R Package Essentials

Table of contents

A	bout	this	7
\mathbf{A}		this v to use this book	7 8
\mathbf{G}	ettin	g course materials	9
\mathbf{G}	ettin Lice	ng course materials	9
Li	cens	e	11
Li	cens	e	11
1	Phi	llosophy	13
2	Inst	tallation	15
	2.1	Overview	15
	2.2	Questions	15
	2.3	Software Setup	15
		2.3.1 Installing R	15
		2.3.2 Installing RStudio	16
		2.3.3 Installing R packages for development	16
		2.3.4 git and github	18
		2.3.5 Installing RTools	20
3	RSt	tudio, What and Why	21
	3.1	Overview	21
	3.2	Questions	21
	3.3	Objectives	21
	3.4	What is RStudio, and why should I use it?	21
	3.5	Learning more	24
4	Wo	rkflow	25
	4.1	Overview	25
	4.2	Questions	26
	4.3	Objectives	26

1	0 (Contents

	4.4	When you start a new project: Open a new RStudio project 4.4.1 So what does this do?
	4.5	What is a file path?
	4.6	Is there an answer to the madness?
	$\frac{4.0}{4.7}$	
	4.7	1 0
	4.0	Remember
5	Sun	amary 33
6	•	y functions? 35
	6.1	Overview
	6.2	Questions
	6.3	Objectives
	6.4	Prior Art
	6.5	Code is for people
	6.6	OK, but what actually is a function?
	6.7	script
	6.8	function
		6.8.1 Functions give ideas a home
	6.9	Anatomy of a function
	6.10	How to think about writing functions
		6.10.1 Identifying the output - what do we need? 44
		6.10.2 Identifying the input
		6.10.3 Managing scope - functions are best (generally) when
		they do one thing
		When to function
	6.12	Naming things is hard
	6.13	The other hard part of writing functions
	6.14	Conclusion
7	Mot	ivation 53
	7.1	Overview
	7.2	Questions
	7.3	Objectives
	7.4	How this works
	7.5	The example: "learned"
	7.6	Discussion of potential problems
	7.7	Identifying the report outputs
	7.8	Plot of proportion of people educated in each age group in
	-	each state
	7.9	box plot of proportion of people educated for each state.
	7.10	Table of The 5 number summary of proportion of people
	0	educated for each state
8		59

0.0 Contents	5
9	61
10	63
11	65
12	67
13	69
14	71
15	73
16	75
17	77
18	79
19	81
20	83
21	85
22	87
23	89
Appendices	91
A	91

About this

This is a book on the essential components of creating an R package. It is aimed at those who want to learn how to make R packages. You probably have written some functions, but if you haven't, we discuss how to do that. I care a lot about writing functions, and have a lot of thoughts and ideas on how to do it.

It was initially developed as a full-day hour workshop, "R package essentials". It is a developed into a resource that will grow and change over time as a living book.

This book aims to teach the following:

- Installation and setup of dependencies
 - git + github
 - R, RStudio
 - package dependencies
- Function essentials
 - DRY;DRY (Don't Repeat Yourself; Don't Reread Yourself)
 - Expression
 - Finding the inputs
- Moving a script to a series of functions
- Create package barebones with create_package()
- How to add dependencies with use_package() DESCRIPTION file
- How to add documentation with roxygen2
- Why you should use R CMD Check
- How to add data to a package
- How to add a README
- How to put your package on github
- How to add vignettes
- Writing tests
- Using a NEWS file
- Adding a website
- Using Continuous Integration to check and test
- Publishing your software on R universe

8 0 About this

How to use this book

This book was written to provide course materials for a 8 hour course on R Packages

We worked through the following sections in the book in 8 hours:

- why R packages
- why functions
- installation
- what is RStudio?
- suggested workflow and hygiene

•

With the remaining sections being used as extra material, or have since been written after the course:

•

$Getting\ course\ materials$

Course materials can be downloaded by using the following command from the usethis package:

usethis::use_course("njtierney/rpkgess-materials")

Licence

This work is licensed under a Creative Commons Attribution-NonCommercial 4.0 International License.

License

This work is licensed under a Creative Commons Attribution-NonCommercial 4.0 International License.

Philosophy

I first learnt to write an R package from [Hilary Parker's famous blog post, "Writing an R package from scratch". Then I consulted Hadley Wickham's "R packages" book (1st edition). I consider the "R packages" book (now in its second edition, by Hadley Wickham and Jenny Bryan), to be the authority on best practices for package development, alongside the rOpenSci guide, "rOpenSci Packages: Development, Maintenance, and Peer Review", by Salmon et al.

These are excellent pieces of reference test, however I think there is a need for a resource that sits somewhere between a blog post on making an R package, and resource. I want something that contains **just enough** information to get you started on the right path to making an R package. This is what that book represents to me. Along the way I'll include breadcrumbs to other resources to look into when you want to learn more.

This book also represents my efforts to explain the key parts of what I think people should know about how to write functions, and also to format this in a teachable way that can be covered in a single workshop.

There are more comprehensive guides, and other guides out there for writing R packages. It has been the fundamental way people have shared code and ideas. So I want to share some resources I really enjoyed and think are great:

So, why write a book?

Similar to my book, "Quarto for Scientists", writing this as a book provides a nice way to structure the content in the form of a workshop, in a way suitable for learning in a day. It is not to say that there aren't already the resources out there; there are. It is instead adding to the list of other (useful, hopefully!) information out there on the internet. To answer a question with another question: "Why NOT write this as a book?"

13

Installation

In this section, the aim is to have everyone setup with R, RStudio, the tools you need to build an R package, and git.

2.1 Overview

• **Duration** 15 minutes

2.2 Questions

- How do I install R?
- How do I install RStudio
 - What about Positron?
- How do I install git?
- How do I install RTools?

