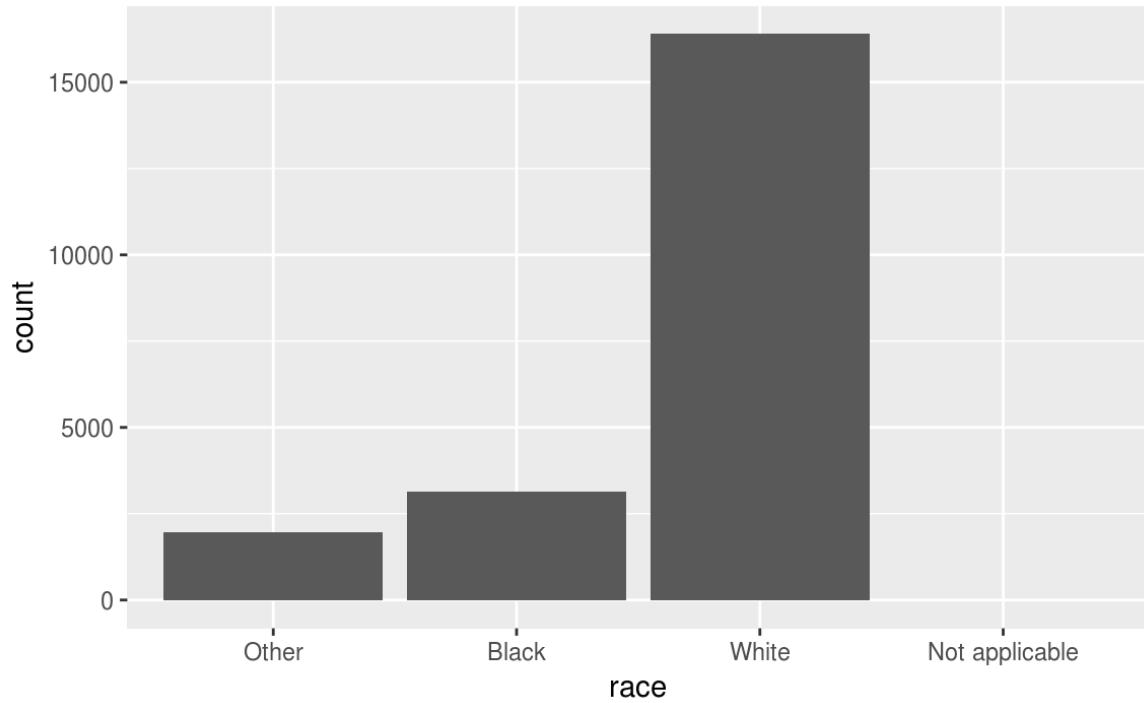
Or with a bar chart:

```
ggplot(gss_cat, aes(race)) +
 geom_bar()
```



By default, ggplot2 will drop levels that don't have any values. You can force them to display with:

```
ggplot(gss_cat, aes(race)) +
 geom_bar() +
 scale_x_discrete(drop = FALSE)
```



These levels represent valid values that simply did not occur in this dataset. Unfortunately, dplyr **doesn't yet have a** drop option, but it will in the future.

When working with factors, the two most common operations are changing the order of the levels, and changing the values of the levels. Those operations are described in the sections below.

### 15.3.1 Exercise

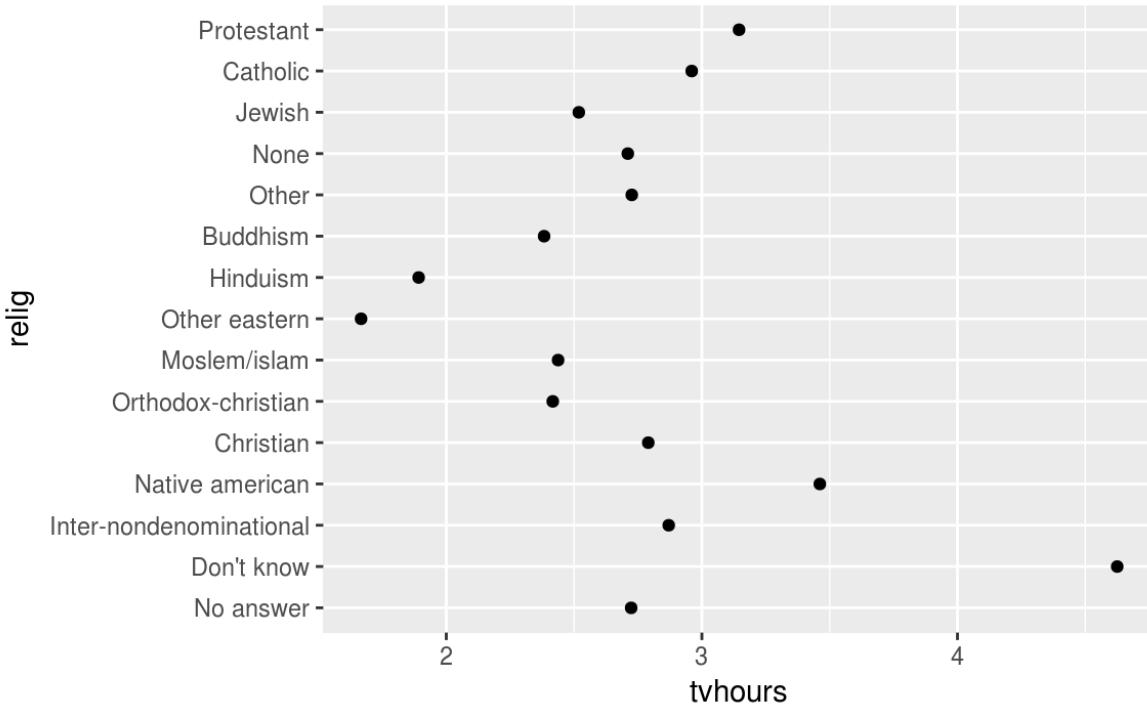
1. Explore the distribution of `rincome` (reported income). What makes the default bar chart hard to understand? How could you improve the plot?
2. What is the most common `relig` in this survey? What's the most common `partyid`?
3. Which `relig` does `denom` (denomination) apply to? How can you find out with a table? How can you find out with a visualisation?

## 15.4 Modifying factor order

It's often useful to change the order of the factor levels in a visualisation. For example, imagine you want to explore the average number of hours spent watching TV per day across religions:

```
relig_summary <- gss_cat %>%
 group_by(relig) %>%
 summarise(
 age = mean(age, na.rm = TRUE),
 tvhours = mean(tvhours, na.rm = TRUE),
 n = n()
)

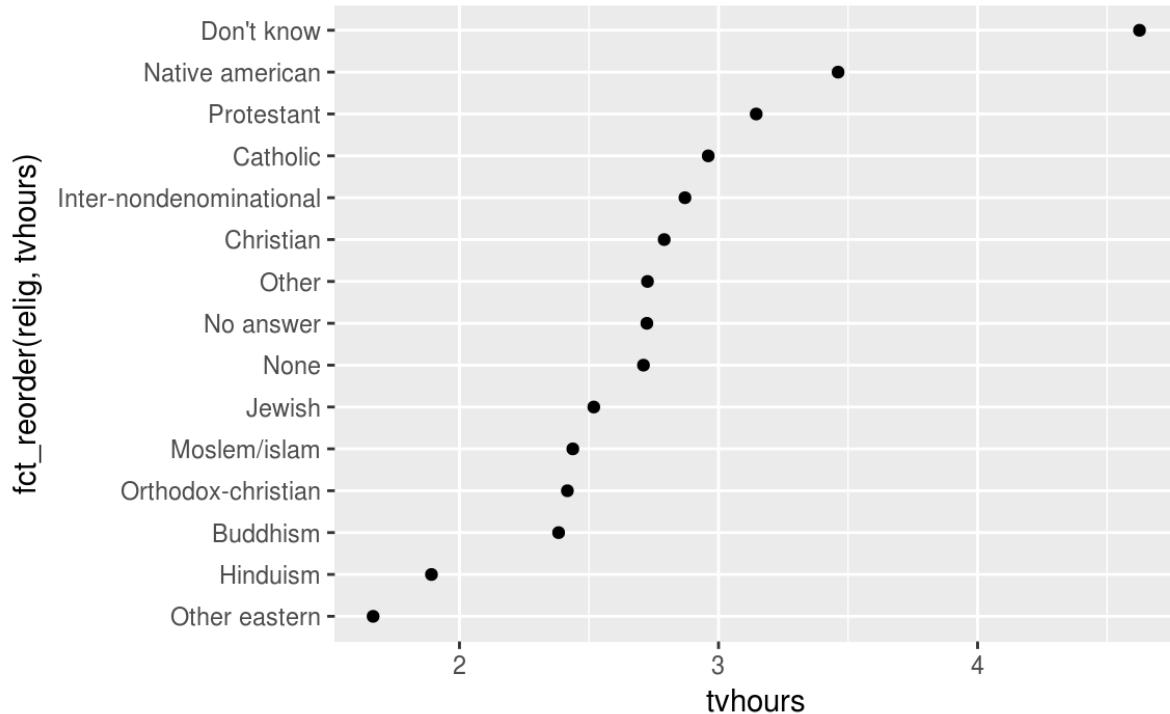
ggplot(relig_summary, aes(tvhours, relig)) + geom_point()
```



It is difficult to interpret this plot because there's no overall pattern. We can improve it by reordering the levels of `relig` using `fct_reorder()`. `fct_reorder()` takes three arguments:

- `f`, the factor whose levels you want to modify.
- `x`, a numeric vector that you want to use to reorder the levels.
- Optionally, `fun`, a function that's used if there are multiple values of `x` for each value of `f`. The default value is `median`.

```
ggplot(relig_summary, aes(tvhours, fct_reorder(relig, tvhours))) +
 geom_point()
```



Reordering religion makes it much easier to see that people in the “Don’t know” category watch much more TV, and Hinduism & Other Eastern religions watch much less.

