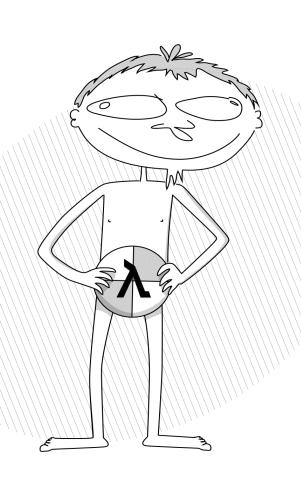
# FUNCTIONAL PROGRAMMING AND UNIT TESTING FOR DATA MUNGING WITH R



# 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 2017-01-07



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:

https://twitter.com/search?q=##fput

# **Contents**

Chapter 1 Why this book?	
1.1 Motivation	
1.2 Who am I?	
1.3 Thanks	
1.4 License	
Chapter 2 Introduction	
2.1 Getting R	
2.2 A short overview of functional programming	
2.3 A short overview of unit testing	
Chapter 3 Functional Programming	
3.1 Introduction	
3.2 Mapping and Reducing: the <i>base</i> way	10
3.3 Mapping and Reducing: the purrr way	
3.4 Anonymous functions	
3.5 Wrap-up	
3.6 Exercises	
Chapter 4 Unit testing	29
4.1 Introduction	
4.2 Unit testing with the testthat package	
4.3 Actually running your tests	
4.4 Wrap-up	
4.5 Exercises	
Chapter 5 Packages	
5.1 Why you need your own packages in your life	37
5.2 R packages: the basics	3'
5.3 Writing documentation for your functions	
5.4 Unit test your package	
5.5 Checking the coverage of your unit tests with covr	
5.6 Wrap-up	
Chapter 6 Putting it all together: writing a package to work on data	a45

#### CONTENTS

6.1 Getting the data	45
6.2 Your first data munging package: prepareData	46

# **Chapter 1 Why this book?**

This short 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.

#### 1.1 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.

#### 1.2 Who am !?

I use R daily at my current job, and discovered R some years ago while I was at the University of Strasbourg<sup>1</sup>. 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<sup>2</sup> or check my blog<sup>3</sup>.

<sup>&</sup>lt;sup>1</sup>http://www.unistra.fr/index.php?id=accueil

<sup>&</sup>lt;sup>2</sup>https://twitter.com/brodriguesco

<sup>&</sup>lt;sup>3</sup>http://brodrigues.co

#### 1.3 Thanks

I'd like to thank Ross Ihaka<sup>4</sup> and Robert Gentleman<sup>5</sup> for developing the R programming language. Many thanks to Hadley Wickham<sup>6</sup> for all the wonderful packages he developed that make R much more pleasant to use. Thanks to Yihui Yie<sup>7</sup> for bookdown without which this book would not exist (at least not in this very nice format).

Thanks to Hans-Martin von Gaudecker<sup>8</sup> 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!).

#### 1.4 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<sup>9</sup>.

#### 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 Grolemund. 2016. *R for Data Science*. 1st ed. O'Reilly. http://r4ds. had.co.nz/.

<sup>&</sup>lt;sup>4</sup>https://www.stat.auckland.ac.nz/~ihaka/

<sup>&</sup>lt;sup>5</sup>https://en.wikipedia.org/wiki/Robert\_Gentleman\_(statistician)

<sup>&</sup>lt;sup>6</sup>http://hadley.nz/

<sup>&</sup>lt;sup>7</sup>http://yihui.name/

<sup>8</sup>https://www.iame.uni-bonn.de/people/hm-gaudecker

<sup>9</sup>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:

Chapter 2 Introduction 4

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<sup>10</sup> 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

<sup>10</sup>https://en.wikipedia.org/wiki/Functional\_programming

Chapter 2 Introduction 5

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*).