2.3 Software Setup

2.3.1 Installing R

2.3.1.1 Windows

https://cloud.r-project.org/bin/windows/

2.3.1.2 MacOS

https://cloud.r-project.org/bin/macosx/

16 2 Installation

2.3.1.3 Linux

https://cloud.r-project.org/bin/linux/

2.3.2 Installing RStudio

https://posit.co/download/rstudio-desktop/#download

2.3.3 Installing R packages for development

To ensure you are up to date, run the following script to install the packages.

```
install.packages(c("devtools", "roxygen2", "testthat", "knitr", "pak"))
```

2.3.3.1 Personalising your R Profile

This is really neat, and I think it's actually worthwhile doing, but it does take up some time, and there are some warnings.

As you develop R packages, you'll need to go through a cycle of restarting R, and loading things up to be ready. One of the issues with this is that you'll find yourself writing code like:

```
library(devtools)
```

A lot. To save you time, we can edit a very special file called "The R profile", which is saved as .RProfile. This code is special, and awesome, because it is run every time you start R. It is also dangerous, for exactly the same reason.

I recommend running the following code from devtools to help set this up:

```
use_devtools()
```

Which will bring up the following message:

```
Include this code in .Rprofile to make devtools
available in all interactive sessions:
if (interactive()) {
   suppressMessages(require(devtools))
}
[Copied to clipboard]
Modify /Users/nick/.Rprofile.
Restart R for changes to take effect.
```

So, copy and paste the above, which I will now explain. There are three parts to this that I will break down:

```
require(devtools)
```

we usually recommend writing library(devtools), but in this instance,

require is what we want, because if the package is not installed, require will throw a warning, rather than an error:

```
# warn
require(whatevenisthis)
```

Loading required package: whatevenisthis

Warning in library(package, lib.loc = lib.loc, character.only = TRUE,
logical.return = TRUE, : there is no package called 'whatevenisthis'
error

```
library(whatevenisthis)

Error in library(whatevenisthis): there is no package called 'whatevenisthis'
```

We do not want an error when we start R, it is annoying.

```
suppressMessages()
```

This code suppresses any messages that appear from running this code, which again, we want, because we don't (generally) want our R session to announce something upon startup.

```
if (interactive()) {
   suppressMessages(require(devtools))
}
```

This means that this code is only run if the R session is interactive. This always felt a bit strange to me - because I had only ever run R interactively. But you don't want to run require(devtools) when we aren't using R interactively, because it means we are potentially changing the state of things. Essentially, it's good practice.

Also, here are a couple of times that you might not realise you are using R non-interactively:

- rendering a document using quarto or rmarkdown
- building an R package (which you'll learn about later)

You also use R non-interactively when you are running Rscript in the command line.

Finally, another bit of useful code in your R profile is something like this:

18 2 Installation

```
person(
    given = "Nicholas",
    family = "Tierney",
    role = c("aut", "cre"),
    email = "nicholas.tierney@gmail.com",
    comment = c(ORCID = "https://orcid.org/0000-0003-1460-8722")
    )
    )',
    License = "MIT + file LICENSE",
    Language = "en-GB",
    Version = "0.0.0.9000"
),
    # set SI to true
    reprex.session_info = TRUE
)
```

This helps when setting up your R package for the first time, to make sure you set up your DESCRIPTION file. It isn't required, but it is neat, and I think worthwhile.

Because I need to set these things up on different laptops sometimes, I actually write all these files to github. They are typically called "dotfiles" - you can see mine at http://github.com/njtierney/dotfiles.

2.3.4 git and github

Very briefly, git is essentially a way of managing versions and changes. You can think of it like a product such as dropbox, but with super powers. You can go back in time, you can make copies for changing, and delicately and precisely mege them back in, or leave them where they are.

Your software needs a home. You'll typically start with your project on your laptop or computer. GitHub is where you can store it online. The benefits to sharing your work on github are many, but my personal top reasons are:

- Build trust in your software. If the community can see your code, they can trust it better.
- Provides a way to log ideas and bugs via issues.
- Provides a way for the community to contribute to your code.

My favourite book on using git and github with R is the book "happy git with R" By Jenny Bryan, Jim Hester, and the Stat 545 TAs. Honestly, it's hard to recommend better installation instructions than their battle tested ones, so I'll point you to this resource in case you run into troubles here.

2.3.4.1 setting up github

Getting set up on github you need an account. It's easy enough to set up - go to https://github.com/ . When picking a username, I recommend the following:

- 1. Keep it short. jsmith is better than jonathansmith.
- 2. Avoid numbers and jokes. jsmith is better than jsmith123
- 3. Keep it professional. jsmithisthebest
- 4. Keep it lowercase

2.3.4.2 installing git

Installing git can sometimes be a challenge. This is largely because there are different ways to install it on windows vs mac vs linux. As states earlier, the best, most battle tested instructions are at https://happygitwithr.com/install-git.

Once you've installed git, I recommend running this:

```
usethis::git_vaccinate()
```

Which ensures that you ignore specific files (specifically, Rproj.user, .Rhistory, .Rdata, .httr-oauth, .DS_Store, and .quarto). This is important because it decreases your chances of leaking credentials or other important details to GitHub.

2.3.4.3 The "git handshake"

In order for your computer to talk to git and github properly, it needs to know three things:

- 1. Name
- 2. Email
- 3. Credentials

git needs to know your name and email - this should be the name and email you used to set up your github account. Set this up with use_git_config()

```
library(usethis)
use_git_config(
  user.name = "Ned Kelly",
  user.email = "ned@example.org"
)
```

github needs a personal access token - this is so you can talk to github from R. This becomes really handy, and dare I say it, nearly magical later on. To get this, run:

20 2 Installation

```
usethis::create_github_token()
```

This will open up GitHub and create a Personal Access Token. If this doesn't work, go to https://github.com/settings/tokens and click "Generate New Token", and select the (classic)."