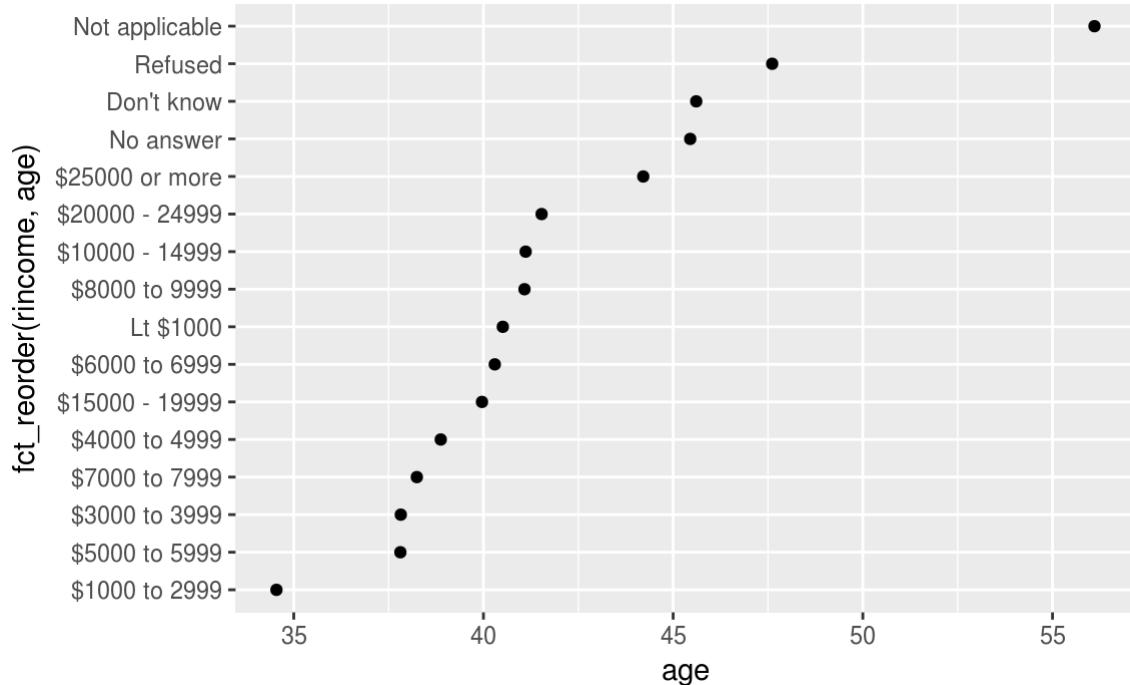
As you start making more complicated transformations, I’d recommend moving them out of `aes()` and into a separate `mutate()` step. For example, you could rewrite the plot above as:

```
relig_summary %>%
 mutate(relig = fct_reorder(relig, tvhours)) %>%
 ggplot(aes(tvhours, relig)) +
 geom_point()
```

What if we create a similar plot looking at how average age varies across reported income level?

```
rincome_summary <- gss_cat %>%
 group_by(rincome) %>%
 summarise(
 age = mean(age, na.rm = TRUE),
 tvhours = mean(tvhours, na.rm = TRUE),
 n = n()
)

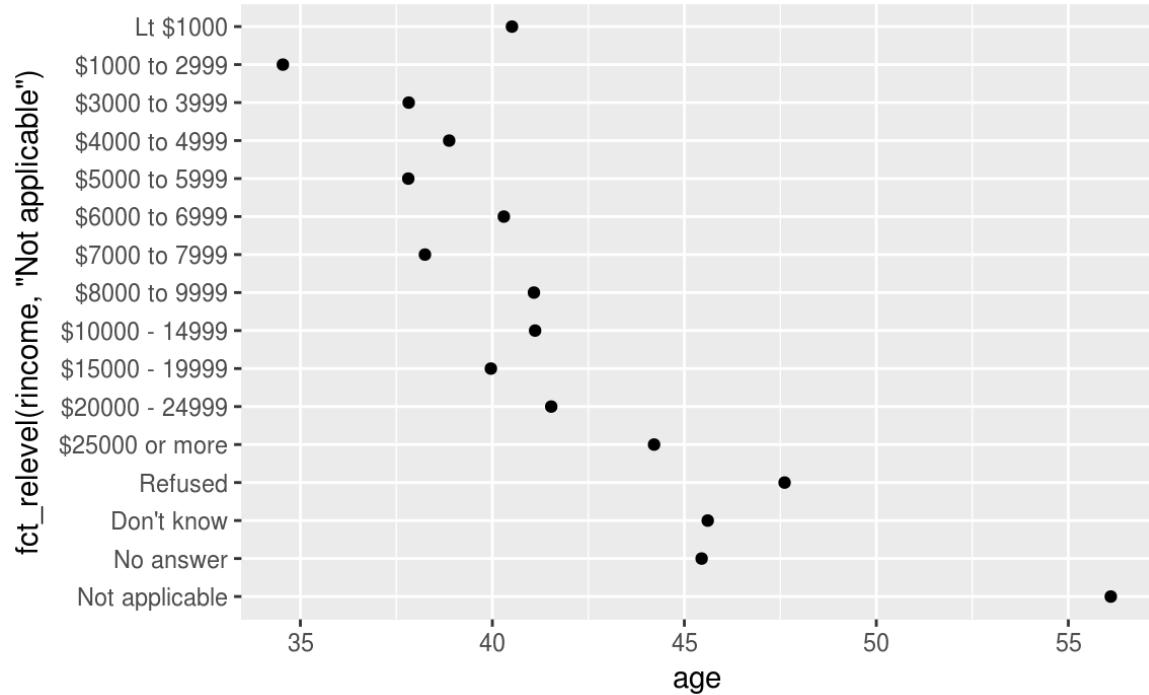
ggplot(rincome_summary, aes(age, fct_reorder(rincome, age))) +
 geom_point()
```



Here, arbitrarily reordering the levels isn't a good idea! That's because `rincome` already has a principled order that we shouldn't mess with. Reserve `fct_reorder()` for factors whose levels are arbitrarily ordered.

However, it does make sense to pull "Not applicable" to the front with the other special levels. You can use `fct_relevel()`. It takes a factor, `f`, and then any number of levels that you want to move to the front of the line.

```
ggplot(rincome_summary, aes(age, fct_relevel(rincome, "Not applicable")) +
 geom_point())
```



Why do you think the average age for “Not applicable” is so high?

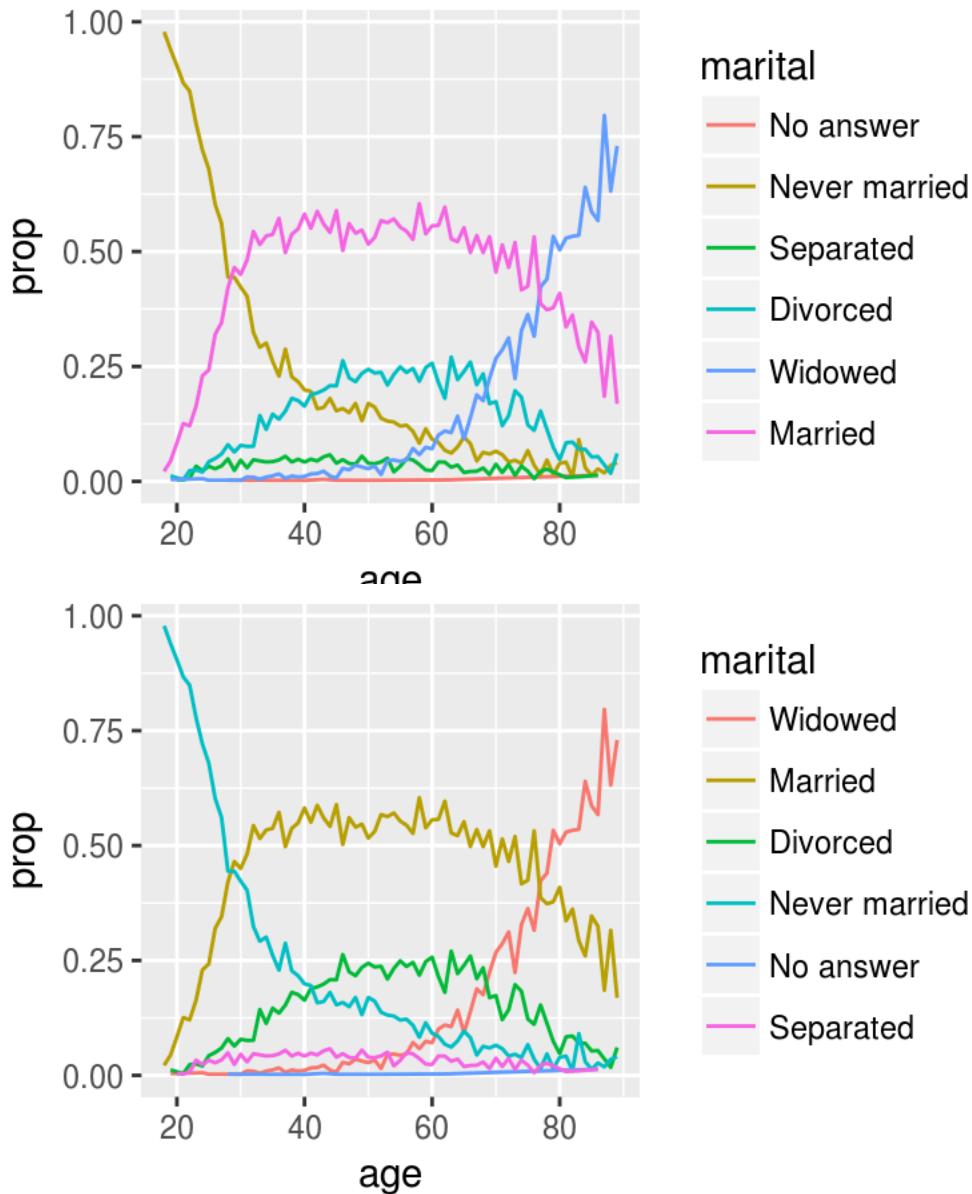
Another type of reordering is useful when you are colouring the lines on a plot.

