

FUNCTIONAL PROGRAMMING AND UNIT TESTING FOR DATA MUNGING WITH R



Bruno
Rodrigues

Functional programming and unit testing for data munging with R

Bruno Rodrigues

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

This version was published on 2017-12-27



Leanpub

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

© 2016 - 2017 Bruno Rodrigues

Tweet This Book!

Please help Bruno Rodrigues by spreading the word about this book on [Twitter](#)!

The suggested hashtag for this book is [#fput](#).

Find out what other people are saying about the book by clicking on this link to search for this hashtag on Twitter:

[#fput](#)

Contents

Chapter 1 Why this book?	1
1.1 Important notice	1
1.2 Motivation	1
1.3 Who am I?	2
1.4 Thanks	2
1.5 License	2
Chapter 2 Introduction	3
2.1 Getting R	3
2.2 A short overview of functional programming	3
2.3 A short overview of unit testing	4
2.4 General recommendations to follow this book	5
Chapter 3 Functional Programming	6
3.1 Introduction	6
3.2 Mapping and Reducing: the <i>base</i> way	10
3.3 Mapping and Reducing: the <i>purrr</i> way	17
3.4 Basic anonymous functions	19
3.5 Wrap-up	22
3.6 Exercises	23
Chapter 4 The tidyverse	25
4.1 Smoking is bad for you, but pipes are your friend	25
4.2 Getting data into R with <i>readr</i> , <i>readxl</i> , <i>haven</i> and what are <i>tibbles</i>	27
4.3 Transforming your data with <i>dplyr</i>	36
4.4 Tidy your data with <i>tidyr</i>	84
4.5 Functional programming with <i>purrr</i> and <i>purrrlyr</i>	84
4.6 Special packages for special kinds of data: <i>forcats</i> , <i>lubridate</i> , and <i>stringr</i>	93
4.7 Exercises	93
Chapter 5 Programming with the tidyverse	95
Chapter 6 Packages	100
6.1 Why you need your own packages in your life	100

CONTENTS

6.2 R packages: the basics	100
6.3 Writing documentation for your functions	102
6.4 Extra files inside your package and dependencies	104
6.5 Unit test your package	105
6.6 Checking the coverage of your unit tests with covr	108
6.7 Wrap-up	110
Chapter 7 Unit testing	111
7.1 Introduction	111
7.2 Unit testing with the testthat package	112
7.3 Actually running your tests	115
7.4 Wrap-up	117
7.5 Exercises	118
Chapter 8 Putting it all together: writing a package to work on data	119
8.1 Getting the data	119
8.2 Your first data munging package: prepareData	120

Chapter 1 Why this book?

1.1 Important notice

This book is still being written, some chapters are not finished yet, and there might be (there are) some typos. Don't hesitate to write to me if you notice something weird.

You can purchase a digital copy of this book at [leanpub](https://www.leanpub.com/fput)¹. The version on Leanpub will not always be up-to-date, I only update it when I made very big changes (new chapters, etc). But once this book will be finished, both version are going to be the same.

This book serves to show how functional programming and unit testing can be useful for the task of data munging. This book is not an in-depth guide to functional programming, nor unit testing with R. If you want to have an in-depth understanding of the concepts presented in these books, I can't but recommend Wickham (2014a), Wickham (2015) and Wickham and Grolemund (2016) enough. Here, I will only briefly present functional programming, unit testing and building your own R packages. Just enough to get you (hopefully) interested and going.

This book is not an introduction to R either. I will assume that you have intermediate knowledge of R.

1.2 Motivation

Functional programming has very nice features that make working on data sets much more pleasant. It is common that you have to repeat the same instructions over and over again for different data sets that look very similar (for example, same, or similar column names). Of course, it is possible to loop over these data sets and repeat a set of instructions that change these data sets. However, we will see why a functional programming approach is to be preferred.

Unit testing then allows you to make sure that the functions you want to apply to your data sets actually do what you really want them to do. Knowing and applying these two concepts together will make you hopefully a better data analyst. Then we will learn to develop our own packages; not with the goal of publishing them in CRAN, but with the goal of making programming more streamlined.

¹<https://www.leanpub.com/fput>

1.3 Who am I?

I use R daily at my current job, and discovered R some years ago while I was at the [University of Strasbourg](#)². I'm not an R developer, and don't have a CS background. Most, if not everything, that I know about R is self-taught. I hope however that you will find this book useful. You can [follow me on twitter](#)³ or check [my blog](#)⁴.

1.4 Thanks

I'd like to thank [Ross Ihaka](#)⁵ and [Robert Gentleman](#)⁶ for developing the R programming language. Many thanks to [Hadley Wickham](#)⁷ for all the wonderful packages he developed that make R much more pleasant to use. Thanks to [Yihui Yie](#)⁸ for bookdown without which this book would not exist (at least not in this very nice format).

Thanks to [Hans-Martin von Gaudecker](#)⁹ for introducing me to unit testing and writing elegant code. The PEP 8 style guidelines will forever remain etched in my brain.

Finally I have to thank my wife for putting up with my endless rants against people not using functional programming nor testing their code (or worse, using proprietary software!).

1.5 License

This book is licensed under the GNU Free Documentation License, version 1.3. A copy of the license is available on the repo, or you can read it [online](#)¹⁰.

References

Wickham, Hadley. 2014a. *Advanced R*. CRC Press.

Wickham, Hadley. 2015. *R Packages*. 1st ed. O'Reilly. <http://r-pkgs.had.co.nz/>.

Wickham, Hadley, and Garrett Grolmund. 2016. *R for Data Science*. 1st ed. O'Reilly. <http://r4ds.had.co.nz/>.

²<http://www.unistra.fr/index.php?id=accueil>

³<https://twitter.com/brodriguesco>

⁴<http://brodrigues.co>

⁵<https://www.stat.auckland.ac.nz/~ihaka/>

⁶[https://en.wikipedia.org/wiki/Robert_Gentleman_\(statistician\)](https://en.wikipedia.org/wiki/Robert_Gentleman_(statistician))

⁷<http://hadley.nz/>

⁸<http://yihui.name/>

⁹<https://www.iame.uni-bonn.de/people/hm-gaudecker>

¹⁰<https://www.gnu.org/licenses/fdl-1.3.txt>

Chapter 2 Introduction

2.1 Getting R

Since I'm assuming you have an intermediate level in R, you already should have R and Rstudio installed on your machine. However, you may lack some of the following packages that are needed to follow the examples in this book:

- `covr`: to check the coverage of your unit tests
- `dplyr`: to clean, transform, prepare data
- `lazyeval`: for lazy evaluation
- `lubridate`: makes working with dates easier
- `memoise`: makes your function remember intermediate results
- `purrr`: extends R's functional programming capabilities
- `readr`: provides alternative functions to `read.csv()` and such
- `roxygen2`: creates documentation files from comments
- `stringr`: makes working with characters easier
- `testthat`: the library we are going to use for unit testing
- `tibble`: provides a nice, cleaner alternative to `data.frame`
- `tidyr`: works hand in hand with `dplyr`

If you're missing some or all of these packages, install them. You'll notice that most, if not all, of these packages were authored or co-authored by Hadley Wickham, currently chief scientist at Rstudio, so you can install most of these packages by installing a single package called `tidyverse`:

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

The `tidyverse` package installs some other useful packages that we will not use, but you should check them out anyways!

2.2 A short overview of functional programming

What is functional programming? Wikipedia tells us the following:

In computer science, functional programming is a programming paradigm —a style of building the structure and elements of computer programs— that treats computation as the evaluation of mathematical functions and avoids changing state and mutable data. It is a declarative programming paradigm, which means programming is done with expressions or declarations instead of statements. In functional code, the output value of a function depends only on the arguments that are input to the function, so calling a function f twice with the same value for an argument x will produce the same result $f(x)$ each time. Eliminating side effects, i.e. changes in state that do not depend on the function inputs, can make it much easier to understand and predict the behavior of a program, which is one of the key motivations for the development of functional programming.

That's the first paragraph of the [Wikipedia page](https://en.wikipedia.org/wiki/Functional_programming)¹¹ and it's quite heavy already!

So let's try to decrypt what is said in this paragraph. Functional programming is a programming paradigm. You may have heard of object oriented programming, or imperative programming before. You actually probably program in an imperative way without knowing it. Imperative programming is usually how programming is taught at universities, and most people then keep on programming in this way, especially in applied sciences like applied econometrics. Usually, people that write code in an imperative way tend to write very long scripts that change the state of the program gradually. In the case of a statistician (I will use the word 'statistician' to mean any person that works with datasets. Be it an economist, biologist, data scientist, etc.) this usually means loading a dataset, doing whatever has to be done by writing each instruction in a file, then running everything. Sometimes this statistician has to save temporary datasets, and then write other scripts that do a series of computations on these temporary datasets and then not forget to delete said temporary datasets. Functional programming is different, in that you write functions that do one single task and then call these functions successively on your data set. These functions can be used for any other project, can be easily documented and tested (more on this below). Because each function performs a single task and is well documented, it is also easier to understand what the program is supposed to do. Comments in a thousand-lines file are actually not that much useful. The file is so long, that even when commented you simply cannot make any sense of what is going on. It is also easier to automate tasks and navigate through the code. Since one function does one single task, if you're looking for the line of code that creates variable X , just look in the function called `create_var_X()`, instead of CTRL-Fing around. 1000 lines long script. You can also be sure that your functions do not do anything else (basically, this is what is meant by "eliminating side effects") than the single task you gave them. You can *trust your functions*.

2.3 A short overview of unit testing

At the end of the last section I wrote that you can *trust your functions*. Is that true though? Functional programming can make your life easier, but it does not prevent you from introducing bugs in your

¹¹https://en.wikipedia.org/wiki/Functional_programming

code. However, what functional programming makes easily possible, is to very easily and effectively test your code thanks to unit testing. You probably already test your code, by hand. You write some loop that is supposed to sum the first 10 integers and then you try it out and check if, indeed, your loop returns 55. Because this is the correct result, you save your work and continue programming something else, and so on. Unit testing is this, but in an automated way. Instead of just trying things out in the interpreter, you write unit tests. You write code that actually checks your functions. You save this unit tests somewhere, and then re-run them whenever you make changes to your code. Even if you don't change some parts of your code, you re-run every unit test. Because you actually never know what may happen. Maybe changing a single line in one of your functions introduced some unforeseen consequences that breaks functionality some place else. When you change code, and *all* your unit tests still pass, then you can be confident that your code is correct (actually, don't be too confident, because maybe you didn't write enough unit tests to cover every case. But we will see how we can be sure there is enough *coverage*).

2.4 General recommendations to follow this book

You should follow the examples in this book as closely as possible. I advise you to use exactly the same names as I do, because I will sometimes refer to previous examples and if you use different names, you will then need to go back and waste time to change the names.

Chapter 3 Functional Programming

3.1 Introduction

3.1.1 Function definitions

As mentioned in the [functional programming overview](#) functional programming is one of the numerous ways to write code. In functional programming, you write functions that do the computations and then as the user, you call these functions to work for you.

You should be familiar with function definitions in R. For example, suppose you want to compute the square root of a number and want to do so using Newton's algorithm:

```
1 sqrt_newton <- function(a, init, eps = 0.01){
2   while(abs(init**2 - a) > eps){
3     init <- 1/2 *(init + a/init)
4   }
5   return(init)
6 }
```

You can then use this function to get the square root of a number:

```
1 sqrt_newton(16, 2)

1 ## [1] 4.00122
```

We are using a `while` loop inside the body. The *body* of a function are the instructions that define the function. You can get the body of a function with `body(some_func)` of the function. In *pure* functional programming languages, like Haskell, you don't have loops. How can you program without loops, you may ask? In functional programming, loops are replaced by recursion. Let's rewrite our little example above with recursion:

```
1 sqrt_newton_recur <- function(a, init, eps = 0.01){
2   if(abs(init**2 - a) < eps){
3     result <- init
4   } else {
5     init <- 1/2 * (init + a/init)
6     result <- sqrt_newton_recur(a, init, eps)
7   }
8   return(result)
9 }
```

```
1 sqrt_newton_recur(16, 2)
```

```
1 ## [1] 4.00122
```

R is not a pure functional programming language though, so we can still use loops (be it while or for loops) in the bodies of our functions. Actually, for R specifically, it is better, performance-wise, to use loops instead of recursion, because R is not tail-call optimized. I won't get into the details of what tail-call optimization is but just remember that if performance is important a loop will be faster. However, sometimes, it is easier to write a function using recursion. I personally tend to avoid loops if performance is not important, because I find that code that avoids loops is easier to read and debug. However, knowing that you have can use loops is reassuring. In the coming sections I will show you some built-in function that make it possible to avoid writing loops and that don't rely on recursion, so performance won't be penalized.

3.1.2 Properties of functions

Mathematical functions have a nice property: we always get the same output for a given input. This is called referential transparency and we should aim to write our R functions in such a way.

For example, the following function:

```
1 increment <- function(x){
2   return(x + 1)
3 }
```

Is a referential transparent function. We always get the same result for any x that we give to this function.

This:

```
1 increment(10)
```

```
1 ## [1] 11
```

will always produce 11.

However, this one:

```
1 increment_opaque <- function(x){  
2   return(x + spam)  
3 }
```

is not a referential transparent function, because its value depends on the global variable spam.

```
1 spam <- 1  
2  
3 increment_opaque(10)
```

```
1 ## [1] 11
```

will only produce 11 if spam = 1. But what if spam = 19?

```
1 spam <- 19  
2  
3 increment_opaque(10)
```

```
1 ## [1] 29
```

To make `increment_opaque()` a referential transparent function, it is enough to make spam an argument:

```
1 increment_not_opaque <- function(x, spam){  
2   return(x + spam)  
3 }
```

Now even if there is a global variable called spam, this will not influence our function:

```

1 spam <- 19
2
3 increment_not_opaque(10, 34)

```

```

1 ## [1] 44

```

This is because the variable `spam` defined in the body of the function is a local variable. It could have been called anything else, really. Avoiding opaque functions makes our life easier.

Another property that adepts of functional programming value is that functions should have no, or very limited, side-effects. This means that functions should not change the state of your program.

For example this function (which is not a referential transparent function):

```

1 count_iter <- 0
2
3 sqrt_newton_side_effect <- function(a, init, eps = 0.01){
4   while(abs(init**2 - a) > eps){
5     init <- 1/2 *(init + a/init)
6     count_iter <- count_iter + 1 # The "<-" symbol means that we assign the
7   }                             # RHS value in a variable in the global en\
8   vironment
9   return(init)
10 }

```

If you look in the environment pane, you will see that `count_iter` equals 0. Now call this function with the following arguments:

```

1 sqrt_newton_side_effect(16000, 2)

```

```

1 ## [1] 126.4911

```

```

1 print(count_iter)

```

```

1 ## [1] 9

```

If you check the value of `count_iter` now, you will see that it increased! This is a side effect, because the function changed something outside its scope. It changed a value in the global environment. In general, it is good practice to avoid side-effects. For example, we could make the above function not have any side effects like this:

```

1 sqrt_newton_count <- function(a, init, count_iter = 0, eps = 0.01){
2   while(abs(init**2 - a) > eps){
3     init <- 1/2 *(init + a/init)
4     count_iter <- count_iter + 1
5   }
6   return(c(init, count_iter))
7 }

```

Now, this function returns a list with two elements, the result, and the number of iterations it took to get the result:

```

1 sqrt_newton_count(16000, 2)

1 ## [1] 126.4911  9.0000

```

Writing to disk is also considered a side effect, because the function changes something (a file) outside its scope. But this cannot be avoided (and it's actually a good thing to have, functions that can write to disk) so just remember: try to avoid having functions changing variables in the global environment unless you have a very good reason of doing so.

Finally, another property of mathematical functions, is that they do one single thing. Functional programming purists also program their functions to do one single task. This has benefits, but can complicate things. The function we wrote previously does two things: it computes the square root of a number and also returns the number of iterations it took to compute the result. However, this is not a bad thing; the function is doing two tasks, but these tasks are related to each other and it makes sense to have them together. My piece of advice: avoid having functions that do too many *unrelated* things. This makes debugging harder.

In conclusion: you should strive for referential transparency, try to avoid side effects unless you have a good reason to have them and try to keep your functions short and do as little tasks as possible. This makes testing and debugging easier, as you will see.

3.2 Mapping and Reducing: the *base* way

No introduction to functional programming would be complete without some discussion about the functions `Map()` (and the associated `*apply()` family of functions) and `Reduce()`. `Map()` allows you to map your function to every element of a list of arguments and is easy to understand, while `Reduce()` (sometimes called `fold()` in other programming languages) *reduces* a list of values to a single value by successively applying a function. It's a bit harder to understand, but with some examples it will become clear soon enough. In this section we will focus on how to do things using base functions. In the next section we will take a look at the `purrr` package which extends R's functional programming capabilities tremendously.

3.2.1 Mapping with `Map()` and the `*apply()` family of functions

Now that we have our nice function that computes square roots using Newton's algorithm, we would like to compute the square root of every element in the following list:

```

1 numbers <- c(16, 25, 36, 49, 64, 81)
2
3 sqrt_newton(numbers, init = rep(1, 6), eps = rep(0.001, 6))

1 ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
2 ## and only the first element will be used
3
4 ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
5 ## and only the first element will be used
6
7 ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
8 ## and only the first element will be used
9
10 ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
11 ## and only the first element will be used
12
13 ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
14 ## and only the first element will be used
15
16 ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
17 ## and only the first element will be used
18
19 ## [1] 4.000001 5.000023 6.000253 7.001406 8.005148 9.014272

```

We get a whole bunch of nasty warning messages, but we do get the expected result. But you should not leave it like this. Who knows what may happen some time down the road, when you try to compose this function with another? Maybe you'll get an error and you won't understand why! Let's rewrite the function properly.

We get these warnings because the condition `(init^2 - a) > eps` does not make sense for vectors. Here, R tells the user that it only uses the first element and then does the computation anyways. I would prefer if R would stop the execution and print an error message. This would force the user to have to rewrite the function to explicitly take vectors into account. And there is a very simple way of doing it, by using the function `Map()`:


```
1 Map(sqrt_newton, numbers, init = 1)
```

```
1 ## [[1]]
2 ## [1] 4.000001
3 ##
4 ## [[2]]
5 ## [1] 5.000023
6 ##
7 ## [[3]]
8 ## [1] 6.000253
9 ##
10 ## [[4]]
11 ## [1] 7
12 ##
13 ## [[5]]
14 ## [1] 8.000002
15 ##
16 ## [[6]]
17 ## [1] 9.000011
```

`Map()` applies a function to every element of a list and returns a list.

We could then write a wrapper around `Map()`:

```
1 sqrt_newton_vec <- function(numbers, init, eps = 0.01){
2   return(Map(sqrt_newton, numbers, init, eps))
3 }
4
5 sqrt_newton_vec(numbers, 1)
```

```
1 ## [[1]]
2 ## [1] 4.000001
3 ##
4 ## [[2]]
5 ## [1] 5.000023
6 ##
7 ## [[3]]
8 ## [1] 6.000253
9 ##
10 ## [[4]]
11 ## [1] 7
```

```

12 ##
13 ## [[5]]
14 ## [1] 8.000002
15 ##
16 ## [[6]]
17 ## [1] 9.000011

```

As you can see, we can give a function as an argument to another function. This makes `Map()` a *higher-order function*. Higher-order functions are functions that take other functions as arguments and return either another function, or a value. This is another important concept in functional programming and encourages modularity. It makes your code easily reusable!

R has other higher-order functions that work like `Map()`, such as `apply()`, `lapply()`, `mapply()`, `sapply()`, `vapply()` and `tapply()`. Depending on what you want to do, you will have to use one or the other. `apply()` and `tapply()` are different from the other `*apply()` functions, because they work on arrays. You can apply a function on the rows or columns of an array, for example if you want a row-wise sum:

```

1 a <- cbind(c(1, 2, 3), c(4, 5, 6), c(7, 8, 9))
2 apply(a, 1, sum)

1 ## [1] 12 15 18

```

We could use `lapply()` instead of `Map()`:

```

1 lapply(numbers, sqrt_newton, init = 1)

1 ## [[1]]
2 ## [1] 4.000001
3 ##
4 ## [[2]]
5 ## [1] 5.000023
6 ##
7 ## [[3]]
8 ## [1] 6.000253
9 ##
10 ## [[4]]
11 ## [1] 7
12 ##
13 ## [[5]]

```

```

14 ## [1] 8.000002
15 ##
16 ## [[6]]
17 ## [1] 9.000011

```

or `sapply()`:

```

1  sapply(numbers, sqrt_newton, init = 1)

1  ## [1] 4.000001 5.000023 6.000253 7.000000 8.000002 9.000011

```

We could rewrite `sqrt_newton_vec()` with `sapply()` which would return a better looking result (a list of numbers instead of a list of lists):

```

1  sqrt_newton_vec <- function(numbers, init, eps = 0.01){
2    return(sapply(numbers, sqrt_newton, init, eps))
3  }
4
5  sqrt_newton_vec(numbers, 1)

1  ## [1] 4.000001 5.000023 6.000253 7.000000 8.000002 9.000011

```

`mapply()` is different from these two:

```

1  inits <- c(100, 20, 3212, 487, 5, 9888)
2  mapply(sqrt_newton, numbers, init = inits)

1  ## [1] 4.000284 5.000001 6.000003 7.000006 8.000129 9.000006

```

What happens here is that `sqrt_newton()` gets called with following arguments:

```

1  sqrt_newton(numbers[1], inits[1])

1  ## [1] 4.000284

```

```
1  sqrt_newton(numbers[2], inits[2])
```

```
1  ## [1] 5.000001
```

```
1  sqrt_newton(numbers[3], inits[3])
```

```
1  ## [1] 6.000003
```

```
1  sqrt_newton(numbers[4], inits[4])
```

```
1  ## [1] 7.000006
```

```
1  sqrt_newton(numbers[5], inits[5])
```

```
1  ## [1] 8.000129
```

```
1  sqrt_newton(numbers[6], inits[6])
```

```
1  ## [1] 9.000006
```

From the `Map()`'s documentation, we learn that:

```
1  `Map()` is wrapper to `mapply()` which does not attempt to simplify the result...
```

All this behaviour can be replicated using loops, but once you get the gist of these functions, you can write code that is shorter and easier to read and unlike in the case of recursion, without any loss in performance (but without any gains either).

3.2.2 Reduce()

`Reduce()` is another very useful higher-order function, especially if you want to avoid loops to make your code easier to read. In some programming languages, `Reduce()` is called `fold()`.

I think that the following example illustrates the power of `Reduce()` well:

```
1 Reduce(`+`, numbers, init = 0)
```

```
1 ## [1] 271
```

Can you guess what happens? `Reduce()` takes a function as an argument, here the function `+¹` and then does the following computation:

```
1 0 + numbers[1] + numbers[2] + numbers[3]...
```

It applies the user supplied function successively but has to start with something, so we give it the argument `init` also. This argument is actually optional, but I show it here because in some cases it might be useful to start the computations at another value than 0. This function generalizes functions that only take two arguments. If you were to write a function that returns the minimum between two numbers:

```
1 my_min <- function(a, b){
2   if(a < b){
3     return(a)
4   } else {
5     return(b)
6   }
7 }
```

You could use `Reduce()` to get the minimum of a list of numbers:

```
1 print(numbers)
```

```
1 ## [1] 16 25 36 49 64 81
```

```
1 Reduce(my_min, numbers)
```

```
1 ## [1] 16
```

Here we don't supply an `init` because there is no need for it. Of course R's built-in `min()` function works on a list of values. But `Reduce()` is a very powerful function that can make our life much easier and most importantly avoid writing clumsy loops.

3.3 Mapping and Reducing: the `purrr` way

Hadley Wickham developed a package called `purrr` which contains a lot of very useful functions. I will show some of them here, but in the next chapter, we are going to study `purrr` in greater depth. Also, take the time to read `purrr`'s documentation, and of course you can read more about `purrr` in Wickham and Grolemund (2016).

3.3.1 The `map*()` family of functions

In the previous section we saw how to map a function to each element of a list. Each version of an `*apply()` function has a different purpose, but it is not very easy to remember which one returns a list, which other one returns an atomic vector and so on. If you're working on data frames you can use `apply()` to sum (for example) over columns or rows, because you can specify which `MARGIN` you want to sum over. But you do not get a data frame back. In the `purrr` package, each of the functions that do mapping have a similar name. The first part of these functions' names all start with `map_` and the second part tells you what this function is going to output. For example, if you want doubles out, you would use `map_dbl()`. If you are working on data frames want a data frame back, you would use `map_df()`. These are much more intuitive and easier to remember and we're going to learn how to use them in the chapter about [The Tidyverse](#)¹². For now, let's just focus on the basic functions, `map()` and `reduce()` (and some variants of `reduce()`). To map a function to every element of a list, simply use `map()`:

```
1 library("purrr")
2 map(numbers, sqrt_newton, init = 1)

1 ## [[1]]
2 ## [1] 4.000001
3 ##
4 ## [[2]]
5 ## [1] 5.000023
6 ##
7 ## [[3]]
8 ## [1] 6.000253
9 ##
10 ## [[4]]
11 ## [1] 7
12 ##
13 ## [[5]]
14 ## [1] 8.000002
```

¹²tidyverse.html#tidyverse

```

15 ##
16 ## [[6]]
17 ## [1] 9.000011

```

Compared to `Map()`, the function and the list are given in the reverse order, but the result is the same, of course.

3.3.2 Reducing with `purrr`

In the `purrr` package, you can find two more functions for folding: `reduce()` and `reduce_right()`. The difference between `reduce()` and `reduce_right()` is pretty obvious: `reduce_right()` starts from the right!

```

1 a <- seq(1, 10)
2
3 reduce(a, `-`)

```

```

1 ## [1] -53

```

```

1 reduce_right(a, `-`)

```

```

1 ## [1] -35

```

For operations that are not commutative, this makes a difference. Other interesting folding functions are `accumulate()` and `accumulate_right()`:

```

1 a <- seq(1, 10)
2
3 accumulate(a, `-`)

```

```

1 ## [1] 1 -1 -4 -8 -13 -19 -26 -34 -43 -53

```

```

1 accumulate_right(a, `-`)

```

```
1 ## [1] -35 -34 -32 -29 -25 -20 -14 -7 1 10
```

These two functions keep the intermediary results.

3.4 Basic anonymous functions

One last very useful concept are anonymous functions. Suppose that you want to apply one of your own functions to a list of datasets. For instance, you want to have a histogram of a variable that is called the same accross a list of datasets. Maybe your datasets are yearly surveys and each year the survey was conducted is another .csv file. For illustration purposes, let us use the `mtcars` dataset with some minor changes:

```
1 data(mtcars)
2
3 mtcars2000 <- mtcars
4 mtcars2001 <- mtcars
5 mtcars2001$cyl <- mtcars2001$cyl+3
6 datasets <- list("mtcars2000" = mtcars2000,
7 "mtcars2001" = mtcars2001)
```

In the next chapters we will learn how to load a lot of datasets at once and store them in a list. So it is important to know how to work with datasets that are stored on lists. Now suppose you want to use `purrr::map()` to plot a histogram of variable `cyl` for each dataset that is contained in your list.

```
1 map(datasets, hist, cyl)

1 Error in hist.default(.x[[i]], ...) : 'x' must be numeric
```

Maybe try this:

```
1 map(datasets, hist(cyl))

1 Error in hist(cyl) : object 'cyl' not found
```

So how can we solve this issue? One way is to use an anonymous function. Anonymous functions are functions that get declared on the fly and do not have names. These are especially useful inside higher order functions such as `purrr::map()`:


```
1 map(datasets, (function(x) hist(x$cyl)))
```



```
1 ## $mtcars2000
2 ## $breaks
3 ## [1] 4.0 4.5 5.0 5.5 6.0 6.5 7.0 7.5 8.0
4 ##
5 ## $counts
6 ## [1] 11 0 0 7 0 0 0 14
7 ##
8 ## $density
9 ## [1] 0.6875 0.0000 0.0000 0.4375 0.0000 0.0000 0.0000 0.8750
10 ##
11 ## $mids
12 ## [1] 4.25 4.75 5.25 5.75 6.25 6.75 7.25 7.75
13 ##
14 ## $xname
15 ## [1] "x$cyl"
16 ##
17 ## $equidist
18 ## [1] TRUE
19 ##
20 ## attr("class")
21 ## [1] "histogram"
22 ##
23 ## $mtcars2001
24 ## $breaks
25 ## [1] 7.0 7.5 8.0 8.5 9.0 9.5 10.0 10.5 11.0
26 ##
27 ## $counts
28 ## [1] 11 0 0 7 0 0 0 14
29 ##
30 ## $density
31 ## [1] 0.6875 0.0000 0.0000 0.4375 0.0000 0.0000 0.0000 0.8750
32 ##
33 ## $mids
34 ## [1] 7.25 7.75 8.25 8.75 9.25 9.75 10.25 10.75
35 ##
36 ## $xname
37 ## [1] "x$cyl"
```

```

38 ##
39 ## $equidist
40 ## [1] TRUE
41 ##
42 ## attr(,"class")
43 ## [1] "histogram"

```

Here the function is enclosed between `()` and is not named. The function has a single argument `x`, which is supposed to be a dataset. We then plot a histogram of the variable `cyl` from this dataset. Then this function is mapped to every dataset contained in the list `datasets` that we created above.