Generally speaking you want the following scopes selected: "repo", "user", and "workflow".

A token will be created - keep this page open, and copy the token to your clipboard.

Then, go to R, and run:

```
gitcreds::gitcreds_set()
```

And paste this PAT code in. Then, verify all of this with:

```
usethis::git_sitrep()
```

2.3.5 Installing RTools

This is actually something that you only need to do if you want to use C or C++ with your R package, which isn't something you need to do for this course. To read more on this, see "The R build toolchain" from the R Packages book.

RStudio, What and Why

(This section is also in my other book, "Quarto for Scientists")

3.1 Overview

- Teaching 5 minutes
- Exercises 2 minutes

3.2 Questions

- What is RStudio?
- Why should I use RStudio?
- What features should I change?

3.3 Objectives

- Get familiarised with RStudio
- Get set up with not storing the RStudio workspace
- Download the course materials for the workshop

3.4 What is RStudio, and why should I use it?

If R is the engine and bare bones of your car, then RS tudio is like the rest of the car. The engine is super critical part of your car. But in order to make things properly functional, you need to have a steering wheel, comfy seats, a radio, rear and side view mirrors, storage, and seatbelts. RStudio is all those niceties

The RStudio layout has the following features:

- On the upper left, the Quarto script
- On the lower left, the R console
- On the lower right, the view for files, plots, packages, help, and viewer.
- On the upper right, the environment / history pane

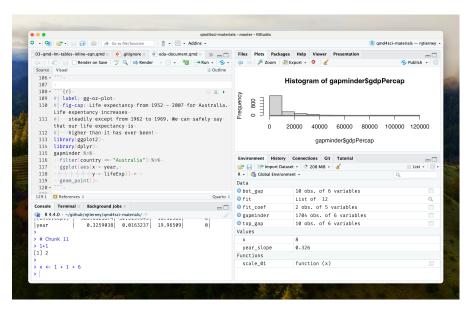


Figure 3.1: A screenshot of the RStudio working environment.

We saw a bit of what an Quarto script does.

- The R console is the bit where you can run your code.
- The file/plot/package viewer is a handy browser for your current files, like Finder, or File Explorer.
- Plots are where your plots appear, you can view packages, see the help files.
- The environment / history pane contains the list of things you have created, and the past commands that you have run.

Your Turn: RStudio default options

To first get set up, I highly recommend changing the following setting Tools > Global Options (or Cmd + , on macOS) Under the **General** tab:

- For workspace:
 - Uncheck restore .RData into workspace at startup.
 - Save workspace to .RData on exit : "Never".
- For **History**:
 - Uncheck "Always save history (even when not saving .RData).
 - Uncheck "Remove duplicate entries in history".

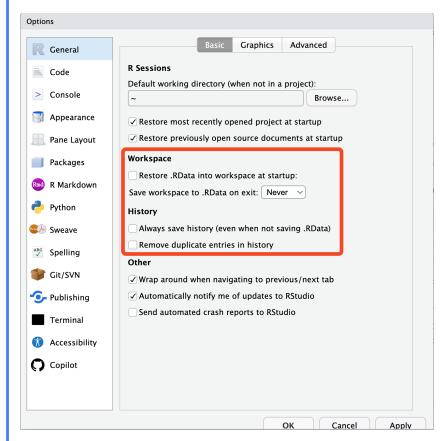


Figure 3.2: Setting the options right for RStudio, so you don't restore previous sessions work, and don't save it either.

This means that you won't save the objects and other things that you create in your R session and reload them. This is important for two reasons

1. **Reproducibility**: you don't want to have objects from last week cluttering your session

2. **Privacy**: you don't want to save private data or other things to your session. You only want to read these in.

Your "history" is the commands that you have entered into R. Additionally, not saving your history means that you won't be relying on things that you typed in the last session, which is a good habit to get into!

3.5 Learning more

• RStudio IDE cheatsheet

4

Workflow

(Note that this section is borrowed from my book, Quarto for Scientists: "work-flow")

Before we start with Quarto, we need to make sure that you understand file storage hygiene.

We can prevent **unexpected problems** if we can maintain an order to your files, paths, and directories. A common problem that arises is R not knowing where a certain file is. For example, we get the error:

```
read.csv("my-very-important-data-file-somewhere.csv")
```

Warning in file(file, "rt"): cannot open file
'my-very-important-data-file-somewhere.csv': No such file or directory

Error in file(file, "rt"): cannot open the connection

Because R doesn't know where "my-very-important-data-file-somewhere.csv" is.

Practicing good file storage hygiene will help maintain an order to files, paths, and directories. This will make you more productive in the future, because you'll spend less time fighting against file paths.

Not sure what a file path is? We explain that as well.

4.1 Overview

- Teaching 10 minutes
- Exercises 10 minutes

26 4 Workflow

4.2 Questions

- Where should I put all my files?
- What is an RStudio project, anyway?
- What is a file path?

4.3 Objectives

- Understand what a file path is
- Set up an RStudio Project to organise your work
- Put some data in your project to set up the next tasks

i Your Turn

In groups of 2-4 discuss:

- 1. What your normal "workflow" is for starting a new project
- 2. Possible challenges that might arise when maintaining your project

4.4 When you start a new project: Open a new RStudio project

This section is heavily influenced by Jenny Bryan's great blog post on project based workflows.

Sometimes this is the first line of an R Script or R markdown file.

setwd("c:/really/long/file/path/to/this/directory")



What do you think the setwd code does?

4.4.1 So what does this do?

This says, "set my working directory to this specific working directory".

It means that you can read in data and other things like this:

```
data <- read_csv("data/mydata.csv")</pre>
```

Instead of

```
data <- read_csv("c:/really/long/file/path/to/this/directory/data/mydata.csv")</pre>
```

So while this has the effect of **making the file paths work in your file**, it is a problem. It is a problem because, among other things, using **setwd()** like this:

- Has 0% chance of working on someone else's machine (this could include you in 6 months!)
- Your file is not self-contained and portable. (Think: "What if this folder moved to /Downloads, or onto another machine?")