`fct_reorder2()` reorders the factor by the y values associated with the largest x values. This makes the plot easier to read because the line colours line up with the legend.

```
by_age <- gss_cat %>%
 filter(!is.na(age)) %>%
 group_by(age, marital) %>%
 count() %>%
 mutate(prop = n / sum(n))

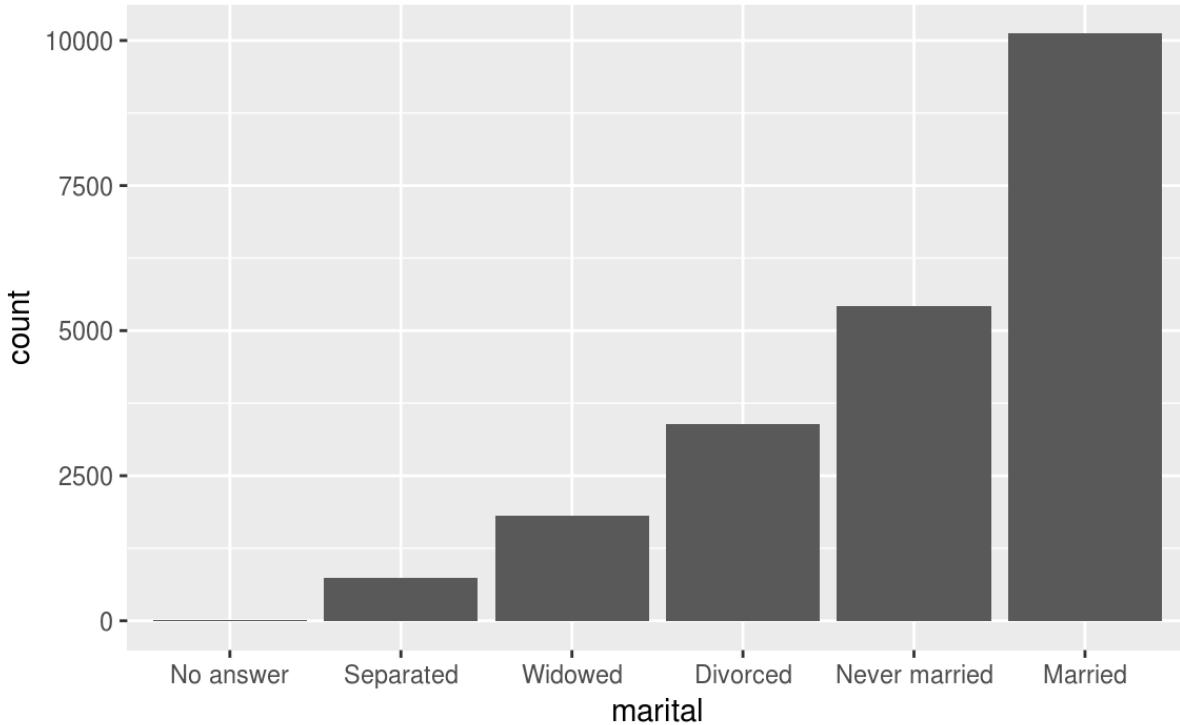
ggplot(by_age, aes(age, prop, colour = marital)) +
 geom_line(na.rm = TRUE)

ggplot(by_age, aes(age, prop, colour = fct_reorder2(marital, age,
prop))) +
 geom_line() +
 labs(colour = "marital")
```



Finally, for bar plots, you can use `fct_infreq()` to order levels in increasing frequency: this is the simplest type of reordering because it doesn't need any extra variables. You may want to combine with `fct_rev()`.

```
gss_cat %>%
 mutate(marital = marital %>% fct_infreq() %>% fct_rev()) %>%
 ggplot(aes(marital)) +
 geom_bar()
```



#### 15.4.1 Exercises

1. There are some suspiciously high numbers in `tvhours`. Is the mean a good summary?
2. For each factor in `gss_cat` identify whether the order of the levels is arbitrary or principled.
3. **Why did moving “Not applicable” to the front of the levels move it to the bottom of the plot?**

## 15.5 Modifying factor levels

More powerful than changing the orders of the levels is changing their values. This allows you to clarify labels for publication, and collapse levels for high-level displays. The most general and powerful tool is `fct_recode()`. It allows you to recode, or change, the value of each level. For example, take the `gss_cat$partyid`:

```
gss_cat %>% count(partyid)
#> # A tibble: 10 × 2
#> partyid n
#> <fctr> <int>
#> 1 No answer 154
#> 2 Don't know 1
#> 3 Other party 393
#> 4 Strong republican 2314
#> 5 Not str republican 3032
#> 6 Ind,near rep 1791
#> # ... with 4 more rows
```

The levels are terse and inconsistent. Let's tweak them to be longer and use a parallel construction.

```
gss_cat %>%
 mutate(partyid = fct_recode(partyid,
 "Republican, strong" = "Strong republican",
 "Republican, weak" = "Not str republican",
 "Independent, near rep" = "Ind,near rep",
 "Independent, near dem" = "Ind,near dem",
 "Democrat, weak" = "Not str democrat",
 "Democrat, strong" = "Strong democrat"
)) %>%
 count(partyid)
#> # A tibble: 10 × 2
#> partyid n
#> <fctr> <int>
#> 1 No answer 154
#> 2 Don't know 1
#> 3 Other party 393
#> 4 Republican, strong 2314
#> 5 Republican, weak 3032
#> 6 Independent, near rep 1791
#> # ... with 4 more rows
```

`fct_recode()` will leave levels that aren't explicitly mentioned as is, and will warn you if you accidentally refer to a level that doesn't exist.

To combine groups, you can assign multiple old levels to the same new level:

```
gss_cat %>%
 mutate(partyid = fct_recode(partyid,
 "Republican, strong" = "Strong republican",
 "Republican, weak" = "Not str republican",
 "Independent, near rep" = "Ind,near rep",
 "Independent, near dem" = "Ind,near dem",
 "Democrat, weak" = "Not str democrat",
 "Democrat, strong" = "Strong democrat",
 "Other" = "No answer",
 "Other" = "Don't know",
 "Other" = "Other party"
)) %>%
 count(partyid)
#> # A tibble: 8 × 2
#> partyid n
#> <fctr> <int>
#> 1 Other 548
#> 2 Republican, strong 2314
#> 3 Republican, weak 3032
#> 4 Independent, near rep 1791
#> 5 Independent, near dem 4119
#> 6 Independent, near dem 2499
#> # ... with 2 more rows
```

You must use this technique with care: if you group together categories that are truly different you will end up with misleading results.

If you want to collapse a lot of levels, `fct_collapse()` is a useful variant of `fct_recode()`. For each new variable, you can provide a vector of old levels:

```
gss_cat %>%
 mutate(partyid = fct_collapse(partyid,
 other = c("No answer", "Don't know", "Other party"),
 rep = c("Strong republican", "Not str republican"),
 ind = c("Ind,near rep", "Independent", "Ind,near dem"),
 dem = c("Not str democrat", "Strong democrat"))
) %>%
 count(partyid)
#> # A tibble: 4 × 2
#> partyid n
#> <fctr> <int>
#> 1 other 548
#> 2 rep 5346
#> 3 ind 8409
#> 4 dem 7180
```

Sometimes you just want to lump together all the small groups to make a plot or table simpler. That's the job of `fct_lump()`:

```
gss_cat %>%
 mutate(relig = fct_lump(relig)) %>%
 count(relig)
#> # A tibble: 2 × 2
#> relig n
#> <fctr> <int>
#> 1 Protestant 10846
#> 2 Other 10637
```

The default behaviour is to progressively lump together the smallest groups, ensuring that the **aggregate is still the smallest group. In this case it's not very helpful: it is true that the majority of Americans in this survey are Protestant, but we've probably over collapsed.**

Instead, we can use the `n` parameter to specify how many groups (excluding other) we want to keep:

```
gss_cat %>%
 mutate(relig = fct_lump(relig, n = 10)) %>%
 count(relig, sort = TRUE) %>%
 print(n = Inf)
#> # A tibble: 10 × 2
#> relig n
#> <fctr> <int>
#> 1 Protestant 10846
#> 2 Catholic 5124
#> 3 None 3523
#> 4 Christian 689
```

```
#> 5 Other 458
#> 6 Jewish 388
#> 7 Buddhism 147
#> 8 Inter-nondenominational 109
#> 9 Moslem/islam 104
#> 10 Orthodox-christian 95
```

### 15.5.1 Exercises

1. How have the proportions of people identifying as Democrat, Republican, and Independent changed over time?
2. How could you collapse `rincome` into a small set of categories?

# 19 Functions

## 19.1 Introduction

One of the best ways to improve your reach as a data scientist is to write functions. Functions allow you to automate common tasks in a more powerful and general way than copy-and-pasting. Writing a function has three big advantages over using copy-and-paste:

1. You can give a function an evocative name that makes your code easier to understand.
2. As requirements change, you only need to update code in one place, instead of many.
3. You eliminate the chance of making incidental mistakes when you copy and paste (i.e. updating a variable name in one place, but not in another).

Writing good functions is a lifetime journey. Even after using R for many years I still learn new techniques and better ways of approaching old problems. The goal of this chapter is not to teach you every esoteric detail of functions but to get you started with some pragmatic advice that you can apply immediately.

As well as practical advice for writing functions, this chapter also gives you some suggestions for how to style your code. Good code style is like correct punctuation. You can manage without it, but it sure makes things easier to read! As with styles of punctuation, there are many possible variations. Here we present the style we use in our code, but the most important thing is to be consistent.

### 19.1.1 Prerequisites

The focus of this chapter is on writing functions in base R, so you won't need any extra packages.