# **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  }</pre>
```

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<sup>1</sup> 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:

```
sqrt_newton_recur <- function(a, init, eps = 0.01){</pre>
1
2
       if(abs(init**2 - a) < eps){
3
           result <- init
4
       } else {
5
            init < -1/2 * (init + a/init)
           result <- sqrt_newton_recur(a, init, eps)
6
       return(result)
8
9
   }
   sqrt_newton_recur(16, 2)
   ## [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 got 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:

```
increment <- function(x){
return(x + 1)
}
</pre>
```

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] 11
   will always produce 11.
   However, this one:
   increment_opaque <- function(x){</pre>
        return(x + spam)
2
   }
3
   is not a referential transparent function, because its value depends on the global variable spam.
   spam <- 1
1
2
   increment_opaque(10)
   ## [1] 11
   will only produce 11 if spam = 1. But what if spam = 19?
   spam <- 19
1
2
   increment_opaque(10)
   ## [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){</pre>
        return(x + spam)
2
   }
```

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</pre>
```

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):

```
count iter <- ∅
 1
 2
    sqrt_newton_side_effect <- function(a, init, eps = 0.01){</pre>
 3
         while(abs(init**2 - a) \rightarrow eps){
 4
 5
              init \langle -1/2 \rangle *(init + a/init)
 6
             count_iter <<- count_iter + 1 # The "<<-" symbol means that we assign the</pre>
                                                # RHS value in a variable in the global en\
         }
    vironment
 8
         return(init)
 9
    }
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 <code>count\_iter</code> 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 }</pre>
```

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:

```
numbers \leftarrow c(16, 25, 36, 49, 64, 81)
 1
 2
 3
    sqrt_newton(numbers, init = rep(1, 6), eps = rep(0.001, 6))
    ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
 1
 2
    ## and only the first element will be used
 3
    ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
 4
    ## and only the first element will be used
 5
 6
    ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
 7
    ## and only the first element will be used
 8
 9
10
    ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
    ## and only the first element will be used
11
12
    ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
13
    ## and only the first element will be used
14
15
    ## Warning in while (abs(init^2 - a) > eps) {: the condition has length > 1
16
    ## and only the first element will be used
17
18
    ## [1] 4.000001 5.000023 6.000253 7.001406 8.005148 9.014272
19
```

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
    ## [1] 4.000001
 2
 3
 4
    ## [[2]]
    ## [1] 5.000023
 5
 6
    ##
    ## [[3]]
 7
    ## [1] 6.000253
 8
 9
    ## [[4]]
10
    ## [1] 7
11
12
    ##
13
    ## [[5]]
    ## [1] 8.000002
14
15
    ##
16
   ## [[6]]
    ## [1] 9.000011
17
    Map() applies a function to every element of a list and returns a list.
    We could then write a wrapper around Map():
    sqrt_newton_vec <- function(numbers, init, eps = 0.01){</pre>
 1
        return(Map(sqrt_newton, numbers, init, eps))
 2
    }
 3
    sqrt_newton_vec(numbers, 1)
    ## [[1]]
 1
    ## [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:

```
a \leftarrow cbind(c(1, 2, 3), c(4, 5, 6), c(7, 8, 9))
    apply(a, 1, sum)
    ## [1] 12 15 18
    We could use lapply() instead of Map():
    lapply(numbers, sqrt_newton, init = 1)
    ## [[1]]
 1
 2
    ## [1] 4.000001
    ##
 3
    ## [[2]]
 4
    ## [1] 5.000023
 5
    ##
 6
 7
    ## [[3]]
    ## [1] 6.000253
 8
 9
    ## [[4]]
10
    ## [1] 7
11
12
13
    ## [[5]]
```

```
14
    ## [1] 8.000002
15
    ##
16 ## [[6]]
17 ## [1] 9.000011
    or sapply():
   sapply(numbers, sqrt_newton, init = 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):
    sqrt_newton_vec <- function(numbers, init, eps = 0.01){</pre>
 1
        return(sapply(numbers, sqrt_newton, init, eps))
 2
    }
 3
 4
    sqrt_newton_vec(numbers, 1)
    ## [1] 4.000001 5.000023 6.000253 7.000000 8.000002 9.000011
    mapply() is different from these two:
  inits <- c(100, 20, 3212, 487, 5, 9888)
    mapply(sqrt_newton, numbers, init = inits)
    ## [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] 4.000284
```

```
sqrt_newton(numbers[2], inits[2])
## [1] 5.000001
sqrt_newton(numbers[3], inits[3])
## [1] 6.000003
sqrt_newton(numbers[4], inits[4])
## [1] 7.000006
sqrt_newton(numbers[5], inits[5])
## [1] 8.000129
sqrt_newton(numbers[6], inits[6])
## [1] 9.000006
From the Map()'s documentation, we learn that:
`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 +<sup>2</sup> 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 Ø. 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 }</pre>
```

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.