So, to get this to work, you need to hand edit the file path to your machine.

This is painful.

When you do this all the time, it gets old, fast.

4.5 What is a file path?

This might all be a bit confusing if you don't know what a file path is. A file path is the machine-readable directions to where files on your computer live. So, the file path:

/Users/njtierney/Desktop/qmd4sci-materials/demo.R

Describes the location of the file "demo.R". This could be visualised as:

users

```
njtierney
Desktop
qmd4sci-materials
demo.R << THIS IS THE FILE HERE
exercises
exploratory-data-analysis
eda-document.qmd
eda-script.R
data
gapminder.csv
```

So, if you want to read in the gapminder.csv file, you might need to write code like this:

28 4 Workflow

gapminder <- read csv("/Users/njtierney/Desktop/qmd4sci-materials/data/gapminder.csv")</pre>

As we now know, this is a problem, because this is not portable code. It is unlikely someone else will have the gapminder.csv data stored under the folders, "Users/njtierney/Desktop".

If you have an RStudio project file inside the qmd4sci-materials folder, you can instead write the following:

gapminder <- read_csv("data/gapminder.csv")</pre>

Your Turn

- (1-2 minutes) Imagine you see the following directory path: "/Users/miles/etc1010/week1/data/health.csv" what are the folders above the file, health.csv?
- What would be the result of using the following code in demo-gapminder.qmd, and then using the code, and then moving this to another location, say inside your C drive?

setwd("Downloads/etc1010/week1/week1.qmd)

4.6 Is there an answer to the madness?

This file path situation is a real pain. Is there an answer to the madness?

The answer is yes!

I highly recommend when you start on a new idea, new research project, paper. Anything that is new. It should start its life as an **rstudio project**.

An rstudio project helps keep related work together in the same place. Amongst other things, they:

- Keep all your files together.
- Set the working directory to the project directory.
- Starts a new session of R.
- Restore previously edited files into the editor tabs.
- Restore other rstudio settings.
- Allow for multiple R projects open at the same time.

This helps keep you sane, because:

- Your projects are each independent.
- You can work on different projects at the same time.

- Objects and functions you create and run from project idea won't impact one another.
- You can refer to your data and other projects in a consistent way.
 And finally, the big one:

RStudio projects help resolve file path problems, because they automatically set the working directory to the location of the rstudio project.

Let's open one together.

Your Turn Use your own rstudio project

1. In RStudio, and run the following code to start a new rstudio project called "qmd4sci-materials".

usethis::use_course("njtierney/qmd4sci-materials")

- 2. Follow the prompts to download this to your desktop and then run the rstudio project. (You can move it later if you like!)
- 3. You are now in an rstudio project!

Your Turn: open the demo.R file

- 1. Run the code inside the demo.R file
- 2. Why does the read_csv code work?
- Run the code inside the exploratory-data-analysis folder - eda-script.R.
- 4. Does the read csv code work?
- 5. Run the code inside the exploratory-data-analysis folder eda-document.qmd, by clicking the "render" button (we'll go into this in more detail soon!)
- 6. Does it work?

4.7 The "here" package

Although RStudio projects help resolve file path problems, in some cases you might have many folders in your r project. To help navigate them appropriately, you can use the here package to provide the full path directory, in a compact way.

30 4 Workflow

here::here("data")

returns

[1] "/Users/nick/github/njtierney/qmd4sci-materials/data"

And

here::here("data", "gapminder.csv")

returns

[1] "/Users/nick/github/njtierney/qmd4sci-materials/data/gapminder.csv"

(Note that these absolute file paths will indeed be different on my computer compared to yours - super neat!)

You can read the above here code as:

In the folder data, there is a file called gapminder.csv, can you please give me the full path to that file?

This is really handy for a few reasons:

- 1. It makes things *completely* portable
- 2. Quarto documents have a special way of looking for files, this helps eliminate file path pain.
- 3. If you decide to not use RStudio projects, you have code that will work on any machine

4.8 Remember

If the first line of your R script is

setwd("C:\Users\jenny\path\that\only\I\have")

I will come into your office and SET YOUR COMPUTER ON FIRE .

– Jenny Bryan

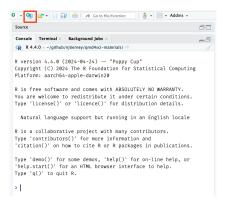
🌢 Aside: Creating an RStudio project

You can create an Rstudio project by going to:

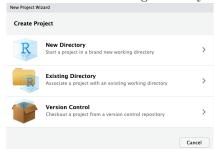
file > new project > new directory > new project > name your project > create project.

You can also click on the create project button in the top left corner

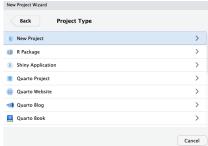
4.8 Remember 31



Then go to new directory, if it is a new folder - otherwise if you have an existing folder you have - click on existing directory.



Then go to new project



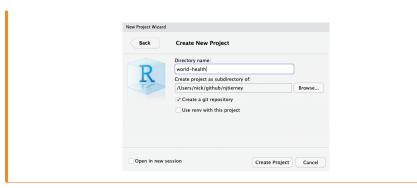
Then write the name of your project. I think it is usually worthwhile spending a bit of time thinking of a name for your project. Even if it is only a few minutes, it can make a difference. You want to think about:

- Keeping it short.
- No spaces.
- Combining words.

For example, I had a project looking at bat calls, so I called it screech, because bats make a screech-y noise. But maybe you're doing some global health analysis so you call it "world-health".

And click "create project".

32 4 Workflow



5

Summary

In this lesson we've:

- $\bullet~$ Learnt what file paths are
- How to setup an rstudio project
- $\bullet\,$ How to construct full file paths with the here package

Why functions?