## 19.2 When should you write a function?

You should consider writing a function whenever you've copied and pasted a block of code more than twice (i.e. you now have three copies of the same code). For example, take a look at this code. What does it do?

```
df <- tibble::tibble(
 a = rnorm(10),
 b = rnorm(10),
 c = rnorm(10),
 d = rnorm(10)
)

df$a <- (df$a - min(df$a, na.rm = TRUE)) /
 (max(df$a, na.rm = TRUE) - min(df$a, na.rm = TRUE))
df$b <- (df$b - min(df$b, na.rm = TRUE)) /
 (max(df$b, na.rm = TRUE) - min(df$a, na.rm = TRUE))
df$c <- (df$c - min(df$c, na.rm = TRUE)) /
 (max(df$c, na.rm = TRUE) - min(df$c, na.rm = TRUE))
df$d <- (df$d - min(df$d, na.rm = TRUE)) /
 (max(df$d, na.rm = TRUE) - min(df$d, na.rm = TRUE))
```

You might be able to puzzle out that this rescales each column to have a range from 0 to 1. But did you spot the mistake? I made an error when copying-and-pasting the code for `df$b`: I forgot to change an `a` to a `b`. Extracting repeated code out into a function is a good idea because it prevents you from making this type of mistake.

To write a function you need to first analyse the code. How many inputs does it have?

```
(df$a - min(df$a, na.rm = TRUE)) /
(max(df$a, na.rm = TRUE) - min(df$a, na.rm = TRUE))
```

This code only has one input: `df$a`. (If you're surprised that `TRUE` is not an input, you can explore why in the exercise below.) To make the inputs more clear, it's a good idea to rewrite the code using temporary variables with general names. Here this code only requires a single numeric vector, so I'll call it `x`:

```
x <- df$a
(x - min(x, na.rm = TRUE)) / (max(x, na.rm = TRUE) - min(x, na.rm =
TRUE))
#> [1] 0.289 0.751 0.000 0.678 0.853 1.000 0.172 0.611 0.612 0.601
```

There is some duplication in this code. We're computing the range of the data three times, but it makes sense to do it in one step:

```
rng <- range(x, na.rm = TRUE)
(x - rng[1]) / (rng[2] - rng[1])
#> [1] 0.289 0.751 0.000 0.678 0.853 1.000 0.172 0.611 0.612 0.601
```

Pulling out intermediate calculations into named variables is a good practice because it makes it more clear what the code is doing. Now that I've simplified the code, and checked that it still works, I can turn it into a function:

```
rescale01 <- function(x) {
 rng <- range(x, na.rm = TRUE)
 (x - rng[1]) / (rng[2] - rng[1])
}
rescale01(c(0, 5, 10))
#> [1] 0.0 0.5 1.0
```

There are three key steps to creating a new function:

1. You need to pick a name **for the function**. Here I've used `rescale01` because this function rescales a vector to lie between 0 and 1.
2. You list the inputs, or arguments, to the function inside `function`. Here we have just one argument. If we had more the call would look like `function(x, y, z)`.
3. You place the code you have developed in body of the function, a `{` block that immediately follows `function(...)`.

Note the overall process: I only made the function after I'd figured out how to make it work with a simple input. It's easier to start with working code and turn it into a function; it's harder to create a function and then try to make it work.

At this point it's a good idea to check your function with a few different inputs:

```
rescale01(c(-10, 0, 10))
#> [1] 0.0 0.5 1.0
rescale01(c(1, 2, 3, NA, 5))
#> [1] 0.00 0.25 0.50 NA 1.00
```

As you write more and more functions you'll eventually want to convert these informal, interactive tests into formal, automated tests. That process is called unit testing. Unfortunately, it's beyond the scope of this book, but you can learn about it in <http://r-pkgs.had.co.nz/tests.html>.

We can simplify the original example now that we have a function:

```
df$a <- rescale01(df$a)
df$b <- rescale01(df$b)
df$c <- rescale01(df$c)
df$d <- rescale01(df$d)
```

Compared to the original, this code is easier to understand and we've eliminated one class of copy-and-paste errors. There is still quite a bit of duplication since we're doing the same thing to multiple columns. We'll learn how to eliminate that duplication in [iteration](#), once you've learned more about R's data structures in [vectors](#).

Another advantage of functions is that if our requirements change, we only need to make the change in one place. For example, we might discover that some of our variables include infinite values, and `rescale01()` fails:

```
x <- c(1:10, Inf)
rescale01(x)
#> [1] 0 0 0 0 0 0 0 0 0 NaN
```

Because we've extracted the code into a function, we only need to make the fix in one place:

```
rescale01 <- function(x) {
 rng <- range(x, na.rm = TRUE, finite = TRUE)
 (x - rng[1]) / (rng[2] - rng[1])
}
rescale01(x)
#> [1] 0.000 0.111 0.222 0.333 0.444 0.556 0.667 0.778 0.889 1.000
Inf
```

This is an important part of the “do not repeat yourself” (or DRY) principle. The more repetition you have in your code, the more places you need to remember to update when things change (and they always do!), and the more likely you are to create bugs over time.

## 19.2.1 Practice

1. Why is `TRUE` not a parameter to `rescale01()`? What would happen if `x` contained a single missing value, and `na.rm` was `FALSE`?
2. In the second variant of `rescale01()`, infinite values are left unchanged. Rewrite `rescale01()` so that `-Inf` is mapped to 0, and `Inf` is mapped to 1.

3. Practice turning the following code snippets into functions. Think about what each function does. What would you call it? How many arguments does it need? Can you rewrite it to be more expressive or less duplicative?
4. `mean(is.na(x))`
- 5.
6. `x / sum(x, na.rm = TRUE)`
7. `sd(x, na.rm = TRUE) / mean(x, na.rm = TRUE)`
8. Follow <http://nicercode.github.io/intro/writing-functions.html> to write your own functions to compute the variance and skew of a numeric vector.
9. Write `both_na()`, a function that takes two vectors of the same length and returns the number of positions that have an NA in both vectors.
10. What do the following functions do? Why are they useful even though they are so short?
11. `is_directory <- function(x) file.info(x)$isdir`  
`is_readable <- function(x) file.access(x, 4) == 0`
12. Read the [complete lyrics](#) to “Little Bunny Foo Foo”. There’s a lot of duplication in this song. Extend the initial piping example to recreate the complete song, and use functions to reduce the duplication.

## 19.3 Functions are for humans and computers

It’s important to remember that functions are not just for the computer, but are also for humans. R doesn’t care what your function is called, or what comments it contains, but these are important for human readers. This section discusses some things that you should bear in mind when writing functions that humans can understand.

The name of a function is important. Ideally, the name of your function will be short, but clearly evoke what the function does. That’s hard! But it’s better to be clear than short, as RStudio’s autocomplete makes it easy to type long names.

Generally, function names should be verbs, and arguments should be nouns. There are some exceptions: nouns are ok if the function computes a very well known noun (i.e. `mean()` is better than `compute_mean()`), or accessing some property of an object (i.e. `coef()` is better than `get_coefficients()`). A good sign that a noun might be a better choice is if you’re using a very broad verb like “get”, “compute”, “calculate”, or “determine”. Use your best judgement and don’t be afraid to rename a function if you figure out a better name later.

```
Too short
f()

Not a verb, or descriptive
my_awesome_function()

Long, but clear
impute_missing()
collapse_years()
```

If your function name is composed of multiple words, I recommend using “snake\_case”, where each lowercase word is separated by an underscore. camelCase is a popular alternative. It

**doesn't** really matter which one you pick, the important thing is to be consistent: pick one or the other and stick with it. R itself is not very consistent, but there's nothing you can do about that. Make sure you don't fall into the same trap by making your code as consistent as possible.

```
Never do this!
col_mins <- function(x, y) {}
rowMaxes <- function(y, x) {}
```

If you have a family of functions that do similar things, make sure they have consistent names and arguments. Use a common prefix to indicate that they are connected. That's better than a common suffix because autocomplete allows you to type the prefix and see all the members of the family.

```
Good
input_select()
input_checkbox()
input_text()

Not so good
select_input()
checkbox_input()
text_input()
```

A good example of this design is the `stringr` package: if you don't remember exactly which function you need, you can type `str_` and jog your memory.

Where possible, avoid overriding existing functions and variables. It's impossible to do in general because so many good names are already taken by other packages, but avoiding the most common names from base R will avoid confusion.

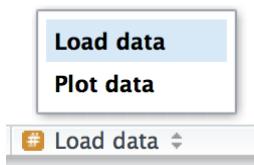
```
Don't do this!
T <- FALSE
c <- 10
mean <- function(x) sum(x)
```

Use comments, lines starting with `#`, to explain the “why” of your code. You generally should avoid comments that explain the “what” or the “how”. If you can't understand what the code does from reading it, you should think about how to rewrite it to be more clear. Do you need to add some intermediate variables with useful names? Do you need to break out a subcomponent of a large function so you can name it? However, your code can never capture the reasoning behind your decisions: why did you choose this approach instead of an alternative? What else did you try that didn't work? It's a great idea to capture that sort of thinking in a comment.

Another important use of comments is to break up your file into easily readable chunks. Use long lines of `-` and `=` to make it easy to spot the breaks.

```
Load data -----
Plot data -----
```

RStudio provides a keyboard shortcut to create these headers (Cmd/Ctrl + Shift + R), and will display them in the code navigation drop-down at the bottom-left of the editor:



### 19.3.1 Exercises

1. Read the source code for each of the following three functions, puzzle out what they do, and then brainstorm better names.
2. 

```
f1 <- function(string, prefix) {
 substr(string, 1, nchar(prefix)) == prefix
}
f2 <- function(x) {
 if (length(x) <= 1) return(NULL)
 x[-length(x)]
}
f3 <- function(x, y) {
 rep(y, length.out = length(x))
}
```
11. Take a function that you've written recently and spend 5 minutes brainstorming a better name for it and its arguments.
12. Compare and contrast `rnorm()` and `MASS::mvrnorm()`. How could you make them more consistent?
13. Make a case for why `norm_r()`, `norm_d()` etc would be better than `rnorm()`, `dnorm()`. Make a case for the opposite.