We can also write anonymous functions that are more complex:

```

1 map2(
2   .x = datasets, .y = names(datasets),
3   (function(.x, .y) hist(.x$cyl, main=paste("Histogram of cyl in", .y)))
4 )

```



```

1 ## $mtcars2000
2 ## $breaks
3 ## [1] 4.0 4.5 5.0 5.5 6.0 6.5 7.0 7.5 8.0
4 ##
5 ## $counts
6 ## [1] 11 0 0 7 0 0 0 14
7 ##
8 ## $density
9 ## [1] 0.6875 0.0000 0.0000 0.4375 0.0000 0.0000 0.0000 0.8750
10 ##
11 ## $mids
12 ## [1] 4.25 4.75 5.25 5.75 6.25 6.75 7.25 7.75
13 ##
14 ## $xname
15 ## [1] ".x$cyl"
16 ##
17 ## $equidist
18 ## [1] TRUE
19 ##
20 ## attr(,"class")
21 ## [1] "histogram"

```

```

22 ##
23 ## $mtcars2001
24 ## $breaks
25 ## [1] 7.0 7.5 8.0 8.5 9.0 9.5 10.0 10.5 11.0
26 ##
27 ## $counts
28 ## [1] 11 0 0 7 0 0 0 14
29 ##
30 ## $density
31 ## [1] 0.6875 0.0000 0.0000 0.4375 0.0000 0.0000 0.0000 0.8750
32 ##
33 ## $mids
34 ## [1] 7.25 7.75 8.25 8.75 9.25 9.75 10.25 10.75
35 ##
36 ## $xname
37 ## [1] ".x$cyl"
38 ##
39 ## $equidist
40 ## [1] TRUE
41 ##
42 ## attr(,"class")
43 ## [1] "histogram"

```

Of course you could have defined the anonymous function as a regular function before using `map()`. But sometimes it is faster to simply use an anonymous function as long as it does not hurt clarity.

This is the end of the introduction to functional programming. Entire books have been written on the subject, such as the upcoming book by Khan (2017) or Lipovaca (2011). If you're curious about functional programming, you should read these books. For our purposes though, knowing how to write functions, and trying to make them referentially transparent as well as knowing about mapping and reducing is enough to get us going.

3.5 Wrap-up

- Make your functions referentially transparent.
- Avoid side effects (if possible).
- Make your functions do one thing (if possible).
- A function that takes another function as an argument is called an higher-order function. You can write your own higher-order functions and this is a way of having short and easily testable functions. Making these functions then work together is trivial and is what makes functional programming very powerful.

3.6 Exercises

For the following exercises, you will have to use any of the functions that we saw in this chapter. `Reduce()`, `Map()` or any function from the `*apply()` family of functions. Do not use loops! If you don't know how to solve these exercises wait for the next section, where we'll learn how to write unit tests. Writing unit tests before the functions they're supposed to test is called test-driven development and can help you write your functions.

1. Create a function that returns the factorial of a number using `Reduce()`. Remember: no recursion nor loops allowed!

```
1 my_fact(5)
2
3 [1] 120
```

1. Suppose you have a list of data set names. Create a function that removes “.csv” from each of these names. Start by creating a function that does so using `stri_split()` from the package `stringi` (you can also use `strsplit()` from base R). Below is an illustration of how it's supposed to work:

```
1 dataset_names <- c("dataset1.csv", "dataset2.csv", "dataset3.csv")
2
3 remove_csv(dataset_names)
4
5 [1] "dataset1" "dataset2" dataset3"
```

1. Create a function that takes a number `a`, and then returns either the sum of the numbers from 1 to this number that are divisible by another number `b` or the product of the numbers from 1 to this number that are divisible by `b`. Your function should be a higher-order function with the following arguments: `a` the number, `divisible_func` the function that checks whether a number is divisible by some number `b` and `reduce_op` the function that either sums or multiplies the numbers from 1 to `a` that are divisible by `b`.

```
1 reduce_some_numbers(a = 10, divisible_func = divisible, b = 2, reduce_op = `*`)
2
3 [1] 3840
```

References

Wickham, Hadley, and Garrett Grolmund. 2016. *R for Data Science*. 1st ed. O'Reilly. <http://r4ds.had.co.nz/>.

Khan, Aslam. 2017. *Grokking Functional Programming*. 1st ed. Manning Publications. <https://www.manning.com/books/grokking-functional-programming>.

Lipovaca, Miran. 2011. *Learn You a Haskell for Great Good!: A Beginner's Guide*. no starch press. <http://learnyouahaskell.com/>.

-
1. This is simply the `+` operator you're used to. Try this out: ``+(1, 5)` and you'll see `+` is a function like any other. You just have to write backticks around the plus symbol to make it work.¹³

¹³[fprog.html#fnref1](#)

Chapter 4 The tidyverse

The tidyverse is the name given to a certain number of packages, most of all (if not all?) developed by, or co-developed by, Hadley Wickham. There's a website that introduces them all: [The tidyverse](http://tidyverse.org/)¹⁴. In this chapter, we are going to learn about some functions of some of these packages. We already know a little bit about purrr; let's discover what these other packages have to offer!

However, before reading everything that follows, I'd suggest you watch Hadley Wickham's talk [Expressing yourself with R](https://www.youtube.com/watch?v=1POb5fx_m3I)¹⁵. R is a computer *language*, and as with any language we really are writing things that are supposed to be read and understood by others, not just the computer. Even if you're working alone, you owe it to your future self to write clean, easy to understand code. Using the tidyverse and adopting the principles presented in the talk will put you in the right mindset for everything that follows!

First of all, let's install the tidyverse packages. You can install them one by one, or you can install the tidyverse meta-package:

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

I suggest you do just that, as we're going to skim over all the packages. To start an analysis, we first have to import data into R.

4.1 Smoking is bad for you, but pipes are your friend

The title of this section might sound weird at first, but by the end of it, you'll get this (terrible) pun. You probably know the following painting by René Magritte, *La trahison des images*:



¹⁴<http://tidyverse.org/>

¹⁵https://www.youtube.com/watch?v=1POb5fx_m3I

It turns out there's an R package from the tidyverse that is called `magrittr`. What does this package do? It brings *pipes* to R. Pipes are a concept from the Unix operating system; if you're using a GNU+Linux distribution or macOS, you're basically using a *modern* unix. (That's an oversimplification, but I'm an economist by training, and outrageously oversimplifying things is what we do, deal with it.)

The idea of pipes is to take the output of a command, and *feed* it as the input of another command. The `magrittr` package brings pipes to R, by using the weird looking `%>%`. Try the following:

```
1 library(magrittr)
```

```
1 16 %>% sqrt
```

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

Super weird right? But you probably understand what happened; 16 got fed as the first argument of the function `sqrt()`. You can chain multiple functions:

```
1 16 %>% sqrt %>% `+`(18)
```

```
1 ## [1] 22
```

The output of 16 (16) got fed to `sqrt()`, and the output of `sqrt(16)` (4) got fed to `+(18)` (22). Without `%>%` you'd write the line just above like this:

```
1 sqrt(16) + 18
```

```
1 ## [1] 22
```

It might not be very clear right now why this is useful, but the `%>%` is probably one of the best things that R has, because when using packages from the tidyverse, you will naturally want to chain a lot of functions together. Without the `%>%` it would become messy very fast.

`%>%` is not the only pipe operator in `magrittr`. There's `%T%`, `%<>%` and `%%$`. All have their uses, but are basically shortcuts to some common tasks with `%>%` plus another function. Which means that you can live without them, and because of this, I will only discuss them briefly once we'll have learned about the other tidyverse packages.

4.2 Getting data into R with `readr`, `readxl`, `haven` and what are *tibbles*

You probably already know how to import data in R, but maybe you are not familiar with these packages. Using them is pretty straightforward, and I will only discuss `haven` a little bit more than `readr` or `readxl`. `readr` allows you to import `*.csv` files as well as other files in plain text. The functions included in `readr` are fairly straightforward, but there is an aspect that I really like about them: if they fail to read your data you can get a report of what went wrong with the `problems()` function. I suggest you read the [Data import chapter](#)¹⁶ of R for Data Science, to get to know `readr` better. But for our purposes, knowing the basic `read_csv()` function is enough.

`readxl` is very similar to `readr` but focuses on importing Excel sheets into R. Read more about it on the [tidyverse website](#)¹⁷.

`haven` imports data from STATA, SAS and SPSS. I'm going into a bit more detail here, by showing an example with a STATA file. STATA files are usually labelled, and I'd like to show how to work with these labels using R. We're going to work with the `mtcars` dataset. I used STATA 14 to label the variables; so the dataset looks like one you could have to work with one day.

```
1 mtcars_stata <- haven::read_dta("images/mtcars.dta")
2 head(mtcars_stata)
```



```
1 ## # A tibble: 6 x 12
2 ##           car      mpg   cyl  disp    hp  drat    wt  qsec    vs    am
3 ##         <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
4 ## 1      Mazda RX4  21.0     6   160   110   3.90  2.620  16.46     0     1
5 ## 2    Mazda RX4 Wag  21.0     6   160   110   3.90  2.875  17.02     0     1
6 ## 3    Datsun 710    22.8     4   108    93   3.85  2.320  18.61     1     1
7 ## 4   Hornet 4 Drive  21.4     6   258   110   3.08  3.215  19.44     1     0
8 ## 5 Hornet Sportabout 18.7     8   360   175   3.15  3.440  17.02     0     0
9 ## 6      Valiant    18.1     6   225   105   2.76  3.460  20.22     1     0
10 ## # ... with 2 more variables: gear <dbl>, carb <dbl>
```

You don't see it here, but the columns are labelled. Try the following:

```
1 str(mtcars_stata$car)
```

¹⁶<http://r4ds.had.co.nz/data-import.html>

¹⁷<http://readxl.tidyverse.org/>


```

1 ## atomic [1:32] Mazda RX4 Mazda RX4 Wag Datsun 710 Hornet 4 Drive ...
2 ## - attr(*, "label")= chr "Make and model of the car"
3 ## - attr(*, "format.stata")= chr "%19s"

```

As you can see, the car column has the `label` attribute, which equals “Make and model of the car”. The other columns are also labelled:

```

1 str(mtcars_stata$cyl)

```

```

1 ## atomic [1:32] 6 6 4 6 8 6 8 4 4 6 ...
2 ## - attr(*, "label")= chr "Number of cylinders"
3 ## - attr(*, "format.stata")= chr "%8.0g"

```

```

1 str(mtcars_stata$am)

```

```

1 ## atomic [1:32] 1 1 1 0 0 0 0 0 0 0 ...
2 ## - attr(*, "label")= chr "Transmission (0 = automatic, 1 = manual)"
3 ## - attr(*, "format.stata")= chr "%8.0g"

```

Another way to get the label is to use the `attr()` function:

```

1 attr(mtcars_stata$cyl, "label")

```

```

1 ## [1] "Number of cylinders"

```

Let’s use what we learned until now to get the labels of all the columns:

```

1 show_labels <- function(dataset){
2   map(dataset, function(col)(attr(col, "label")))
3 }
4
5 show_labels(mtcars_stata)

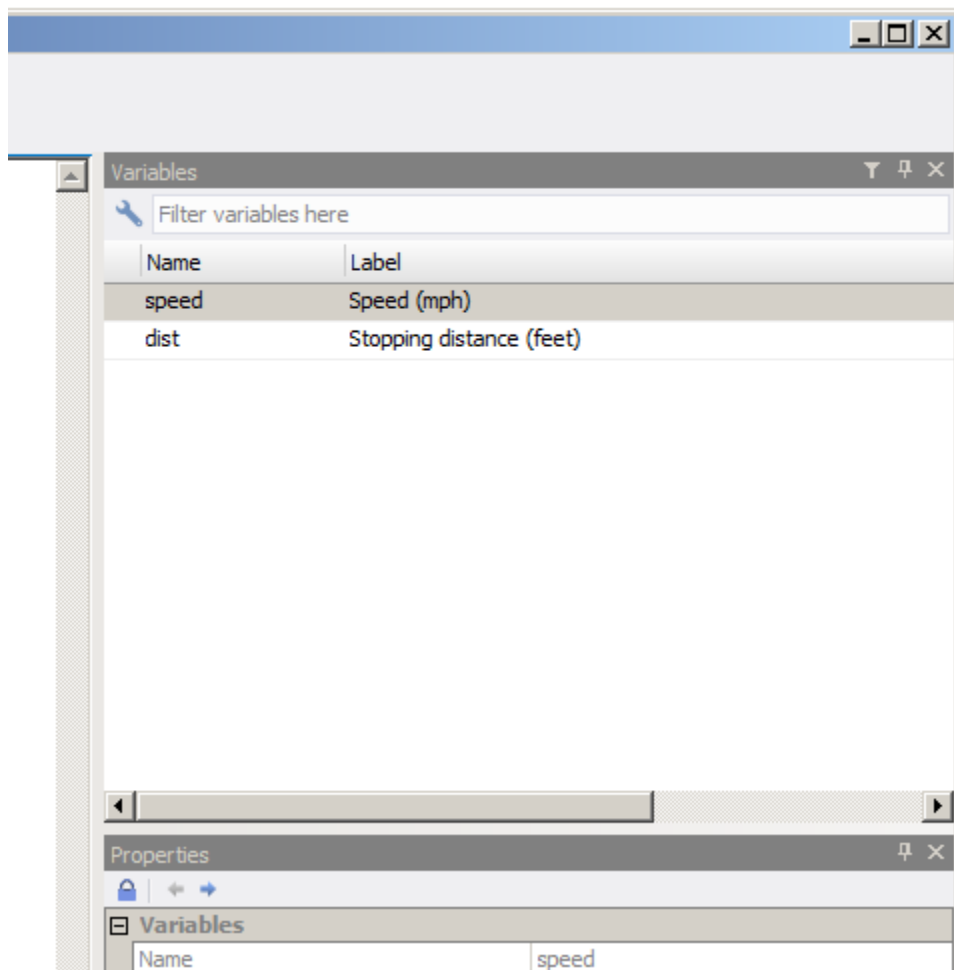
```

```
1  ## $car
2  ## [1] "Make and model of the car"
3  ##
4  ## $mpg
5  ## [1] "Miles/(US) gallon"
6  ##
7  ## $cyl
8  ## [1] "Number of cylinders"
9  ##
10 ## $disp
11 ## [1] "Displacement (cu.in.)"
12 ##
13 ## $hp
14 ## [1] "Gross horsepower"
15 ##
16 ## $drat
17 ## [1] "Rear axle ratio"
18 ##
19 ## $wt
20 ## [1] "Weight (1000 lbs)"
21 ##
22 ## $qsec
23 ## [1] "1/4 mile time"
24 ##
25 ## $vs
26 ## [1] "V/S"
27 ##
28 ## $am
29 ## [1] "Transmission (0 = automatic, 1 = manual)"
30 ##
31 ## $gear
32 ## [1] "Number of forward gears"
33 ##
34 ## $carb
35 ## [1] "Number of carburetors"
```

Could we label any dataset and then export it to a .dta file and have the labels in STATA? Let's find out with the cars dataset:

```
1 data(cars)
2
3 attr(cars$speed, "label") <- "Speed (mph)"
4 attr(cars$dist, "label") <- "Stopping distance (feet)"
5
6 haven::write_dta(cars, "images/cars.dta")
```

Below you see that cars.dta file opened in STATA:



When you use any of the discussed packages to import data, the resulting object is a `tibble`. `tibbles` are modern day “`data.frame`”s. The first thing you might have noticed is when you print a `tibble` vs a ‘`data.frame`’:

```

1 data(mtcars)
2
3 print(mtcars)

```

```

1 ##           mpg  cyl  disp  hp drat   wt  qsec vs am gear carb
2 ## Mazda RX4      21.0    6 160.0 110 3.90 2.620 16.46  0  1    4    4
3 ## Mazda RX4 Wag  21.0    6 160.0 110 3.90 2.875 17.02  0  1    4    4
4 ## Datsun 710     22.8    4 108.0  93 3.85 2.320 18.61  1  1    4    1
5 ## Hornet 4 Drive  21.4    6 258.0 110 3.08 3.215 19.44  1  0    3    1
6 ## Hornet Sportabout 18.7    8 360.0 175 3.15 3.440 17.02  0  0    3    2
7 ## Valiant        18.1    6 225.0 105 2.76 3.460 20.22  1  0    3    1
8 ## Duster 360     14.3    8 360.0 245 3.21 3.570 15.84  0  0    3    4
9 ## Merc 240D      24.4    4 146.7  62 3.69 3.190 20.00  1  0    4    2
10 ## Merc 230       22.8    4 140.8  95 3.92 3.150 22.90  1  0    4    2
11 ## Merc 280       19.2    6 167.6 123 3.92 3.440 18.30  1  0    4    4
12 ## Merc 280C      17.8    6 167.6 123 3.92 3.440 18.90  1  0    4    4
13 ## Merc 450SE     16.4    8 275.8 180 3.07 4.070 17.40  0  0    3    3
14 ## Merc 450SL     17.3    8 275.8 180 3.07 3.730 17.60  0  0    3    3
15 ## Merc 450SLC    15.2    8 275.8 180 3.07 3.780 18.00  0  0    3    3
16 ## Cadillac Fleetwood 10.4    8 472.0 205 2.93 5.250 17.98  0  0    3    4
17 ## Lincoln Continental 10.4    8 460.0 215 3.00 5.424 17.82  0  0    3    4
18 ## Chrysler Imperial 14.7    8 440.0 230 3.23 5.345 17.42  0  0    3    4
19 ## Fiat 128       32.4    4  78.7  66 4.08 2.200 19.47  1  1    4    1
20 ## Honda Civic    30.4    4  75.7  52 4.93 1.615 18.52  1  1    4    2
21 ## Toyota Corolla 33.9    4  71.1  65 4.22 1.835 19.90  1  1    4    1
22 ## Toyota Corona  21.5    4 120.1  97 3.70 2.465 20.01  1  0    3    1
23 ## Dodge Challenger 15.5    8 318.0 150 2.76 3.520 16.87  0  0    3    2
24 ## AMC Javelin    15.2    8 304.0 150 3.15 3.435 17.30  0  0    3    2
25 ## Camaro Z28     13.3    8 350.0 245 3.73 3.840 15.41  0  0    3    4
26 ## Pontiac Firebird 19.2    8 400.0 175 3.08 3.845 17.05  0  0    3    2
27 ## Fiat X1-9      27.3    4  79.0  66 4.08 1.935 18.90  1  1    4    1
28 ## Porsche 914-2  26.0    4 120.3  91 4.43 2.140 16.70  0  1    5    2
29 ## Lotus Europa   30.4    4  95.1 113 3.77 1.513 16.90  1  1    5    2
30 ## Ford Pantera L  15.8    8 351.0 264 4.22 3.170 14.50  0  1    5    4
31 ## Ferrari Dino   19.7    6 145.0 175 3.62 2.770 15.50  0  1    5    6
32 ## Maserati Bora   15.0    8 301.0 335 3.54 3.570 14.60  0  1    5    8
33 ## Volvo 142E     21.4    4 121.0 109 4.11 2.780 18.60  1  1    4    2

```

```

1 print(mtcars_stata)

```

```

1 ## # A tibble: 32 x 12
2 ##           car      mpg    cyl  disp    hp  drat    wt   qsec    vs    am
3 ##           <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
4 ## 1      Mazda RX4  21.0     6 160.0   110  3.90 2.620 16.46     0     1
5 ## 2    Mazda RX4 Wag  21.0     6 160.0   110  3.90 2.875 17.02     0     1
6 ## 3    Datsun 710    22.8     4 108.0    93  3.85 2.320 18.61     1     1
7 ## 4  Hornet 4 Drive  21.4     6 258.0   110  3.08 3.215 19.44     1     0
8 ## 5 Hornet Sportabout 18.7     8 360.0   175  3.15 3.440 17.02     0     0
9 ## 6      Valiant    18.1     6 225.0   105  2.76 3.460 20.22     1     0
10 ## 7     Duster 360   14.3     8 360.0   245  3.21 3.570 15.84     0     0
11 ## 8     Merc 240D   24.4     4 146.7    62  3.69 3.190 20.00     1     0
12 ## 9       Merc 230  22.8     4 140.8    95  3.92 3.150 22.90     1     0
13 ## 10      Merc 280  19.2     6 167.6   123  3.92 3.440 18.30     1     0
14 ## # ... with 22 more rows, and 2 more variables: gear <dbl>, carb <dbl>

```

Only the first 10 lines of the tibble get printed, but the number of remaining lines and the names of the columns that didn't find are shown as well as the types of the columns.

You can easily create a tibble from vectors:

```

1 library(tibble)
2
3 set.seed(123)
4 example <- tibble(a = seq(1,5), b = rnorm(5), c = rpois(5, 3))
5
6 print(example)

```

```

1 ## # A tibble: 5 x 3
2 ##       a         b         c
3 ##   <int>     <dbl> <int>
4 ## 1     1 -0.56047565     6
5 ## 2     2 -0.23017749     3
6 ## 3     3  1.55870831     4
7 ## 4     4  0.07050839     3
8 ## 5     5  0.12928774     1

```

Even better than `print()`, there's `glimpse()`:

```

1 glimpse(mtcars_stata)

```

```

1  ## Observations: 32
2  ## Variables: 12
3  ## $ car   <chr> "Mazda RX4", "Mazda RX4 Wag", "Datsun 710", "Hornet 4 Dri...
4  ## $ mpg   <dbl> 21.0, 21.0, 22.8, 21.4, 18.7, 18.1, 14.3, 24.4, 22.8, 19....
5  ## $ cyl   <dbl> 6, 6, 4, 6, 8, 6, 8, 4, 4, 6, 6, 8, 8, 8, 8, 8, 8, 4, 4, ...
6  ## $ disp  <dbl> 160.0, 160.0, 108.0, 258.0, 360.0, 225.0, 360.0, 146.7, 1...
7  ## $ hp    <dbl> 110, 110, 93, 110, 175, 105, 245, 62, 95, 123, 123, 180, ...
8  ## $ drat  <dbl> 3.90, 3.90, 3.85, 3.08, 3.15, 2.76, 3.21, 3.69, 3.92, 3.9...
9  ## $ wt    <dbl> 2.620, 2.875, 2.320, 3.215, 3.440, 3.460, 3.570, 3.190, 3...
10 ## $ qsec  <dbl> 16.46, 17.02, 18.61, 19.44, 17.02, 20.22, 15.84, 20.00, 2...
11 ## $ vs    <dbl> 0, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, ...
12 ## $ am    <dbl> 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, ...
13 ## $ gear  <dbl> 4, 4, 4, 3, 3, 3, 3, 4, 4, 4, 4, 3, 3, 3, 3, 3, 4, 4, ...
14 ## $ carb  <dbl> 4, 4, 1, 1, 2, 1, 4, 2, 2, 4, 4, 3, 3, 3, 4, 4, 4, 1, 2, ...

```

tibbles are lazy, which means that something like this is valid:

```

1  set.seed(123)
2  example <- tibble(a = seq(1,5), b = rnorm(5), c = 10 * b)
3
4  glimpse(example)

```

```

1  ## Observations: 5
2  ## Variables: 3
3  ## $ a <int> 1, 2, 3, 4, 5
4  ## $ b <dbl> -0.56047565, -0.23017749, 1.55870831, 0.07050839, 0.12928774
5  ## $ c <dbl> -5.6047565, -2.3017749, 15.5870831, 0.7050839, 1.2928774

```

The tibble package contains some other useful functions, such as `tribble()`, which allows you to create a tibble row by row:

```

1  set.seed(123)
2
3  example <- tribble(
4    ~a, ~b, ~c,
5    1, 2, "spam",
6    3, 4, "eggs",
7    5, 6, "bacon"
8  )
9
10 glimpse(example)

```

```

1 ## Observations: 3
2 ## Variables: 3
3 ## $ a <dbl> 1, 3, 5
4 ## $ b <dbl> 2, 4, 6
5 ## $ c <chr> "spam", "eggs", "bacon"

```

Another thing I find very useful is the following:

```

1 mtcars$m

1 ## [1] 21.0 21.0 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 17.8 16.4 17.3 15.2
2 ## [15] 10.4 10.4 14.7 32.4 30.4 33.9 21.5 15.5 15.2 13.3 19.2 27.3 26.0 30.4
3 ## [29] 15.8 19.7 15.0 21.4

```

```

1 mtcars_stata$m

```

```

1 ## Warning: Unknown or uninitialised column: 'm'.
2
3 ## NULL

```

`mtcars$m` shows the `mpg` column for some reason. There might be a good reason for this, but I prefer tibble's behaviour of notifying the user that this column does not exist.

It is possible to convert a lot of objects into tibbles:

```

1 example <- matrix(rnorm(36), nrow = 6)
2
3 as_tibble(example)

```

```

1 ## # A tibble: 6 x 6
2 ##           V1           V2           V3           V4           V5           V6
3 ##       <dbl>       <dbl>       <dbl>       <dbl>       <dbl>       <dbl>
4 ## 1 -0.56047565  0.4609162  0.4007715  0.7013559 -0.6250393  0.4264642
5 ## 2 -0.23017749 -1.2650612  0.1106827 -0.4727914 -1.6866933 -0.2950715
6 ## 3  1.55870831 -0.6868529 -0.5558411 -1.0678237  0.8377870  0.8951257
7 ## 4  0.07050839 -0.4456620  1.7869131 -0.2179749  0.1533731  0.8781335
8 ## 5  0.12928774  1.2240818  0.4978505 -1.0260044 -1.1381369  0.8215811
9 ## 6  1.71506499  0.3598138 -1.9666172 -0.7288912  1.2538149  0.6886403

```

```

1 example_df <- as.data.frame(example)
2
3 as_tibble(example_df)

```

```

1 ## # A tibble: 6 x 6
2 ##           V1           V2           V3           V4           V5           V6
3 ##       <dbl>       <dbl>       <dbl>       <dbl>       <dbl>       <dbl>
4 ## 1 -0.56047565  0.4609162  0.4007715  0.7013559 -0.6250393  0.4264642
5 ## 2 -0.23017749 -1.2650612  0.1106827 -0.4727914 -1.6866933 -0.2950715
6 ## 3  1.55870831 -0.6868529 -0.5558411 -1.0678237  0.8377870  0.8951257
7 ## 4  0.07050839 -0.4456620  1.7869131 -0.2179749  0.1533731  0.8781335
8 ## 5  0.12928774  1.2240818  0.4978505 -1.0260044 -1.1381369  0.8215811
9 ## 6  1.71506499  0.3598138 -1.9666172 -0.7288912  1.2538149  0.6886403

```

```

1 example_list <- list(a = seq(1,5), b = seq(6, 10))
2
3 as_tibble(example_list)

```

```

1 ## # A tibble: 5 x 2
2 ##       a     b
3 ##   <int> <int>
4 ## 1     1     6
5 ## 2     2     7
6 ## 3     3     8
7 ## 4     4     9
8 ## 5     5    10

```

You can also convert named vectors to tibbles with `enframe`:


```

1 recipe <- c("spam" = 1, "eggs" = 3, "bacon" = 10)
2
3 enframe(recipe, "ingredients", "quantity")

```

```

1 ## # A tibble: 3 x 2
2 ##   ingredients quantity
3 ##   <chr>      <dbl>
4 ## 1      spam        1
5 ## 2      eggs        3
6 ## 3      bacon       10

```

Contrast this to `as_tibble()` or `as.data.frame()`:

```

1 as.data.frame(recipe)

```

```

1 ##      recipe
2 ## spam      1
3 ## eggs      3
4 ## bacon     10

```

```

1 as_tibble(recipe)

```

```

1 ## # A tibble: 3 x 1
2 ##   value
3 ## * <dbl>
4 ## 1     1
5 ## 2     3
6 ## 3    10

```

There are a lot of other functions in the `tibble` package that you might find useful. I suggest you take a look at all of them and see what you can integrate in your workflow!

4.3 Transforming your data with `dplyr`

You may have never heard of the tidyverse, but you most certainly heard about `dplyr` and `tidyr`. Both these packages are probably the most popular packages of the tidyverse. Even if you know these packages already, you might not be using some more advanced functions, I'm talking about the *scoped* version of the usual `dplyr verbs` (`dplyr verbs` is how Hadley Wickham refers to the functions included in the package: `group_by()`, `select()`, etc).

This is going to be long, so prepare some coffee, lock the door to your study, turn off your phone and buckle up.