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 Map() and Reduce() is enough to get us going.

# 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, but will only scratch the surface. Take the time to read purrr's documentation! 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. There are also other interesting variants, such as map\_if():

```
library("purrr")
 2
    a < - seq(1,10)
 3
    is_multiple_of_two <- function(x){</pre>
 4
         ifelse(x %% 2 == 0, TRUE, FALSE)
 5
 6
    }
 7
    map_if(a, is_multiple_of_two, sqrt)
    ## [[1]]
 1
 2
    ## [1] 1
 3
 4
    ## [[2]]
    ## [1] 1.414214
 5
 6
    ##
 7
    ## [[3]]
    ## [1] 3
 8
 9
    ## [[4]]
10
11
    ## [1] 2
12
13
    ## [[5]]
```

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

What happened in this snippet of code? First I wrote a function that returns TRUE if a number is a multiple of 2, and FALSE otherwise. Then, I used map\_if() to take the square root of only those numbers in vector a that are divisble by 2.

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

```
map2(numbers, inits, sqrt_newton)
    ## [[1]]
 1
    ## [1] 4.000284
 2
 3
    ## [[2]]
 4
 5
    ## [1] 5.000001
    ##
 7
    ## [[3]]
    ## [1] 6.000003
 8
 9
    ##
10
    ## [[4]]
11
    ## [1] 7.000006
12
    ##
13
    ## [[5]]
    ## [1] 8.000129
14
15
    ##
    ## [[6]]
16
    ## [1] 9.000006
```

#### 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</pre>
```

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</pre>
```

These two functions keep the intermediary results.

#### 3.3.3 Other useful functions from purrr

There are a lot of other useful functions in purrr. For example safely() and possibly() are great:

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

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

Using map() or Map() will result in a similar error. However, using safely() will work for the numbers contained in a and show an error for the first element of a which is a character:

```
a <- list("a", 4, 5)
 1
 2
 3 safe_sqrt <- safely(sqrt)</pre>
 4
   map(a, safe_sqrt)
   ## [[1]]
 1
 2 ## [[1]]$result
 3 ## NULL
 4
   ##
   ## [[1]]$error
   ## \langle simpleError in .f(...) : non-numeric argument to mathematical function \rangle
 6
 7
    ##
 8
   ##
   ## [[2]]
 9
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
```

And possibly() allows you to specify a return value in case of an error:

```
1
   possible_sqrt <- possibly(sqrt, otherwise = NA_real_)</pre>
2
   map(a, possible_sqrt)
3
   ## [[1]]
2
   ## [1] NA
3
   ##
4
   ## [[2]]
   ## [1] 2
5
6
   ##
7
   ## [[3]]
   ## [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.

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)</pre>
```

```
1
    ## [[1]]
    ## [[1]]$result
 2
    ## [1] NA
 3
    ##
 4
    ## [[1]]$error
 5
    ## \langle simpleError in .f(...) \rangle: non-numeric argument to mathematical function
 6
    ##
 8
    ##
   ## [[2]]
10 ## [[2]]$result
   ## [1] 2
11
    ##
12
   ## [[2]]$error
13
   ## NULL
14
   ##
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))
```

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

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

## 3.4 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)</pre>
```

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.

```
map(datasets, hist, cyl)

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

Maybe try this:
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)))