## 19.4 Conditional execution

An `if` statement allows you to conditionally execute code. It looks like this:

```
if (condition) {
 # code executed when condition is TRUE
} else {
 # code executed when condition is FALSE
}
```

To get help on `if` you need to surround it in backticks: `?`if``. The help isn't particularly helpful if you're not already an experienced programmer, but at least you know how to get to it!

Here's a simple function that uses an `if` statement. The goal of this function is to return a logical vector describing whether or not each element of a vector is named.

```
has_name <- function(x) {
 nms <- names(x)
 if (is.null(nms)) {
 rep(FALSE, length(x))
 } else {
 !is.na(nms) & nms != ""
```

```
 }
}
```

This function takes advantage of the standard return rule: a function returns the last value that it computed. Here that is either one of the two branches of the `if` statement.

### 19.4.1 Conditions

The condition must evaluate to either TRUE or FALSE. If it's a vector, you'll get a warning message; if it's an NA, you'll get an error. Watch out for these messages in your own code:

```
if (c(TRUE, FALSE)) {}
#> Warning in if (c(TRUE, FALSE)) {: the condition has length > 1 and
only the
#> first element will be used
#> NULL

if (NA) {}
#> Error in if (NA) {: missing value where TRUE/FALSE needed
```

You can use `||` (or) and `&&` (and) to combine multiple logical expressions. These operators are **“short-circuiting”**: as soon as `||` sees the first TRUE it returns TRUE without computing anything else. As soon as `&&` sees the first FALSE it returns FALSE. You should never use `|` or `&` in an `if` statement: these are vectorised operations that apply to **multiple values (that's why you use them in `filter()`)**. If you do have a logical vector, you can use `any()` or `all()` to collapse it to a single value.

Be careful when testing for equality. `==` is **vectorised**, which means that it's easy to get more than one output. Either check the length is already 1, collapse with `all()` or `any()`, or use the non-vectorised `identical()`. `identical()` is very strict: it always returns either a single TRUE or a single FALSE, and doesn't coerce types. This means that you need to be careful when comparing integers and doubles:

```
identical(0L, 0)
#> [1] FALSE
```

You also need to be wary of floating point numbers:

```
x <- sqrt(2) ^ 2
x
#> [1] 2
x == 2
#> [1] FALSE
x - 2
#> [1] 4.44e-16
```

Instead use `dplyr::near()` for comparisons, as described in [comparisons](#).

And remember, `x == NA` doesn't do anything useful!

## 19.4.2 Multiple conditions

You can chain multiple if statements together:

```
if (this) {
 # do that
} else if (that) {
 # do something else
} else {
 #
}
```

But if you end up with a very long series of chained `if` statements, you should consider rewriting. One useful technique is the `switch()` function. It allows you to evaluate selected code based on position or name.

```
#> function(x, y, op) {
#> switch(op,
#> plus = x + y,
#> minus = x - y,
#> times = x * y,
#> divide = x / y,
#> stop("Unknown op!")
#>)
#> }
```

Another useful function that can often eliminate long chains of `if` statements is `cut()`. It's used to discretise continuous variables.

## 19.4.3 Code style

Both `if` and `function` should (almost) always be followed by squiggly brackets (`{ }` ), and the contents should be indented by two spaces. This makes it easier to see the hierarchy in your code by skimming the left-hand margin.

An opening curly brace should never go on its own line and should always be followed by a new line. A closing curly brace should always go on its own line, unless it's followed by `else`. Always indent the code inside curly braces.

```
Good
if (y < 0 && debug) {
 message("Y is negative")
}

if (y == 0) {
 log(x)
} else {
 y ^ x
}

Bad
```

```

if (y < 0 && debug)
message("Y is negative")

if (y == 0) {
 log(x)
}
else {
 y ^ x
}

```

It's ok to drop the curly braces if you have a very short `if` statement that can fit on one line:

```

y <- 10
x <- if (y < 20) "Too low" else "Too high"

```

I recommend this only for very brief `if` statements. Otherwise, the full form is easier to read:

```

if (y < 20) {
 x <- "Too low"
} else {
 x <- "Too high"
}

```

#### 19.4.4 Exercises

1. What's the difference between `if` and `ifelse()`? Carefully read the help and construct three examples that illustrate the key differences.
2. Write a greeting function that says “good morning”, “good afternoon”, or “good evening”, depending on the time of day. (Hint: use a time argument that defaults to `lubridate::now()`. That will make it easier to test your function.)
3. Implement a `fizzbuzz` function. It takes a single number as input. If the number is divisible by three, it returns “fizz”. If it's divisible by five it returns “buzz”. If it's divisible by three and five, it returns “fizzbuzz”. Otherwise, it returns the number. Make sure you first write working code before you create the function.
4. How could you use `cut()` to simplify this set of nested if-else statements?
  5. if (temp <= 0) {
  6. "freezing"
  7. } else if (temp <= 10) {
  8. "cold"
  9. } else if (temp <= 20) {
  10. "cool"
  11. } else if (temp <= 30) {
  12. "warm"
  13. } else {
  14. "hot"

How would you change the call to `cut()` if I'd used `<` instead of `<=?` What is the other chief advantage of `cut()` for this problem? (Hint: what happens if you have many values in `temp`?)

```

15. What happens if you use switch() with numeric values?
16. What does this switch() call do? What happens if x is "e"?
17. switch(x,
18. a = ,
19. b = "ab",
20. c = ,
21. d = "cd"
)

```

Experiment, then carefully read the documentation.

## 19.5 Function arguments

The arguments to a function typically fall into two broad sets: one set supplies the data to compute on, and the other supplies arguments that control the details of the computation. For example:

- In `log()`, the data is `x`, and the detail is the `base` of the logarithm.
- In `mean()`, the data is `x`, and the details are how much data to trim from the ends (`trim`) and how to handle missing values (`na.rm`).
- In `t.test()`, the data are `x` and `y`, and the details of the test are `alternative`, `mu`, `paired`, `var.equal`, and `conf.level`.
- In `str_c()` you can supply any number of strings to `...`, and the details of the concatenation are controlled by `sep` and `collapse`.

Generally, data arguments should come first. Detail arguments should go on the end, and usually should have default values. You specify a default value in the same way you call a function with a named argument:

```

Compute confidence interval around mean using normal approximation
mean_ci <- function(x, conf = 0.95) {
 se <- sd(x) / sqrt(length(x))
 alpha <- 1 - conf
 mean(x) + se * qnorm(c(alpha / 2, 1 - alpha / 2))
}

x <- runif(100)
mean_ci(x)
#> [1] 0.498 0.610
mean_ci(x, conf = 0.99)
#> [1] 0.480 0.628

```

The default value should almost always be the most common value. The few exceptions to this rule are to do with safety. For example, it makes sense for `na.rm` to default to `FALSE` because missing values are important. Even though `na.rm = TRUE` is what you usually put in your code, **it's a bad idea to silently ignore** missing values by default.

When you call a function, you typically omit the names of the data arguments, because they are used so commonly. If you override the default value of a detail argument, you should use the full name:

```

Good
mean(1:10, na.rm = TRUE)

Bad
mean(x = 1:10, , FALSE)
mean(, TRUE, x = c(1:10, NA))

```

You can refer to an argument by its unique prefix (e.g. `mean(x, n = TRUE)`), but this is generally best avoided given the possibilities for confusion.

Notice that when you call a function, you should place a space around `=` in function calls, and always put a space after a comma, not before (just like in regular English). Using whitespace makes it easier to skim the function for the important components.

```

Good
average <- mean(feet / 12 + inches, na.rm = TRUE)

Bad
average<-mean(feet/12+inches,na.rm=TRUE)

```

### 19.5.1 Choosing names

**The names of the arguments are also important. R doesn't care, but the readers of your code (including future-you!) will.** Generally you should prefer longer, more descriptive names, but there are a handful of very common, very short names. It's worth memorising these:

- `x, y, z`: vectors.
- `w`: a vector of weights.
- `df`: a data frame.
- `i, j`: numeric indices (typically rows and columns).
- `n`: length, or number of rows.
- `p`: number of columns.

Otherwise, consider matching names of arguments in existing R functions. For example, use `na.rm` to determine if missing values should be removed.

### 19.5.2 Checking values

As you start to write more functions, you'll eventually get to the point where you don't remember exactly how your function works. At this point it's easy to call your function with invalid inputs. To avoid this problem, it's often useful to make constraints explicit. For example, imagine you've written some functions for computing weighted summary statistics:

```

wt_mean <- function(x, w) {
 sum(x * w) / sum(w)
}
wt_var <- function(x, w) {
 mu <- wt_mean(x, w)
 sum(w * (x - mu) ^ 2) / sum(w)
}

```

```
wt_sd <- function(x, w) {
 sqrt(wt_var(x, w))
}
```

What happens if `x` and `w` are not the same length?

```
wt_mean(1:6, 1:3)
#> [1] 7.67
```

In this case, because of R's vector recycling rules, we don't get an error.

It's good practice to check important preconditions, and throw an error (with `stop()`), if they are not true:

```
wt_mean <- function(x, w) {
 if (length(x) != length(w)) {
 stop("`x` and `w` must be the same length", call. = FALSE)
 }
 sum(w * x) / sum(w)
}
```

Be careful not to take this too far. There's a tradeoff between how much time you spend making your function robust, versus how long you spend writing it. For example, if you also added a `na.rm` argument, I probably wouldn't check it carefully:

```
wt_mean <- function(x, w, na.rm = FALSE) {
 if (!is.logical(na.rm)) {
 stop("`na.rm` must be logical")
 }
 if (length(na.rm) != 1) {
 stop("`na.rm` must be length 1")
 }
 if (length(x) != length(w)) {
 stop("`x` and `w` must be the same length", call. = FALSE)
 }

 if (na.rm) {
 miss <- is.na(x) | is.na(w)
 x <- x[!miss]
 w <- w[!miss]
 }
 sum(w * x) / sum(w)
}
```

This is a lot of extra work for little additional gain. A useful compromise is the built-in `stopifnot()`: it checks that each argument is TRUE, and produces a generic error message if not.

```
wt_mean <- function(x, w, na.rm = FALSE) {
 stopifnot(is.logical(na.rm), length(na.rm) == 1)
 stopifnot(length(x) == length(w))
```

```

if (na.rm) {
 miss <- is.na(x) | is.na(w)
 x <- x[!miss]
 w <- w[!miss]
}
sum(w * x) / sum(w)
}
wt_mean(1:6, 6:1, na.rm = "foo")
#> Error: is.logical(na.rm) is not TRUE

```

Note that when using `stopifnot()` you assert what should be true rather than checking for what might be wrong.