4.3.1 `filter()` and friends

We're going to use the `Gasoline` dataset from the `plm` package, so install that first:

```
1 install.packages("plm")
```

Then load the required data:

```
1 data(Gasoline, package = "plm")
```

and load `dplyr`:

```
1 library(dplyr)
```

```
1 ##
2 ## Attaching package: 'dplyr'
3
4 ## The following objects are masked from 'package:stats':
5 ##
6 ##   filter, lag
7
8 ## The following objects are masked from 'package:base':
9 ##
10 ##   intersect, setdiff, setequal, union
```

This dataset gives the consumption of gasoline for 18 countries from 1960 to 1978. When you load the data like this, it is a standard `data.frame`. `dplyr` functions can be used on standard `data.frame` objects, but just because we learned about `tibble`'s, let's convert the data to a `tibble` and change its name:

```
1 gasoline <- as_tibble(Gasoline)
```

`filter()` is pretty straightforward. What if you would like to subset the data to focus on the year 1969? Simple:

```
1 filter(gasoline, year == 1969)
```

```

1 ## # A tibble: 18 x 6
2 ##   country year lgaspcar lincomep   lrpmg   lcarpcap
3 ##   <fctr> <int>   <dbl>   <dbl>   <dbl>   <dbl>
4 ## 1 AUSTRIA 1969 4.046355 -6.153140 -0.5591105 -8.788686
5 ## 2 BELGIUM 1969 3.854601 -5.857532 -0.3548085 -8.521453
6 ## 3 CANADA 1969 4.864433 -5.560853 -1.0368639 -8.095113
7 ## 4 DENMARK 1969 4.173561 -5.722769 -0.4068792 -8.470459
8 ## 5 FRANCE 1969 3.773460 -5.840774 -0.3151909 -8.369136
9 ## 6 GERMANY 1969 3.899185 -5.829641 -0.5892314 -8.438061
10 ## 7 GREECE 1969 4.894773 -6.591104 -0.1798700 -10.713848
11 ## 8 IRELAND 1969 4.208613 -6.379743 -0.2716284 -8.947265
12 ## 9 ITALY 1969 3.737389 -6.282857 -0.2475668 -8.666004
13 ## 10 JAPAN 1969 4.518290 -6.159308 -0.4168502 -9.607600
14 ## 11 NETHERLA 1969 3.987689 -5.880556 -0.4169496 -8.634102
15 ## 12 NORWAY 1969 4.086823 -5.735319 -0.3382305 -8.694593
16 ## 13 SPAIN 1969 3.994103 -5.601046 0.6694895 -9.720425
17 ## 14 SWEDEN 1969 3.991715 -7.771081 -2.7319041 -8.197462
18 ## 15 SWITZERL 1969 4.211290 -5.912172 -0.9181216 -8.473379
19 ## 16 TURKEY 1969 5.720705 -7.388646 -0.2984542 -12.518545
20 ## 17 U.K. 1969 3.948058 -6.031953 -0.3833246 -8.468119
21 ## 18 U.S.A. 1969 4.841383 -5.414374 -1.2231427 -7.792706

```

Remember the pipe operator, `%>%` from the start of this chapter? Here's how this would work with it:

```
1 gasoline %>% filter(year == 1969)
```

```

1 ## # A tibble: 18 x 6
2 ##   country year lgaspcar lincomep   lrpmg   lcarpcap
3 ##   <fctr> <int>   <dbl>   <dbl>   <dbl>   <dbl>
4 ## 1 AUSTRIA 1969 4.046355 -6.153140 -0.5591105 -8.788686
5 ## 2 BELGIUM 1969 3.854601 -5.857532 -0.3548085 -8.521453
6 ## 3 CANADA 1969 4.864433 -5.560853 -1.0368639 -8.095113
7 ## 4 DENMARK 1969 4.173561 -5.722769 -0.4068792 -8.470459
8 ## 5 FRANCE 1969 3.773460 -5.840774 -0.3151909 -8.369136
9 ## 6 GERMANY 1969 3.899185 -5.829641 -0.5892314 -8.438061
10 ## 7 GREECE 1969 4.894773 -6.591104 -0.1798700 -10.713848
11 ## 8 IRELAND 1969 4.208613 -6.379743 -0.2716284 -8.947265
12 ## 9 ITALY 1969 3.737389 -6.282857 -0.2475668 -8.666004
13 ## 10 JAPAN 1969 4.518290 -6.159308 -0.4168502 -9.607600
14 ## 11 NETHERLA 1969 3.987689 -5.880556 -0.4169496 -8.634102

```

```

15 ## 12 NORWAY 1969 4.086823 -5.735319 -0.3382305 -8.694593
16 ## 13 SPAIN 1969 3.994103 -5.601046 0.6694895 -9.720425
17 ## 14 SWEDEN 1969 3.991715 -7.771081 -2.7319041 -8.197462
18 ## 15 SWITZERL 1969 4.211290 -5.912172 -0.9181216 -8.473379
19 ## 16 TURKEY 1969 5.720705 -7.388646 -0.2984542 -12.518545
20 ## 17 U.K. 1969 3.948058 -6.031953 -0.3833246 -8.468119
21 ## 18 U.S.A. 1969 4.841383 -5.414374 -1.2231427 -7.792706

```

So gasoline, which is a tibble object, is passed as the first argument of the `filter()` function. Starting now, we're only going to use these pipes. You will see why soon enough, so bear with me.

You can also filter more than just one year, by using the `%in%` operator:

```
1 gasoline %>% filter(year %in% seq(1969, 1973))
```

```

1 ## # A tibble: 90 x 6
2 ##   country year lgaspcar lincomep      lrpmg lcarpcap
3 ##   <fctr> <int>   <dbl>   <dbl>   <dbl>   <dbl>
4 ## 1 AUSTRIA 1969 4.046355 -6.153140 -0.5591105 -8.788686
5 ## 2 AUSTRIA 1970 4.080888 -6.081712 -0.5965612 -8.728200
6 ## 3 AUSTRIA 1971 4.106720 -6.043626 -0.6544591 -8.635898
7 ## 4 AUSTRIA 1972 4.128018 -5.981052 -0.5963318 -8.538338
8 ## 5 AUSTRIA 1973 4.199381 -5.895153 -0.5944468 -8.487289
9 ## 6 BELGIUM 1969 3.854601 -5.857532 -0.3548085 -8.521453
10 ## 7 BELGIUM 1970 3.870392 -5.797201 -0.3779404 -8.453043
11 ## 8 BELGIUM 1971 3.872245 -5.761050 -0.3992299 -8.409457
12 ## 9 BELGIUM 1972 3.905402 -5.710230 -0.3106458 -8.362588
13 ## 10 BELGIUM 1973 3.895996 -5.644145 -0.3730919 -8.314447
14 ## # ... with 80 more rows

```

or even non-consecutive years:

```
1 gasoline %>% filter(year %in% c(1969, 1973, 1977))
```

```

1 ## # A tibble: 54 x 6
2 ##   country year lgaspcar lincomep   lrpmg   lcarpcap
3 ##   <fctr> <int>   <dbl>   <dbl>   <dbl>   <dbl>
4 ## 1 AUSTRIA 1969 4.046355 -6.153140 -0.5591105 -8.788686
5 ## 2 AUSTRIA 1973 4.199381 -5.895153 -0.5944468 -8.487289
6 ## 3 AUSTRIA 1977 3.931676 -5.833288 -0.4219156 -8.249563
7 ## 4 BELGIUM 1969 3.854601 -5.857532 -0.3548085 -8.521453
8 ## 5 BELGIUM 1973 3.895996 -5.644145 -0.3730919 -8.314447
9 ## 6 BELGIUM 1977 3.854311 -5.556697 -0.4316413 -8.138534
10 ## 7 CANADA 1969 4.864433 -5.560853 -1.0368639 -8.095113
11 ## 8 CANADA 1973 4.899694 -5.414753 -1.1331614 -7.942140
12 ## 9 CANADA 1977 4.810992 -5.336967 -1.0708445 -7.768793
13 ## 10 DENMARK 1969 4.173561 -5.722769 -0.4068792 -8.470459
14 ## # ... with 44 more rows

```

`%in%` tests if an object is part of a set.

`filter()` is not the only *filtering* verb there is. Suppose that we have a condition that we want to use to filter out a lot of columns at once. For example, for every column that is of type `numeric`, keep only the lines where the condition `value > -8` is satisfied. The next line does that:

```
1 gasoline %>% filter_if( ~all(is.numeric(.)), all_vars(. > -8))
```

```

1 ## # A tibble: 30 x 6
2 ##   country year lgaspcar lincomep   lrpmg   lcarpcap
3 ##   <fctr> <int>   <dbl>   <dbl>   <dbl>   <dbl>
4 ## 1 CANADA 1972 4.889302 -5.436603 -1.0996670 -7.989531
5 ## 2 CANADA 1973 4.899694 -5.414753 -1.1331614 -7.942140
6 ## 3 CANADA 1974 4.891591 -5.418456 -1.1238000 -7.900758
7 ## 4 CANADA 1975 4.888471 -5.379097 -1.1856843 -7.873313
8 ## 5 CANADA 1976 4.837359 -5.361285 -1.0617966 -7.808425
9 ## 6 CANADA 1977 4.810992 -5.336967 -1.0708445 -7.768793
10 ## 7 CANADA 1978 4.855846 -5.311272 -1.0749507 -7.788061
11 ## 8 GERMANY 1978 3.883879 -5.561733 -0.6281728 -7.950079
12 ## 9 SWEDEN 1975 3.973840 -7.679557 -2.7673146 -7.994217
13 ## 10 SWEDEN 1976 3.983997 -7.672043 -2.8229448 -7.956066
14 ## # ... with 20 more rows

```

It's a bit more complicated than before. `filter_if()` needs 3 arguments to work; the data, a predicate function (a function that returns `TRUE`, or `FALSE`) which will select the columns we want to work on, and then the condition. The condition can be applied to *all* the columns that were selected by the

predicate function (hence the `all_vars()`) or only to at least one (you'd use `any_vars()` then). Try to change the condition, or the predicate function, to figure out how `filter_if()` works. The dot is a placeholder that stands for whatever columns were selected.

`filter_at()` works differently; it allows the user to filter columns by position:

```
1 gasoline %>% filter_at(vars(ends_with("p")), all_vars(. > -8))

2 ## # A tibble: 30 x 6
3 ##   country year lgaspcar lincomep      lrpmg lcarpcap
4 ##   <fctr> <int>   <dbl>   <dbl>   <dbl>   <dbl>
5 ## 1 CANADA  1972  4.889302 -5.436603 -1.0996670 -7.989531
6 ## 2 CANADA  1973  4.899694 -5.414753 -1.1331614 -7.942140
7 ## 3 CANADA  1974  4.891591 -5.418456 -1.1238000 -7.900758
8 ## 4 CANADA  1975  4.888471 -5.379097 -1.1856843 -7.873313
9 ## 5 CANADA  1976  4.837359 -5.361285 -1.0617966 -7.808425
10 ## 6 CANADA  1977  4.810992 -5.336967 -1.0708445 -7.768793
11 ## 7 CANADA  1978  4.855846 -5.311272 -1.0749507 -7.788061
12 ## 8 GERMANY 1978  3.883879 -5.561733 -0.6281728 -7.950079
13 ## 9 SWEDEN  1975  3.973840 -7.679557 -2.7673146 -7.994217
14 ## 10 SWEDEN 1976  3.983997 -7.672043 -2.8229448 -7.956066
15 ## # ... with 20 more rows
```

`end_with()` is a helper function that we are going to use a lot (as well as `starts_with()` and some others, you'll see.). So the above line means “for the columns whose name end with a ‘p’ only keep the lines where, for all the selected columns, the values are strictly superior to -8”. Again, this is not very easy the first time you deal with that, so play around with it for a bit.

`filter_all()`, as the name implies, considers all variables for the filtering step.

`filter_if()` and `filter_at()` are very useful when you have very large datasets with a lot of variables and you want to apply a filtering function only to a subset of them. `filter_all()` is useful if, for example, you only want to keep the positive values for all the columns.

4.3.2 `select()` and its helpers

While `filter()` and its scoped versions allow you to keep or discard rows of data, `select()` (and its scoped versions) allow you to keep or discard entire columns. To keep columns:

```
1 gasoline %>% select(country, year, lrpmg)
```

```

1  ## # A tibble: 342 x 3
2  ##   country year      lrpmpg
3  ## *   <fctr> <int>      <dbl>
4  ## 1 AUSTRIA  1960 -0.3345476
5  ## 2 AUSTRIA  1961 -0.3513276
6  ## 3 AUSTRIA  1962 -0.3795177
7  ## 4 AUSTRIA  1963 -0.4142514
8  ## 5 AUSTRIA  1964 -0.4453354
9  ## 6 AUSTRIA  1965 -0.4970607
10 ## 7 AUSTRIA  1966 -0.4668377
11 ## 8 AUSTRIA  1967 -0.5058834
12 ## 9 AUSTRIA  1968 -0.5224125
13 ## 10 AUSTRIA 1969 -0.5591105
14 ## # ... with 332 more rows

```

To discard them:

```
1 gasoline %>% select(-country, -year, -lrpmpg)
```

```

1  ## # A tibble: 342 x 3
2  ##   lgaspcar lincomep lcarpcap
3  ## *   <dbl>      <dbl>      <dbl>
4  ## 1 4.173244 -6.474277 -9.766840
5  ## 2 4.100989 -6.426006 -9.608622
6  ## 3 4.073177 -6.407308 -9.457257
7  ## 4 4.059509 -6.370679 -9.343155
8  ## 5 4.037689 -6.322247 -9.237739
9  ## 6 4.033983 -6.294668 -9.123903
10 ## 7 4.047537 -6.252545 -9.019822
11 ## 8 4.052911 -6.234581 -8.934403
12 ## 9 4.045507 -6.206894 -8.847967
13 ## 10 4.046355 -6.153140 -8.788686
14 ## # ... with 332 more rows

```

To rename them:

```
1 gasoline %>% select(country, date = year, lrpmpg)
```

```

1 ## # A tibble: 342 x 3
2 ##   country date      lrpmg
3 ## *   <fctr> <int>      <dbl>
4 ## 1 AUSTRIA 1960 -0.3345476
5 ## 2 AUSTRIA 1961 -0.3513276
6 ## 3 AUSTRIA 1962 -0.3795177
7 ## 4 AUSTRIA 1963 -0.4142514
8 ## 5 AUSTRIA 1964 -0.4453354
9 ## 6 AUSTRIA 1965 -0.4970607
10 ## 7 AUSTRIA 1966 -0.4668377
11 ## 8 AUSTRIA 1967 -0.5058834
12 ## 9 AUSTRIA 1968 -0.5224125
13 ## 10 AUSTRIA 1969 -0.5591105
14 ## # ... with 332 more rows

```

There's also `rename()`, but it works a bit differently:

```

1 gasoline %>% rename(date = year)

1 ## # A tibble: 342 x 6
2 ##   country date lgaspcar lincomep      lrpmg lcarpcap
3 ## *   <fctr> <int>      <dbl>      <dbl>      <dbl>      <dbl>
4 ## 1 AUSTRIA 1960 4.173244 -6.474277 -0.3345476 -9.766840
5 ## 2 AUSTRIA 1961 4.100989 -6.426006 -0.3513276 -9.608622
6 ## 3 AUSTRIA 1962 4.073177 -6.407308 -0.3795177 -9.457257
7 ## 4 AUSTRIA 1963 4.059509 -6.370679 -0.4142514 -9.343155
8 ## 5 AUSTRIA 1964 4.037689 -6.322247 -0.4453354 -9.237739
9 ## 6 AUSTRIA 1965 4.033983 -6.294668 -0.4970607 -9.123903
10 ## 7 AUSTRIA 1966 4.047537 -6.252545 -0.4668377 -9.019822
11 ## 8 AUSTRIA 1967 4.052911 -6.234581 -0.5058834 -8.934403
12 ## 9 AUSTRIA 1968 4.045507 -6.206894 -0.5224125 -8.847967
13 ## 10 AUSTRIA 1969 4.046355 -6.153140 -0.5591105 -8.788686
14 ## # ... with 332 more rows

```

`rename()` does not do any kind of selection, but just renames.

To re-order them:

```

1 gasoline %>% select(year, country, lrpmg, everything())

```



```

1 ## # A tibble: 342 x 6
2 ##   year country      lrpmpg lgaspcar  lincomep  lcarpcap
3 ## * <int>  <fctr>      <dbl>    <dbl>    <dbl>    <dbl>
4 ## 1  1960 AUSTRIA -0.3345476  4.173244 -6.474277 -9.766840
5 ## 2  1961 AUSTRIA -0.3513276  4.100989 -6.426006 -9.608622
6 ## 3  1962 AUSTRIA -0.3795177  4.073177 -6.407308 -9.457257
7 ## 4  1963 AUSTRIA -0.4142514  4.059509 -6.370679 -9.343155
8 ## 5  1964 AUSTRIA -0.4453354  4.037689 -6.322247 -9.237739
9 ## 6  1965 AUSTRIA -0.4970607  4.033983 -6.294668 -9.123903
10 ## 7  1966 AUSTRIA -0.4668377  4.047537 -6.252545 -9.019822
11 ## 8  1967 AUSTRIA -0.5058834  4.052911 -6.234581 -8.934403
12 ## 9  1968 AUSTRIA -0.5224125  4.045507 -6.206894 -8.847967
13 ## 10 1969 AUSTRIA -0.5591105  4.046355 -6.153140 -8.788686
14 ## # ... with 332 more rows

```

everything() is another of those helper functions (like starts_with(), and ends_with()). What if we are only interested in columns whose name start with "l"?

```
1 gasoline %>% select(starts_with("l"))
```

```

1 ## # A tibble: 342 x 4
2 ##   lgaspcar  lincomep      lrpmpg  lcarpcap
3 ## * <dbl>    <dbl>      <dbl>    <dbl>
4 ## 1  4.173244 -6.474277 -0.3345476 -9.766840
5 ## 2  4.100989 -6.426006 -0.3513276 -9.608622
6 ## 3  4.073177 -6.407308 -0.3795177 -9.457257
7 ## 4  4.059509 -6.370679 -0.4142514 -9.343155
8 ## 5  4.037689 -6.322247 -0.4453354 -9.237739
9 ## 6  4.033983 -6.294668 -0.4970607 -9.123903
10 ## 7  4.047537 -6.252545 -0.4668377 -9.019822
11 ## 8  4.052911 -6.234581 -0.5058834 -8.934403
12 ## 9  4.045507 -6.206894 -0.5224125 -8.847967
13 ## 10 4.046355 -6.153140 -0.5591105 -8.788686
14 ## # ... with 332 more rows

```

The same can be achieved with select_at():

```
1 gasoline %>% select_at(vars(starts_with("l")))
```

```

1  ## # A tibble: 342 x 4
2  ##   lgaspcar  lincomep    lrpmg  lcarpcap
3  ##   *    <dbl>      <dbl>    <dbl>    <dbl>
4  ## 1 4.173244 -6.474277 -0.3345476 -9.766840
5  ## 2 4.100989 -6.426006 -0.3513276 -9.608622
6  ## 3 4.073177 -6.407308 -0.3795177 -9.457257
7  ## 4 4.059509 -6.370679 -0.4142514 -9.343155
8  ## 5 4.037689 -6.322247 -0.4453354 -9.237739
9  ## 6 4.033983 -6.294668 -0.4970607 -9.123903
10 ## 7 4.047537 -6.252545 -0.4668377 -9.019822
11 ## 8 4.052911 -6.234581 -0.5058834 -8.934403
12 ## 9 4.045507 -6.206894 -0.5224125 -8.847967
13 ## 10 4.046355 -6.153140 -0.5591105 -8.788686
14 ## # ... with 332 more rows

```

`select_at()` can be quite useful if you know the position of the columns you're interested in:

```
1 gasoline %>% select_at(vars(c(1,2,5)))
```

```

1  ## # A tibble: 342 x 3
2  ##   country  year    lrpmg
3  ##   <fctr> <int>    <dbl>
4  ## 1 AUSTRIA 1960 -0.3345476
5  ## 2 AUSTRIA 1961 -0.3513276
6  ## 3 AUSTRIA 1962 -0.3795177
7  ## 4 AUSTRIA 1963 -0.4142514
8  ## 5 AUSTRIA 1964 -0.4453354
9  ## 6 AUSTRIA 1965 -0.4970607
10 ## 7 AUSTRIA 1966 -0.4668377
11 ## 8 AUSTRIA 1967 -0.5058834
12 ## 9 AUSTRIA 1968 -0.5224125
13 ## 10 AUSTRIA 1969 -0.5591105
14 ## # ... with 332 more rows

```

This also works with `filter_at()` by the way.

`select_if()` makes it easy to select columns that satisfy a criterium:

```
1 gasoline %>% select_if(is.numeric)
```

```

1  ## # A tibble: 342 x 5
2  ##   year lgaspcar  lincomep      lrpmg  lcarpcap
3  ## * <int>     <dbl>     <dbl>     <dbl>     <dbl>
4  ## 1  1960 4.173244 -6.474277 -0.3345476 -9.766840
5  ## 2  1961 4.100989 -6.426006 -0.3513276 -9.608622
6  ## 3  1962 4.073177 -6.407308 -0.3795177 -9.457257
7  ## 4  1963 4.059509 -6.370679 -0.4142514 -9.343155
8  ## 5  1964 4.037689 -6.322247 -0.4453354 -9.237739
9  ## 6  1965 4.033983 -6.294668 -0.4970607 -9.123903
10 ## 7  1966 4.047537 -6.252545 -0.4668377 -9.019822
11 ## 8  1967 4.052911 -6.234581 -0.5058834 -8.934403
12 ## 9  1968 4.045507 -6.206894 -0.5224125 -8.847967
13 ## 10 1969 4.046355 -6.153140 -0.5591105 -8.788686
14 ## # ... with 332 more rows

```

You can even pass a further function to `select_if()` that will be applied to the selected columns:

```
1 gasoline %>% select_if(is.numeric, toupper)
```

```

1  ## # A tibble: 342 x 5
2  ##   YEAR LGASPCAR  LINCOME  LRPMG  LCARPCAP
3  ## * <int>     <dbl>     <dbl>     <dbl>     <dbl>
4  ## 1  1960 4.173244 -6.474277 -0.3345476 -9.766840
5  ## 2  1961 4.100989 -6.426006 -0.3513276 -9.608622
6  ## 3  1962 4.073177 -6.407308 -0.3795177 -9.457257
7  ## 4  1963 4.059509 -6.370679 -0.4142514 -9.343155
8  ## 5  1964 4.037689 -6.322247 -0.4453354 -9.237739
9  ## 6  1965 4.033983 -6.294668 -0.4970607 -9.123903
10 ## 7  1966 4.047537 -6.252545 -0.4668377 -9.019822
11 ## 8  1967 4.052911 -6.234581 -0.5058834 -8.934403
12 ## 9  1968 4.045507 -6.206894 -0.5224125 -8.847967
13 ## 10 1969 4.046355 -6.153140 -0.5591105 -8.788686
14 ## # ... with 332 more rows

```

Another verb, similar to `select()`, is `pull()`. Let's compare the two:

```
1 gasoline %>% select(lrpmg)
```

```

1  ## # A tibble: 342 x 1
2  ##       lrpmpg
3  ##   *       <dbl>
4  ##   1 -0.3345476
5  ##   2 -0.3513276
6  ##   3 -0.3795177
7  ##   4 -0.4142514
8  ##   5 -0.4453354
9  ##   6 -0.4970607
10 ##   7 -0.4668377
11 ##   8 -0.5058834
12 ##   9 -0.5224125
13 ##  10 -0.5591105
14 ## # ... with 332 more rows

```

```

1  gasoline %>% pull(lrpmpg)

```

```

1  ##   [1] -0.33454761 -0.35132761 -0.37951769 -0.41425139 -0.44533536
2  ##   [6] -0.49706066 -0.46683773 -0.50588340 -0.52241255 -0.55911051
3  ##  [11] -0.59656122 -0.65445914 -0.59633184 -0.59444681 -0.46602693
4  ##  [16] -0.45414221 -0.50008372 -0.42191563 -0.46960312 -0.16570961
5  ##  [21] -0.17173098 -0.22229138 -0.25046225 -0.27591057 -0.34493695
6  ##  [26] -0.23639770 -0.26699499 -0.31116076 -0.35480852 -0.37794044
7  ##  [31] -0.39922992 -0.31064584 -0.37309192 -0.36223563 -0.36430848
8  ##  [36] -0.37896584 -0.43164133 -0.59094964 -0.97210650 -0.97229024
9  ##  [41] -0.97860756 -1.01904791 -1.00285696 -1.01712549 -1.01694436
10 ##  [46] -1.02359713 -1.01984524 -1.03686389 -1.06733308 -1.05803676
11 ##  [51] -1.09966703 -1.13316142 -1.12379997 -1.18568427 -1.06179659
12 ##  [56] -1.07084448 -1.07495073 -0.19570260 -0.25361844 -0.21875400
13 ##  [61] -0.24800936 -0.30654923 -0.32701542 -0.39618846 -0.44257369
14 ##  [66] -0.35204752 -0.40687922 -0.44046082 -0.45473954 -0.49918863
15 ##  [71] -0.43257185 -0.42517720 -0.39395431 -0.35361534 -0.35690917
16 ##  [76] -0.29068135 -0.01959833 -0.02386000 -0.06892022 -0.13792900
17 ##  [81] -0.19784646 -0.23365325 -0.26427164 -0.29405795 -0.32316179
18 ##  [86] -0.31519087 -0.33384616 -0.37945667 -0.40781642 -0.47503429
19 ##  [91] -0.21698191 -0.25838174 -0.24651309 -0.22550681 -0.38075942
20 ##  [96] -0.18591078 -0.23095384 -0.34384171 -0.37464672 -0.39965256
21 ## [101] -0.43987825 -0.54000197 -0.54998139 -0.43824222 -0.58923137
22 ## [106] -0.63329520 -0.67176311 -0.71797458 -0.72587521 -0.56982876
23 ## [111] -0.56482380 -0.62481298 -0.59761210 -0.62817279 -0.08354740
24 ## [116] -0.10421997 -0.13320751 -0.15653576 -0.18051772 -0.07793999

```

```

25 ## [121] -0.11491900 -0.13775849 -0.15375883 -0.17986997 -0.20252426
26 ## [126] -0.06761078 -0.11973059 -0.05191029  0.31625351  0.20631574
27 ## [131]  0.19319312  0.23502961  0.16896037 -0.07648118 -0.12040874
28 ## [136] -0.14160039 -0.15232915 -0.24428212 -0.16899366 -0.21071901
29 ## [141] -0.17383533 -0.21339314 -0.27162842 -0.32069023 -0.36041067
30 ## [146] -0.42393131 -0.64567297 -0.55343875 -0.64126416 -0.66134256
31 ## [151] -0.56011483 -0.66277808  0.16507708 -0.08559038 -0.18351291
32 ## [156] -0.26541405 -0.42609643 -0.32712637 -0.24887418 -0.19160048
33 ## [161] -0.20616656 -0.24756681 -0.23271512 -0.14822267 -0.21508857
34 ## [166] -0.32508487 -0.22290860 -0.03270913  0.10292798  0.16418805
35 ## [171]  0.03482212 -0.14532271 -0.14874940 -0.18731459 -0.19996473
36 ## [176] -0.20386433 -0.23786571 -0.27411537 -0.33167240 -0.35126918
37 ## [181] -0.41685019 -0.46203546 -0.43941354 -0.52100094 -0.46270739
38 ## [186] -0.19090636 -0.15948473 -0.20726559 -0.21904447 -0.28707638
39 ## [191] -0.20148480 -0.21599265 -0.25968008 -0.29718661 -0.36929389
40 ## [196] -0.34197503 -0.34809007 -0.31232019 -0.44450431 -0.41694955
41 ## [201] -0.39954544 -0.43393029 -0.31903240 -0.42728193 -0.35253685
42 ## [206] -0.43426178 -0.42908393 -0.46474195 -0.55791459 -0.13968957
43 ## [211] -0.15790514 -0.19908809 -0.23263318 -0.26374731 -0.31593124
44 ## [216] -0.25011726 -0.26555763 -0.30036775 -0.33823045 -0.39072560
45 ## [221] -0.30127223 -0.26023925 -0.33880765 -0.15100924 -0.32726757
46 ## [226] -0.35308752 -0.38255762 -0.30765935  1.12531070  1.10956235
47 ## [231]  1.05700394  0.97683534  0.91532254  0.81666055  0.75671751
48 ## [236]  0.74130811  0.70386453  0.66948950  0.61217208  0.60699563
49 ## [241]  0.53716844  0.43377166  0.52492096  0.62955545  0.68385409
50 ## [246]  0.52627167  0.62141374 -2.52041588 -2.57148340 -2.53448158
51 ## [251] -2.60511224 -2.65801626 -2.64476790 -2.63901460 -2.65609762
52 ## [256] -2.67918662 -2.73190414 -2.73359211 -2.77884554 -2.77467537
53 ## [261] -2.84142900 -2.79840677 -2.76731461 -2.82294480 -2.82005896
54 ## [266] -2.89649671 -0.82321833 -0.86558473 -0.82218510 -0.86012004
55 ## [271] -0.86767682 -0.90528668 -0.85956665 -0.90656671 -0.87232520
56 ## [276] -0.91812162 -0.96344188 -1.03746081 -0.94015345 -0.86722756
57 ## [281] -0.88692306 -0.88475790 -0.90736205 -0.91147285 -1.03208811
58 ## [286] -0.25340821 -0.34252375 -0.40820484 -0.22499174 -0.25219448
59 ## [291] -0.29347614 -0.35640491 -0.33515022 -0.36507386 -0.29845417
60 ## [296] -0.39882648 -0.30461880 -0.54637424 -0.69162023 -0.33965308
61 ## [301] -0.53794675 -0.75141027 -0.95552413 -0.35290961 -0.39108581
62 ## [306] -0.45185308 -0.42287690 -0.46335147 -0.49577430 -0.42654915
63 ## [311] -0.47068145 -0.44118786 -0.46245080 -0.38332457 -0.41899030
64 ## [316] -0.46135978 -0.52777246 -0.56529718 -0.56641296 -0.20867428
65 ## [321] -0.27354010 -0.50886285 -0.78652911 -1.12111489 -1.14624034
66 ## [326] -1.16187449 -1.17991524 -1.20026222 -1.19428750 -1.19026054

```

```

67 ## [331] -1.18991215 -1.20730059 -1.22314272 -1.25176347 -1.28131560
68 ## [336] -1.33116930 -1.29066967 -1.23146686 -1.20037697 -1.15468197
69 ## [341] -1.17590974 -1.21206183

```

`pull()`, unlike `select()`, does not return a tibble, but only the atomic vector.