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

```
32
   ##
   ## $mids
33
34
   ## [1] 7.25 7.75 8.25 8.75 9.25 9.75 10.25 10.75
   ##
35
36
   ## $xname
   ## [1] "x$cyl"
37
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 )
```



```
## $mtcars2000
 1
    ## $breaks
 3
   ## [1] 4.0 4.5 5.0 5.5 6.0 6.5 7.0 7.5 8.0
 4
   ##
    ## $counts
 5
    ## [1] 11 0 0 7 0 0 0 14
 6
 7
    ##
   ## $density
 8
   ## [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
    ##
   ## $equidist
17
   ## [1] TRUE
18
19
   ##
    ## attr(,"class")
20
   ## [1] "histogram"
21
   ##
22
23
   ## $mtcars2001
24
   ## $breaks
   ## [1] 7.0 7.5 8.0 8.5 9.0 9.5 10.0 10.5 11.0
25
   ##
26
27
    ## $counts
   ## [1] 11 0 0 7 0 0 0 14
28
29
30
   ## $density
31
   ## [1] 0.6875 0.0000 0.0000 0.4375 0.0000 0.0000 0.0000 0.8750
32
   ##
    ## $mids
33
   ## [1] 7.25 7.75 8.25 8.75 9.25 9.75 10.25 10.75
34
35
   ## $xname
36
   ## [1] ".x$cy1"
37
38
   ##
   ## $equidist
39
   ## [1] TRUE
40
   ##
41
   ## attr(,"class")
42
   ## [1] "histogram"
43
```

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.

## 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:

```
dataset_names <- c("dataset1.csv", "dataset2.csv", "dataset3.csv")
remove_csv(dataset_names)

[1] "dataset1" "dataset2" dataset3"</pre>
```

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.

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

#### **References**

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/.

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

<sup>1.</sup> The *body* of a function are the instructions that define the function. You can get the body of a function with body(some\_func) $\boxtimes^{11}$ 

<sup>2.</sup> 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. $\boxtimes^{12}$ 

<sup>&</sup>lt;sup>11</sup>fprog.html#fnref1

<sup>&</sup>lt;sup>12</sup>fprog.html#fnref2

# **Chapter 4 Unit testing**

#### 4.1 Introduction

Let's take a look at Wikipedia's definition<sup>13</sup> 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.

<sup>13</sup>https://en.wikipedia.org/wiki/Unit\_testing

Chapter 4 Unit testing 30

# 4.2 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:

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")
```

Chapter 4 Unit testing 31

```
##
1
2
   ## Attaching package: 'testthat'
3
4
   ## The following object is masked from 'package:purrr':
5
   ##
   ##
6
          is_null
   test_that("Test sqrt_newton: positive numeric",{
1
2
                  eps <- 0.001
3
                  expected <- 2
                  actual <- sqrt_newton(4, 1, eps = eps)</pre>
4
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!

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:

```
sqrt_newton <- function(a, init, eps = 0.01){</pre>
1
2
       stopifnot(a >= ∅)
       while(abs(init**2 - a) > eps){
3
            init < -1/2 *(init + a/init)
4
5
       }
       return(init)
6
  }
   Now let's run our test again:
   library("testthat")
1
2
   test_that("Test sqrt_newton: negative numeric",{
                  expect_error(sqrt_newton(-4, 1))
4
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.

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

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

We can now write some more tests:

# 4.3 Actually running your tests

One of the easiest ways to run your tests is when your 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:

```
library("testthat")
 1
 2
 3
    test_that("Test sqrt_newton: positive numeric",{
 4
        eps <- 0.001
 5
        expected <- 2
        actual <- sqrt_newton(4, 1, eps = eps)</pre>
 6
        expect_lt(abs(expected - actual), eps)
    })
 8
 9
10
    test_that("Test sqrt_newton: negative numeric",{
        expect_error(sqrt_newton(-4, 1))
11
12
    })
13
    test_that("Test sqrt_newton: not enough iterations",{
14
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:

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:

```
test_that("Test sqrt_newton: wrong on purpose",{
    eps <- 0.001
    expected <- 12
    actual <- sqrt_newton(4, 1, eps = eps)
    expect_lt(abs(expected - actual), eps)
}</pre>
```

This is the output we get now:

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.