### 19.5.3 Dot-dot-dot (...)

Many functions in R take an arbitrary number of inputs:

```

sum(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)
#> [1] 55
stringr::str_c("a", "b", "c", "d", "e", "f")
#> [1] "abcdef"

```

How do these functions work? They rely on a special argument: `...` (pronounced dot-dot-dot). **This special argument captures any number of arguments that aren't otherwise matched.**

**It's useful because you can then send those `...` on to another function.** This is a useful catch-all if your function primarily wraps another function. For example, I commonly create these helper functions that wrap around `str_c()`:

```

commas <- function(...) stringr::str_c(..., collapse = ", ")
commas(letters[1:10])
#> [1] "a, b, c, d, e, f, g, h, i, j"

rule <- function(..., pad = "-") {
 title <- paste0(...)
 width <- getOption("width") - nchar(title) - 5
 cat(title, " ", stringr::str_dup(pad, width), "\n", sep = "")
}
rule("Important output")
#> Important output -----
--
```

Here `...` lets me forward on any arguments that I don't want to deal with to `str_c()`. It's a very convenient technique. But it does come at a price: any misspelled arguments will not raise an error. This makes it easy for typos to go unnoticed:

```

x <- c(1, 2)
sum(x, na.mr = TRUE)
#> [1] 4

```

If you just want to capture the values of the `...`, use `list(...)`.

## 19.5.4 Lazy evaluation

Arguments in R are lazily evaluated: they're not computed until they're needed. That means if they're never used, they're never called. This is an important property of R as a programming language, but is generally not important when you're writing your own functions for data analysis. You can read more about lazy evaluation at <http://adv-r.had.co.nz/Functions.html#lazy-evaluation>.

## 19.5.5 Exercises

1. What does `commas(letters, collapse = "-")` do? Why?
2. It'd be nice if you could supply multiple characters to the `pad` argument, e.g. `rule("Title", pad = "-+").` Why doesn't this currently work? How could you fix it?
3. What does the `trim` argument to `mean()` do? When might you use it?
4. The default value for the `method` argument to `cor()` is `c("pearson", "kendall", "spearmann")`. What does that mean? What value is used by default?

## 19.6 Return values

Figuring out what your function should return is usually straightforward: it's why you created the function in the first place! There are two things you should consider when returning a value:

1. Does returning early make your function easier to read?
2. Can you make your function pipeable?

### 19.6.1 Explicit return statements

The value returned by the function is usually the last statement it evaluates, but you can choose to return early by using `return()`. I think it's best to save the use of `return()` to signal that you can return early with a simpler solution. A common reason to do this is because the inputs are empty:

```
complicated_function <- function(x, y, z) {
 if (length(x) == 0 || length(y) == 0) {
 return(0)
 }

 # Complicated code here
}
```

Another reason is because you have a `if` statement with one complex block and one simple block. For example, you might write an `if` statement like this:

```
f <- function() {
 if (x) {
 # Do
 # something
 # that
 # takes
```

```

many
lines
to
express
} else {
 # return something short
}
}

```

But if the first block is very long, by the time you get to the `else`, **you've forgotten the condition**. One way to rewrite it is to use an early return for the simple case:

```

f <- function() {
 if (!x) {
 return(something_short)
 }

 # Do
 # something
 # that
 # takes
 # many
 # lines
 # to
 # express
}

```

This tends to make the code easier to understand, because you don't need quite so much context to understand it.

## 19.6.2 Writing pipeable functions

If you want to write your own pipeable functions, it's important to think about the return value. Knowing the return value's object type will mean that your pipeline will "just work". For example, with `dplyr` and `tidyverse` the object type is the data frame.

There are two basic types of pipeable functions: transformations and side-effects. With transformations, **an object is passed to the function's first argument and a modified object is returned**. With side-effects, the passed object is not transformed. Instead, the function performs an action on the object, like drawing a plot or saving a file. Side-effects functions should **"invisibly" return the first argument, so that while they're not printed they can still be used in a pipeline**. For example, this simple function prints the number of missing values in a data frame:

```

show_missings <- function(df) {
 n <- sum(is.na(df))
 cat("Missing values: ", n, "\n", sep = "")

 invisible(df)
}

```

If we call it interactively, the `invisible()` means that the input `df` **doesn't get printed out**:

```
show_missings(mtcars)
#> Missing values: 0
```

But it's still there, it's just not printed by default:

```
x <- show_missings(mtcars)
#> Missing values: 0
class(x)
#> [1] "data.frame"
dim(x)
#> [1] 32 11
```

And we can still use it in a pipe:

```
mtcars %>%
 show_missings() %>%
 mutate(mpg = ifelse(mpg < 20, NA, mpg)) %>%
 show_missings()
#> Missing values: 0
#> Missing values: 18
```

## 19.7 Environment

The last component of a function is its environment. This is not something you need to **understand deeply when you first start writing functions. However, it's important to know a little bit** about environments because they are crucial to how functions work. The environment of a function controls how R finds the value associated with a name. For example, take this function:

```
f <- function(x) {
 x + y
}
```

In many programming languages, this would be an error, because `y` is not defined inside the function. In R, this is valid code because R uses rules called lexical scoping to find the value associated with a name. Since `y` is not defined inside the function, R will look in the environment where the function was defined:

```
y <- 100
f(10)
#> [1] 110

y <- 1000
f(10)
#> [1] 1010
```

This behaviour seems like a recipe for bugs, and indeed you should avoid creating functions like **this deliberately, but by and large it doesn't cause too many problems (especially if you regularly restart R to get to a clean slate)**.

The advantage of this behaviour is that from a language standpoint it allows R to be very consistent. Every name is looked up using the same set of rules. For `f()` that includes the

behavior of two things that you might not expect: `{` and `+`. This allows you to do devious things like:

```
`+` <- function(x, y) {
 if (runif(1) < 0.1) {
 sum(x, y)
 } else {
 sum(x, y) * 1.1
 }
}
table(replicate(1000, 1 + 2))
#>
#> 3 3.3
#> 100 900
rm(`+`)
```

This is a common phenomenon in R. R places few limits on your power. You can do many things **that you can't do in other** programming languages. You can do many things that 99% of the time are extremely ill-advised (like overriding how addition works!). But this power and flexibility is what makes tools like ggplot2 and dplyr possible. Learning how to make best use of this flexibility is beyond the scope of this book, but you can read about in [Advanced R](#).

# 27 R Markdown

## 27.1 Introduction

R Markdown provides an unified authoring framework for data science, combining your code, its results, and your prose commentary. R Markdown documents are fully reproducible and support dozens of output formats, like PDFs, Word files, slideshows, and more.

R Markdown files are designed to be used in three ways:

1. For communicating to decision makers, who want to focus on the conclusions, not the code behind the analysis.
2. For collaborating with other data scientists (including future you!), who are interested in both your conclusions, and how you reached them (i.e. the code).
3. As an environment in which to do data science, as a modern day lab notebook where you can capture not only what you did, but also what you were thinking.

R Markdown integrates a number of R packages and external tools. This means that help is, by-and-large, not available through ?. Instead, as you work through this chapter, and use R Markdown in the future, keep these resources close to hand:

- R Markdown Cheat Sheet: `Help > Cheatsheets > R Markdown CheatSheet`
- R Markdown Reference Guide: `Help > Cheatsheets > R Markdown Reference Guide`.

Both cheatsheets are also available at <http://rstudio.com/cheatsheets>.

### 27.1.1 Prerequisites

You need the rmarkdown **package**, but you don't need to explicitly install it or load it, as RStudio automatically does both when needed.

## 27.2 R Markdown basics

This is an R Markdown file, a plain text file that has the extension `.Rmd`:

```

title: "Diamond sizes"
date: 2016-08-25
output: html_document

```{r setup, include = FALSE}
library(ggplot2)
library(dplyr)

smaller <- diamonds %>%
  filter(carat <= 2.5)
...``
```

We have data about `r nrow(diamonds)` diamonds. Only `r nrow(diamonds) - nrow(smaller)` are larger than 2.5 carats. The distribution of the remainder is shown below:

```
```{r, echo = FALSE}
smaller %>%
 ggplot(aes(carat)) +
 geom_freqpoly(binwidth = 0.01)
```

```

It contains three important types of content:

1. An (optional) YAML header surrounded by ---s.
2. Chunks of R code surrounded by ```.
3. Text mixed with simple text formatting like # heading and italics.

When you open an .Rmd, you get a notebook interface where code and output are interleaved. You can run each code chunk by clicking the Run icon (it looks like a play button at the top of the chunk), or by pressing Cmd/Ctrl + Shift + Enter. RStudio executes the code and displays the results inline with the code:



To produce a complete report containing all text, code, and results, click “Knit” or press Cmd/Ctrl + Shift + K. You can also do this programmatically with `rmarkdown::render("1-example.Rmd")`. This will display the report in the viewer pane, and create a self-contained HTML file that you can share with others.

The screenshot shows the RStudio interface with two main panes. The left pane is the "Code" editor for the file "diamond-sizes.Rmd". It contains R code for reading diamonds data, filtering for carat <= 2.5, and creating a geom_freqpoly plot. The right pane is the "Viewer" pane, which displays the generated HTML report. The title is "Diamond sizes" (2016-08-25). The report states there are 53940 diamonds, with 126 larger than 2.5 carats. It includes a histogram of diamond carat distribution.