4.3.3 `group_by()`

`group_by()` is a very useful verb; as the name implies, it allows you to create groups and then, for example, compute descriptive statistics by groups. For example, let's group our data by country:

```

1 gasoline %>% group_by(country)

1 ## # A tibble: 342 x 6
2 ## # Groups:   country [18]
3 ##   country year lgaspcar lincomep      lrpmpg lcarpcap
4 ## *   <fctr> <int>   <dbl>     <dbl>     <dbl>     <dbl>
5 ## 1 AUSTRIA  1960  4.173244 -6.474277 -0.3345476 -9.766840
6 ## 2 AUSTRIA  1961  4.100989 -6.426006 -0.3513276 -9.608622
7 ## 3 AUSTRIA  1962  4.073177 -6.407308 -0.3795177 -9.457257
8 ## 4 AUSTRIA  1963  4.059509 -6.370679 -0.4142514 -9.343155
9 ## 5 AUSTRIA  1964  4.037689 -6.322247 -0.4453354 -9.237739
10 ## 6 AUSTRIA  1965  4.033983 -6.294668 -0.4970607 -9.123903
11 ## 7 AUSTRIA  1966  4.047537 -6.252545 -0.4668377 -9.019822
12 ## 8 AUSTRIA  1967  4.052911 -6.234581 -0.5058834 -8.934403
13 ## 9 AUSTRIA  1968  4.045507 -6.206894 -0.5224125 -8.847967
14 ## 10 AUSTRIA 1969  4.046355 -6.153140 -0.5591105 -8.788686
15 ## # ... with 332 more rows

```

It looks like nothing much happened, but if you look at the second line of the output you can read the following:

```

1 ## # Groups:   country [18]

```

this means that the data is grouped, and every computation you will do now will take these groups into account. This will be clearer in the next subsection.

It is also possible to group according to various variables:

```

1 gasoline %>% group_by(country, year)

1 ## # A tibble: 342 x 6
2 ## # Groups:   country, year [342]
3 ##   country year lgaspcar lincomep      lrpmpg lcarpcap
4 ## *   <fctr> <int>   <dbl>   <dbl>   <dbl>   <dbl>
5 ## 1 AUSTRIA  1960  4.173244 -6.474277 -0.3345476 -9.766840
6 ## 2 AUSTRIA  1961  4.100989 -6.426006 -0.3513276 -9.608622
7 ## 3 AUSTRIA  1962  4.073177 -6.407308 -0.3795177 -9.457257
8 ## 4 AUSTRIA  1963  4.059509 -6.370679 -0.4142514 -9.343155
9 ## 5 AUSTRIA  1964  4.037689 -6.322247 -0.4453354 -9.237739
10 ## 6 AUSTRIA  1965  4.033983 -6.294668 -0.4970607 -9.123903
11 ## 7 AUSTRIA  1966  4.047537 -6.252545 -0.4668377 -9.019822
12 ## 8 AUSTRIA  1967  4.052911 -6.234581 -0.5058834 -8.934403
13 ## 9 AUSTRIA  1968  4.045507 -6.206894 -0.5224125 -8.847967
14 ## 10 AUSTRIA 1969  4.046355 -6.153140 -0.5591105 -8.788686
15 ## # ... with 332 more rows

```

and so on. You can then also ungroup:

```

1 gasoline %>% group_by(country, year) %>% ungroup()

1 ## # A tibble: 342 x 6
2 ##   country year lgaspcar lincomep      lrpmpg lcarpcap
3 ## *   <fctr> <int>   <dbl>   <dbl>   <dbl>   <dbl>
4 ## 1 AUSTRIA  1960  4.173244 -6.474277 -0.3345476 -9.766840
5 ## 2 AUSTRIA  1961  4.100989 -6.426006 -0.3513276 -9.608622
6 ## 3 AUSTRIA  1962  4.073177 -6.407308 -0.3795177 -9.457257
7 ## 4 AUSTRIA  1963  4.059509 -6.370679 -0.4142514 -9.343155
8 ## 5 AUSTRIA  1964  4.037689 -6.322247 -0.4453354 -9.237739
9 ## 6 AUSTRIA  1965  4.033983 -6.294668 -0.4970607 -9.123903
10 ## 7 AUSTRIA  1966  4.047537 -6.252545 -0.4668377 -9.019822
11 ## 8 AUSTRIA  1967  4.052911 -6.234581 -0.5058834 -8.934403
12 ## 9 AUSTRIA  1968  4.045507 -6.206894 -0.5224125 -8.847967
13 ## 10 AUSTRIA 1969  4.046355 -6.153140 -0.5591105 -8.788686
14 ## # ... with 332 more rows

```

4.3.4 summarise()

Ok, now that we have learned the basic verbs, we can start to do more interesting stuff. For example, one might want to compute the average gasoline consumption in each country, for the whole period:

```

1 gasoline %>%
2   group_by(country) %>%
3   summarise(mean(lgascar))

```

```

1 ## # A tibble: 18 x 2
2 ##   country `mean(lgascar)`
3 ##   <fctr>      <dbl>
4 ## 1 AUSTRIA      4.056487
5 ## 2 BELGIUM      3.922286
6 ## 3 CANADA      4.862402
7 ## 4 DENMARK     4.189886
8 ## 5 FRANCE      3.815198
9 ## 6 GERMANY     3.893389
10 ## 7 GREECE      4.878679
11 ## 8 IRELAND     4.225560
12 ## 9 ITALY       3.729646
13 ## 10 JAPAN      4.699642
14 ## 11 NETHERLA  4.080338
15 ## 12 NORWAY     4.109773
16 ## 13 SPAIN      4.055314
17 ## 14 SWEDEN     4.006055
18 ## 15 SWITZERL  4.237586
19 ## 16 TURKEY     5.766355
20 ## 17 U.K.       3.984685
21 ## 18 U.S.A.     4.819075

```

`mean()` was given as an argument to `summarise()`, which is a `dplyr` verb. What we get is another tibble, that contains the variable we used to group, as well as the average per country. We can also rename this column:

```

1 gasoline %>%
2   group_by(country) %>%
3   summarise(mean_gascar = mean(lgascar))

```



```

1  ## # A tibble: 18 x 2
2  ##   country mean_gaspcar
3  ##   <fctr>      <dbl>
4  ## 1 AUSTRIA      4.056487
5  ## 2 BELGIUM      3.922286
6  ## 3 CANADA       4.862402
7  ## 4 DENMARK      4.189886
8  ## 5 FRANCE       3.815198
9  ## 6 GERMANY      3.893389
10 ## 7 GREECE       4.878679
11 ## 8 IRELAND      4.225560
12 ## 9 ITALY        3.729646
13 ## 10 JAPAN        4.699642
14 ## 11 NETHERLA    4.080338
15 ## 12 NORWAY      4.109773
16 ## 13 SPAIN       4.055314
17 ## 14 SWEDEN      4.006055
18 ## 15 SWITZERL    4.237586
19 ## 16 TURKEY      5.766355
20 ## 17 U.K.        3.984685
21 ## 18 U.S.A.      4.819075

```

and because the output is a tibble, we can continue to use `dplyr` verbs on it:

```

1 gasoline %>%
2   group_by(country) %>%
3   summarise(mean_gaspcar = mean(lgaspcar)) %>%
4   filter(country == "FRANCE")

```

```

1  ## # A tibble: 1 x 2
2  ##   country mean_gaspcar
3  ##   <fctr>      <dbl>
4  ## 1 FRANCE      3.815198

```

Ok, let's pause here. See what I did in the last example? I chained 3 `dplyr` verbs together with `%>%`. Without using `%>%` I would have written:

```

1 filter(
2   summarise(
3     group_by(gasoline, country),
4     mean_gaspcar = mean(lgaspcar)),
5   country == "FRANCE")

```

```

1 ## # A tibble: 1 x 2
2 ##   country mean_gaspcar
3 ##   <fctr>      <dbl>
4 ## 1  FRANCE      3.815198

```

I don't know about you, but this is much more difficult to read than the version with `%>%`. It is possible to work like that, of course, but personally, I would advise you bite the bullet and learn to love the pipe. It won't give you cancer.

Ok, back to `summarise()`. We can really do a lot of stuff with this verb. For example, we can compute several descriptive statistics at once:

```

1 gasoline %>%
2   group_by(country) %>%
3   summarise(mean_gaspcar = mean(lgaspcar),
4             sd_gaspcar = sd(lgaspcar),
5             max_gaspcar = max(lgaspcar),
6             min_gaspcar = min(lgaspcar))

```

```

1 ## # A tibble: 18 x 5
2 ##   country mean_gaspcar sd_gaspcar max_gaspcar min_gaspcar
3 ##   <fctr>      <dbl>      <dbl>      <dbl>      <dbl>
4 ## 1  AUSTRIA      4.056487 0.06929942  4.199381  3.922750
5 ## 2  BELGIUM      3.922286 0.10339189  4.164016  3.818230
6 ## 3  CANADA      4.862402 0.02618377  4.899694  4.810992
7 ## 4  DENMARK      4.189886 0.15819728  4.501986  4.000461
8 ## 5  FRANCE      3.815198 0.04986425  3.908116  3.749535
9 ## 6  GERMANY      3.893389 0.02389849  3.932402  3.848782
10 ## 7  GREECE      4.878679 0.25467445  5.381495  4.479956
11 ## 8  IRELAND      4.225560 0.04369894  4.325585  4.164896
12 ## 9  ITALY      3.729646 0.22001527  4.050728  3.380209
13 ## 10 JAPAN      4.699642 0.68411717  5.995287  3.948746
14 ## 11 NETHERLA    4.080338 0.28642682  4.646268  3.711384
15 ## 12 NORWAY      4.109773 0.12306866  4.435041  3.960331
16 ## 13 SPAIN      4.055314 0.31696784  4.749409  3.620444

```

```

17 ## 14 SWEDEN      4.006055 0.03639626    4.067373    3.913159
18 ## 15 SWITZERL    4.237586 0.10178743    4.441330    4.050048
19 ## 16 TURKEY      5.766355 0.32901391    6.156644    5.141255
20 ## 17 U.K.        3.984685 0.04787887    4.100244    3.912584
21 ## 18 U.S.A.      4.819075 0.02189802    4.860286    4.787895

```

Because the output is a tibble, you can save it in a variable of course:

```

1 desc_gasoline <- gasoline %>%
2   group_by(country) %>%
3   summarise(mean_gaspcar = mean(lgaspcar),
4             sd_gaspcar = sd(lgaspcar),
5             max_gaspcar = max(lgaspcar),
6             min_gaspcar = min(lgaspcar))

```

And then you can answer questions such as, *which country has the maximum average gasoline consumption?*:

```

1 desc_gasoline %>%
2   filter(max(mean_gaspcar) == mean_gaspcar)

1 ## # A tibble: 1 x 5
2 ##   country mean_gaspcar sd_gaspcar max_gaspcar min_gaspcar
3 ##   <fctr>      <dbl>      <dbl>      <dbl>      <dbl>
4 ## 1 TURKEY      5.766355 0.3290139 6.156644 5.141255

```

Turns out it's Turkey. What about the minimum consumption?

```

1 desc_gasoline %>%
2   filter(min(mean_gaspcar) == mean_gaspcar)

1 ## # A tibble: 1 x 5
2 ##   country mean_gaspcar sd_gaspcar max_gaspcar min_gaspcar
3 ##   <fctr>      <dbl>      <dbl>      <dbl>      <dbl>
4 ## 1 ITALY      3.729646 0.2200153 4.050728 3.380209

```

Just like for `filter()` and `select()`, `summarise()` comes with scoped versions:

```

1 gasoline %>%
2   group_by(country) %>%
3   summarise_at(vars(starts_with("l")), mean)

```

```

1 ## # A tibble: 18 x 5
2 ##   country lgaspcar  lincomep      lrpmg  lcarpcap
3 ##   <fctr>    <dbl>    <dbl>    <dbl>    <dbl>
4 ## 1 AUSTRIA 4.056487 -6.119613 -0.48578185 -8.848114
5 ## 2 BELGIUM 3.922286 -5.852297 -0.32575856 -8.630392
6 ## 3 CANADA 4.862402 -5.576948 -1.04918735 -8.081975
7 ## 4 DENMARK 4.189886 -5.756725 -0.35761243 -8.583795
8 ## 5 FRANCE 3.815198 -5.866167 -0.25277821 -8.452957
9 ## 6 GERMANY 3.893389 -5.845314 -0.51718417 -8.506392
10 ## 7 GREECE 4.878679 -6.606373 -0.03391043 -10.782066
11 ## 8 IRELAND 4.225560 -6.441571 -0.34754288 -9.035927
12 ## 9 ITALY 3.729646 -6.350354 -0.15219273 -8.827179
13 ## 10 JAPAN 4.699642 -6.248629 -0.28662755 -9.945087
14 ## 11 NETHERLA 4.080338 -5.920132 -0.36977928 -8.817087
15 ## 12 NORWAY 4.109773 -5.753316 -0.27767861 -8.765066
16 ## 13 SPAIN 4.055314 -5.627756 0.73937888 -9.896247
17 ## 14 SWEDEN 4.006055 -7.816214 -2.70917074 -8.250729
18 ## 15 SWITZERL 4.237586 -5.927320 -0.90165998 -8.541029
19 ## 16 TURKEY 5.766355 -7.336992 -0.42151399 -12.458858
20 ## 17 U.K. 3.984685 -6.015377 -0.45929339 -8.548493
21 ## 18 U.S.A. 4.819075 -5.448560 -1.20756453 -7.781090

```

See how I managed to summarise every variable in one simple call to `summarise_at()`? Simply by using `vars()` and specifying that I was interested in the ones that started with `l` and then I specified the function I wanted. But what if I wanted to use more than one function to summarise the data? Very easy:

```

1 gasoline %>%
2   group_by(country) %>%
3   summarise_at(vars(starts_with("l")), funs(mean, sd, max, min))

```

```

1 ## # A tibble: 18 x 17
2 ##   country lgaspcar_mean lincomep_mean lrpmg_mean lcarpcap_mean
3 ##   <fctr>      <dbl>      <dbl>      <dbl>      <dbl>
4 ## 1 AUSTRIA      4.056487      -6.119613 -0.48578185      -8.848114
5 ## 2 BELGIUM      3.922286      -5.852297 -0.32575856      -8.630392
6 ## 3 CANADA       4.862402      -5.576948 -1.04918735      -8.081975
7 ## 4 DENMARK      4.189886      -5.756725 -0.35761243      -8.583795
8 ## 5 FRANCE       3.815198      -5.866167 -0.25277821      -8.452957
9 ## 6 GERMANY      3.893389      -5.845314 -0.51718417      -8.506392
10 ## 7 GREECE      4.878679      -6.606373 -0.03391043     -10.782066
11 ## 8 IRELAND      4.225560      -6.441571 -0.34754288      -9.035927
12 ## 9 ITALY       3.729646      -6.350354 -0.15219273      -8.827179
13 ## 10 JAPAN      4.699642      -6.248629 -0.28662755      -9.945087
14 ## 11 NETHERLA   4.080338      -5.920132 -0.36977928      -8.817087
15 ## 12 NORWAY     4.109773      -5.753316 -0.27767861      -8.765066
16 ## 13 SPAIN      4.055314      -5.627756  0.73937888      -9.896247
17 ## 14 SWEDEN     4.006055      -7.816214 -2.70917074      -8.250729
18 ## 15 SWITZERL   4.237586      -5.927320 -0.90165998      -8.541029
19 ## 16 TURKEY     5.766355      -7.336992 -0.42151399     -12.458858
20 ## 17 U.K.       3.984685      -6.015377 -0.45929339      -8.548493
21 ## 18 U.S.A.     4.819075      -5.448560 -1.20756453      -7.781090
22 ## # ... with 12 more variables: lgaspcar_sd <dbl>, lincomep_sd <dbl>,
23 ## #   lrpmg_sd <dbl>, lcarpcap_sd <dbl>, lgaspcar_max <dbl>,
24 ## #   lincomep_max <dbl>, lrpmg_max <dbl>, lcarpcap_max <dbl>,
25 ## #   lgaspcar_min <dbl>, lincomep_min <dbl>, lrpmg_min <dbl>,
26 ## #   lcarpcap_min <dbl>

```

But maybe you're just interested in descriptive statistics for some variables, but not all those that start with "car"? What if you want to use another pattern? Easy to do with the `contains()` helper:

```

1 gasoline %>%
2   group_by(country) %>%
3   summarise_at(vars(dplyr::contains("car")), funs(mean, sd, max, min))

```

```

1 ## # A tibble: 18 x 9
2 ##   country lgaspcar_mean lcarpcap_mean lgaspcar_sd lcarpcap_sd
3 ##   <fctr>      <dbl>      <dbl>      <dbl>      <dbl>
4 ## 1 AUSTRIA      4.056487      -8.848114  0.06929942  0.4728231
5 ## 2 BELGIUM      3.922286      -8.630392  0.10339189  0.4171514
6 ## 3 CANADA       4.862402      -8.081975  0.02618377  0.1953069
7 ## 4 DENMARK      4.189886      -8.583795  0.15819728  0.3486135
8 ## 5 FRANCE       3.815198      -8.452957  0.04986425  0.3436969
9 ## 6 GERMANY      3.893389      -8.506392  0.02389849  0.4060370
10 ## 7 GREECE      4.878679     -10.782066  0.25467445  0.8388589
11 ## 8 IRELAND      4.225560      -9.035927  0.04369894  0.3452272
12 ## 9 ITALY       3.729646      -8.827179  0.22001527  0.6389769
13 ## 10 JAPAN      4.699642      -9.945087  0.68411717  1.1969275
14 ## 11 NETHERLA   4.080338      -8.817087  0.28642682  0.6173209
15 ## 12 NORWAY     4.109773      -8.765066  0.12306866  0.4382484
16 ## 13 SPAIN      4.055314      -9.896247  0.31696784  0.9596034
17 ## 14 SWEDEN     4.006055      -8.250729  0.03639626  0.2422792
18 ## 15 SWITZERL   4.237586      -8.541029  0.10178743  0.3775211
19 ## 16 TURKEY     5.766355     -12.458858  0.32901391  0.7512506
20 ## 17 U.K.       3.984685      -8.548493  0.04787887  0.2812851
21 ## 18 U.S.A.     4.819075      -7.781090  0.02189802  0.1617998
22 ## # ... with 4 more variables: lgaspcar_max <dbl>, lcarpcap_max <dbl>,
23 ## #   lgaspcar_min <dbl>, lcarpcap_min <dbl>

```

I used `dplyr::contains()` instead of simply `contains()` because there's also a `purrr::contains()`. If you load `purrr` after `dplyr`, `contains()` will actually be `purrr::contains()` and not `dplyr::contains()` which causes the above code to fail.

There's also `summarise_if()`:

```

1 gasoline %>%
2   group_by(country) %>%
3   summarise_if(is.double, funs(mean, sd, min, max))

```

```

1 ## # A tibble: 18 x 17
2 ##   country lgaspcar_mean lincomep_mean lrpmg_mean lcarpcap_mean
3 ##   <fctr>      <dbl>      <dbl>      <dbl>      <dbl>
4 ## 1 AUSTRIA      4.056487    -6.119613 -0.48578185    -8.848114
5 ## 2 BELGIUM      3.922286    -5.852297 -0.32575856    -8.630392
6 ## 3 CANADA       4.862402    -5.576948 -1.04918735    -8.081975
7 ## 4 DENMARK      4.189886    -5.756725 -0.35761243    -8.583795
8 ## 5 FRANCE       3.815198    -5.866167 -0.25277821    -8.452957
9 ## 6 GERMANY      3.893389    -5.845314 -0.51718417    -8.506392
10 ## 7 GREECE      4.878679    -6.606373 -0.03391043   -10.782066
11 ## 8 IRELAND      4.225560    -6.441571 -0.34754288    -9.035927
12 ## 9 ITALY       3.729646    -6.350354 -0.15219273    -8.827179
13 ## 10 JAPAN      4.699642    -6.248629 -0.28662755    -9.945087
14 ## 11 NETHERLA   4.080338    -5.920132 -0.36977928    -8.817087
15 ## 12 NORWAY     4.109773    -5.753316 -0.27767861    -8.765066
16 ## 13 SPAIN      4.055314    -5.627756  0.73937888    -9.896247
17 ## 14 SWEDEN     4.006055    -7.816214 -2.70917074    -8.250729
18 ## 15 SWITZERL   4.237586    -5.927320 -0.90165998    -8.541029
19 ## 16 TURKEY     5.766355    -7.336992 -0.42151399   -12.458858
20 ## 17 U.K.       3.984685    -6.015377 -0.45929339    -8.548493
21 ## 18 U.S.A.     4.819075    -5.448560 -1.20756453    -7.781090
22 ## # ... with 12 more variables: lgaspcar_sd <dbl>, lincomep_sd <dbl>,
23 ## #   lrpmg_sd <dbl>, lcarpcap_sd <dbl>, lgaspcar_min <dbl>,
24 ## #   lincomep_min <dbl>, lrpmg_min <dbl>, lcarpcap_min <dbl>,
25 ## #   lgaspcar_max <dbl>, lincomep_max <dbl>, lrpmg_max <dbl>,
26 ## #   lcarpcap_max <dbl>

```

This allows you to summarise every column that contain real numbers (if you use `is.numeric()` instead, year will also be summarised, which is not really interesting).

To go faster, you can also use `summarise_all()`:

```

1 gasoline %>%
2   select(-year) %>%
3   group_by(country) %>%
4   summarise_all(funs(mean, sd, min, max))

```

```

1 ## # A tibble: 18 x 17
2 ##   country lgaspcar_mean lincomep_mean lrpmg_mean lcarpcap_mean
3 ##   <fctr>      <dbl>      <dbl>      <dbl>      <dbl>
4 ## 1 AUSTRIA      4.056487      -6.119613 -0.48578185      -8.848114
5 ## 2 BELGIUM      3.922286      -5.852297 -0.32575856      -8.630392
6 ## 3 CANADA       4.862402      -5.576948 -1.04918735      -8.081975
7 ## 4 DENMARK      4.189886      -5.756725 -0.35761243      -8.583795
8 ## 5 FRANCE       3.815198      -5.866167 -0.25277821      -8.452957
9 ## 6 GERMANY      3.893389      -5.845314 -0.51718417      -8.506392
10 ## 7 GREECE      4.878679      -6.606373 -0.03391043     -10.782066
11 ## 8 IRELAND      4.225560      -6.441571 -0.34754288      -9.035927
12 ## 9 ITALY       3.729646      -6.350354 -0.15219273      -8.827179
13 ## 10 JAPAN      4.699642      -6.248629 -0.28662755      -9.945087
14 ## 11 NETHERLA   4.080338      -5.920132 -0.36977928      -8.817087
15 ## 12 NORWAY     4.109773      -5.753316 -0.27767861      -8.765066
16 ## 13 SPAIN      4.055314      -5.627756  0.73937888      -9.896247
17 ## 14 SWEDEN     4.006055      -7.816214 -2.70917074      -8.250729
18 ## 15 SWITZERL   4.237586      -5.927320 -0.90165998      -8.541029
19 ## 16 TURKEY     5.766355      -7.336992 -0.42151399     -12.458858
20 ## 17 U.K.       3.984685      -6.015377 -0.45929339      -8.548493
21 ## 18 U.S.A.     4.819075      -5.448560 -1.20756453      -7.781090
22 ## # ... with 12 more variables: lgaspcar_sd <dbl>, lincomep_sd <dbl>,
23 ## #   lrpmg_sd <dbl>, lcarpcap_sd <dbl>, lgaspcar_min <dbl>,
24 ## #   lincomep_min <dbl>, lrpmg_min <dbl>, lcarpcap_min <dbl>,
25 ## #   lgaspcar_max <dbl>, lincomep_max <dbl>, lrpmg_max <dbl>,
26 ## #   lcarpcap_max <dbl>

```

I removed the year variable because it's not a variable for which we want to have descriptive statistics.

4.3.5 mutate() and transmute()

mutate() adds a column to the tibble, which can contain any transformation of any other variable:

```

1 gasoline %>%
2   group_by(country) %>%
3   mutate(n())

```



```

1 ## # A tibble: 342 x 7
2 ## # Groups:   country [18]
3 ##   country year lgaspcar lincomep   lrpmg   lcarpcap `n()`
4 ##   <fctr> <int>   <dbl>   <dbl>   <dbl>   <dbl> <int>
5 ## 1 AUSTRIA 1960 4.173244 -6.474277 -0.3345476 -9.766840 19
6 ## 2 AUSTRIA 1961 4.100989 -6.426006 -0.3513276 -9.608622 19
7 ## 3 AUSTRIA 1962 4.073177 -6.407308 -0.3795177 -9.457257 19
8 ## 4 AUSTRIA 1963 4.059509 -6.370679 -0.4142514 -9.343155 19
9 ## 5 AUSTRIA 1964 4.037689 -6.322247 -0.4453354 -9.237739 19
10 ## 6 AUSTRIA 1965 4.033983 -6.294668 -0.4970607 -9.123903 19
11 ## 7 AUSTRIA 1966 4.047537 -6.252545 -0.4668377 -9.019822 19
12 ## 8 AUSTRIA 1967 4.052911 -6.234581 -0.5058834 -8.934403 19
13 ## 9 AUSTRIA 1968 4.045507 -6.206894 -0.5224125 -8.847967 19
14 ## 10 AUSTRIA 1969 4.046355 -6.153140 -0.5591105 -8.788686 19
15 ## # ... with 332 more rows

```

Using `mutate()` I've added a column that counts how many times the country appears in the tibble, using `n()`, another dplyr function. There's also `count()` and `tally()`, which we are going to see further down. It is also possible to rename the column on the fly:

```

1 gasoline %>%
2   group_by(country) %>%
3   mutate(freq = n())

1 ## # A tibble: 342 x 7
2 ## # Groups:   country [18]
3 ##   country year lgaspcar lincomep   lrpmg   lcarpcap freq
4 ##   <fctr> <int>   <dbl>   <dbl>   <dbl>   <dbl> <int>
5 ## 1 AUSTRIA 1960 4.173244 -6.474277 -0.3345476 -9.766840 19
6 ## 2 AUSTRIA 1961 4.100989 -6.426006 -0.3513276 -9.608622 19
7 ## 3 AUSTRIA 1962 4.073177 -6.407308 -0.3795177 -9.457257 19
8 ## 4 AUSTRIA 1963 4.059509 -6.370679 -0.4142514 -9.343155 19
9 ## 5 AUSTRIA 1964 4.037689 -6.322247 -0.4453354 -9.237739 19
10 ## 6 AUSTRIA 1965 4.033983 -6.294668 -0.4970607 -9.123903 19
11 ## 7 AUSTRIA 1966 4.047537 -6.252545 -0.4668377 -9.019822 19
12 ## 8 AUSTRIA 1967 4.052911 -6.234581 -0.5058834 -8.934403 19
13 ## 9 AUSTRIA 1968 4.045507 -6.206894 -0.5224125 -8.847967 19
14 ## 10 AUSTRIA 1969 4.046355 -6.153140 -0.5591105 -8.788686 19
15 ## # ... with 332 more rows

```

It is possible to do any arbitrary operation:

```

1 gasoline %>%
2   group_by(country) %>%
3   mutate(spam = exp(lgaspcar + lincomep))

```

```

1 ## # A tibble: 342 x 7
2 ## # Groups:   country [18]
3 ##   country year lgaspcar lincomep   lrpmg   lcarpcap   spam
4 ##   <fctr> <int>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
5 ## 1 AUSTRIA 1960 4.173244 -6.474277 -0.3345476 -9.766840 0.10015533
6 ## 2 AUSTRIA 1961 4.100989 -6.426006 -0.3513276 -9.608622 0.09778181
7 ## 3 AUSTRIA 1962 4.073177 -6.407308 -0.3795177 -9.457257 0.09689458
8 ## 4 AUSTRIA 1963 4.059509 -6.370679 -0.4142514 -9.343155 0.09914524
9 ## 5 AUSTRIA 1964 4.037689 -6.322247 -0.4453354 -9.237739 0.10181905
10 ## 6 AUSTRIA 1965 4.033983 -6.294668 -0.4970607 -9.123903 0.10427907
11 ## 7 AUSTRIA 1966 4.047537 -6.252545 -0.4668377 -9.019822 0.11024954
12 ## 8 AUSTRIA 1967 4.052911 -6.234581 -0.5058834 -8.934403 0.11285291
13 ## 9 AUSTRIA 1968 4.045507 -6.206894 -0.5224125 -8.847967 0.11516524
14 ## 10 AUSTRIA 1969 4.046355 -6.153140 -0.5591105 -8.788686 0.12162839
15 ## # ... with 332 more rows

```

`transmute()` is the same as `mutate()`, but only returns the created variable:

```

1 gasoline %>%
2   group_by(country) %>%
3   transmute(spam = exp(lgaspcar + lincomep))