#### 4.4 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.

# 4.5 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.

#### 5.1 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 it 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.

#### 5.2 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:

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

## 5.3 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:

```
sqrt_newton <- function(a, init, eps = 0.01, iter = 100){</pre>
 1
 2
        stopifnot(a >= ∅)
        i <- 1
 3
 4
        while(abs(init**2 - a) > eps){
             init < -1/2 *(init + a/init)
 5
            i < -i + 1
 6
            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:

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

The first difference with standard comments is that roxygen2 type comments start with the #' symbol 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 automagically 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.

#### 5.4 Unit test your package

Now that we know the basics of creating a package, we move on to unit testing your package. Unit testing is very useful, but require some work, especially because you have to run them 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:

```
library("testthat")
1
2
  library("myFirstPackage")
3
  test_that("Test sqrt_newton: positive numeric",{
4
      eps <- 0.001
5
      expected <- 2
6
      actual <- sqrt_newton(4, 1, eps = eps)</pre>
      expect_lt(abs(expected - actual), eps)
8
  })
   Save this file and use the following keyboard shortcut: CTRL-SHIFT-T to run your unit test. You will
   see the following output:
  ==> devtools::test()
1
2
3 Loading myFirstPackage
4 Loading required package: testthat
  Testing myFirstPackage
6
  You can of course add more unit tests inside the same file. Add the following code to test_sqrt_-
   newton.R:
```

You will now see the following output:

2
3 })

test\_that("Test sqrt\_newton: negative numeric",{

expect\_error(sqrt\_newton(-4, 1))

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:

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

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

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

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:

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

## 5.5 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)</pre>
```

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



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:



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!

### 5.6 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 6 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.

## 6.1 Getting the data

You can download the data from the github repository<sup>14</sup> 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:

<sup>14</sup>https://github.com/b-rodrigues/functional\_programming\_and\_unit\_testing\_for\_data\_munging

```
data_2000 <- read.csv("/path/to/data/data_2000.csv", header = T)
data_2001 <- read.csv("/path/to/data/data_2001.csv", header = T)
data_2002 <- read.csv("/path/to/data/data_2002.csv", header = T)
data_2003 <- read.csv("/path/to/data/data_2003.csv", header = T)
data_2004 <- read.csv("/path/to/data/data_2004.csv", header = T)</pre>
```

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.

## 6.2 Your first data munging package: prepareData

#### 6.2.1 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 you 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<sup>3</sup> 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:

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

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

```
First time using roxygen2. Upgrading automatically...

Updating roxygen version in /home/bro/Dropbox/prepareData/DESCRIPTION

Writing NAMESPACE

Writing read_list.Rd

devtools::use_testthat()

* Adding testthat to Suggests

* Creating `tests/testthat`.

* 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)</pre>
```

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 stopi fnot(length(list\_of\_datasets)>0) works. Because you ran detools::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 were you will save your unit tests, and any file needed for the tests to run (for example, mock datasets that are used by tests).

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

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

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 highlightened 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())</pre>
```

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)
```

```
## $`assets/data 2000`
1
   ## # A tibble: 6 × 6
2
   ##
            id Variable1 other2000 gender2000 eggs2000
3
                                                               spam2000
   ##
        <int>
                   <int>
                              <int>
                                          <chr>
                                                   <int>
                                                                  <chr>>
4
   ## 1
5
            1
                      32
                                  3
                                              F
                                                      80 -1.5035369157
   ## 2
            2
                      28
                                  2
                                              F
                                                      20 -0.1836726393
6
   ## 3
             3
                      36
                                  4
                                              М
                                                      58 -0.6851988608
7
   ## 4
             4
                                              F
                      28
                                  1
                                                      30 1.9900760191
8
                                              F
   ## 5
            5
                      34
                                                      14 0.4324725273
```

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

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

The datasets we will work with all have the same variables and the same inviduals. 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.