At their core, an R package is a way to share code. The way we share that code is primarily through R functions. There is a lot about the mechanics, and the tools to create and write R packages, but what I want to communicate here is the **what**, **why**, **when**, **and how** of using functions.

6.1 Overview

- Teaching 20 minutes
- Exercises 15 minutes

6.2 Questions

- What is a function?
- Why should I use a function?
- When should I use a function?
- How do I create a function?

6.3 Objectives

- Understand why functions should be used
- Understand when do use functions
- Understand how to write functions

6.4 Prior Art

There's a lot of work and thought that's gone into writing functions. A lot of my own understanding of this has been informed by others, and I want to make sure I properly acknowledge them:

- Joe Cheng: You have to be able to reason about it
- Hadley Wickham's 'Many Models' talk
- Hadley Wickham's 'The design of everyday functions'
- Miles Mcbain's 'Our colour of magic'
- Jenny Bryan's 'Code Smells and Feels'
- Roger Peng's 'From tapply to Tidyverse'
- Advanced R: Functions
- Tidy Design Principles
- Lexical Scope and Statistical Computing
- stat545 chapter on functions

These are all well worth the time reading, but if I had to pick two, I would say that Hadley Wickham's "Many Models" talk, and Jenny Bryan's "Code Smells and Feels" have been two of the most influential on me.

6.5 Code is for people

If I could have you walk away with one key idea, it would be this:

Functions are tools to manage complexity that allow us to reason with and understand our code.

In essence, **code** is for **people**. This stems from a famous (well, I think it's famous), quote:

[W]e want to establish the idea that a computer language is not just a way of getting a computer to perform operations but rather that it is a novel formal medium for expressing ideas about methodology. Thus, **programs** must be written for people to read, and only incidentally for machines to execute.

— Structure and Interpretation of Computer Programs. Abelson, Sussman, and Sussman, 1984.

6.6 OK, but what actually is a function?

Going back to my quote:

Functions are tools to manage complexity that allow us to reason with and understand our code.

I actually think before we talk about the anatomy, the **what**. We first must discuss **why** functions.

A function is something that helps us manage complexity. You can think about this as something that allows us to repeat certain tasks. Kind of like how a robot, or a manufacturing line can repeat manual tasks.

Let's say we had some data on age groups - the number of contacts these people record on a given day.

250

library(tidyverse)

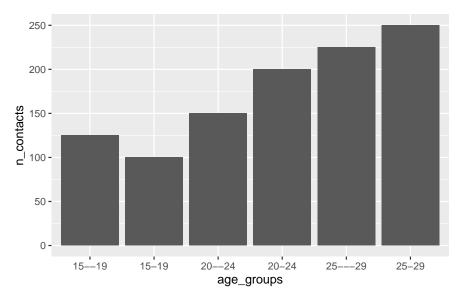
6 NSW

25-29

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr
            1.1.4
                      v readr
                                   2.1.5
v forcats
            1.0.0
                      v stringr
                                   1.5.1
v ggplot2
            3.5.1
                      v tibble
                                   3.2.1
                                   1.3.1
v lubridate 1.9.4
                      v tidyr
v purrr
            1.0.4
-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()
                  masks stats::lag()
i Use the conflicted package (<a href="http://conflicted.r-lib.org/">http://conflicted.r-lib.org/</a>) to force all conflicts to be
contact <- tibble(</pre>
  location = rep(c("QLD", "NSW"), 3),
  age_groups = c("15-19", "15--19", "20--24", "20-24", "25---29", "25-29"),
  n_{contacts} = c(100, 125, 150, 200, 225, 250)
contact
# A tibble: 6 x 3
  location age_groups n_contacts
  <chr>
           <chr>>
                            <dbl>
1 QLD
           15-19
                              100
2 NSW
           15--19
                              125
3 QLD
           20--24
                              150
4 NSW
           20 - 24
                              200
5 QLD
           25---29
                              225
```

We want to produce a plot of age groups and the number of contacts, but we can't do this, because there are all these different ways of representing "age_group".

```
ggplot(
  contact,
  aes(x = age_groups,
     y = n_contacts)) +
  geom_col()
```



Well rather, we CAN do this, but we want to get the totals of each age group.

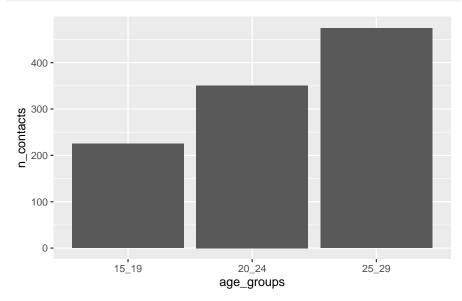
What we want out of this is them all to be separated out by an underscore "_", and turned into a factor:

```
library(stringr)
tidy_contact <- contact |>
  mutate(
    age_groups = str_replace_all(
    string = age_groups,
    pattern = "---|--|-",
    replacement = "_"
    ),
    age_groups = as.factor(age_groups)
)
tidy_contact
```

```
# A tibble: 6 x 3
  location age_groups n_contacts
  <chr>>
                             <dbl>
            <fct>
1 QLD
            15_19
                               100
2 NSW
            15_19
                               125
3 QLD
            20_24
                               150
4 NSW
            20_24
                               200
5 QLD
            25_29
                               225
6 NSW
            25_29
                               250
```

And then we can plot this:

```
ggplot(
  tidy_contact,
  aes(x = age_groups,
      y = n_contacts)) +
  geom_col()
```



Sure, job done.

But now we have some new data, this one contains similar information, but it has population data that we need to join onto it so we can get proportion information.