```

```

1 ## Adding missing grouping variables: `country`
2
3 ## # A tibble: 342 x 2
4 ## # Groups:   country [18]
5 ##   country   spam
6 ##   <fctr>   <dbl>
7 ## 1 AUSTRIA 0.10015533
8 ## 2 AUSTRIA 0.09778181
9 ## 3 AUSTRIA 0.09689458
10 ## 4 AUSTRIA 0.09914524
11 ## 5 AUSTRIA 0.10181905
12 ## 6 AUSTRIA 0.10427907
13 ## 7 AUSTRIA 0.11024954
14 ## 8 AUSTRIA 0.11285291

```

```

15 ## 9 AUSTRIA 0.11516524
16 ## 10 AUSTRIA 0.12162839
17 ## # ... with 332 more rows

```

`mutate()` and `transmute()` also come with scoped version:

```

1 gasoline %>%
2   mutate_if(is.double, exp)

1 ## # A tibble: 342 x 6
2 ##   country year lgaspcar   lincomep   lrpmg   lcarpcap
3 ##   <fctr> <int>   <dbl>     <dbl>   <dbl>     <dbl>
4 ## 1 AUSTRIA 1960 64.92574 0.001542614 0.7156618 5.732123e-05
5 ## 2 AUSTRIA 1961 60.40000 0.001618904 0.7037532 6.714730e-05
6 ## 3 AUSTRIA 1962 58.74327 0.001649458 0.6841913 7.812062e-05
7 ## 4 AUSTRIA 1963 57.94586 0.001710998 0.6608348 8.756274e-05
8 ## 5 AUSTRIA 1964 56.69516 0.001795904 0.6406094 9.729730e-05
9 ## 6 AUSTRIA 1965 56.48546 0.001846122 0.6083161 1.090283e-04
10 ## 7 AUSTRIA 1966 57.25624 0.001925547 0.6269818 1.209877e-04
11 ## 8 AUSTRIA 1967 57.56477 0.001960451 0.6029727 1.317766e-04
12 ## 9 AUSTRIA 1968 57.14015 0.002015487 0.5930880 1.436735e-04
13 ## 10 AUSTRIA 1969 57.18861 0.002126794 0.5717174 1.524481e-04
14 ## # ... with 332 more rows

1 gasoline %>%
2   mutate_at(vars(starts_with("l")), exp)

1 ## # A tibble: 342 x 6
2 ##   country year lgaspcar   lincomep   lrpmg   lcarpcap
3 ##   <fctr> <int>   <dbl>     <dbl>   <dbl>     <dbl>
4 ## 1 AUSTRIA 1960 64.92574 0.001542614 0.7156618 5.732123e-05
5 ## 2 AUSTRIA 1961 60.40000 0.001618904 0.7037532 6.714730e-05
6 ## 3 AUSTRIA 1962 58.74327 0.001649458 0.6841913 7.812062e-05
7 ## 4 AUSTRIA 1963 57.94586 0.001710998 0.6608348 8.756274e-05
8 ## 5 AUSTRIA 1964 56.69516 0.001795904 0.6406094 9.729730e-05
9 ## 6 AUSTRIA 1965 56.48546 0.001846122 0.6083161 1.090283e-04
10 ## 7 AUSTRIA 1966 57.25624 0.001925547 0.6269818 1.209877e-04
11 ## 8 AUSTRIA 1967 57.56477 0.001960451 0.6029727 1.317766e-04
12 ## 9 AUSTRIA 1968 57.14015 0.002015487 0.5930880 1.436735e-04
13 ## 10 AUSTRIA 1969 57.18861 0.002126794 0.5717174 1.524481e-04
14 ## # ... with 332 more rows

```

```

1 gasoline %>%
2   mutate_all(as.character)

1 ## # A tibble: 342 x 6
2 ##   country year   lgaspcar   lincomep   lrpmg   lcarpcap
3 ##   <chr> <chr>     <chr>     <chr>     <chr>     <chr>
4 ## 1 AUSTRIA 1960  4.173244195 -6.474277179 -0.334547613 -9.766839569
5 ## 2 AUSTRIA 1961  4.1009891049 -6.426005835 -0.351327614 -9.608621845
6 ## 3 AUSTRIA 1962  4.0731765511 -6.407308295 -0.379517692 -9.457256552
7 ## 4 AUSTRIA 1963  4.0595091239 -6.370678539 -0.414251392 -9.343154947
8 ## 5 AUSTRIA 1964  4.037688787 -6.322246805 -0.445335362 -9.237739346
9 ## 6 AUSTRIA 1965  4.033983285 -6.294667914 -0.497060662 -9.123903477
10 ## 7 AUSTRIA 1966  4.0475365589 -6.252545451 -0.466837731 -9.019822048
11 ## 8 AUSTRIA 1967  4.0529106939 -6.234580709 -0.505883405 -8.934402537
12 ## 9 AUSTRIA 1968  4.045507048 -6.206894403 -0.522412545 -8.847967407
13 ## 10 AUSTRIA 1969  4.0463547891 -6.153139668 -0.559110514 -8.788686207
14 ## # ... with 332 more rows

```

and there are a lot of useful functions that you can use within `mutate()` or `transmute()`: `lead()`, `lag()`, `dense_rank()`, `ntile()`, `cumsum()`, `cummax()`, `coalesce()`, `na_if()`, etc. We are going to study some of them in the next section.

4.3.5.1 `if_else()`, `case_when()` and `recode()`

The two helper functions I use the most are probably `if_else()` and `case_when()`. These two functions, combined with `mutate()` make it easy to create a new variable conditionally on the values of other variables. For instance, we might want to have a dummy that equals 1 if a country in the European Union (to simplify, say as of 2017) and 0 if not. First let's create a list of countries that are in the EU:

```

1 eu_countries <- c("austria", "belgium", "bulgaria", "croatia", "republic of cypr\
2 us",
3               "czech republic", "denmark", "estonia", "finland", "france", "\
4 germany",
5               "greece", "hungary", "ireland", "italy", "latvia", "lithuania"\
6 , "luxembourg",
7               "malta", "netherla", "poland", "portugal", "romania", "slovaki\
8 a", "slovenia",
9               "spain", "sweden", "u.k.")

```

I've had to change "netherlands" to "netherla" because that's how the country is called in the data. Now let's create a dummy variable that equals 1 for EU countries, and 0 for the others:

```

1 gasoline %>%
2   mutate(country = tolower(country)) %>%
3   mutate(in_eu = if_else(country %in% eu_countries, 1, 0))

1 ## # A tibble: 342 x 7
2 ##   country year lgaspcar lincomep      lrpmg  lcarpcap in_eu
3 ##   <chr> <int>   <dbl>    <dbl>    <dbl>    <dbl> <dbl>
4 ## 1 austria 1960 4.173244 -6.474277 -0.3345476 -9.766840 1
5 ## 2 austria 1961 4.100989 -6.426006 -0.3513276 -9.608622 1
6 ## 3 austria 1962 4.073177 -6.407308 -0.3795177 -9.457257 1
7 ## 4 austria 1963 4.059509 -6.370679 -0.4142514 -9.343155 1
8 ## 5 austria 1964 4.037689 -6.322247 -0.4453354 -9.237739 1
9 ## 6 austria 1965 4.033983 -6.294668 -0.4970607 -9.123903 1
10 ## 7 austria 1966 4.047537 -6.252545 -0.4668377 -9.019822 1
11 ## 8 austria 1967 4.052911 -6.234581 -0.5058834 -8.934403 1
12 ## 9 austria 1968 4.045507 -6.206894 -0.5224125 -8.847967 1
13 ## 10 austria 1969 4.046355 -6.153140 -0.5591105 -8.788686 1
14 ## # ... with 332 more rows

```

Instead of 1 and 0, we can of course use strings (I add `filter(year == 1960)` at the end to have a better view of what happened):

```

1 gasoline %>%
2   mutate(country = tolower(country)) %>%
3   mutate(in_eu = if_else(country %in% eu_countries, "yes", "no")) %>%
4   filter(year == 1960)

1 ## # A tibble: 18 x 7
2 ##   country year lgaspcar lincomep      lrpmg  lcarpcap in_eu
3 ##   <chr> <int>   <dbl>    <dbl>    <dbl>    <dbl> <chr>
4 ## 1 austria 1960 4.173244 -6.474277 -0.33454761 -9.766840 yes
5 ## 2 belgium 1960 4.164016 -6.215091 -0.16570961 -9.405527 yes
6 ## 3 canada 1960 4.855238 -5.889713 -0.97210650 -8.378917 no
7 ## 4 denmark 1960 4.501986 -6.061726 -0.19570260 -9.326161 yes
8 ## 5 france 1960 3.907704 -6.264363 -0.01959833 -9.145706 yes
9 ## 6 germany 1960 3.916953 -6.159837 -0.18591078 -9.342481 yes
10 ## 7 greece 1960 5.037406 -7.164861 -0.08354740 -12.173814 yes
11 ## 8 ireland 1960 4.270421 -6.722466 -0.07648118 -9.698144 yes
12 ## 9 italy 1960 4.050728 -6.727487 0.16507708 -10.142098 yes
13 ## 10 japan 1960 5.995287 -6.986196 -0.14532271 -12.235079 no

```

```

14 ## 11 netherla 1960 4.646268 -6.216365 -0.20148480 -9.998449 yes
15 ## 12  norway 1960 4.435041 -6.090356 -0.13968957 -9.675052 no
16 ## 13   spain 1960 4.749409 -6.166085  1.12531070 -11.588403 yes
17 ## 14  sweden 1960 4.063010 -8.072524 -2.52041588 -8.742679 yes
18 ## 15 switzerl 1960 4.397621 -6.156074 -0.82321833 -9.262400 no
19 ## 16  turkey 1960 6.129553 -7.801144 -0.25340821 -13.475185 no
20 ## 17    u.k. 1960 4.100244 -6.186849 -0.39108581 -9.117623 yes
21 ## 18   u.s.a. 1960 4.823965 -5.698374 -1.12111489 -8.019458 no

```

I think that `if_else()` is fairly straightforward, especially if you know `ifelse()` already. You might be wondering what is the difference between these two. `if_else()` is stricter than `ifelse()` and does not do type conversion. Compare the two next lines:

```

1 ifelse(1 == 1, "0", 1)

1 ## [1] "0"

1 if_else(1 == 1, "0", 1)

1 Error: `false` must be type string, not double

```

Type conversion, especially without a warning is very dangerous. `if_else()`'s behaviour which consists in failing as soon as possible avoids a lot of pain and suffering, especially when programming non-interactively.

`if_else()` also accepts an optional argument, that allows you to specify what should be returned in case of NA:

```

1 if_else(1 == NA, 0, 1, 999)

1 ## [1] 999

1 # Or
2 if_else(1 == NA, 0, 1, NA_real_)

```

```
1 ## [1] NA
```

`case_when()` can be seen as a generalization of `if_else()`. Whenever you want to use multiple `if_else()`s, that's when you know you should use `case_when()` (I'm adding the filter at the end for the same reason as before, to see the output better):

```
1 gasoline %>%
2   mutate(country = tolower(country)) %>%
3   mutate(region = case_when(
4     country %in% c("france", "italy", "turkey", "greece", "spain") ~ "med\
5 iterranean",
6     country %in% c("germany", "austria", "switzerl", "belgium", "netherla\
7 ") ~ "central europe",
8     country %in% c("canada", "u.s.a.", "u.k.", "ireland") ~ "anglosphere",
9     country %in% c("denmark", "norway", "sweden") ~ "nordic",
10    country %in% c("japan") ~ "asia")) %>%
11   filter(year == 1960)
```

```
1 ## # A tibble: 18 x 7
2 ##   country year lgaspcar lincomep   lrpmg   lcarpcap   region
3 ##   <chr> <int>   <dbl>   <dbl>   <dbl>   <dbl>   <chr>
4 ## 1 austria 1960 4.173244 -6.474277 -0.33454761 -9.766840 central europe
5 ## 2 belgium 1960 4.164016 -6.215091 -0.16570961 -9.405527 central europe
6 ## 3 canada 1960 4.855238 -5.889713 -0.97210650 -8.378917 anglosphere
7 ## 4 denmark 1960 4.501986 -6.061726 -0.19570260 -9.326161 nordic
8 ## 5 france 1960 3.907704 -6.264363 -0.01959833 -9.145706 mediterranean
9 ## 6 germany 1960 3.916953 -6.159837 -0.18591078 -9.342481 central europe
10 ## 7 greece 1960 5.037406 -7.164861 -0.08354740 -12.173814 mediterranean
11 ## 8 ireland 1960 4.270421 -6.722466 -0.07648118 -9.698144 anglosphere
12 ## 9 italy 1960 4.050728 -6.727487 0.16507708 -10.142098 mediterranean
13 ## 10 japan 1960 5.995287 -6.986196 -0.14532271 -12.235079 asia
14 ## 11 netherla 1960 4.646268 -6.216365 -0.20148480 -9.998449 central europe
15 ## 12 norway 1960 4.435041 -6.090356 -0.13968957 -9.675052 nordic
16 ## 13 spain 1960 4.749409 -6.166085 1.12531070 -11.588403 mediterranean
17 ## 14 sweden 1960 4.063010 -8.072524 -2.52041588 -8.742679 nordic
18 ## 15 switzerl 1960 4.397621 -6.156074 -0.82321833 -9.262400 central europe
19 ## 16 turkey 1960 6.129553 -7.801144 -0.25340821 -13.475185 mediterranean
20 ## 17 u.k. 1960 4.100244 -6.186849 -0.39108581 -9.117623 anglosphere
21 ## 18 u.s.a. 1960 4.823965 -5.698374 -1.12111489 -8.019458 anglosphere
```

If all you want is to recode values, you can use `recode()`. For example, the Netherlands is written as `“NETHERLA”` in the which is quite ugly. Same for Switzerland:

```

1 gasoline <- gasoline %>%
2   mutate(country = tolower(country)) %>%
3   mutate(country = recode(country, "netherla" = "netherlands", "switzerl" = "swi\
4   tzerland"))

```

I saved the data with these changes as they will become useful in the future. Let's take a look at the data:

```

1 gasoline %>%
2   filter(country %in% c("netherlands", "switzerland"), year == 1960)

1 ## # A tibble: 2 x 6
2 ##       country year lgaspcar lincomep      lrpmpg lcarpcap
3 ##       <chr> <int>   <dbl>    <dbl>    <dbl>    <dbl>
4 ## 1 netherlands 1960 4.646268 -6.216365 -0.2014848 -9.998449
5 ## 2 switzerland 1960 4.397621 -6.156074 -0.8232183 -9.262400

```

4.3.5.2 lead() and lag()

lead() and lag() are especially useful in econometrics. When I was doing my masters, in 4 B.d. (*Before dplyr*) lagging variables in panel data was quite tricky. Now, with dplyr it's really very easy:

```

1 gasoline %>%
2   group_by(country) %>%
3   mutate(lag_lgaspcar = lag(lgaspcar)) %>%
4   mutate(lead_lgaspcar = lead(lgaspcar)) %>%
5   filter(year %in% seq(1960, 1963))

1 ## # A tibble: 72 x 8
2 ## # Groups:   country [18]
3 ##   country year lgaspcar lincomep      lrpmpg lcarpcap lag_lgaspcar
4 ##   <chr> <int>   <dbl>    <dbl>    <dbl>    <dbl>    <dbl>
5 ## 1 austria 1960 4.173244 -6.474277 -0.3345476 -9.766840      NA
6 ## 2 austria 1961 4.100989 -6.426006 -0.3513276 -9.608622  4.173244
7 ## 3 austria 1962 4.073177 -6.407308 -0.3795177 -9.457257  4.100989
8 ## 4 austria 1963 4.059509 -6.370679 -0.4142514 -9.343155  4.073177
9 ## 5 belgium 1960 4.164016 -6.215091 -0.1657096 -9.405527      NA
10 ## 6 belgium 1961 4.124356 -6.176843 -0.1717310 -9.303149  4.164016
11 ## 7 belgium 1962 4.075962 -6.129638 -0.2222914 -9.218070  4.124356
12 ## 8 belgium 1963 4.001266 -6.094019 -0.2504623 -9.114932  4.075962
13 ## 9 canada 1960 4.855238 -5.889713 -0.9721065 -8.378917      NA
14 ## 10 canada 1961 4.826555 -5.884344 -0.9722902 -8.346729  4.855238
15 ## # ... with 62 more rows, and 1 more variables: lead_lgaspcar <dbl>

```


To lag every variable, remember that you can use `mutate_if()`:

```
1 gasoline %>%
2   group_by(country) %>%
3   mutate_if(is.double, lag) %>%
4   filter(year %in% seq(1960, 1963))

1 ## # A tibble: 72 x 6
2 ## # Groups:   country [18]
3 ##   country year lgaspcar lincomep   lrpmg   lcarpcap
4 ##   <chr> <int>   <dbl>   <dbl>   <dbl>   <dbl>
5 ## 1 austria 1960      NA      NA      NA      NA
6 ## 2 austria 1961 4.173244 -6.474277 -0.3345476 -9.766840
7 ## 3 austria 1962 4.100989 -6.426006 -0.3513276 -9.608622
8 ## 4 austria 1963 4.073177 -6.407308 -0.3795177 -9.457257
9 ## 5 belgium 1960      NA      NA      NA      NA
10 ## 6 belgium 1961 4.164016 -6.215091 -0.1657096 -9.405527
11 ## 7 belgium 1962 4.124356 -6.176843 -0.1717310 -9.303149
12 ## 8 belgium 1963 4.075962 -6.129638 -0.2222914 -9.218070
13 ## 9 canada 1960      NA      NA      NA      NA
14 ## 10 canada 1961 4.855238 -5.889713 -0.9721065 -8.378917
15 ## # ... with 62 more rows
```

you can replace `lag()` with `lead()`, but just keep in mind that the columns get transformed in place.

4.3.5.3 `ntile()`

The last helper function I will discuss is `ntile()`. There are some other, so do read `mutate()`'s documentation with `help(mutate)`!

If you need quantiles, you need `ntile()`. Let's see how it works:

```
1 gasoline %>%
2   mutate(quintile = ntile(lgaspcar, 5)) %>%
3   mutate(decile = ntile(lgaspcar, 10)) %>%
4   select(country, year, lgaspcar, quintile, decile)
```

```

1 ## # A tibble: 342 x 5
2 ##   country year lgaspcar quintile decile
3 ##   <chr> <int>   <dbl>   <int> <int>
4 ## 1 austria 1960 4.173244     3     6
5 ## 2 austria 1961 4.100989     3     6
6 ## 3 austria 1962 4.073177     3     5
7 ## 4 austria 1963 4.059509     3     5
8 ## 5 austria 1964 4.037689     3     5
9 ## 6 austria 1965 4.033983     3     5
10 ## 7 austria 1966 4.047537     3     5
11 ## 8 austria 1967 4.052911     3     5
12 ## 9 austria 1968 4.045507     3     5
13 ## 10 austria 1969 4.046355     3     5
14 ## # ... with 332 more rows

```

quintile and decile do not hold the values but the quantile the value lies in. If you want to have a column that contains the median for instance, you can use good ol' `quantile()`:

```

1 gasoline %>%
2   group_by(country) %>%
3   mutate(median = quantile(lgaspcar, 0.5)) %>% # quantile(x, 0.5) is equivalent \
4   to median(x)
5   filter(year == 1960) %>%
6   select(country, year, median)

```

```

1 ## # A tibble: 18 x 3
2 ## # Groups:   country [18]
3 ##   country year median
4 ##   <chr> <int>   <dbl>
5 ## 1 austria 1960 4.047537
6 ## 2 belgium 1960 3.877778
7 ## 3 canada 1960 4.855846
8 ## 4 denmark 1960 4.161687
9 ## 5 france 1960 3.807995
10 ## 6 germany 1960 3.889362
11 ## 7 greece 1960 4.894773
12 ## 8 ireland 1960 4.221146
13 ## 9 italy 1960 3.737389
14 ## 10 japan 1960 4.518290
15 ## 11 netherlands 1960 3.987689
16 ## 12 norway 1960 4.084607

```

```

17 ## 13      spain  1960 3.994103
18 ## 14      sweden 1960 4.002557
19 ## 15  switzerland 1960 4.259159
20 ## 16      turkey  1960 5.722105
21 ## 17        u.k.  1960 3.976781
22 ## 18      u.s.a.  1960 4.811032

```

4.3.5.4 arrange()

arrange() re-orders the whole tibble according to values of the supplied variable:

```

1 gasoline %>%
2   arrange(lgaspcar)

1 ## # A tibble: 342 x 6
2 ##   country year lgaspcar lincomep      lrpmg  lcarpcap
3 ##   <chr> <int>   <dbl>   <dbl>     <dbl>   <dbl>
4 ## 1  italy  1977  3.380209 -6.104541  0.16418805 -8.145660
5 ## 2  italy  1978  3.394504 -6.083685  0.03482212 -8.112852
6 ## 3  italy  1976  3.427629 -6.119222  0.10292798 -8.167756
7 ## 4  italy  1974  3.499470 -6.125844 -0.22290860 -8.262256
8 ## 5  italy  1975  3.515680 -6.170100 -0.03270913 -8.217666
9 ## 6  spain  1978  3.620444 -5.285959  0.62141374 -8.634697
10 ## 7  italy  1972  3.629243 -6.205654 -0.21508857 -8.380176
11 ## 8  italy  1971  3.648205 -6.222602 -0.14822267 -8.470968
12 ## 9  spain  1977  3.650735 -5.299089  0.52627167 -8.727218
13 ## 10 italy  1973  3.652263 -6.157297 -0.32508487 -8.316705
14 ## # ... with 332 more rows

```

If you want to re-order the tibble in descending order of the variable:

```

1 gasoline %>%
2   arrange(desc(lgaspcar))

```

```

1 ## # A tibble: 342 x 6
2 ##   country year lgaspcar lincomep   lrpmg   lcarpcap
3 ##   <chr> <int>   <dbl>   <dbl>   <dbl>   <dbl>
4 ## 1 turkey  1966  6.156644 -7.510895 -0.3564049 -12.95114
5 ## 2 turkey  1960  6.129553 -7.801144 -0.2534082 -13.47518
6 ## 3 turkey  1961  6.106213 -7.786727 -0.3425237 -13.38473
7 ## 4 turkey  1962  6.084587 -7.836272 -0.4082048 -13.24594
8 ## 5 turkey  1968  6.076595 -7.421201 -0.3650739 -12.80735
9 ## 6 turkey  1963  6.075129 -7.631193 -0.2249917 -13.25505
10 ## 7 turkey  1964  6.064601 -7.626898 -0.2521945 -13.21031
11 ## 8 turkey  1967  6.044479 -7.460810 -0.3351502 -12.80185
12 ## 9  japan  1960  5.995287 -6.986196 -0.1453227 -12.23508
13 ## 10 turkey  1965  5.823046 -7.622027 -0.2934761 -12.87934
14 ## # ... with 332 more rows

```

arrange's documentation alerts the user that re-ordering by group is only possible by explicitly specifying an option:

```

1 gasoline %>%
2   filter(year %in% seq(1960, 1963)) %>%
3   group_by(country) %>%
4   arrange(desc(lgaspcar), .by_group = TRUE)

1 ## # A tibble: 72 x 6
2 ## # Groups:   country [18]
3 ##   country year lgaspcar lincomep   lrpmg   lcarpcap
4 ##   <chr> <int>   <dbl>   <dbl>   <dbl>   <dbl>
5 ## 1 austria  1960  4.173244 -6.474277 -0.3345476 -9.766840
6 ## 2 austria  1961  4.100989 -6.426006 -0.3513276 -9.608622
7 ## 3 austria  1962  4.073177 -6.407308 -0.3795177 -9.457257
8 ## 4 austria  1963  4.059509 -6.370679 -0.4142514 -9.343155
9 ## 5 belgium  1960  4.164016 -6.215091 -0.1657096 -9.405527
10 ## 6 belgium  1961  4.124356 -6.176843 -0.1717310 -9.303149
11 ## 7 belgium  1962  4.075962 -6.129638 -0.2222914 -9.218070
12 ## 8 belgium  1963  4.001266 -6.094019 -0.2504623 -9.114932
13 ## 9  canada  1960  4.855238 -5.889713 -0.9721065 -8.378917
14 ## 10 canada  1962  4.850533 -5.844552 -0.9786076 -8.320512
15 ## # ... with 62 more rows

```

4.3.6 tally() and count()

tally() and count() count the number of observations in your data. I believe count() is the more useful of the two, as it counts the number of observations within a group that you can provide:

```
1 gasoline %>%
2   count(country)
```

```
1 ## # A tibble: 18 x 2
2 ##       country      n
3 ##       <chr> <int>
4 ## 1    austria    19
5 ## 2    belgium    19
6 ## 3    canada     19
7 ## 4    denmark    19
8 ## 5    france     19
9 ## 6    germany     19
10 ## 7    greece     19
11 ## 8    ireland     19
12 ## 9     italy     19
13 ## 10   japan      19
14 ## 11 netherlands  19
15 ## 12   norway     19
16 ## 13   spain      19
17 ## 14   sweden     19
18 ## 15 switzerland  19
19 ## 16   turkey     19
20 ## 17     u.k.      19
21 ## 18    u.s.a.     19
```

There's also `add_count()` which adds the column to the data:

```
1 gasoline %>%
2   add_count(country)
```

```
1 ## # A tibble: 342 x 7
2 ##   country year lgaspcar lincomep   lrpmg   lcarpcap     n
3 ##   <chr> <int>   <dbl>   <dbl>   <dbl>   <dbl> <int>
4 ## 1 austria 1960 4.173244 -6.474277 -0.3345476 -9.766840    19
5 ## 2 austria 1961 4.100989 -6.426006 -0.3513276 -9.608622    19
6 ## 3 austria 1962 4.073177 -6.407308 -0.3795177 -9.457257    19
7 ## 4 austria 1963 4.059509 -6.370679 -0.4142514 -9.343155    19
8 ## 5 austria 1964 4.037689 -6.322247 -0.4453354 -9.237739    19
9 ## 6 austria 1965 4.033983 -6.294668 -0.4970607 -9.123903    19
10 ## 7 austria 1966 4.047537 -6.252545 -0.4668377 -9.019822    19
```

```

11 ## 8 austria 1967 4.052911 -6.234581 -0.5058834 -8.934403 19
12 ## 9 austria 1968 4.045507 -6.206894 -0.5224125 -8.847967 19
13 ## 10 austria 1969 4.046355 -6.153140 -0.5591105 -8.788686 19
14 ## # ... with 332 more rows

```

`add_count()` is a shortcut for the following code:

```

1 gasoline %>%
2   group_by(country) %>%
3   mutate(n = n())

```

```

1 ## # A tibble: 342 x 7
2 ## # Groups:   country [18]
3 ##   country year lgaspcar lincomep   lrpmg lcarpcap     n
4 ##   <chr> <int>   <dbl>   <dbl>   <dbl>   <dbl> <int>
5 ## 1 austria 1960 4.173244 -6.474277 -0.3345476 -9.766840 19
6 ## 2 austria 1961 4.100989 -6.426006 -0.3513276 -9.608622 19
7 ## 3 austria 1962 4.073177 -6.407308 -0.3795177 -9.457257 19
8 ## 4 austria 1963 4.059509 -6.370679 -0.4142514 -9.343155 19
9 ## 5 austria 1964 4.037689 -6.322247 -0.4453354 -9.237739 19
10 ## 6 austria 1965 4.033983 -6.294668 -0.4970607 -9.123903 19
11 ## 7 austria 1966 4.047537 -6.252545 -0.4668377 -9.019822 19
12 ## 8 austria 1967 4.052911 -6.234581 -0.5058834 -8.934403 19
13 ## 9 austria 1968 4.045507 -6.206894 -0.5224125 -8.847967 19
14 ## 10 austria 1969 4.046355 -6.153140 -0.5591105 -8.788686 19
15 ## # ... with 332 more rows

```

where `n()` is a `dplyr` function that can only be used within `summarise()`, `mutate()` and `filter()`.

4.3.7 Joining tibbles with `full_join()`, `left_join()`, `right_join()` and all the others

I will end this section on `dplyr` with the very useful verbs: the `*_join()` verbs. Let's first start by loading another dataset from the `plm` package. `SumHes` and let's convert it to tibble and rename it:

```

1 data(SumHes, package = "plm")
2
3 pwt <- SumHes %>%
4   as_tibble() %>%
5   mutate(country = tolower(country))

```

Let's take a quick look at the data:

```

1 glimpse(pwt)

1 ## Observations: 3,250
2 ## Variables: 7
3 ## $ year      <int> 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968, ...
4 ## $ country   <chr> "algeria", "algeria", "algeria", "algeria", "algeria", ...
5 ## $ opec      <fctr> no, no, no, no, no, no, no, no, no, no, no, no, no, n...
6 ## $ com       <fctr> no, no, no, no, no, no, no, no, no, no, no, no, no, n...
7 ## $ pop       <int> 10800, 11016, 11236, 11460, 11690, 11923, 12267, 12622...
8 ## $ gdp       <int> 1723, 1599, 1275, 1517, 1589, 1584, 1548, 1600, 1758, ...
9 ## $ sr        <dbl> 19.9, 21.1, 15.0, 13.9, 10.6, 11.0, 8.3, 11.3, 15.1, 1...