```

1 ---  
2 title: "Diamond sizes"  
3 date: 2016-08-25  
4 output: html_document  
5 ---  
7 `r setup, include = FALSE}  
8 library(ggplot2)  
9 library(dplyr)  
10 smaller <- diamonds %>%  
11   filter(carat <= 2.5)  
12 `r  
13 We have data about `r nrow(diamonds)` diamonds. Only  
14 `r nrow(diamonds) - nrow(smaller)` are larger than  
15 2.5 carats. The distribution of the remainder is shown  
16 below:  
17`r, echo = FALSE}  
18 smaller %>%  
19   ggplot(aes(carat)) +  
20     geom_freqpoly(binwidth = 0.01)  
21 `r  
22 8:17 | Chunk 1: setup | R Markdown
```

Console output:

```

~/Documents/r4ds/r4ds/rmarkdown/ ↵
platform: x86_64-apple-darwin13.4.0 (64-bit)

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

Natural language support but running in an English locale

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

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

When you knit the document, R Markdown sends the .Rmd file to knitr, <http://yihui.name/knitr/>, which executes all of the code chunks and creates a new markdown (.md) document which includes the code and its output. The markdown file generated by knitr is then processed by pandoc, <http://pandoc.org/>, which is responsible for creating the finished file. The advantage of this two step workflow is that you can create a very wide range of output formats, as you'll learn about in [R markdown formats](#).



To get started with your own `.Rmd` file, select ***File > New File > R Markdown...*** in the menubar. RStudio will launch a wizard that you can use to pre-populate your file with useful content that reminds you how the key features of R Markdown work.

The following sections dive into the three components of an R Markdown document in more details: the markdown text, the code chunks, and the YAML header.

27.2.1 Exercises

1. Create a new notebook using ***File > New File > R Notebook***. Read the instructions. Practice running the chunks. Verify that you can modify the code, re-run it, and see modified output.
2. Create a new R Markdown document with ***File > New File > R Markdown...*** Knit it by clicking the appropriate button. Knit it by using the appropriate keyboard short cut. Verify that you can modify the input and see the output update.
3. Compare and contrast the R notebook and R markdown files you created above. How are the outputs similar? How are they different? How are the inputs similar? How are they different? What happens if you copy the YAML header from one to the other?
4. Create one new R Markdown document for each of the three built-in formats: HTML, PDF and Word. Knit each of the three documents. How does the output differ? How does the input differ? (You may need to install LaTeX in order to build the PDF output — RStudio will prompt you if this is necessary.)

27.3 Text formatting with Markdown

Prose in `.Rmd` files is written in Markdown, a lightweight set of conventions for formatting plain text files. Markdown is designed to be easy to read and easy to write. It is also very easy to learn. **The guide below shows how to use Pandoc's Markdown, a slightly extended version of Markdown that R Markdown understands.**

Text formatting

```
*italic* or _italic_
**bold** __bold__
`code`
superscript^2^ and subscript~2~
```

Headings

```
# 1st Level Header
## 2nd Level Header
### 3rd Level Header
```

Lists

```
* Bulleted list item 1
```

```
* Item 2
  * Item 2a
  * Item 2b
1. Numbered list item 1
1. Item 2. The numbers are incremented automatically in the output.
```

Links and images

```
<http://example.com>
[linked phrase](http://example.com)
![optional caption text](path/to/img.png)
```

Tables

| First Header | Second Header |
|--------------|---------------|
| Content Cell | Content Cell |
| Content Cell | Content Cell |

The best way to learn these is simply to try them out. It will take a few days, but soon they will **become second nature, and you won't need to think about them. If you forget, you can get to a handy reference sheet with Help > Markdown Quick Reference.**

27.3.1 Exercises

1. Practice what you've learned by creating a brief CV. The title should be your name, and you should include headings for (at least) education or employment. Each of the sections should include a bulleted list of jobs/degrees. Highlight the year in bold.
2. Using the R Markdown quick reference, figure out how to:
 1. Add a footnote.
 2. Add a horizontal rule.
 3. Add a block quote
3. Copy and paste the contents of `diamond-sizes.Rmd` from <https://github.com/hadley/r4ds/tree/master/rmarkdown> in to a local R markdown document. Check that you can run it, then add text after the frequency polygon that describes its most striking features..

27.4 Code chunks

To run code inside an R Markdown document, you need to insert a chunk. There are three ways to do so:

1. The keyboard shortcut Cmd/Ctrl + Alt + I
2. **The “Insert” button icon in the editor toolbar.**
3. By manually typing the chunk delimiters ` `` { r } and `` `.

Obviously, I’d recommend you learn the keyboard shortcut. It will save you a lot of time in the long run!

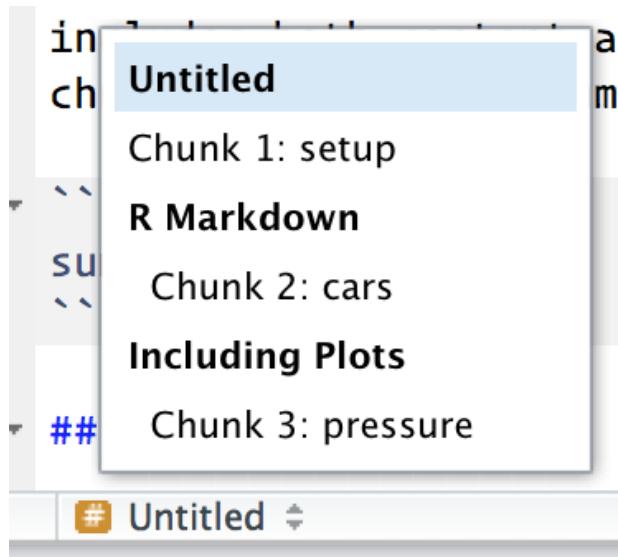
You can continue to run the code using the keyboard shortcut that by now (I hope!) you know and love: Cmd/Ctrl + Enter. However, chunks get a new keyboard shortcut: Cmd/Ctrl + Shift + Enter, which runs all the code in the chunk. Think of a chunk like a function. A chunk should be relatively self-contained, and focussed around a single task.

The following sections describe the chunk header which consists of `` ` { r , followed by an optional chunk name, followed by comma separated options, followed by } . Next comes your R code and the chunk end is indicated by a final `` ` .

27.4.1 Chunk name

Chunks can be given an optional name: `` ` { r by-name } . This has three advantages:

1. You can more easily navigate to specific chunks using the drop-down code navigator in the bottom-left of the script editor:



2. Graphics produced by the chunks will have useful names that make them easier to use elsewhere. More on that in [other important options](#).
3. You can set up networks of cached chunks to avoid re-performing expensive computations on every run. More on that below.

There is one chunk name that imbues special behaviour: `setup`. When you’re in a notebook mode, the chunk named `setup` will be run automatically once, before any other code is run.

27.4.2 Chunk options

Chunk output can be customised with options, arguments supplied to chunk header. Knitr provides almost 60 options that you can use to customize your code chunks. Here we'll cover the most important chunk options that you'll use frequently. You can see the full list at <http://yihui.name/knitr/options/>.

The most important set of options controls if your code block is executed and what results are inserted in the finished report:

- `eval = FALSE` prevents code from being evaluated. (And obviously if the code is not run, no results will be generated). This is useful for displaying example code, or for disabling a large block of code without commenting each line.
- `include = FALSE` runs the code, but doesn't show the code or results in the final document. Use this for setup code that you don't want cluttering your report.
- `echo = FALSE` prevents code, but not the results from appearing in the finished file. Use this when writing reports aimed at people who don't want to see the underlying R code.
- `message = FALSE` or `warning = FALSE` prevents messages or warnings from appearing in the finished file.
- `results = 'hide'` hides printed output; `fig.show = 'hide'` hides plots.
- `error = TRUE` causes the render to continue even if code returns an error. This is rarely something you'll want to include in the final version of your report, but can be very useful if you need to debug exactly what is going on inside your .Rmd. It's also useful if you're teaching R and want to deliberately include an error. The default, `error = FALSE` causes knitting to fail if there is a single error in the document.

The following table summarises which types of output each option suppresses:

| Option | Run code | Show code | Output | Plots | Messages | Warnings |
|--------------------------------|----------|-----------|--------|-------|----------|----------|
| <code>eval = FALSE</code> | - | | - | - | - | - |
| <code>include = FALSE</code> | | - | - | - | - | - |
| <code>echo = FALSE</code> | | - | | | | |
| <code>results = "hide"</code> | | | - | | | |
| <code>fig.show = "hide"</code> | | | | - | | |
| <code>message = FALSE</code> | | | | | - | |
| <code>warning = FALSE</code> | | | | | | - |

27.4.3 Table

By default, R Markdown prints data frames and matrices as you'd see them in the console:

```
mtcars[1:5, ]
#>          mpg cyl disp  hp drat    wt  qsec vs am gear carb
#> Mazda RX4     21.0   6 160 110 3.90 2.62 16.5  0  1     4     4
#> Mazda RX4 Wag  21.0   6 160 110 3.90 2.88 17.0  0  1     4     4
#> Datsun 710    22.8   4 108  93 3.85 2.32 18.6  1  1     4     1
```

```
#> Hornet 4 Drive     21.4   6 258 110 3.08 3.21 19.4 1 0 3 1
#> Hornet Sportabout 18.7   8 360 175 3.15 3.44 17.0 0 0 3 2
```

If you prefer that data be displayed with additional formatting you can use the `knitr::kable` function. The code below generates Table 27.1.

```
knitr::kable(
  mtcars[1:5, ],
  caption = "A knitr kable."
)
```

Table 27.1: A knitr kable.

| | mpg | cyl | disp | hp | drat | wt | qsec | vs | am | gear | carb |
|-------------------|------|-----|------|-----|------|------|------|----|----|------|------|
| Mazda RX4 | 21.0 | 6 | 160 | 110 | 3.90 | 2.62 | 16.5 | 0 | 1 | 4 | 4 |
| Mazda RX4 Wag | 21.0 | 6 | 160 | 110 | 3.90 | 2.88 | 17.0 | 0 | 1 | 4 | 4 |
| Datsun 710 | 22.8 | 4 | 108 | 93 | 3.85 | 2.32 | 18.6 | 1 | 1 | 4 | 1 |
| Hornet 4 Drive | 21.4 | 6 | 258 | 110 | 3.08 | 3.21 | 19.4 | 1 | 0 | 3 | 1 |
| Hornet Sportabout | 18.7 | 8 | 360 | 175 | 3.15 | 3.44 | 17.0 | 0 | 0 | 3 | 2 |

Read the documentation for `?knitr::kable` to see the other ways in which you can customise the table. For even deeper customisation, consider the `xtable`, `stargazer`, `pander`, `tables`, and `ascii` packages. Each provides a set of tools for returning formatted tables from R code.