#### 6.2.2 Treating the columns of your datasets

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

```
map(datasets, colnames)
    ## $`assets/data_2000`
 1
    ## [1] "id"
 2
                         "Variable1"
                                       "other2000"
                                                     "gender2000" "eggs2000"
    ## [6] "spam2000"
 3
 4
    ##
 5
    ## $`assets/data_2001`
       [1] "id"
                                                     "Gender2001" "eggs2001"
                         "VARIABLE1"
 6
                                       "other2001"
 7
    ## [6] "spam2001"
 8
    ##
    ## $`assets/data_2002`
 9
    ## [1] "ID"
                         "variable1"
                                                     "gender2002" "eggs2002"
10
                                       "Other2002"
    ## [6] "Spam2002"
11
    ##
12
    ## $`assets/data_2003`
13
    ## [1] "id"
                         "variable1" "other2003"
                                                     "gender2003" "EGGS2003"
14
15
    ## [6] "spam2003"
    ##
16
    ## $`assets/data 2004`
17
    ## [1] "Id"
                         "Variable1"
                                       "Other2004"
                                                     "Gender2004" "Eggs2004"
18
    ## [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:

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

This is much better. If <code>clean\_names()</code> 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 <code>clean\_names()</code> does a great job, so there's no need to reivent 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
   #' @description This function removes year strings from column names, meaning th\
 5
 6 at a column called
   #' "eggs9000" gets renamed into "eggs"
 7
 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){</pre>
18
19
      for_one_dataset <- function(dataset){</pre>
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
        colnames(dataset) <- unlist(strsplit(colnames(dataset), "\\d{4}\", perl = TRU\</pre>
24
25 E))
        return(dataset)
26
27
      }
28
29
      output <- purrr::map(list_of_datasets, for_one_dataset)</pre>
30
      return(output)
31
32
    }
```

and here is the accompanying unit test:

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

```
datasets <- remove_years_from_strings(datasets)</pre>
 1
 2
 3
    map(datasets, colnames)
    ## $`assets/data_2000`
 1
    ## [1] "id"
                        "variable1" "other"
 2
                                                  "gender"
                                                               "eggs"
                                                                            "spam"
 3
    ##
 4
    ## $`assets/data_2001`
    ## [1] "id"
                        "variable1" "other"
 5
                                                  "gender"
                                                               "eggs"
                                                                            "spam"
 6
    ##
 7
    ## $`assets/data_2002`
    ## [1] "id"
                        "variable1" "other"
                                                  "gender"
 8
                                                               "eggs"
                                                                            "spam"
 9
    ## $`assets/data_2003`
10
    ## [1] "id"
                        "variable1" "other"
                                                  "gender"
11
                                                               "eggs"
                                                                            "spam"
12
    ##
13
    ## $`assets/data_2004`
    ## [1] "id"
                        "variable1" "other"
14
                                                  "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:

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

And its unit test:

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

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

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

```
data_sets <- list.files(pattern = "wrong")

data_list <- read_list(data_sets, read_csv, col_types = cols())

test_that("Test add year column: wrong name",{
    data_list_result <- purrr::map(data_list, janitor::clean_names)
    expect_error(add_year_column(data_list_result))

})</pre>
```

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:

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

```
39
    ## 2
             2
                       28
                                            20 -0.1836726
                                                            2003
    ## 3
              3
                       36
                                      М
                                            58 -0.6851989
                                                            2003
40
                                      F
41
    ## 4
              4
                       28
                                               1.9900760
                                                            2003
    ## 5
             5
                                      F
42
                       34
                               3
                                            14
                                                0.4324725
                                                            2003
                                            40 -0.7900185
    ## 6
             6
                       30
43
                                                            2003
44
    ## $`assets/data_2004`
45
    ## # A tibble: 6 × 7
46
47
    ##
             id variable1 other gender
                                          eggs
                                                            year
                                                     spam
48
    ##
         <int>
                    <int> <int>
                                  <chr> <int>
                                                     <dbl> <dbl>
                                      F
    ## 1
             1
                       32
                               3
                                            80 -1.5035369
                                                            2004
49
    ## 2
                       28
                               2
                                      F
50
             2
                                            20 -0.1836726
                                                            2004
51
    ## 3
                       36
                               4
                                            58 -0.6851989
              3
                                                            2004
    ## 4
             4
                       28
                                      F
                                               1.9900760
52
                               1
                                            30
                                                            2004
53
    ## 5
             5
                       34
                               3
                                      F
                                            14 0.4324725
                                                            2004
54
    ## 6
             6
                       30
                               3
                                            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<sup>15</sup>.

1. It stands for Bruno Rodrigues' Tools. I'm still working on releasing the package on Github, and maybe CRAN.  $\square^{16}$ 

<sup>&</sup>lt;sup>15</sup>https://doi.org/10.18637/jss.v059.i10

 $<sup>^{\</sup>bf 16} putting-it-all-together-writing-a-package-to-work-on-data.html\#fnref3$