```

We can merge both gasoline and pwt by country and year, as these two variables are common to both datasets. There are more countries and years in the pwt dataset, so when merging both, and depending on which function you use, you will either have NA's for the variables where there is no match, or rows that will be dropped. Let's start with `full_join`:

```

1 gas_pwt_full <- gasoline %>%
2   full_join(pwt, by = c("country", "year"))

```

Let's see which countries and years are included:

```

1 gas_pwt_full %>%
2   count(country, year)

```

```

1 ## # A tibble: 3,269 x 3
2 ##   country year    n
3 ##   <chr> <int> <int>
4 ## 1 algeria 1960     1
5 ## 2 algeria 1961     1
6 ## 3 algeria 1962     1
7 ## 4 algeria 1963     1
8 ## 5 algeria 1964     1
9 ## 6 algeria 1965     1
10 ## 7 algeria 1966     1
11 ## 8 algeria 1967     1
12 ## 9 algeria 1968     1
13 ## 10 algeria 1969     1
14 ## # ... with 3,259 more rows

```

As you see, every country and year was included, but what happened for, say, the U.S.S.R? This country is in pwt but not in gasoline at all:

```

1 gas_pwt_full %>%
2   filter(country == "u.s.s.r.")

```

```

1 ## # A tibble: 26 x 11
2 ##   country year lgaspcar lincomep lrpmpg lcarpcap opec com  pop
3 ##   <chr> <int>   <dbl>   <dbl> <dbl>   <dbl> <fctr> <fctr> <int>
4 ## 1 u.s.s.r. 1960      NA      NA    NA      NA    no    yes 214400
5 ## 2 u.s.s.r. 1961      NA      NA    NA      NA    no    yes 217896
6 ## 3 u.s.s.r. 1962      NA      NA    NA      NA    no    yes 221449
7 ## 4 u.s.s.r. 1963      NA      NA    NA      NA    no    yes 225060
8 ## 5 u.s.s.r. 1964      NA      NA    NA      NA    no    yes 227571
9 ## 6 u.s.s.r. 1965      NA      NA    NA      NA    no    yes 230109
10 ## 7 u.s.s.r. 1966      NA      NA    NA      NA    no    yes 232676
11 ## 8 u.s.s.r. 1967      NA      NA    NA      NA    no    yes 235272
12 ## 9 u.s.s.r. 1968      NA      NA    NA      NA    no    yes 237896
13 ## 10 u.s.s.r. 1969      NA      NA    NA      NA    no    yes 240550
14 ## # ... with 16 more rows, and 2 more variables: gdp <int>, sr <dbl>

```

As you probably guessed, the variables from gasoline that are not included in pwt are filled with NAs. One could remove all these lines and only keep countries for which these variables are not NA everywhere with `filter()`, but there is a simpler solution:


```
1 gas_pwt_inner <- gasoline %>%
2   inner_join(pwt, by = c("country", "year"))
```

Let's use the `tabyl()` from the `janitor` packages which is a very nice alternative to the `table()` function from base R:

```
1 library(janitor)
2
3 gas_pwt_inner %>%
4   tabyl(country)
```

```
1 ##      country n    percent
2 ## 1      austria 19 0.05882353
3 ## 2      belgium 19 0.05882353
4 ## 3       canada 19 0.05882353
5 ## 4      denmark 19 0.05882353
6 ## 5       france 19 0.05882353
7 ## 6       greece 19 0.05882353
8 ## 7      ireland 19 0.05882353
9 ## 8        italy 19 0.05882353
10 ## 9        japan 19 0.05882353
11 ## 10 netherlands 19 0.05882353
12 ## 11       norway 19 0.05882353
13 ## 12        spain 19 0.05882353
14 ## 13        sweden 19 0.05882353
15 ## 14  switzerland 19 0.05882353
16 ## 15        turkey 19 0.05882353
17 ## 16          u.k. 19 0.05882353
18 ## 17        u.s.a. 19 0.05882353
```

Only countries with values in both datasets were returned. It's almost every country from `gasoline`, apart from Germany (called "germany west" in `pwt` and "germany" in `gasoline`). I left it as is to provide an example of a country not in `pwt`). Let's also look at the variables:

```
1 glimpse(gas_pwt_inner)
```

```

1 ## Observations: 323
2 ## Variables: 11
3 ## $ country <chr> "austria", "austria", "austria", "austria", "austria"...
4 ## $ year <int> 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968,...
5 ## $ lgaspcar <dbl> 4.173244, 4.100989, 4.073177, 4.059509, 4.037689, 4.0...
6 ## $ lincomep <dbl> -6.474277, -6.426006, -6.407308, -6.370679, -6.322247...
7 ## $ lrpmg <dbl> -0.3345476, -0.3513276, -0.3795177, -0.4142514, -0.44...
8 ## $ lcarpcap <dbl> -9.766840, -9.608622, -9.457257, -9.343155, -9.237739...
9 ## $ opec <fctr> no, no, no, no, no, no, no, no, no, no, no, no, ...
10 ## $ com <fctr> no, no, no, no, no, no, no, no, no, no, no, no, no, ...
11 ## $ pop <int> 7048, 7087, 7130, 7172, 7215, 7255, 7308, 7338, 7362,...
12 ## $ gdp <int> 5143, 5388, 5481, 5688, 5978, 6144, 6437, 6596, 6847,...
13 ## $ sr <dbl> 24.3, 24.5, 23.3, 22.9, 25.2, 25.2, 26.7, 25.6, 25.7,...

```

The variables from both datasets are in the joined data.

Contrast this to `semi_join()`:

```

1 gas_pwt_semi <- gasoline %>%
2   semi_join(pwt, by = c("country", "year"))
3
4 glimpse(gas_pwt_semi)

```

```

1 ## Observations: 323
2 ## Variables: 6
3 ## $ country <chr> "austria", "austria", "austria", "austria", "austria"...
4 ## $ year <int> 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968,...
5 ## $ lgaspcar <dbl> 4.173244, 4.100989, 4.073177, 4.059509, 4.037689, 4.0...
6 ## $ lincomep <dbl> -6.474277, -6.426006, -6.407308, -6.370679, -6.322247...
7 ## $ lrpmg <dbl> -0.3345476, -0.3513276, -0.3795177, -0.4142514, -0.44...
8 ## $ lcarpcap <dbl> -9.766840, -9.608622, -9.457257, -9.343155, -9.237739...

```

```

1 gas_pwt_semi %>%
2   tabyl(country)

```

```

1  ##      country n    percent
2  ## 1      austria 19 0.05882353
3  ## 2      belgium 19 0.05882353
4  ## 3      canada 19 0.05882353
5  ## 4      denmark 19 0.05882353
6  ## 5      france 19 0.05882353
7  ## 6      greece 19 0.05882353
8  ## 7      ireland 19 0.05882353
9  ## 8      italy 19 0.05882353
10 ## 9      japan 19 0.05882353
11 ## 10     netherlands 19 0.05882353
12 ## 11     norway 19 0.05882353
13 ## 12     spain 19 0.05882353
14 ## 13     sweden 19 0.05882353
15 ## 14     switzerland 19 0.05882353
16 ## 15     turkey 19 0.05882353
17 ## 16      u.k. 19 0.05882353
18 ## 17     u.s.a. 19 0.05882353

```

Only columns of gasoline are returned, and only rows of gasoline that were matched with rows from pwt. `semi_join()` is not a commutative operation:

```

1  pwt_gas_semi <- pwt %>%
2    semi_join(gasoline, by = c("country", "year"))
3
4  glimpse(pwt_gas_semi)

```

```

1  ## Observations: 323
2  ## Variables: 7
3  ## $ year      <int> 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968, ...
4  ## $ country   <chr> "canada", "canada", "canada", "canada", "canada", "can...
5  ## $ opec      <fctr> no, no, no, no, no, no, no, no, no, no, no, no, n...
6  ## $ com       <fctr> no, no, no, no, no, no, no, no, no, no, no, no, n...
7  ## $ pop       <int> 17910, 18270, 18614, 18963, 19326, 19678, 20049, 20411...
8  ## $ gdp       <int> 7258, 7261, 7605, 7876, 8244, 8664, 9093, 9231, 9582, ...
9  ## $ sr       <dbl> 22.7, 21.5, 22.1, 21.9, 22.9, 24.8, 25.4, 23.1, 22.6, ...

```

```

1 gas_pwt_semi %>%
2   tabyl(country)

1 ##      country  n    percent
2 ## 1    austria 19 0.05882353
3 ## 2    belgium 19 0.05882353
4 ## 3    canada  19 0.05882353
5 ## 4    denmark 19 0.05882353
6 ## 5    france  19 0.05882353
7 ## 6    greece  19 0.05882353
8 ## 7    ireland 19 0.05882353
9 ## 8    italy   19 0.05882353
10 ## 9    japan  19 0.05882353
11 ## 10 netherlands 19 0.05882353
12 ## 11    norway 19 0.05882353
13 ## 12    spain  19 0.05882353
14 ## 13    sweden 19 0.05882353
15 ## 14    switzerland 19 0.05882353
16 ## 15    turkey 19 0.05882353
17 ## 16      u.k.  19 0.05882353
18 ## 17    u.s.a. 19 0.05882353

```

The rows are the same, but not the columns.

`left_join()` and `right_join()` return all the rows from either the dataset that is on the “left” (the first argument of the fonction) or on the “right” (the second argument of the function) but all columns from both datasets. So depending on which countries you’re interested in, you’re going to use either one of these functions:

```

1 gas_pwt_left <- gasoline %>%
2   left_join(pwt, by = c("country", "year"))
3
4 gas_pwt_left %>%
5   tabyl(country)

```

```

1  ##          country  n    percent
2  ## 1          austria 19 0.05555556
3  ## 2          belgium 19 0.05555556
4  ## 3          canada  19 0.05555556
5  ## 4          denmark 19 0.05555556
6  ## 5          france  19 0.05555556
7  ## 6          germany 19 0.05555556
8  ## 7          greece  19 0.05555556
9  ## 8          ireland 19 0.05555556
10 ## 9          italy   19 0.05555556
11 ## 10         japan   19 0.05555556
12 ## 11 netherlands 19 0.05555556
13 ## 12         norway  19 0.05555556
14 ## 13         spain   19 0.05555556
15 ## 14         sweden  19 0.05555556
16 ## 15 switzerland 19 0.05555556
17 ## 16         turkey  19 0.05555556
18 ## 17          u.k.   19 0.05555556
19 ## 18         u.s.a.  19 0.05555556

```

```

1  gas_pwt_right <- gasoline %>%
2    right_join(pwt, by = c("country", "year"))
3
4  gas_pwt_right %>%
5    tabyl(country)

```

```

1  ##          country  n percent
2  ## 1          algeria 26  0.008
3  ## 2          angola  26  0.008
4  ## 3          argentina 26  0.008
5  ## 4          australia 26  0.008
6  ## 5          austria  26  0.008
7  ## 6          bangladesh 26  0.008
8  ## 7          barbados 26  0.008
9  ## 8          belgium  26  0.008
10 ## 9          benin    26  0.008
11 ## 10         bolivia  26  0.008
12 ## 11         botswana 26  0.008
13 ## 12         brazil   26  0.008
14 ## 13         burkina faso 26  0.008
15 ## 14         burundi  26  0.008

```

16	##	15	cameroon	26	0.008
17	##	16	canada	26	0.008
18	##	17	cape verde is.	26	0.008
19	##	18	central afr.r.	26	0.008
20	##	19	chad	26	0.008
21	##	20	chile	26	0.008
22	##	21	china	26	0.008
23	##	22	colombia	26	0.008
24	##	23	comoros	26	0.008
25	##	24	congo	26	0.008
26	##	25	costa rica	26	0.008
27	##	26	cyprus	26	0.008
28	##	27	czechoslovakia	26	0.008
29	##	28	denmark	26	0.008
30	##	29	dominican rep.	26	0.008
31	##	30	ecuador	26	0.008
32	##	31	egypt	26	0.008
33	##	32	el salvador	26	0.008
34	##	33	ethiopia	26	0.008
35	##	34	fiji	26	0.008
36	##	35	finland	26	0.008
37	##	36	france	26	0.008
38	##	37	gabon	26	0.008
39	##	38	gambia	26	0.008
40	##	39	germany west	26	0.008
41	##	40	ghana	26	0.008
42	##	41	greece	26	0.008
43	##	42	guatemala	26	0.008
44	##	43	guinea	26	0.008
45	##	44	guinea-biss	26	0.008
46	##	45	guyana	26	0.008
47	##	46	haiti	26	0.008
48	##	47	honduras	26	0.008
49	##	48	hong kong	26	0.008
50	##	49	iceland	26	0.008
51	##	50	india	26	0.008
52	##	51	indonesia	26	0.008
53	##	52	iran	26	0.008
54	##	53	iraq	26	0.008
55	##	54	ireland	26	0.008
56	##	55	israel	26	0.008
57	##	56	italy	26	0.008

```

58 ## 57      ivory coast 26  0.008
59 ## 58      jamaica 26  0.008
60 ## 59      japan 26  0.008
61 ## 60      jordan 26  0.008
62 ## 61      kenya 26  0.008
63 ## 62      korea 26  0.008
64 ## 63      lesotho 26  0.008
65 ## 64      liberia 26  0.008
66 ## 65      luxembourg 26  0.008
67 ## 66      madagascar 26  0.008
68 ## 67      malawi 26  0.008
69 ## 68      malaysia 26  0.008
70 ## 69      mali 26  0.008
71 ## 70      malta 26  0.008
72 ## 71      mauritania 26  0.008
73 ## 72      mauritius 26  0.008
74 ## 73      mexico 26  0.008
75 ## 74      morocco 26  0.008
76 ## 75      mozambique 26  0.008
77 ## 76      myanmar 26  0.008
78 ## 77      namibia 26  0.008
79 ## 78      nepal 26  0.008
80 ## 79      netherlands 26  0.008
81 ## 80      new zealand 26  0.008
82 ## 81      nicaragua 26  0.008
83 ## 82      niger 26  0.008
84 ## 83      nigeria 26  0.008
85 ## 84      norway 26  0.008
86 ## 85      pakistan 26  0.008
87 ## 86      panama 26  0.008
88 ## 87      papua n.guinea 26  0.008
89 ## 88      paraguay 26  0.008
90 ## 89      peru 26  0.008
91 ## 90      philippines 26  0.008
92 ## 91      portugal 26  0.008
93 ## 92      puerto rico 26  0.008
94 ## 93      reunion 26  0.008
95 ## 94      romania 26  0.008
96 ## 95      rwanda 26  0.008
97 ## 96      saudi arabia 26  0.008
98 ## 97      senegal 26  0.008
99 ## 98      seychelles 26  0.008

```

```

100 ## 99      singapore 26  0.008
101 ## 100     somalia 26  0.008
102 ## 101    south africa 26  0.008
103 ## 102      spain 26  0.008
104 ## 103    sri lanka 26  0.008
105 ## 104     suriname 26  0.008
106 ## 105     swaziland 26  0.008
107 ## 106      sweden 26  0.008
108 ## 107    switzerland 26  0.008
109 ## 108      syria 26  0.008
110 ## 109      taiwan 26  0.008
111 ## 110     tanzania 26  0.008
112 ## 111     thailand 26  0.008
113 ## 112      togo 26  0.008
114 ## 113 trinidad&tobago 26  0.008
115 ## 114      tunisia 26  0.008
116 ## 115      turkey 26  0.008
117 ## 116       u.k. 26  0.008
118 ## 117      u.s.a. 26  0.008
119 ## 118      u.s.s.r. 26  0.008
120 ## 119      uganda 26  0.008
121 ## 120      uruguay 26  0.008
122 ## 121     venezuela 26  0.008
123 ## 122     yugoslavia 26  0.008
124 ## 123      zaire 26  0.008
125 ## 124      zambia 26  0.008
126 ## 125     zimbabwe 26  0.008

```

The last merge function is `anti_join()`:

```

1 gas_pwt_anti <- gasoline %>%
2   anti_join(pwt, by = c("country", "year"))
3
4 glimpse(gas_pwt_anti)

```



```

1 ## Observations: 19
2 ## Variables: 6
3 ## $ country <chr> "germany", "germany", "germany", "germany", "germany"...
4 ## $ year <int> 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968,...
5 ## $ lgaspcar <dbl> 3.916953, 3.885345, 3.871484, 3.848782, 3.868993, 3.8...
6 ## $ lincomep <dbl> -6.159837, -6.120923, -6.094258, -6.068361, -6.013442...
7 ## $ lrpmpg <dbl> -0.1859108, -0.2309538, -0.3438417, -0.3746467, -0.39...
8 ## $ lcarpcap <dbl> -9.342481, -9.183841, -9.037280, -8.913630, -8.811013...

```

```

1 gas_pwt_anti %>%
2   tabyl(country)

```

```

1 ##   country  n percent
2 ## 1 germany 19      1

```

`gas_pwt_anti` has the columns the `gasoline` dataset as well as the only country from `gasoline` that is not in `pwt`: “germany”.

4.4 Tidy your data with `tidyr`

4.5 Functional programming with `purrr` and `purrrlyr`

4.5.1 Mapping and reducing with `purrr`, continued

We’ve already seen how to map functions to elements of lists and how to reduce lists to a single value in [Mapping and Reducing: the `purrr` way](#)¹⁸, and in this section we’re going to dig a little deeper into the `purrr` package. We know the standard `map()` function, which returns a list, but there are a number of variants of this function. `map_dbl()` returns an atomic vector of doubles:

```

1 map_dbl(numbers, sqrt_newton, init = 1)

1 ## [1] 4.000001 5.000023 6.000253 7.000000 8.000002 9.000011

```

`map_chr()` returns an atomic vector of strings:

¹⁸[fprog.html#map_reduce_purrr](#)

```
1 map_chr(numbers, sqrt_newton, init = 1)
```

```
1 ## [1] "4.000001" "5.000023" "6.000253" "7.000000" "8.000002" "9.000011"
```

`map_lgl()` returns an atomic vector of TRUE or FALSE:

```
1 divisible <- function(x, y){
2   if_else(x %% y == 0, TRUE, FALSE)
3 }
4
5 map_lgl(seq(1:100), divisible, 3)
```

```
1 ## [1] FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE
2 ## [12] TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE
3 ## [23] FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE
4 ## [34] FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE
5 ## [45] TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE
6 ## [56] FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE
7 ## [67] FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE
8 ## [78] TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE
9 ## [89] FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE FALSE FALSE TRUE
10 ## [100] FALSE
```

There are also other interesting variants, such as `map_if()`:

```
1 a <- seq(1,10)
2
3 map_if(a, (function(x) divisible(x, 2)), sqrt)
```

```
1 ## [[1]]
2 ## [1] 1
3 ##
4 ## [[2]]
5 ## [1] 1.414214
6 ##
7 ## [[3]]
8 ## [1] 3
9 ##
```

```

10 ## [[4]]
11 ## [1] 2
12 ##
13 ## [[5]]
14 ## [1] 5
15 ##
16 ## [[6]]
17 ## [1] 2.44949
18 ##
19 ## [[7]]
20 ## [1] 7
21 ##
22 ## [[8]]
23 ## [1] 2.828427
24 ##
25 ## [[9]]
26 ## [1] 9
27 ##
28 ## [[10]]
29 ## [1] 3.162278

```

I used `map_if()` to take the square root of only those numbers in vector `a` that are divisible by 2, by using an anonymous function that checks if a number is divisible by 2 (by wrapping `divisible()`).

`map_at()` is similar to `map_if()` but maps the function at a position specified by the user:

```

1 map_at(numbers, c(1, 3), sqrt)

1 ## [[1]]
2 ## [1] 4
3 ##
4 ## [[2]]
5 ## [1] 25
6 ##
7 ## [[3]]
8 ## [1] 6
9 ##
10 ## [[4]]
11 ## [1] 49
12 ##
13 ## [[5]]
14 ## [1] 64

```

```

15 ##
16 ## [[6]]
17 ## [1] 81

```

or if you have a named list:

```

1 recipe <- list("spam" = 1, "eggs" = 3, "bacon" = 10)
2
3 map_at(recipe, "bacon", `*`, 2)

```

```

1 ## $spam
2 ## [1] 1
3 ##
4 ## $eggs
5 ## [1] 3
6 ##
7 ## $bacon
8 ## [1] 20

```

I used `map_at()` to double the quantity of bacon in the recipe (by using the `*` function, and specifying its second argument, 2. Try the following in the command prompt: ``*`(3, 4)`).

`map2()` is the equivalent of `mapply()` and `pmap()` is the generalisation of `map2()` for more than 2 arguments:

```

1 print(a)

```

```

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

```

```

1 b <- seq(1, 2, length.out = 10)
2
3 print(b)

```

```

1 ## [1] 1.000000 1.111111 1.222222 1.333333 1.444444 1.555556 1.666667
2 ## [8] 1.777778 1.888889 2.000000

```

```
1 map2(a, b, `*`)
```

```
1 ## [[1]]
2 ## [1] 1
3 ##
4 ## [[2]]
5 ## [1] 2.222222
6 ##
7 ## [[3]]
8 ## [1] 3.666667
9 ##
10 ## [[4]]
11 ## [1] 5.333333
12 ##
13 ## [[5]]
14 ## [1] 7.222222
15 ##
16 ## [[6]]
17 ## [1] 9.333333
18 ##
19 ## [[7]]
20 ## [1] 11.66667
21 ##
22 ## [[8]]
23 ## [1] 14.22222
24 ##
25 ## [[9]]
26 ## [1] 17
27 ##
28 ## [[10]]
29 ## [1] 20
```

```
1 n <- seq(1:10)
2
3 pmap(list(a, b, n), rnorm)
```

```

1  ## [[1]]
2  ## [1] 1.553918
3  ##
4  ## [[2]]
5  ## [1] 0.9872877 0.4991858
6  ##
7  ## [[3]]
8  ## [1] 0.08080922 -0.86189871 0.59847039
9  ##
10 ## [[4]]
11 ## [1] -3.728252 10.009157 6.165181 -3.159101
12 ##
13 ## [[5]]
14 ## [1] -0.5699797 -0.8888323 5.3442700 1.0275991 2.7110370
15 ##
16 ## [[6]]
17 ## [1] 1.3842750 1.2983328 9.7671693 0.2009296 10.6543792 -7.7369613
18 ##
19 ## [[7]]
20 ## [1] 5.758963 2.533646 3.178258 4.324143 -1.849598 -0.665785 -5.463361
21 ##
22 ## [[8]]
23 ## [1] -6.796552 4.206007 5.363456 2.201812 9.155918 18.178455
24 ## [7] -2.150472 -16.695573
25 ##
26 ## [[9]]
27 ## [1] 10.9405356 -4.4939180 -4.3031887 11.1190312 -0.6740682 -9.0975705
28 ## [7] 3.5206202 0.6388666 1.9407666
29 ##
30 ## [[10]]
31 ## [1] 5.8528040 -1.7066003 8.4437655 -0.2048656 5.3178196 12.9683901
32 ## [7] 6.3518149 -1.2593159 13.4880762 11.9350386

```

4.5.2 safely() and possibly()

safely() and possibly() are very useful functions. Consider the following situation:

```

1  a <- list("a", 4, 5)
2
3  sqrt(a)

```

```
1 Error in sqrt(a) : non-numeric argument to mathematical function
```

Using `map()` or `Map()` will result in a similar error. `safely()` is an higher-order function that takes one function as an argument and executes it *safely*, meaning the execution of the function will not stop if there is an error. The error message gets captured alongside valid results.

```
1 a <- list("a", 4, 5)
2
3 safe_sqrt <- safely(sqrt)
4
5 map(a, safe_sqrt)

1 ## [[1]]
2 ## [[1]]$result
3 ## NULL
4 ##
5 ## [[1]]$error
6 ## <simpleError in sqrt(x = x): non-numeric argument to mathematical function>
7 ##
8 ##
9 ## [[2]]
10 ## [[2]]$result
11 ## [1] 2
12 ##
13 ## [[2]]$error
14 ## NULL
15 ##
16 ##
17 ## [[3]]
18 ## [[3]]$result
19 ## [1] 2.236068
20 ##
21 ## [[3]]$error
22 ## NULL
```

`possibly()` works similarly, but also allows you to specify a return value in case of an error:

```

1 possible_sqrt <- possibly(sqrt, otherwise = NA_real_)
2
3 map(a, possible_sqrt)

```

```

1 ## [[1]]
2 ## [1] NA
3 ##
4 ## [[2]]
5 ## [1] 2
6 ##
7 ## [[3]]
8 ## [1] 2.236068

```

Of course, in this particular example, the same effect could be obtained way more easily:

```

1 sqrt(as.numeric(a))

1 ## Warning: NAs introduced by coercion
2
3 ## [1] NA 2.000000 2.236068

```

However, in some situations, this trick does not work as intended (or at all), so `possibly()` and `safely()` are the way to go.

4.5.3 «Transposing lists»

Another interesting function is `transpose()`. It is not an alternative to the function `t()` from base but, has a similar effect. `transpose()` works on lists. Let's take a look at the example from before:

```

1 safe_sqrt <- safely(sqrt, otherwise = NA_real_)
2
3 map(a, safe_sqrt)

```



```

1  ## [[1]]
2  ## [[1]]$result
3  ## [1] NA
4  ##
5  ## [[1]]$error
6  ## <simpleError in sqrt(x = x): non-numeric argument to mathematical function>
7  ##
8  ##
9  ## [[2]]
10 ## [[2]]$result
11 ## [1] 2
12 ##
13 ## [[2]]$error
14 ## NULL
15 ##
16 ##
17 ## [[3]]
18 ## [[3]]$result
19 ## [1] 2.236068
20 ##
21 ## [[3]]$error
22 ## NULL

```

The output is a list with the first element being a list with a result and an error message. One might want to have all the results in a single list, and all the error messages in another list. This is safe with transpose:

```

1  transpose(map(a, safe_sqrt))

1  ## $result
2  ## $result[[1]]
3  ## [1] NA
4  ##
5  ## $result[[2]]
6  ## [1] 2
7  ##
8  ## $result[[3]]
9  ## [1] 2.236068
10 ##
11 ##
12 ## $error

```

```

13 ## $error[[1]]
14 ## <simpleError in sqrt(x = x): non-numeric argument to mathematical function>
15 ##
16 ## $error[[2]]
17 ## NULL
18 ##
19 ## $error[[3]]
20 ## NULL

```

4.6 Special packages for special kinds of data: forcats, lubridate, and stringr

4.6.1 Factor variables

Factor variables are very useful but not very easy to manipulate. `forcats` contains very useful functions that make working on factor variables painless. I use mainly three functions in my work, `fct_recode()`, `fct_relevel()` and `fct_relabel()`, so that's what I'll be showing.

4.6.1.1 `fct_recode()`

4.7 Exercises

1. Suppose you have an Excel workbook that contains data on three sheets. Create a function that reads entire workbooks, and that returns a list of tibbles, where each tibble is the data of one sheet (download the example Excel workbook, `example_workbook.xlsx`, from the `assets` folder on the books github).
2. Use one of the `map()` functions to combine two lists into one. Consider the following two lists:

```

1 mediterranean <- list("starters" = list("humous", "lasagna"), "dishes" = list("s\
2 ardines", "olives"))
3
4 continental <- list("starters" = list("pea soup", "terrine"), "dishes" = list("f\
5 rikadelle", "sauerkraut"))

```

The result we'd like to have would look like this:

```
1 $starters
2 $starters[[1]]
3 [1] "humous"
4
5 $starters[[2]]
6 [1] "olives"
7
8 $starters[[3]]
9 [1] "pea soup"
10
11 $starters[[4]]
12 [1] "terriner"
13
14
15 $dishes
16 $dishes[[1]]
17 [1] "sardines"
18
19 $dishes[[2]]
20 [1] "lasagna"
21
22 $dishes[[3]]
23 [1] "frikadelle"
24
25 $dishes[[4]]
26 [1] "sauerkraut"
```

Chapter 5 Programming with the tidyverse

Functions are very powerful because by using them, we avoid repetition. This means that we must be able to write functions that allow the user to abstract over certain things, such as columns names of datasets. So for example, one would like to write a function that would look like that:

```
1 my_function(my_data, column)
```

and in this chapter we will learn how to do that using `dplyr` (version 0.70 or above).

I advise you to also read the “Programming with `dplyr`” vignette [here](#)¹⁹, which explains with great detail the concept I will only skim in this chapter!

Consider the following code:

```
1 data(mtcars)
2 simple_function = function(dataset, col_name){
3   dataset %>%
4     group_by(col_name) %>%
5     summarise(mean_mpg = mean(mpg)) -> dataset
6   return(dataset)
7 }
```

When you try to run this:

```
1 simple_function(mtcars, "cyl")
```

This is the error you get:

```
1 Error in grouped_df_impl(data, unname(vars), drop) :
2   Column `col_name` is unknown
```

R is *literally* looking for a column called `col_name` in the `mtcars` dataset. How to solve this issue and make R understand to not take “cyl” literally as a string, but to interpret it?