```
population <- tibble(
  location = rep(c("QLD", "NSW"), 3),
  age_groups = c("15--19", "15-19", "20---24", "20-24", "25-29", "25--29"),
  population = c(319014, 468550, 338824, 540233, 370468, 607891)
)</pre>
```

population

```
# A tibble: 6 x 3
  location age_groups population
  <chr>>
            <chr>
                             <dbl>
1 QLD
            15--19
                            319014
2 NSW
            15-19
                            468550
           20---24
                            338824
3 QLD
4 NSW
           20 - 24
                            540233
5 QLD
            25-29
                            370468
6 NSW
           25--29
                            607891
```

And this is why I think you should write a function. We want to *encapsulate the idea* of cleaning up age group. That is: "clean age groups". So let's write a function that captures this idea.

```
clean_age_groups <- function(age_groups){

  age_underscore <- str_replace_all(
      string = age_groups,
      pattern = "---|--|-",
      replacement = "_"
      )
  as.factor(age_underscore)
}</pre>
```

And this is the difference in the worflow, for each of these tidying up processes:

6.7 script

```
tidy_contact <- contact |>
  mutate(
    age_groups = str_replace_all(
        string = age_groups,
        pattern = "---|--|-",
        replacement = "_"
        ),
        age_groups = as.factor(age_groups)
)

tidy_population <- population |>
```

6.8 function 41

```
mutate(
    age_groups = str_replace_all(
     string = age_groups,
     pattern = "---|--|-",
     replacement = "_"
    age_groups = as.factor(age_groups)
tidy_proportion <- tidy_contact |>
 left_join(tidy_population,
            by = c("location", "age_groups")) |>
 mutate(proportion = n_contacts / population)
tidy_proportion
# A tibble: 6 x 5
 location age_groups n_contacts population proportion
  <chr> <fct>
                        <dbl>
                                    <dbl>
1 QLD
         15_19
                            100
                                    319014 0.000313
2 NSW
          15_19
                            125
                                    468550
                                           0.000267
                                    338824 0.000443
3 QLD
          20_24
                            150
4 NSW
          20_24
                            200
                                    540233 0.000370
          25_29
                            225
5 QLD
                                    370468
                                           0.000607
6 NSW
          25_29
                            250
                                    607891
                                           0.000411
```

6.8 function

```
clean_age_groups <- function(age_groups){
    age_underscore <- str_replace_all(
        string = age_groups,
        pattern = "---|--|-",
        replacement = "_"
        )
    as.factor(age_underscore)
}
tidy_contact <- contact |>
```

6.8.1 Functions give ideas a home

Functions provide a way to **express** the idea of what we want to do. They also provide your ideas a home. What if the data changes? Do you want to go back and change each line of code? No! You can update the function in one place, and then repeat it again.

Once you start writing functions to do things, they will start to be little repositories of knowledge. Little shortcuts that you can use to just remember the most important part.

Now, on to the anatomy of functions

6.9 Anatomy of a function

Now, to speak about the mechanics of writing functions: a function is composed of three parts:

- 1. Name
- 2. Arguments
- 3. Body

To look at our clean_age_groups function again, we can see the following:

```
# The name of the function
clean_age_groups <- function(age_groups){ # The argument - age_groups

# The body of the function
age_underscore <- str_replace_all(</pre>
```

```
string = age_groups,
   pattern = "---|--|-",
   replacement = "_"
)

# The last thing you do with the function is what it returns
as.factor(age_underscore)
}
```

⚠ The last thing you do shouldn't be assignment <-

The last thing that a function does is what it returns. If we take our example above and change the last line to assign to some variable, then the function will not return anything!

This is a pretty common mistake, one I still make. Just something to be aware of! The way to fix this is to make sure that the last thing you do isn't assigned. So, our example above should look like so:

```
# The name of the function
clean_age_groups <- function(age_groups){ # The argument - age_groups</pre>
  # The body of the function
  age_underscore <- str_replace_all(</pre>
      string = age_groups,
      pattern = "---|--|-",
      replacement = "_"
  # The last thing you do with the function is what it returns
  # NOT THIS
  # factored <- as.factor(age underscore)</pre>
  # THIS
  as.factor(age_underscore)
}
clean_age_groups("10--11")
[1] 10_11
Levels: 10_11
```

6.10 How to think about writing functions

There are many ways to start writing functions. Fundamentally, it is about identifying inputs and outputs. One useful approach, I think, is to identify the outputs before the inputs:

- 1. The output. What **one thing** do you want this function to return?
- 2. The input. What (potentially many things) goes in to this.

This "gestalt", or top-down approach isn't how it always needs to be done. But I think it helps you identify **the thing you need** first, which can help guide you.

6.10.1 Identifying the output - what do we need?

It might feel a bit like putting the cart before the horse, but I think there is a nice advantage to thinking about the output first: you focus on what you want the function to do.

In the case of our clean_age_groups function, we want to get values like "15_19" that are factors.

6.10.2 Identifying the input

So now we have a clear idea of what we need - we can now clarify what we have, which in our case earlier, was some contact data

contact

```
# A tibble: 6 x 3
  location age_groups n_contacts
  <chr>
            <chr>
                             <dbl>
            15-19
1 QLD
                               100
2 NSW
            15--19
                                125
3 QLD
            20--24
                               150
4 NSW
            20 - 24
                                200
5 QLD
            25---29
                                225
6 NSW
            25-29
                                250
```

Where we want to focus on age groups, and take inputs like

```
c("15-19", "15--19")
```

```
[1] "15-19" "15--19"
```

And then turn them into:

```
c("15_19", "15_19")
```

```
[1] "15_19" "15_19"
```

Breaking things down like this means we can focus on a really small example of the thing we want, which makes the problem easier to solve.

There are many ways to manage turning strings into other strings, and I like to use the stringr package to do this. We can use the str_replace_all function. So I'll start by scratching up some inputs like so, and seeing if this works

```
ages <- c("15-19", "15--19")
str_replace_all(
   string = ages,
   pattern = "-",
   replacement = "_"
)</pre>
```