There is also a rich set of options for controlling how figures are embedded. You'll learn about these in [saving your plots](#).

27.4.4 Caching

Normally, each knit of a document starts from a completely clean slate. This is great for **reproducibility, because it ensures that you've captured every important computation in code**. However, it can be painful if you have some computations that take a long time. The solution is `cache = TRUE`. When set, this will save the output of the chunk to a specially named file on disk. On subsequent runs, knitr will check to see if the code has changed, and if it hasn't, it will reuse the cached results.

The caching system must be used with care, because by default it is based on the code only, not its dependencies. For example, here the `processed_data` chunk depends on the `raw_data` chunk:

```
```{r raw_data}
rawdata <- readr::read_csv("a_very_large_file.csv")
```

```{r processed_data, cache = TRUE}
processed_data <- rawdata %>%
 filter(!is.na(import_var)) %>%
```

```
mutate(new_variable = complicated_transformation(x, y, z))
```
```

Caching the `processed_data` chunk means that it will get re-run if the dplyr pipeline is **changed, but it won't get rerun if the `read_csv()` call changes**. You can avoid that problem with the `dependson` chunk option:

```
```{r processed_data, cache = TRUE, dependson = "raw_data"}  
processed_data <- rawdata %>%
 filter(!is.na(import_var)) %>%
 mutate(new_variable = complicated_transformation(x, y, z))
```
```

`dependson` should contain a character vector of **every** chunk that the cached chunk depends on. Knitr will update the results for the cached chunk whenever it detects that one of its dependencies have changed.

Note that the chunks won't update if `a_very_large_file.csv` changes, because knitr caching only tracks changes within the `.Rmd` file. If you want to also track changes to that file you can use the `cache.extra` option. This is an arbitrary R expression that will invalidate the cache whenever it changes. A good function to use is `file.info()`: it returns a bunch of information about the file including when it was last modified. Then you can write:

```
```{r raw_data, cache.extra = file.info("a_very_large_file.csv")}  
rawdata <- readr::read_csv("a_very_large_file.csv")
```
```

As your caching strategies get progressively more complicated, it's a good idea to regularly clear out all your caches with `knitr::clean_cache()`.

I've used the advice of [David Robinson](#) to name these chunks: each chunk is named after the primary object that it creates. This makes it easier to understand the `dependson` specification.

27.4.5 Global options

As you work more with knitr, you will discover that some of the default **chunk options don't fit your needs** and you want to change them. You can do this by calling `knitr::opts_chunk$set()` in a code chunk. For example, when writing books and tutorials I set:

```
knitr::opts_chunk$set(  
  comment = "#>",  
  collapse = TRUE  
)
```

This uses my preferred comment formatting, and ensures that the code and output are kept closely entwined. On the other hand, if you were preparing a report, you might set:

```
knitr::opts_chunk$set(  
  echo = FALSE  
)
```

That will hide the code by default, so only showing the chunks you deliberately choose to show (with `echo = TRUE`). You might consider setting `message = FALSE` and `warning = FALSE`, **but that would make it harder to debug problems because you wouldn't see any messages in the final document.**

27.4.6 Inline code

There is one other way to embed R code into an R Markdown document: directly into the text, with: ``r``. This can be very useful if you mention properties of your data in the text. For example, in the example document I used at the start of the chapter I had:

We have data about ``r nrow(diamonds)`` diamonds. Only ``r nrow(diamonds) - nrow(smaller)`` are larger than 2.5 carats. The distribution of the remainder is shown below:

When the report is knit, the results of these computations are inserted into the text:

We have data about 53940 diamonds. Only 126 are larger than 2.5 carats. The distribution of the remainder is shown below:

When inserting numbers into text, `format()` is your friend. It allows you to set the number of digits **so you don't print to a ridiculous degree** of accuracy, and a `big.mark` to make numbers easier to read. I'll often combine these into a helper function:

```
comma <- function(x) format(x, digits = 2, big.mark = ",")  
comma(3452345)  
#> [1] "3,452,345"  
comma(.12358124331)  
#> [1] "0.12"
```

27.4.7 Exercises

1. Add a section that explores how diamond sizes vary by cut, colour, and clarity. Assume **you're writing a report for someone who doesn't know R, and instead of setting `echo = FALSE` on each chunk, set a global option.**
2. Download `diamond-sizes.Rmd` from <https://github.com/hadley/r4ds/tree/master/rmarkdown>. Add a section that describes the largest 20 diamonds, including a table that displays their most important attributes.
3. Modify `diamonds-sizes.Rmd` to use `comma()` to produce nicely formatted output. Also include the percentage of diamonds that are larger than 2.5 carats.
4. Set up a network of chunks where `a` depends on `c` and `b`, and both `b` and `c` depend on `a`. Have each chunk print `lubridate::now()`, set `cache = TRUE`, then verify your understanding of caching.

27.5 Troubleshooting

Troubleshooting R Markdown documents can be challenging because you are no longer in an interactive R environment, and you will need to learn some new tricks. The first thing you should **always try is to recreate the problem in an interactive session. Restart R, then “Run all chunks”**

(either from Code menu, under Run region), or with the keyboard shortcut Ctrl + Alt + R. If you're lucky, that will recreate the problem, and you can figure out what's going on interactively.

If that doesn't help, there must be something different between your interactive environment and the R markdown environment. You're going to need to systematically explore the options. The most common difference is the working directory: the working directory of an R Markdown is the directory in which it lives. Check the working directory is what you expect by including `getwd()` in a chunk.

Next, brainstorm all the things that might cause the bug. You'll need to systematically check that they're the same in your R session and your R markdown session. The easiest way to do that is to set `error = TRUE` on the chunk causing the problem, then use `print()` and `str()` to check that settings are as you expect.

27.6 YAML header

You can control many other “whole document” settings by tweaking the parameters of the YAML header. You might wonder what YAML stands for: it's “yet another markup language”, which is designed for representing hierarchical data in a way that's easy for humans to read and write. R Markdown uses it to control many details of the output. Here we'll discuss two: document parameters and bibliographies.

27.6.1 Parameters

R Markdown documents can include one or more parameters whose values can be set when you render the report. Parameters are useful when you want to re-render the same report with distinct values for various key inputs. For example, you might be producing sales reports per branch, exam results by student, or demographic summaries by country. To declare one or more parameters, use the `params` field.

This example use a `my_class` parameter to determines which class of cars to display:

```
---
```

```
output: html_document
```

```
params:
```

```
  my_class: "suv"
```

```
---
```

```
```{r setup, include = FALSE}
```

```
library(ggplot2)
```

```
library(dplyr)
```

```
class <- mpg %>% filter(class == params$my_class)
```

```
```
```

```
# Fuel economy for `r params$my_class`s
```

```
```{r, message = FALSE}
```

```
ggplot(class, aes(displ, hwy)) +
```

```
 geom_point() +
```

```
 geom_smooth(se = FALSE)
```

~ ~ ~

As you can see, parameters are available within the code chunks as a read-only list named `params`.

You can write atomic vectors directly into the YAML header. You can also run arbitrary R expressions by prefacing the parameter value with `!r`. This is a good way to specify date/time parameters.

```
params:
 start: !r lubridate::ymd("2015-01-01")
 snapshot: !r lubridate::ymd_hms("2015-01-01 12:30:00")
```

In RStudio, you can click the “Knit with Parameters” option in the Knit dropdown menu to set parameters, render, and preview the report in a single user friendly step. You can customise the dialog by setting other options in the header. See [http://rmarkdown.rstudio.com/developer\\_parameterized\\_reports.html#parameter\\_user\\_interfaces](http://rmarkdown.rstudio.com/developer_parameterized_reports.html#parameter_user_interfaces) for more details.

Alternatively, if you need to produce many such paramterised reports, you can call `rmarkdown::render()` with a list of `params`:

```
rmarkdown::render("fuel-economy.Rmd", params = list(my_class = "suv"))
```

This is particularly powerful in conjunction with `purrr::pwalk()`. The following example creates a report for each value of `class` found in `mpg`. First we create a data frame that has one row for each class, giving the `filename` of the report and the `params`:

```
reports <- tibble(
 class = unique(mpg$class),
 filename = stringr::str_c("fuel-economy-", class, ".html"),
 params = purrr::map(class, ~ list(my_class = .))
)
reports
#> # A tibble: 7 × 3
#> class filename params
#> <chr> <chr> <list>
#> 1 compact fuel-economy-compact.html <list [1]>
#> 2 midsize fuel-economy-midsiz.html <list [1]>
#> 3 suv fuel-economy-suv.html <list [1]>
#> 4 2seater fuel-economy-2seater.html <list [1]>
#> 5 minivan fuel-economy-minivan.html <list [1]>
#> 6 pickup fuel-economy-pickup.html <list [1]>
#> # ... with 1 more rows
```

Then we match the column names to the argument names of `render()`, and use `purrr`'s parallel walk to call `render()` once for each row:

```
reports %>%
 select(output_file = filename, params) %>%
 purrr::pwalk(rmarkdown::render, input = "fuel-economy.Rmd")
```