One way is to use the `enquo()` function (or `quo()` if working interactively), in conjunction with the `!!` operator introduced with `dplyr` 0.7 (but actually part of the `rlang` package, which gets used by `dplyr` seamlessly). First let’s look at the solution and then I’ll explain how it works:

¹⁹<https://cran.r-project.org/web/packages/dplyr/vignettes/programming.html>

```

1 library(dplyr)
2
3 simple_function = function(dataset, col_name){
4   col_name = enquo(col_name)
5   dataset = dataset %>%
6     group_by(!!col_name) %>%
7     summarise(mean_mpg = mean(mpg))
8   return(dataset)
9 }
10
11
12 simple_function(mtcars, cyl)

```

```

1 ## # A tibble: 3 x 2
2 ##   cyl mean_mpg
3 ##   <dbl>   <dbl>
4 ## 1     4 26.66364
5 ## 2     6 19.74286
6 ## 3     8 15.10000

```

In the above example, I wanted the mean of `mpg` but first by grouping by `cyl`. The `enquo()` to quote the input. This tells your function that the variable `col_name` should be quoted. However then, `filter()` (and the other `dplyr` functions) need to have an unquoted variable name. So `!!()` does this and evaluates its argument.

If you want to use `filter()` with a value that the user has to provide, you can also do that:

```

1 simpleFunction = function(dataset, col_name, value){
2   col_name = enquo(col_name)
3   dataset = dataset %>%
4     filter(!!col_name == value) %>%
5     summarise(mean_cyl = mean(cyl))
6   return(dataset)
7 }
8
9
10 simpleFunction(mtcars, am, 1)

```

```
1 ##    mean_cyl
2 ## 1 5.076923
```

There is something that you must pay attention to in the above example. Notice that I've written:

```
1 filter((!!col_name) == value)
```

and not:

```
1 filter(!col_name == value)
```

I have enclosed `!!col_name` inside parentheses, because `==` has precedence over `!!`.

Let's make this function a bit more general. I hard-coded the variable `cyl` inside the body of the function, but what if you need more flexibility and let the user provide the variable to group by to?

```
1 simpleFunction = function(dataset, group_col, mean_col, value){
2   group_col = enquo(group_col)
3   mean_col = enquo(mean_col)
4   dataset = dataset %>%
5     filter((!!group_col) == value) %>%
6     summarise(mean(!!mean_col))
7   return(dataset)
8 }
9
10
11 simpleFunction(mtcars, am, cyl, 1)
```

```
1 ##    mean((cyl))
2 ## 1      5.076923
```

It is possible to set the name of the column in the output using `:=` instead of `=`:

```

1  simpleFunction = function(dataset, group_col, mean_col, value){
2    group_col = enquo(group_col)
3    mean_col = enquo(mean_col)
4    mean_name = paste0("mean_", mean_col)[2]
5    dataset %>%
6      filter((!!group_col) == value) %>%
7      summarise(!!mean_name := mean(!!mean_col))) -> dataset
8    return(dataset)
9  }
10
11
12  simpleFunction(mtcars, am, cyl, 1)

```

```

1  ##   mean_cyl
2  ## 1 5.076923

```

To get the name of the column I added this line:

```

1  mean_name = paste0("mean_", mean_col)[2]

```

To see what it does, try the following inside an R interpreter (remember to use `quo()` instead of `enquo()` outside functions!):

```

1  paste0("mean_", quo(cyl))

1  ## [1] "mean_~"   "mean_cyl"

```

`enquo()` quotes the input, and with `paste0()` it gets converted to a string that can be used as a column name. However, the `~` is in the way and the output of `paste0()` is a vector of two strings: the correct name is contained in the second element, hence the `[2]`. There might be a more elegant way of doing that, but for now this has been working well for me.

It is also possible to write functions that take any amount of variables the user wants:

```

1 nice_function = function(dataset, ...){
2   variables = quos(...)
3   dataset %>%
4     summarise_at(vars(!!!variables), mean)
5 }
6
7 nice_function(mtcars, mpg, cyl)

```

```

1 ##      mpg      cyl
2 ## 1 20.09062 6.1875

```

You can even be more flexible, and let the user define the functions that will be used by `summarised_at()`:

```

1 nice_function = function(dataset, cols, funcs){
2   dataset %>%
3     summarise_at(vars(!!!cols), funs(!!!funcs))
4 }
5
6
7 list_cols = quos(mpg, cyl)
8
9
10 list_funs = quos(mean, sd, sum)
11
12
13 nice_function(mtcars, list_cols, list_funs)

```

```

1 ##      mpg_mean cyl_mean   mpg_sd   cyl_sd mpg_sum cyl_sum
2 ## 1 20.09062   6.1875 6.026948 1.785922  642.9   198

```

In this last example however, the user has to provide a list of quoted variables and functions.

Chapter 6 Packages

6.1 Why you need your own packages in your life

One of the reasons you might have tried R in the first place is the abundance of packages. As I'm writing these lines (in August 2016), 8922 packages are available on CRAN. That's almost over 9000. This is an absolutely crazy amount of packages! Chances are that if you want to do something, there's a package for that (I'll stop here with the lame references, promise!).

So why the heck should you write your own packages? After all, with 8922 packages you're sure to find something that suits your needs, right? No. Simply because the data sets that you're working with are probably unique to your workplace or maybe what you want to do with them is unique to your needs. You won't find a package that will take care of cleaning *your* data for you.

Ok, but is it necessary to write a package? Why not just write functions inside some scripts and then simply run these scripts? This seems like a valid solution at first. However, it quickly becomes tedious, especially if you have multiple scripts scattered around your computer or inside different subfolders. You'll also have to write the documentation on separate files and these can easily get lost or become outdated.

Having everything inside a package takes care of these headaches for you. And code that is inside packages is very easy to test, especially if you're using Rstudio. It also makes it possible to use the wonderful `covr` package, which tells you which lines in which functions are called by your tests. If some lines are missing, write tests that invoke them and increase the coverage of your tests!

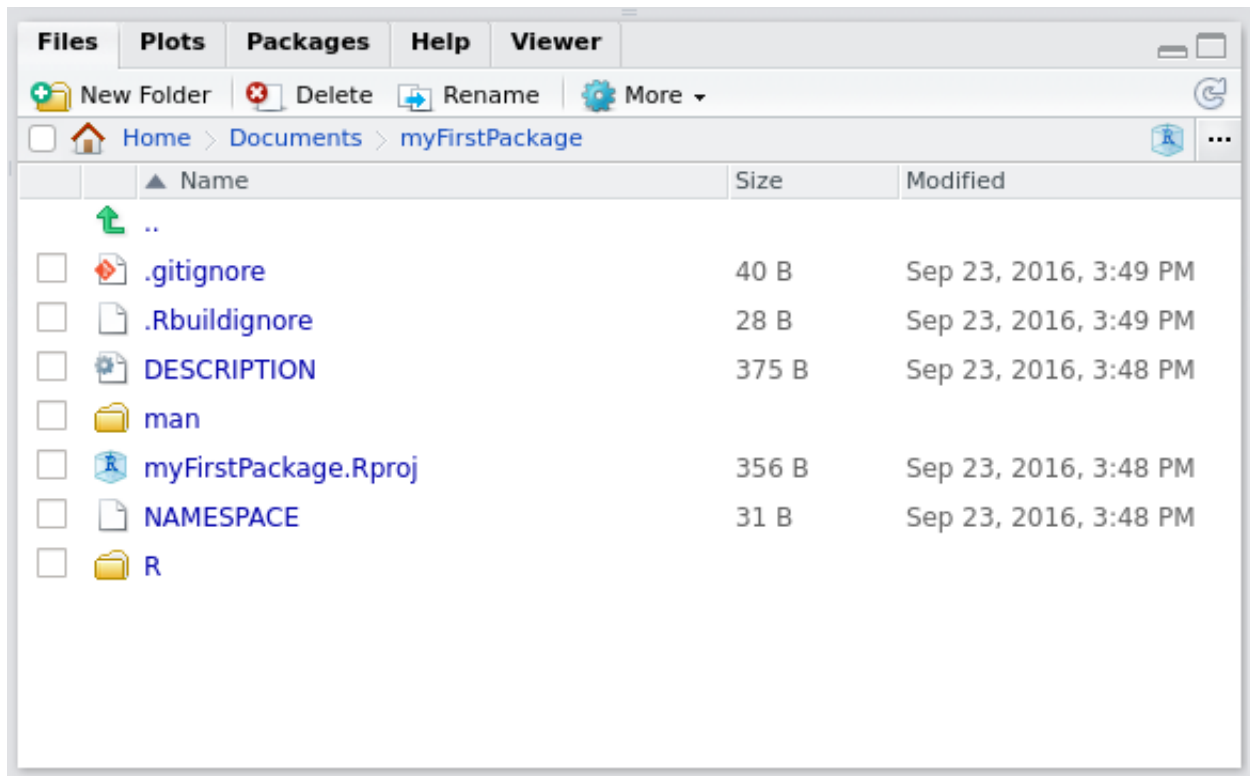
As I mentioned in the introduction, if you want to learn much more than I'll show about packages read Wickham (2014a). I will only show you the basics, but it should be enough to get you productive.

One last thing: if you don't know git, you really should learn git. I won't talk about it here, because there's a ton of books on git, such as Silverman (2013). I learned by reading it and googling whenever I had a problem. Learning git is really worth it, especially if you're collaborating with some colleagues on your packages.

6.2 R packages: the basics

To start writing a package, the easiest way is to load up Rstudio and start a new project, under the *File* menu. If you're starting from scratch, just choose the first option, *New Directory* and then *R package*. Give a new to your package, for example `myFirstPackage` and you can also choose to use git for version control. Now if you check the folder where you chose to save your package, you will see a folder with the same name as your package, and inside this folder a lot of new files and other

folders. The most important folder for now is the `R` folder. This is the folder that will hold your `.R` source code files. You can also see these files and folders inside the *Files* panel from within Rstudio. Rstudio will also have `hello.R` opened, which is a single demo source file inside the `R` folder. You can get rid of this file.



The picture above shows the basic structure of your package. As a first step, create a script called `square_root_loop.R` and put the following code in it:

```

1  sqrt_newton <- function(a, init, eps = 0.01, iter = 100){
2    stopifnot(a >= 0)
3    i <- 1
4    while(abs(init**2 - a) > eps){
5      init <- 1/2 *(init + a/init)
6      i <- i + 1
7      if(i > iter) stop("Maximum number of iterations reached")
8    }
9    return(init)
10 }
```

Then save this script. You can now test your package by building your package, either by clicking on the button named *Build and Reload* button which you can find inside the *Build* pane or by using the following keyboard shortcut: CTRL-SHIFT-B. You will use *Build and Reload* quite often, so I advise

you remember this shortcut! In the next section we will see how we can add documentation to our functions.

6.3 Writing documentation for your functions

Writing documentation for your functions is very streamlined, thanks to the `roxygen2` package. Suppose we want to write documentation for our square root function:

```

1  sqrt_newton <- function(a, init, eps = 0.01, iter = 100){
2    stopifnot(a >= 0)
3    i <- 1
4    while(abs(init**2 - a) > eps){
5      init <- 1/2 *(init + a/init)
6      i <- i + 1
7      if(i > iter) stop("Maximum number of iterations reached")
8    }
9    return(init)
10 }
```

Usually, you would write comments to describe what your function does, what are its inputs and outputs. ‘`roxygen2`’ is a package that turns these comments into documentation. Here is what our function would look like with `roxygen2` type comments:

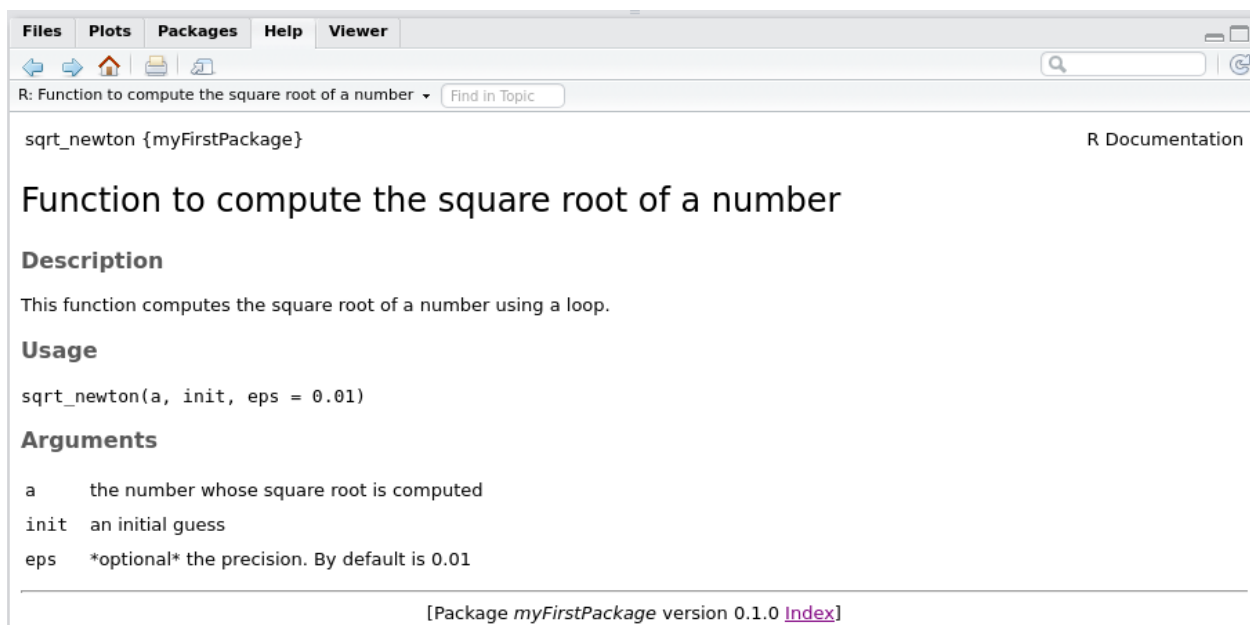
```

1  #' Function to compute the square root of a number
2  #' @param a the number whose square root is computed
3  #' @param init an initial guess
4  #' @param eps *optional* the precision. Default value: 0.01
5  #' @param iter *optional* the number of iteration. Default value: 100
6  #' @description This function computes the square root of a number using a loop.
7  #' @export
8  sqrt_newton <- function(a, init, eps = 0.01, iter = 100){
9    stopifnot(a >= 0)
10   i <- 1
11   while(abs(init**2 - a) > eps){
12     init <- 1/2 *(init + a/init)
13     i <- i + 1
14     if(i > iter) stop("Maximum number of iterations reached")
15   }
16   return(init)
17 }
```

The first difference with standard comments is that roxygen2 type comments start with two symbols, `#'` instead of simply the `#` symbol. Then, after `#'` you can supply different keywords such as `@param`, `@description`, `@export`. These keywords are then used by the `roxygenise()` function from the roxygen package to create the documentation files inside your package. Before roxygen, these documentation files were written in the `.Rd` format by hand. Now these files get created automatically by simply formatting your comments with this specific syntax and then running `roxygen2::roxygenise()` in the command prompt. Try it, you should see the following in the command prompt:

1 Writing `sqrt_newton.Rd`

then you can *Build and Reload* your package again using CTRL-SHIFT-B. If you go check the documentation of your function inside your package, this is what you should see:



There is still a keyword that I did not mention: the `@export` keyword. This keyword is needed if you want your function to be accessible by the user without prepending the package name, like this:

1 `my_package::my_function`

Not using `@export` can be useful though, if you want to have helper functions that are used by your other functions inside your package, and if you wish to not make these functions accessible to the users. In the next subsection, I will mention two files that got created with your package, `NAMESPACE` and `DESCRIPTION`.

6.4 Extra files inside your package and dependencies

6.4.1 The NAMESPACE file

The NAMESPACE file gets generated automatically by roxygen2. You do not have to worry much about it; it lists the functions that are made available to the user when your package gets loaded. If you have helper functions that you do not want to make available to the user, remove the @export keyword from the function comments and the function will not get listed in the NAMESPACE file. There is still a way for the user to access these helper functions like so:

```
1 package::helper_function
```

This is somewhat similar to the concept of private/public methods in object oriented programming.

6.4.2 How can you use functions from other packages inside your package?

There are two solutions for this. Without going into much details, this is mostly handled by the DESCRIPTION file. This file is very important for two reasons: if you want to publish a package on CRAN, this is the file where you specify your name, contact information and the license of the package. But this does not mean that this file is not important if you want to write your own personal package to clean data: this is also the file where you specify the version of your package, but also where you specify the dependencies of your package. You have two fields for this: Imports and Depends. If you are lazy and want to use functions from other packages inside your package, you can list them in the Depends field. However, you should be careful with this approach; indeed, if you named some of your functions the same as functions from other packages, by loading your package after another, you will “overwrite” these functions. For example, imagine you have function `f` in package A but also in package B. If you first load A, and then B, `f` will refer to `B::f`. It’s basically the same problem when attaching datasets, but with functions, so it’s even worse! A cleaner way of solving the problem of using functions from other packages inside your package is to list them in the Imports field, and then use `::` inside your functions. For example:

```
1 my_nice_function <- function(x){  
2   readr::read_csv( bla bla )  
3   bla bla  
4 }
```

This is much cleaner and avoids the problem described above. If you have to use a lot of functions from the same package, having always to use `::` can get annoying quite fast. In these cases, you can use the @import keyword when you document your function:

```

1  #' Function to compute the square root of a number
2  #' @param a the number whose square root is computed
3  #' @param init an initial guess
4  #' @param eps *optional* the precision. Default value: 0.01
5  #' @param iter *optional* the number of iteration. Default value: 100
6  #' @description This function computes the square root of a number using a loop.
7  #' @export
8  #` @import dplyr
9  sqrt_newton <- function(a, init, eps = 0.01, iter = 100){
10     stopifnot(a >= 0)
11     i <- 1
12     while(abs(init**2 - a) > eps){
13         init <- 1/2 *(init + a/init)
14         i <- i + 1
15         if(i > iter) stop("Maximum number of iterations reached")
16     }
17     return(init)
18 }

```

Above, I've added the line:

```

1  #` @import dplyr

```

which allows me to call dplyr functions inside my functions. This is cleaner and does not pollute your session with attached packages. If you check the `NAMESPACE` you will see the following line:

```

1  import(dplyr)

```

6.5 Unit test your package

Now that we know the basics of package creation, we move on to unit testing your package. Unit testing is very useful, but requires some work, especially because you have to run your unit tests often to make them truly worth your time. However running them often can be painful because you have to be careful with the current working directory. The simplest way to do unit testing is to put your functions inside a package and write unit tests for these functions and use Rstudio's keyboard shortcuts to run your tests. First of all, create a folder called `tests` in the root of your package and inside this `tests` folder create another folder, called `testthat`. The `testthat` folder will hold your unit tests. Inside the `tests` folder, create a script called `test_sqrt_newton.R` and put the following code in it:

```

1 library("testthat")
2 library("myFirstPackage")
3
4 test_that("Test sqrt_newton: positive numeric",{
5   eps <- 0.001
6   expected <- 2
7   actual <- sqrt_newton(4, 1, eps = eps)
8   expect_lt(abs(expected - actual), eps)
9 })

```

Save this file and use the following keyboard shortcut: CTRL-SHIFT-T to run your unit test. You will see the following output:

```

1 ==> devtools::test()
2
3 Loading myFirstPackage
4 Loading required package: testthat
5 Testing myFirstPackage
6 .
7 DONE =====

```

You can of course add more unit tests inside the same file. Add the following code to `test_sqrt_newton.R`:

```

1 test_that("Test sqrt_newton: negative numeric",{
2   expect_error(sqrt_newton(-4, 1))
3 })

```

You will now see the following output:

```

1 ==> devtools::test()
2
3 Loading myFirstPackage
4 Loading required package: testthat
5 Testing myFirstPackage
6 ..
7 DONE =====

```

Notice the two `.` above `DONE`. This means that two unit tests passed. If a unit test does not pass, you will of course get notified. For example, add the following test to `test_sqrt_newton.R`:

```

1 test_that("Test sqrt_newton: with a string!",{
2     expect_equal(4, sqrt_newton("WontWork", 1))
3 })

```

and if you try running your tests this is what you will see:

```

1 ==> devtools::test()
2
3 Loading myFirstPackage
4 Loading required package: testthat
5 Testing myFirstPackage
6 ..1
7 Failed -----
8 1. Error: Test sqrt_newton: with a string! (@test_sqrt_newton.R#15) -----
9 non-numeric argument to binary operator
10 1: expect_equal(4, sqrt_newton("WontWork", 1)) at /home/bro/Documents/myFirstPac\
11 kage/inst/tests/test_sqrt_newton.R:15
12 2: compare(object, expected, ...)
13 3: compare.numeric(object, expected, ...)
14 4: all.equal(x, y, tolerance = tolerance, ...)
15 5: all.equal.numeric(x, y, tolerance = tolerance, ...)
16 6: attr.all.equal(target, current, tolerance = tolerance, scale = scale, ...)
17 7: mode(current)
18 8: sqrt_newton("WontWork", 1)
19
20 DONE =====

```

You can then either modify the test if you made a mistake writing the test, or amend your function if your test is correct and needs to pass, but does not because there is an error in your function. For now, simply remove these lines for your `test_sqrt_newton.R` script.

Another interesting feature you should use once in a while, is the *Check Package* command using CTRL-SHIFT-E. This command will find errors and other mistakes and warns you. For example, when I ran this command I got the following report:


```
1  checking DESCRIPTION meta-information ... WARNING
2  Non-standard license specification:
3    What license is it under?
4  Standardizable: FALSE
5
6  checking for code/documentation mismatches ... WARNING
7  Codoc mismatches from documentation object 'sqrt_newton':
8  sqrt_newton
9    Code: function(a, init, eps = 0.01, iter = 100)
10   Docs: function(a, init, eps = 0.01)
11   Argument names in code not in docs:
12     iter
```

Check Package is telling me that I did not specify a license for my package, and that I did not document the `iter` parameter. This command takes some time to run, so do not run it as often as your unit tests, but do not forget about it either!

6.6 Checking the coverage of your unit tests with `covr`

To check the coverage of your package run the following code:

```
1  library("covr")
2
3  cov <- package_coverage()
4
5  shine(cov)
```

The line `shine(cov)` launches an interactive shiny app inside your viewer pane with the following:

```

1 #' Function to compute the square root of a number
2 #' @param a the number whose square root is computed
3 #' @param init an initial guess
4 #' @param eps *optional* the precision. By default is 0.01
5 #' @description This function computes the square root of a number using a loop.
6 #' @export
7 sqrt_newton <- function(a, init, eps = 0.01, iter = 100){
8   stopifnot(a >= 0)
9   i <- 1
10  while(abs(init**2 - a) > eps){
11    init <- 1/2 *(init + a/init)
12    i <- i + 1
13    if(i > iter) stop("Maximum number of iterations reached")
14  }
15  return(init)
16 }

```

We see that no unit test executes the highlighted line. So let's write a unit test to test this line and increase the coverage of our package! Add the following test to `test_sqrt_newton.R`:

```

1 test_that("Test maximum number of iterations",{
2   expect_error(sqrt_newton(10, 1E10, eps=1E-10, 5))
3 })

```

Now if you look at the coverage of the package:

```

1 #' Function to compute the square root of a number
2 #' @param a the number whose square root is computed
3 #' @param init an initial guess
4 #' @param eps *optional* the precision. By default is 0.01
5 #' @description This function computes the square root of a number using a loop.
6 #' @export
7 sqrt_newton <- function(a, init, eps = 0.01, iter = 100){
8   stopifnot(a >= 0)
9   i <- 1
10  while(abs(init**2 - a) > eps){
11    init <- 1/2 *(init + a/init)
12    i <- i + 1
13    if(i > iter) stop("Maximum number of iterations reached")
14  }
15  return(init)
16 }

```

In this example, we used `package_coverage()`, but if you are interested in the coverage of a single function you can use `function_coverage()`, or even `file_coverage()` to get the coverage of a single file. However, I suggest to always run `package_coverage()` since we are working inside a package. There are other functions in the `covr` package that might be useful depending on your needs, so do not hesitate to explore `covr` documentation!

6.7 Wrap-up

- Packages are the easiest way to organize, document and test your code.
- You do not need to take care of paths anymore.
- You do not need to write documentation “by hand”.
- If you use Rstudio, the workflow is very streamlined and you can use version control to keep track of your changes.
- Developing a package is also the easiest way to share your code with colleagues at your company or online.

References

Wickham, Hadley. 2014a. *Advanced R*. CRC Press.

Silverman, Richard E. 2013. *Git Pocket Guide*. “O’Reilly Media, Inc.”

Chapter 7 Unit testing

7.1 Introduction

Let's take a look at [Wikipedia's definition](https://en.wikipedia.org/wiki/Unit_testing)²⁰ of unit testing:

In computer programming, unit testing is a software testing method by which individual units of source code, sets of one or more computer program modules together with associated control data, usage procedures, and operating procedures, are tested to determine whether they are fit for use. Intuitively, one can view a unit as the smallest testable part of an application. In procedural programming, a unit could be an entire module, but it is more commonly an individual function or procedure. In object-oriented programming, a unit is often an entire interface, such as a class, but could be an individual method. Unit tests are short code fragments created by programmers or occasionally by white box testers during the development process. It forms the basis for component testing.

So unit tests are small pieces of code that test your code. They're called *unit* tests, because they test the smallest unit composing your code, in the case of functional programming, the smallest units are functions. You've probably been testing your code *manually* since you've started programming. For example, you would simply do something like this:

```
1  sqrt_newton(4, 1)

1  ## [1] 2.00061
```

and check if the result is equal to 2 and stop there. Usually you would probably write this in the console and then forget about it. If you need to check again, you would write this small test again in the console. But what if some of your functions have to work together with other functions? Maybe changing something in these other functions will indirectly break in other functions. You would have to retest everything together again! In this chapter you will learn the basics of unit testing, which is simply writing these tests in a file, and running this file each time you change your code. If all your unit tests still pass, you can be more confident that your code works as intended.

Unit tests can also be useful to guide you as you program. Some programmers do test-driven development. These programmers start by writing the unit tests first, and then the code to make them pass. This can be useful sometimes, if you don't really know where you should start but know what you want.

²⁰https://en.wikipedia.org/wiki/Unit_testing

7.2 Unit testing with the `testthat` package

We are going to test the function we wrote in the previous chapter, `sqrt_newton()`. The basic steps are:

1. Write a file containing your tests
2. Run the tests

It's very simple! You only need to install the `testthat` package for this. In this section I'll only show you how to write tests and try to illustrate their usefulness. In the next section, we'll see how we can run the tests.

Below is the code that we are going to put in the file `test_my_functions.R`:

```
1 library("testthat")
2
3 test_that("Test sqrt_newton: positive numeric",{
4     expected <- 2
5     actual <- sqrt_newton(4, 1)
6     expect_equal(expected, actual)
7 })
```

The syntax of the test is pretty straightforward. We start with a short description of what the test is about, and then we define two variables: the result we expect, and the actual result that is returned by the function we wish to test. When we run this test (we'll discuss running tests in the next section), this is what we get:

```
1 Error: Test failed: 'Test sqrt_newton: positive numeric'
2 * `expected` not equal to `actual`.
3 1/1 mismatches
4 [1] 2 - 2 == -0.00061
```

This is because the value that `sqrt_newton()` returns is not exactly equal to 2. How to solve this? We could simply check if the difference of the value expected and the value returned is smaller than `eps` (which is actually how the function works):

```
1 library("testthat")
```

```

1  ##
2  ## Attaching package: 'testthat'
3
4  ## The following object is masked from 'package:dplyr':
5  ##
6  ##      matches
7
8  ## The following objects are masked from 'package:magrittr':
9  ##
10 ##      equals, is_less_than, not
11
12 ## The following object is masked from 'package:purrr':
13 ##
14 ##      is_null

1  test_that("Test sqrt_newton: positive numeric",{
2      eps <- 0.001
3      expected <- 2
4      actual <- sqrt_newton(4, 1, eps = eps)
5      expect_lt(abs(expected - actual), eps)
6  })

```

There's no visible output, meaning that the test passes. Don't worry, we'll see how to run these tests in the next section, and we'll get a nice output confirming that tests did, indeed, pass.