```
[1] "15_19" "15__19"
```

6.10.2.1 Iteration: Writing functions is writing

I didn't get this right the first time - and I rarely do! The point I want to make here is:

Writing functions is just like writing. It takes iteration.

We have incidentally replaced every "-" with "_-", which means "–" becomes "___".

Let's change that by using \mid in the "pattern" argument, which allows us to specify "- \mid -", which means, "-" OR "-":

```
ages <- c("15-19", "15--19")
str_replace_all(
    string = ages,
    pattern = "-|--",
    replacement = "_"
)</pre>
```

```
[1] "15_19" "15__19"
```

OK, the same problem. We actually need to flip the order here, so we change "-" first:

```
ages <- c("15-19", "15--19")
str_replace_all(
    string = ages,
    pattern = "--|-",
    replacement = "_"
)</pre>
```

```
[1] "15_19" "15_19"
```

Great! Now let's put that into the body of the function, and give the function a good name.

```
clean_age_groups <- function(age_groups){
   str_replace_all(
   string = ages,
   pattern = "--|-",
   replacement = "_"
)
}
clean_age_groups(ages)</pre>
```

```
[1] "15_19" "15_19"
```

It's a useful process to scratch out a function like this. As you get more

confident with this, you will start to be able to write the code as a function first, and then iterate in that way.

▲ beware copying and pasting into functions

The process of writing a function out in scratchings as we've done, is that we can leave some scraps in the code. In this case, I've actually left the ages object in the function, but the argument is age_groups:

```
clean_age_groups <- function(age_groups){
   str_replace_all(
   string = ages,
   pattern = "--|-",
   replacement = "_"
)
}
clean_age_groups(ages)</pre>
```

[1] "15 19" "15 19"

Notice that this still works! This is because the ages object still exists as a variable I've created. But if we try another input, we'll get some strange output:

```
clean_age_groups(c("10-12", "10--12"))
```

[1] "15_19" "15_19"

So, make sure to clean up after you've copied and pasted - remember to check the arguments match how they are used in the function.

And on that note, let's redefine clean_age_groups correctly so we don't get an error later on (which happened during the development of the book)

```
clean_age_groups <- function(age_groups){
   str_replace_all(
   string = age_groups,
   pattern = "--|-",
   replacement = "_"
)
}
clean_age_groups(ages)
[1] "15_19" "15_19"</pre>
```

6.10.3 Managing scope - functions are best (generally) when they do one thing.

Also, note that we wrote clean_age_groups to just focus on converting input like "10–12" into "10_12". We could have instead focussed on cleaning up the data frame, like so:

contact

```
# A tibble: 6 x 3
  location age_groups n_contacts
  <chr>
            <chr>>
                             <dbl>
                               100
1 QLD
            15-19
2 NSW
            15--19
                               125
           20--24
3 QLD
                               150
4 NSW
           20-24
                               200
            25---29
                               225
5 QLD
6 NSW
            25-29
                               250
```

```
clean_age_groups_data <- function(data){
  tidy_contact <- data |>
  mutate(
    age_groups = str_replace_all(
        string = age_groups,
        pattern = "---|--|-",
        replacement = "_"
        ),
        age_groups = as.factor(age_groups)
  )
  tidy_contact
}
clean_age_groups_data(contact)
```

A tibble: 6 x 3 location age_groups n_contacts <chr> <fct> <dbl> 1 QLD 15_19 100 2 NSW 15_19 125 3 QLD 20_24 150 4 NSW 20_24 200 5 QLD 25_29 225 6 NSW 250 25_29

I think there are a couple of issues with this:

- 1. We assume the age groups column is always age_groups
- 2. The scope is now larger we are always working with data and returning data
- 3. We haven't necessarily made the expression easier.

It is fine to wrap up the existing function into another function that cleans the data - to me this better encapsulates and expresses the ideas:

contact

```
# A tibble: 6 x 3
  location age_groups n_contacts
  <chr>
           <chr>
                             <dbl>
1 QLD
           15-19
                               100
2 NSW
           15--19
                               125
3 QLD
           20--24
                               150
4 NSW
           20-24
                               200
5 QLD
           25---29
                               225
6 NSW
           25-29
                               250
```

```
clean_contacts <- function(data){
  data |>
  mutate(
    age_groups = clean_age_groups(age_groups)
  )
}
clean_contacts(contact)
```

A tibble: 6 x 3

```
location age_groups n_{contacts}
  <chr>
            <chr>>
                              <dbl>
1 QLD
            15_19
                                100
2 NSW
            15_19
                                125
3 QLD
            20_24
                                150
4 NSW
            20_24
                                200
            25__29
                                225
5 QLD
6 NSW
            25_29
                                250
```

Some of the improvements I notice

- We are just focussing on cleaning up the age group column.
- We have given it a name that refers to cleaning up the data, which might also give us some space and room to add more cleaning function here.

6.11 When to function

One of my overall points with functions is:

functions help you express your intention.

However, there are some generally good heuristics to follow to help guide you towards writing a function. Here are some of these.

Generally, it is time to write a function if:

- 1. You've copied and pasted the code 3 or more times.
- 2. You've re-read your code more than 3 times.

6.12 Naming things is hard

There are only two hard things in Computer Science: cache invalidation and naming things.

- Phil Karlton

What does this function **do**?

```
myfun <- function(x){
  (x * 9/5) + 32
}</pre>
```

Converting temperature?

```
temperature_conversion <- function(x){
  (x * 9/5) + 32
}</pre>
```

Clearly state input_to_output()

```
celcius_to_fahrenheit <- function(x){
  (x * 9/5) + 32
}</pre>
```

Name argument and intermediate variables

```
celcius_to_fahrenheit <- function(celcius){
  fahrenheit <- (celcius * 9/5) + 32
  fahrenheit
}</pre>
```

What, what does make functions hard?