I didn't talk about the functions `expect_equal()` and `expect_lt()`, but now is the moment. These functions are part of the `testthat` package and these are what allow you to test your functions. There's a number of them that allow you to test for a variety of situations. Check the documentation of `testthat` for more info. Let's continue to write more tests!

```

1  library("testthat")
2
3  test_that("Test sqrt_newton: negative numeric",{
4      expect_error(sqrt_newton(-4, 1))
5  })

```

We would like our function to return an error message if the user tries to get the square root of a negative number (let's say we don't want to generalize our function to complex numbers). But what happens here is that the function runs forever! This is because we are using a while loop whose condition is never fulfilled. This test basically allowed us to find two problems with our function:

- it doesn't deal with negative numbers

- the while loop may run forever if the condition is never fulfilled (for example if eps is too small)

Let's rewrite our function to take care of this, one problem at a time:

```
1 sqrt_newton <- function(a, init, eps = 0.01){
2   stopifnot(a >= 0)
3   while(abs(init**2 - a) > eps){
4     init <- 1/2 *(init + a/init)
5   }
6   return(init)
7 }
```

Now let's run our test again:

```
1 library("testthat")
2
3 test_that("Test sqrt_newton: negative numeric",{
4   expect_error(sqrt_newton(-4, 1))
5 })
```

Again no output, so things are good. Now to the next issue: we need to write a safeguard in the function to avoid having the while loop running for too long. For example if you try to run this:

```
1 sqrt_newton(49, 1E100000, 1E-100000)
```

You will see that it takes an awful lot of time! Let's limit the number of iterations to 100.

```
1 sqrt_newton <- function(a, init, eps = 0.01){
2   stopifnot(a >= 0)
3   i <- 1
4   while(abs(init**2 - a) > eps){
5     init <- 1/2 *(init + a/init)
6     i <- i + 1
7     if(i > 100) stop("Maximum number of iterations reached")
8   }
9   return(init)
10 }
```

Now when we try to run the following expression we get an error message:

```
1 sqrt_newton(49, 1E100, 1E-100)
```

```
1 Error in sqrt_newton(49, 1e+100, 1e-100) :  
2   Maximum number of iterations reached
```

But wouldn't it be better if the user could change the number of iterations himself?

```
1 sqrt_newton <- function(a, init, eps = 0.01, iter = 100){  
2   stopifnot(a >= 0)  
3   i <- 1  
4   while(abs(init**2 - a) > eps){  
5     init <- 1/2 *(init + a/init)  
6     i <- i + 1  
7     if(i > iter) stop("Maximum number of iterations reached")  
8   }  
9   return(init)  
10 }
```

We can now write some more tests:

```
1 library("testthat")  
2  
3 test_that("Test sqrt_newton: not enough iterations",{  
4   expect_error(sqrt_newton(4, 1E100, 1E-100, iter = 100))  
5 })
```

7.3 Actually running your tests

One of the easiest ways to run your tests is when you're developing a package. We are going to see this in the next chapter, but for now, let's suppose that we have a folder called `my_project` with the code inside of it. There's a file called `my_functions.R` and another file called `test_my_functions.R` which contain the functions you programmed and the unit tests that go with it respectively.

The file `test_my_functions.R` contains the following source code:


```

1  library("testthat")
2
3  test_that("Test sqrt_newton: positive numeric",{
4    eps <- 0.001
5    expected <- 2
6    actual <- sqrt_newton(4, 1, eps = eps)
7    expect_lt(abs(expected - actual), eps)
8  })
9
10 test_that("Test sqrt_newton: negative numeric",{
11   expect_error(sqrt_newton(-4, 1))
12 })
13
14 test_that("Test sqrt_newton: not enough iterations",{
15   expect_error(sqrt_newton(4, 1E100, 1E-100, iter = 100))
16 })

```

Then you simply run the following in the console:

```
1 test_file("test_my_functions.R")
```

of course you have to make sure that you are in the correct working directory. This can be tricky, and is one of the reasons why it's easier to run your tests when you're developing a package.

This is the output we get:

```

1  ...
2  DONE =====\
3  =====

```

See the three dots on the first line? Each dot represents a test that passed successfully. Let's add a test that will not pass on purpose, just to see what happens:

```

1  test_that("Test sqrt_newton: wrong on purpose",{
2    eps <- 0.001
3    expected <- 12
4    actual <- sqrt_newton(4, 1, eps = eps)
5    expect_lt(abs(expected - actual), eps)
6  })

```

This is the output we get now:

```

1  ...1
2  Failed -----\
3  -----
4  1. Failure: Test sqrt_newton: wrong on purpose (@test_my_functions.R#22) -----\
5  -----
6  abs(expected - actual) is not strictly less than `eps`. Difference: 10
7
8
9  DONE =====\
10 =====

```

You can then go back to the file that contains the tests and correct them. If all your tests are in a separate folder, you can use the function `test_dir()` to test all the functions in a given folder. The files containing your tests should all start with the string `test`. You could have a file called `run_tests.R` on the root of the directory and this file could contain the following:

```

1  library("testthat")
2
3  test_dir("tests")

```

You could then run your tests by running this file. You might also be tempted to write a bash script on GNU/Linux distributions or on macOS:

```

1  #!/bin/sh
2
3  Rscript -e "testthat::test_that('/whole/path/to/your/tests')"

```

but you'll probably only get burned because when you run this script, a new R session is started which does not know anything about your functions in your file `my_functions.R`. Managing the working directory is quite a pain. This is why in the next chapter we are going to start learning about packages and why writing our own packages to clean datasets is the best possible way to write your code.

7.4 Wrap-up

- Unit tests are a way of testing your code, and more specifically your functions.
- The basic workflow is to write your code, write tests, and check if your tests pass.
- You can also start with the tests and then write or modify your code to make them pass.
- We didn't talk about *coverage* yet. Are you sure that you test every line of your function? No you're not. In the next chapter I'll show you how can be sure to test each line of your function with the `covr` package.

7.5 Exercises

1. Write unit tests for the functions you wrote in the previous chapter. Just play around a little bit, and get a feeling for unit tests.

[\]\(packages.html\)](#) [²¹

²¹[putting-it-all-together-writing-a-package-to-work-on-data.html](#)

Chapter 8 Putting it all together: writing a package to work on data

Everything we have seen until now allows us to develop our own packages with the goal of *working* on data. By *working* on data I mean any operation that involves cleaning, transforming, analyzing or plotting data. I will summarize why everything we have seen until now helps us in this task:

1. Functional programming makes our code easier to test
2. Unit tests make sure our code is correct
3. Packages allows us to forget about paths, so unit tests are easier to run, makes writing documentation easier and makes sharing our code easier

For the rest of this chapter we are going to work with mock datasets that I created. The data is completely random but for our purposes it does not matter. In this chapter, we are going to write a number of functions with the goal of going from these awful, badly formatted datasets to a nice longitudinal data set.

8.1 Getting the data

You can download the data from the [github repository](https://github.com/b-rodrigues/functional_programming_and_unit_testing_for_data_munging)²² of the book. There are 5 .csv files that comprise the data sets we are going to work with:

- data_2000.csv
- data_2001.csv
- data_2002.csv
- data_2003.csv
- data_2004.csv

The first step, of course, is to load these datasets into R. For 5 datasets, I assume that you would simply write the following into Rstudio:

²²https://github.com/b-rodrigues/functional_programming_and_unit_testing_for_data_munging

```

1 data_2000 <- read.csv("/path/to/data/data_2000.csv", header = T)
2 data_2001 <- read.csv("/path/to/data/data_2001.csv", header = T)
3 data_2002 <- read.csv("/path/to/data/data_2002.csv", header = T)
4 data_2003 <- read.csv("/path/to/data/data_2003.csv", header = T)
5 data_2004 <- read.csv("/path/to/data/data_2004.csv", header = T)

```

This might be ok for 5 datasets which are named very similarly, especially since you can do block editing in Rstudio. However, imagine that you have hundreds, thousands, of datasets? And image that their names are not so well formatted as here? We will start our package by writing a function that reads a lot of datasets at once.

8.2 Your first data munging package: prepareData

8.2.1 Reading a lot of datasets at once

Using Rstudio, create a new project like shown in the previous chapter, and select *R package*. Give it a name, for example `prepareData`. If you are working with datasets that have a name, for example the *Penn World Tables*, you could call your package `preparePWT`, or something similar. By the way, we are going to work on some test data sets that I created for illustration purposes. When you will develop your own package to work on your own data, you do not have to write unit tests that use your original data. A subset can be enough, or taking the time to create a small test dataset might be preferable. It depends on what features of your functions you want to test. The first function I will show you is actually very general and could work with any datasets. This means that I created a package called `broTools` that contains all the little functions that I use daily. But for illustration purposes, we will put this function inside `prepareData`, even if it does not have anything directly to do with it. I have called this function `read_list()` and here is the source code:

```

1 #' Reads a list of datasets
2 #' @param list_of_datasets A list of datasets (names of datasets are strings)
3 #' @param read_func A function, the read function to use to read the data
4 #' @return Returns a list of the datasets
5 #' @export
6 #' @examples
7 #' \dontrun{
8 #'   setwd("/path/to/datasets/")
9 #'   list_of_datasets <- list.files(pattern = "*.csv")
10 #'   list_of_loaded_datasets <- read_list(list_of_datasets, read_func = read.csv)
11 #' }
12 read_list <- function(list_of_datasets, read_func, ...){
13
14   stopifnot(length(list_of_datasets)>0)

```

```

15
16   read_and_assign <- function(dataset, read_func){
17     dataset_name <- as.name(dataset)
18     dataset_name <- read_func(dataset, ...)
19   }
20
21   # invisible is used to suppress the unneeded output
22   output <- invisible(
23     purrr::map(list_of_datasets,
24               read_and_assign,
25               read_func = read_func)
26   )
27
28   # Remove the ".csv" at the end of the data set names
29   names_of_datasets <- c(unlist(strsplit(list_of_datasets, "[.]"))[c(T, F)])
30   names(output) <- names_of_datasets
31   return(output)
32 }

```

The basic idea of `read_list()` is that it takes a list of datasets as the first argument, then a function to read in the datasets as a second argument and as a third argument the famous `...`, which allows the user to specify further options to other functions that are contained in the body of the main function. In this case, further arguments are passed to the `read_func` function, for example if your data does not contain headers, you could pass the option `header = FALSE` to `read_list()` which would then get passed to `read_func`. I use `purrr::map()` to apply `read_and_assign()`; a helper function whose role is to read in a dataset and save it with its name, to the whole list of datasets. This step is wrapped inside `invisible()` as to remove unnecessary output. Finally I use `strsplit()` with a regular expression to remove the extension of the dataset from its name. The output is thus a list of datasets where each dataset is named as it is on your hard drive. Save this function in a script called `read_list.R` and save it in the `R` folder of your package. Now you need to invoke `roxygen2::roxygenise()` to create the documentation of your function. I suggest you also run `devtools::use_testthat`. This creates the necessary folder to hold your tests as well as creating a small `testthat.R` file with the code that gets called to run your tests. Without this, you might encounter weird issues (for example, `covr` not finding your tests!).

```
1 roxygen2::roxygenise()
```

```

1 First time using roxygen2. Upgrading automatically...
2 Updating roxygen version in /home/bro/Dropbox/prepareData/DESCRIPTION
3 Writing NAMESPACE
4 Writing read_list.Rd

```

```

1 devtools::use_testthat()

```

```

1 * Adding testthat to Suggests
2 * Creating `tests/testthat`.
3 * Creating `tests/testthat.R` from template.

```

Now let us check the coverage of our package:

```

1 library("covr")
2
3 cov <- package_coverage()
4
5 shine(cov)

```

Unsurprisingly we get a coverage of 0% for our package. We will now write a unit test for this function. For example, let us see if the condition `stopifnot(length(list_of_datasets)>0)` works. Because you ran `devtools::use_testthat()` you should have a folder called `tests` on the root of your project directory. In it, there is a folder called `testthat`. This is where you will save your unit tests, and any file needed for the tests to run (for example, mock datasets that are used by tests).

```

1 library("testthat")
2 library("prepareData")
3
4 test_that("Try to import empty list of datasets: this may be caused because
5           the path to the datasets is wrong for instance",{
6
7     list_datasets <- NULL
8
9     expect_error(read_list(list_datasets, read_csv, col_types = cols()))
10 })

```

Run the test using CTRL-SHIFT-T if you are on Rstudio.

```

1 ==> devtools::test()
2
3 Loading prepareData
4 Loading required package: testthat
5 Testing prepareData
6 .
7 DONE =====

```

This is the output you should see. If you check the coverage of your package, you should see that the line `stopifnot(length(list_of_datasets)>0)` is highlighted in green and you should have around 9% of coverage for your package. You can spend some to to get the coverage as high as possible, but you have to take into account the time it will take you to write tests vs the benefits you are going to get from them. In the case of this function, I do not really see what more you could test.

Let us use this function to read in the datasets:

```

1 library("readr")
2 library("purrr")
3 library("tibble")
4
5 list_of_data <- Sys.glob("assets/*.csv")
6
7 datasets <- read_list(list_of_data, read_csv, col_type = cols())

```

`list_of_data` is a variable that contains the path to the datasets. I used `Sys.glob("assets/*.csv")` to find the datasets. The datasets are saved in the `assets` folder of the book and end with the `.csv` extension. You could also use `list.files("*.csv")` to achieve the same. Let's take a look inside this list using `head()`. Since `head()` only works on single data frames or tibbles, we use `map()` to apply `head()` to each data frame on the list.

```

1 map(datasets, head)

```

```

1 ## $`assets/data_2000`
2 ## # A tibble: 6 x 6
3 ##       id Variable1 other2000 gender2000 eggs2000      spam2000
4 ##   <int>      <int>      <int>      <chr>      <int>      <chr>
5 ## 1     1         32         3         F         80 -1.5035369157
6 ## 2     2         28         2         F         20 -0.1836726393
7 ## 3     3         36         4         M         58 -0.6851988608
8 ## 4     4         28         1         F         30  1.9900760191
9 ## 5     5         34         3         F         14  0.4324725273

```



```

10 ## 6      6      30      3      F      40    -0.79001853
11 ##
12 ## `$assets/data_2001`
13 ## # A tibble: 6 x 6
14 ##       id VARIABLE1 other2001 Gender2001 eggs2001  spam2001
15 ##   <int>     <int>     <int>     <chr>     <int>     <dbl>
16 ## 1       1       1      32       3       F       80 -1.5035369
17 ## 2       2       2      28       2       F       20 -0.1836726
18 ## 3       3       3      36       4       M       58 -0.6851989
19 ## 4       4       4      28       1       F       30  1.9900760
20 ## 5       5       5      34       3       F       14  0.4324725
21 ## 6       6       6      30       3       F       40 -0.7900185
22 ##
23 ## `$assets/data_2002`
24 ## # A tibble: 6 x 6
25 ##       ID variable1 Other2002 gender2002 eggs2002  Spam2002
26 ##   <int>     <int>     <int>     <chr>     <int>     <dbl>
27 ## 1       1       1      32       3       F       80 -1.5035369
28 ## 2       2       2      28       2       F       20 -0.1836726
29 ## 3       3       3      36       4       M       58 -0.6851989
30 ## 4       4       4      28       1       F       30  1.9900760
31 ## 5       5       5      34       3       F       14  0.4324725
32 ## 6       6       6      30       3       F       40 -0.7900185
33 ##
34 ## `$assets/data_2003`
35 ## # A tibble: 6 x 6
36 ##       id variable1 other2003 gender2003 EGGS2003  spam2003
37 ##   <int>     <int>     <int>     <chr>     <int>     <dbl>
38 ## 1       1       1      32       3       F       80 -1.5035369
39 ## 2       2       2      28       2       F       20 -0.1836726
40 ## 3       3       3      36       4       M       58 -0.6851989
41 ## 4       4       4      28       1       F       30  1.9900760
42 ## 5       5       5      34       3       F       14  0.4324725
43 ## 6       6       6      30       3       F       40 -0.7900185
44 ##
45 ## `$assets/data_2004`
46 ## # A tibble: 6 x 6
47 ##       Id Variable1 Other2004 Gender2004 Eggs2004  Spam2004
48 ##   <int>     <int>     <int>     <chr>     <int>     <dbl>
49 ## 1       1       1      32       3       F       80 -1.5035369
50 ## 2       2       2      28       2       F       20 -0.1836726
51 ## 3       3       3      36       4       M       58 -0.6851989

```

```

52 ## 4      4      28      1      F      30  1.9900760
53 ## 5      5      34      3      F      14  0.4324725
54 ## 6      6      30      3      F      40 -0.7900185

```

The datasets we will work with all have the the same variables and the same individuals. We have datasets for the years 2000 to 2004. It would be much better for analysis if we could have clean variable names and merge every datasets together in a single, longitudinal dataset. In short, what we need:

- Have nice names for the columns.
- Remove the year from the name of the columns and add a column containing the year.
- Merge every dataset together.

This is to make the dataset tidy, as explained Wickham (2014b). Of course, depending on your needs, you might need to add further operations, for example creating new variables etc. For now, we are going to focus on these three steps.

8.2.2 Treating the columns of your datasets

Let us take a look at the column names of the datasets:

```

1  map(datasets, colnames)

1  ## `$assets/data_2000`
2  ## [1] "id"          "Variable1"  "other2000"  "gender2000" "eggs2000"
3  ## [6] "spam2000"
4  ##
5  ## `$assets/data_2001`
6  ## [1] "id"          "VARIABLE1"  "other2001"  "Gender2001" "eggs2001"
7  ## [6] "spam2001"
8  ##
9  ## `$assets/data_2002`
10 ## [1] "ID"          "variable1"  "Other2002"  "gender2002" "eggs2002"
11 ## [6] "Spam2002"
12 ##
13 ## `$assets/data_2003`
14 ## [1] "id"          "variable1"  "other2003"  "gender2003" "EGGS2003"
15 ## [6] "spam2003"
16 ##
17 ## `$assets/data_2004`
18 ## [1] "Id"          "Variable1"  "Other2004"  "Gender2004" "Eggs2004"
19 ## [6] "Spam2004"

```

This is very messy, we would need to have a function that would clean all this mess and “normalize” these column names. Turns out that we’re lucky, and there is exactly what we are looking for in the `janitor` package. The function `janitor::clean_names()` does exactly this. Let’s use it and see the output:

```
1 library("janitor")
2
3 datasets <- map(datasets, clean_names)
4
5 map(datasets, colnames)

1 ## $`assets/data_2000`
2 ## [1] "id"          "variable1"  "other2000"  "gender2000" "eggs2000"
3 ## [6] "spam2000"
4 ##
5 ## $`assets/data_2001`
6 ## [1] "id"          "variable1"  "other2001"  "gender2001" "eggs2001"
7 ## [6] "spam2001"
8 ##
9 ## $`assets/data_2002`
10 ## [1] "id"          "variable1"  "other2002"  "gender2002" "eggs2002"
11 ## [6] "spam2002"
12 ##
13 ## $`assets/data_2003`
14 ## [1] "id"          "variable1"  "other2003"  "gender2003" "eggs2003"
15 ## [6] "spam2003"
16 ##
17 ## $`assets/data_2004`
18 ## [1] "id"          "variable1"  "other2004"  "gender2004" "eggs2004"
19 ## [6] "spam2004"
```

This is much better. If `clean_names()` didn’t exist, you would have to have written your own function for this. This could have been a complicated exercise, depending on how messy and heterogenous the variable names would have been in your data. However `clean_names()` does a great job, so there’s no need to reinvent the wheel!

Now we would like to remove the years from the column names and add a column with the name of each dataset. Let us start by removing the years from the column names by writing a function. For this function, a little regular expression knowledge will not hurt. Here is what the function looks like:

```

1  #' Remove year strings from column names
2  #' @param list_of_datasets A list containing named datasets
3  #' @return A list of datasets with the supplied string prepended to the column n\
4  ames
5  #' @description This function removes year strings from column names, meaning th\
6  at a column called
7  #' "eggs9000" gets renamed into "eggs"
8  #' @export
9  #' @examples
10 #' \dontrun{
11 #' #`list_of_data_sets` is a list containing named data sets
12 #' # For example, to access the first data set, called dataset_1 you would
13 #' # write
14 #' list_of_data_sets$dataset_1
15 #' remove_years_from_strings(list_of_data_sets)
16 #' }
17 remove_years_from_strings <- function(list_of_datasets){
18
19   for_one_dataset <- function(dataset){
20     # strsplit() accepts regular expressions, so it's easy to get rid of a numbe\
21 r made up of
22     # *exactly* 4 digits
23
24     colnames(dataset) <- unlist(strsplit(colnames(dataset), "\\d{4}", perl = TRUE\
25 E))
26     return(dataset)
27   }
28
29   output <- purrr::map(list_of_datasets, for_one_dataset)
30
31   return(output)
32 }

```

and here is the accompanying unit test:

```

1 library("testthat")
2 library("prepareData")
3 library("readr")
4
5 data_sets <- list.files(pattern = "2001")
6
7 data_list <- read_list(data_sets, read_csv, col_types = cols())
8
9 test_that("Test remove years from strings",{
10   data_list_result <- purr::map(data_list, janitor::clean_names)
11   data_list_result <- remove_years_from_strings(data_list_result)
12   expect <- c("id", "year_", "variable1", "other", "gender", "eggs", "spam")
13   actual <- colnames(data_list_result[[1]])
14   expect_equal(expect, actual)
15 })

```

For the unit test to work, I had to add the dataset for the year 2001 in the tests/testthat directory. Again, this dataset does not have to be the real dataset you will ultimately be working on. A mock dataset with simulated data on 10 rows and with the same column names works exactly the same!

Let's take a look at the output:

```

1 datasets <- remove_years_from_strings(datasets)
2
3 map(datasets, colnames)

```

```

1 ## $`assets/data_2000`
2 ## [1] "id"          "variable1" "other"      "gender"     "eggs"       "spam"
3 ##
4 ## $`assets/data_2001`
5 ## [1] "id"          "variable1" "other"      "gender"     "eggs"       "spam"
6 ##
7 ## $`assets/data_2002`
8 ## [1] "id"          "variable1" "other"      "gender"     "eggs"       "spam"
9 ##
10 ## $`assets/data_2003`
11 ## [1] "id"          "variable1" "other"      "gender"     "eggs"       "spam"
12 ##
13 ## $`assets/data_2004`
14 ## [1] "id"          "variable1" "other"      "gender"     "eggs"       "spam"

```

This is starting to look like something!

Now, since we removed the years from the column names, we need to add a column containing the year to our datasets. And now to add the year column:

```

1  #' Adds the year column
2  #' @param list_of_datasets A list containing named datasets
3  #' @return A list of datasets with the year column
4  #' @description This function works by extracting the year string contained in
5  #' the data set name and appending a new column to the data set with the numeric
6  #' value of the year. This means that the data sets have to have a name of the
7  #' form data_set_2001 or data_2001_europe, etc
8  #' @export
9  #' @examples
10 #' \dontrun{
11 #' # `list_of_data_sets` is a list containing named data sets
12 #' # For example, to access the first data set, called dataset_1 you would
13 #' # write
14 #' list_of_data_sets$dataset_1
15 #' add_year_column(list_of_data_sets)
16 #' }
17 add_year_column <- function(list_of_datasets){
18
19   for_one_dataset <- function(dataset, dataset_name){
20
21     # Split the name of the dataset at the "_". The datasets must have a name of\
22     the
23     # form "data_2000" (notice the underscore).
24     name_year <- unlist(strsplit(dataset_name, "[_.]"))
25     # Get the index of the string that contains digits
26     index <- grep("\\d+", name_year)
27
28     # Get the year
29     year <- as.numeric(name_year[index])
30
31     # Add it to the data set
32     dataset$year <- year
33     return(dataset)
34   }
35
36   output <- purrr::map2(list_of_datasets, names(list_of_datasets), for_one_data_s\
37   et)
38   return(output)
39 }

```

And its unit test:

```

1  library("testthat")
2  library("prepareData")
3  library("readr")
4
5
6  data_sets <- list.files(pattern = "data")
7
8  data_list <- read_list(data_sets, read_csv, col_types = cols())
9
10 test_that("Test add year column",{
11   data_list_result <- purrr::map(data_list, janitor::clean_names)
12   data_list_result <- add_year_column(data_list_result)
13   expect <- list(rep(2001, 1000), rep(2002, 1000))
14   actual <- list(data_list_result[[1]]$year, data_list_result[[2]]$year)
15   expect_equal(expect, actual)
16 })

```

This function does not work if the names of the datasets are not of the form “data_2000”. This means that this function should have either an additional argument, where you specify the separator (for example “_” or “.” or even “-”) or fail if the name does not contain an “_”. I like the second solution better:

```

1  #' Adds the year column
2  #' @param list_of_datasets A list containing named datasets
3  #' @return A list of datasets with the year column
4  #' @description This function works by extracting the year string contained in
5  #' the data set name and appending a new column to the data set with the numeric
6  #' value of the year. This means that the data sets have to have a name of the
7  #' form data_set_2001 or data_2001_europe, etc
8  #' @export
9  #' @examples
10 #' \dontrun{
11 #' #`list_of_data_sets` is a list containing named data sets
12 #' # For example, to access the first data set, called dataset_1 you would
13 #' # write
14 #' list_of_data_sets$dataset_1
15 #' add_year_column(list_of_data_sets)
16 #' }
17 add_year_column <- function(list_of_datasets){
18
19   for_one_dataset <- function(dataset, dataset_name){
20

```

```

21   if(!("_" %in% unlist(strsplit(dataset_name, split = "")))){
22     stop("Make sure that your datasets are named like
23         `data_2000.csv` or similar. The `_` between `data`
24         and `2000` is what matters")
25   }
26   # Split the name of the dataset at the "_". The datasets must have a name of\
27   the
28   # form "data_2000" (notice the underscore).
29   name_year <- unlist(strsplit(dataset_name, split = "[_.]"))
30   # Get the index of the string that contains digits
31   index <- grep("\\d+", name_year)
32
33   # Get the year
34   year <- as.numeric(name_year[index])
35
36   # Add it to the data set
37   dataset$year <- year
38   return(dataset)
39 }
40
41 output <- purrr::map2(list_of_datasets, names(list_of_datasets), for_one_data\
42 et)
43 return(output)
44 }

```

If you check the coverage of this function, you will see that the lines that test if the datasets are correctly named do not get called. Let's add a unit test that does this, but first, we need to create *wrong* datasets. Just copy the datasets you have in your tests folder, and rename them to `wrongdata2001.csv` and `wrongdata2002.csv`. We expect our function to stop with an error message if it tries anything on these datasets:

```

1 data_sets <- list.files(pattern = "wrong")
2
3 data_list <- read_list(data_sets, read_csv, col_types = cols())
4
5 test_that("Test add year column: wrong name",{
6   data_list_result <- purrr::map(data_list, janitor::clean_names)
7   expect_error(add_year_column(data_list_result))
8 })

```

Now have fully covered your function, and you also know when the function breaks. With the informative error message, future you or your coworkers will know how to correctly name the datasets. Let's try `add_year_column()` to see how it behaves on our data:


```

1 datasets <- add_year_column(datasets)
2
3 map(datasets, head)

1 ## $`assets/data_2000`
2 ## # A tibble: 6 x 7
3 ##       id variable1 other gender  eggs      spam  year
4 ##   <int>      <int> <int>  <chr> <int>      <chr> <dbl>
5 ## 1     1         32     3    F     80 -1.5035369157 2000
6 ## 2     2         28     2    F     20 -0.1836726393 2000
7 ## 3     3         36     4    M     58 -0.6851988608 2000
8 ## 4     4         28     1    F     30  1.9900760191 2000
9 ## 5     5         34     3    F     14  0.4324725273 2000
10 ## 6     6         30     3    F     40  -0.79001853 2000
11 ##
12 ## $`assets/data_2001`
13 ## # A tibble: 6 x 7
14 ##       id variable1 other gender  eggs      spam  year
15 ##   <int>      <int> <int>  <chr> <int>      <dbl> <dbl>
16 ## 1     1         32     3    F     80 -1.5035369 2001
17 ## 2     2         28     2    F     20 -0.1836726 2001
18 ## 3     3         36     4    M     58 -0.6851989 2001
19 ## 4     4         28     1    F     30  1.9900760 2001
20 ## 5     5         34     3    F     14  0.4324725 2001
21 ## 6     6         30     3    F     40 -0.7900185 2001
22 ##
23 ## $`assets/data_2002`
24 ## # A tibble: 6 x 7
25 ##       id variable1 other gender  eggs      spam  year
26 ##   <int>      <int> <int>  <chr> <int>      <dbl> <dbl>
27 ## 1     1         32     3    F     80 -1.5035369 2002
28 ## 2     2         28     2    F     20 -0.1836726 2002
29 ## 3     3         36     4    M     58 -0.6851989 2002
30 ## 4     4         28     1    F     30  1.9900760 2002
31 ## 5     5         34     3    F     14  0.4324725 2002
32 ## 6     6         30     3    F     40 -0.7900185 2002
33 ##
34 ## $`assets/data_2003`
35 ## # A tibble: 6 x 7
36 ##       id variable1 other gender  eggs      spam  year
37 ##   <int>      <int> <int>  <chr> <int>      <dbl> <dbl>
38 ## 1     1         32     3    F     80 -1.5035369 2003

```

```

39 ## 2      2      28      2      F      20 -0.1836726 2003
40 ## 3      3      36      4      M      58 -0.6851989 2003
41 ## 4      4      28      1      F      30  1.9900760 2003
42 ## 5      5      34      3      F      14  0.4324725 2003
43 ## 6      6      30      3      F      40 -0.7900185 2003
44 ##
45 ## `$assets/data_2004`
46 ## # A tibble: 6 x 7
47 ##       id variable1 other gender  eggs      spam  year
48 ##   <int>      <int> <int>   <chr> <int>    <dbl> <dbl>
49 ## 1      1        32     3     F     80 -1.5035369 2004
50 ## 2      2        28     2     F     20 -0.1836726 2004
51 ## 3      3        36     4     M     58 -0.6851989 2004
52 ## 4      4        28     1     F     30  1.9900760 2004
53 ## 5      5        34     3     F     14  0.4324725 2004
54 ## 6      6        30     3     F     40 -0.7900185 2004

```

Just as expected!

TBC...

References

Wickham, Hadley. 2014b. “Tidy Data.” *Journal of Statistical Software* 59 (1): 1–23. doi:[10.18637/jss.v059.i10](https://doi.org/10.18637/jss.v059.i10)²³.

-
1. It stands for Bruno Rodrigues' Tools. I'm still working on releasing the package on Github, and maybe CRAN.[✉](#)²⁴

²³<https://doi.org/10.18637/jss.v059.i10>

²⁴[putting-it-all-together-writing-a-package-to-work-on-data.html#fnref2](#)