```
celcius_to_fahrenheit <- function(celcius){
  (celcius * 9/5) + 32
}</pre>
```

Identifying inputs and outputs is hard.

But what is hard it taking code, (like the code in a data analysis) and finding the parts that need to change

There's a level of "I got it to work" and there's a level of "It works, and I can reason about it"

 Joe Cheng You have to be able to reason about it | Data Science Hangout

I can **reason** about it

...how do you take all this complexity and break it down into smaller pieces...each of which you can **reason about**...each of which you can **hold** in your head...each of which you can look at and be like "yup, I can fully ingest this entire function definition, I can read it line by line and prove to myself this is definitely correct...So software engineering... is a lot about this: How do you break up inherently complicated things that we are trying to do into small pieces that are individually easy to reason about. That's half the battle...The other half of the battle is how do we combine them in ways that can be reliable and also easy to reason about

6.13 The other hard part of writing functions

i Practice naming things

Practice naming these functions, going through

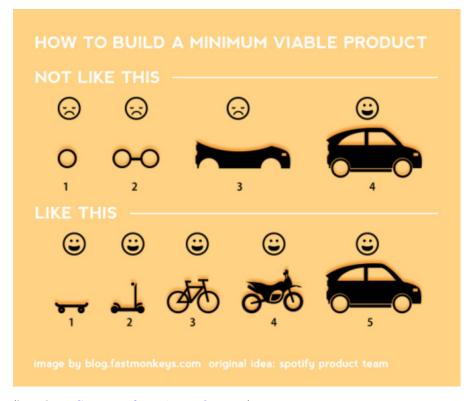
6.14 Conclusion

The process of writing a function is:

- Identify outputs and inputs
- Identify the complexity to abstract away
- Writing functions is iterative, Just like regular writing

• Naming things is hard

On a final note, I think it's worthwhile thinking about the iteration - and the idea of moving from a skateboard to a car, rather than building the car:



(heard via Stat545 functions chapter)

7

Motivation

We've gone through a lot of setup, and now we're going to start building an R package. Soon. But we need to have some motivation, first. It involves a bit of a story, and a bit of imagination.

7.1 Overview

- Teaching 30 minutes
- Exercises 10 minutes

7.2 Questions

.

7.3 Objectives

• Start to wrangle with a script to turn it into functions

7.4 How this works

One of my big goals with teaching functions, and with teaching R packages, is that I want the examples to be somewhat rooted in the familiar and the real. There are really useful toy examples of writing packages that deliver praise (e.g., ones I've used to teach R packages in the past:

54 7 Motivation

https://github.com/njtierney/praiseme), or do simple conversions between units (celcius to farenheit being a very common example).

These examples are useful because they teach you the tools, and the process. However in this course, I want to focus on a bit more than this and incorporate the process of turning code into functions. I think this is important, because it more closely represents other examples we come across in using R, and presents a bit of a richer learning journey, because in addition to learning about the tools and the process of R package building, you will also learn:

- How to think about converting scripts to functions
- How to write better functions

I have written up some example code, which starts as a quarto document. We are going to take this document, and then eventually turn it into an R package.

The structure of this exercise has taken inspiration from "The package within" chapter from R Packages.

7.5 The example: "learned"

We are going to be looking at a role-play situation where we imagine we are at some fictional workplace, where part of our job is to look at education data that we have acquired from some source. The overall goal of our job here is to produce some **key outputs** from some data.

You can see this example at: https://github.com/njtierney/learned

To download it, run the following code

```
library(usethis)
use_course("njtierney/learned")
```

i Your turn

- 1. Download the repository using the code above
- 2. Render the document
- 3. Read over the document, thinking about what we discussed in Why functions.
- 4. Identify some potential problems with the code
- 5. Think about what might happen if we want to read in data from 2015 (or later years), how would you like to do this?

7.6 Discussion of potential problems

After you've taken from time to think about some of the potential problems, open the box below

Some of the potential problems

- Copying and pasting a document could lead to errors
- What if the data changes?
- What if other people collaborate on this project? How do they have the source of truth?
- Is there a way to formalise this all?

7.7 Identifying the report outputs

To get us started with some key things, let's think about what the key outputs of this report are.

i Your turn

- 1. Identify the **key outputs** of the report
- 2. Pick one of those key outputs and start to write out a function for it

♦ Key outputs

They key outputs are related to the "Produce a ..." steps of the document:

- 1. Produce a plot of the **proportion of people educated** in each age group in each state
- Produce a box plot of proportion of people educated for each state.
- 3. Produce a table of The 5 number summary (min, 1st quantile, median, 3rd quantile, max) of **proportion of people** educated for each state.

So now we know where we are headed - we want to write some functions that produce these plots and tables.

56 7 Motivation

However, the main problem that we encountered was that there was actually a bit of data cleaning that needed to happen before we did this. Let's focus on cleaning up and rearranging the quarto document first to identify the data cleaning steps required.

Your turn

- 1. Open, "alpha-analysis2014.qmd"
- 2. Move all the "data quality" checks into a new section called "data quality"
- 3. Move all of the data cleaning code up to the top, so we just work with one data set, named tidy_age_state_education_2014
- 4. Create two functions to clean the data:
- 5. tidy the age groups
- 6. remove the missing values
- 7. Put these two functions into another function that does the data cleaning

7.8 Plot of proportion of people educated in each age group in each state

The key code here is this:

```
ggplot(
  fixed_age_group_education_2014,
  aes(
    x = prop_studying,
    y = age_group
)
) +
  geom_col() +
  facet_wrap(~state_territory, ncol = 2)
```

However, in order to write this, we need to identify the steps that happened earlier to the fixed_age_group_education_2014 data.

i Your turn

What were the fixes that were required for this data?

- 7.9 box plot of proportion of people educated for each state.
- 7.10 Table of The 5 number summary of proportion of people educated for each state.

Α

Bibliography