

Teach Yourself or Others R

Adam Rawles

Contents

1	teacheR	5
1.1	Overview	5
1.2	About Me	5
2	Introduction to R	7
2.1	What is R and RStudio?	7
2.2	Operators	9
2.3	Variable assignment	11
2.4	Data types	13
2.5	Data structures (and subsetting)	19
2.6	Functions	20
3	Data Analysis in R	21
3.1	Installing packages	21
3.2	Loading data	21
3.3	Cleaning data	21
3.4	Summary statistics	21
3.5	Plots	22
4	Programming in R	23
4.1	User-defined functions	23
4.2	For loops	23
4.3	If/else statements	23
5	Tidyverse	25
5.1	tidyr	25
5.2	dplyr	25
5.3	stringr	25
5.4	ggplot2	25
6	Modelling	27
6.1	Linear modelling	27
6.2	Clustering	27
6.3	Programming and modelling	27

Chapter 1

teacheR

1.1 Overview

This book is a collection of training materials for an introduction to the R statistical computing programming language. Broken down into chapters, I've aimed to cover most of the basics. Alongside each chapter is a xaringan presentation. This can help for those looking to learn, but can also function as a first step for those looking to begin teaching R but who don't have the time to fully develop their own training modules. Hopefully, all the topics covered in the text version will be covered in the presentation and vice versa, so if you learn visually then you can rest easy knowing that you're not missing out! You can find a link to each presentation on the first page of each chapter.

This is a work in progress, and so I would greatly appreciate any feedback. Anything from typos to content suggestions, feel free to raise a GitHub issue if you feel something should be changed.

1.2 About Me

I began using R in my second year of university, during an internship looking at publication bias correction methods. I was under the tutorship of a member of staff who helped me immensely, but I must confess that I have never taken an official course in R, online or in person. I like to think, however, that this is not always a bad thing. Learning from the bottom up and struggling along the way is a fantastic way to acquire knowledge and instills a very important lesson:

You're not going to know everything there is to know about R. Ever. But that's okay.

I'm now 4 years into my R career and I use R every day. With that in mind, I don't think there has ever been a day where I haven't referred to a tutorial, or Stack Overflow, or even just Googled the name of a function that I've used 1000 times before. There is a great repository of knowledge for R and it's one of the things I love most about the R community. So please never feel as though you're an impostor in a world of R gurus. In reality, everyone else is just as lost as you. But if you keep ticking along and never feel that learning something new in R isn't worth your time, you'll end up doing some great things.

And in a roundabout way, that is part of the reason I decided to develop these materials. I don't pretend to be the ultimate R programmer, because I still know what it's like to learn something from the start. And everyone has to start somewhere. So I hope that I can help impart some of the lessons that I've learnt over the 4 years to anyone who's looking to learn R in a way that won't leave you feeling lost.

The only final note I have before we start learning how to use R is another bit of advice:

Don't believe everything you read

Whilst this is probably a good thing to keep in mind for any type of training, I feel it's particularly relevant with R for two reasons. Firstly, when it comes to programming languages, lots of people have opinions. Some are true, most are not. Most things you read are a mix between fact and opinion, so take everything with a pinch of salt. For example, the developers of the `ggplot2` package are fervently against arbitrary second axes and so support for them in `ggplot2` is limited. I also share this view, but that doesn't mean that I'm right - read, learn, but question and make your own mind up. Secondly, R and particularly all of its packages are prone to change. For this reason, people may make statements relative to one version of R that aren't necessarily true in the future. Things have changed over the years, and so answers from a 10 year-old Stack Overflow question may not still be true when you come across them. A microcosmic version of this are some recent changes in the `tidyr` package. Historically, converting data from/to long and short formats was done using the `spread()` and `gather()` functions. However, in newer releases, these functions are deprecated in favour of `pivot_wider()` and `pivot_longer()`, which provide the same functionality but also some extra bits. The practical implication of this suggestion is don't always read one tutorial on a subject before you dive in.

Chapter 2

Introduction to R

Chapter 2 is a general overview of R and its basics. The xaringan presentation for this module is [here](#).

2.1 What is R and RStudio?

2.1.1 R

R is a public license programming language. More specifically, it's a statistical programming language meaning that it's often used for statistical analysis rather than software development. R is also a functional programming, rather than an object-oriented programming language like Python. This means that operation in R are primarily performed by functions (input, do something, output), but more about that later.

Strictly speaking, R is not just a functional programming language. In reality, a language is never purely one type and R is no exception. There are object-oriented systems in R (three main ones), meaning that object-oriented programming is possible and relatively straightforward in R. Having said that, I feel now is a good time to refer back to the mantra from the Overview chapter:

Don't believe everything you read

So basically, R is a functional programming language with some object-oriented systems. If that means very little to you, don't worry. For the vast majority of users, this is a purely academic definition.

One important attribute about R however that may affect you, is that R is an interpreted language. This essentially means that when you send someone some R code, they need R installed to be able to run it. This means that making full

programs is difficult. Later on, we'll look at the `shiny` package, which can be used to quickly make web apps based on R code. These apps are no different in the sense that they also need R installed to be able to run, but because they are web-based, they are significantly easier to share.

For the most part however, if you want to share R code with colleagues, they'll need to have R installed as well.

2.1.1.1 Should I use R over another language?

People have and will continue to argue about which is better, R or Python or Java or C or writing down mathematical equations on a piece of paper and handing it to a monkey to solve. I imagine you're reading this because you heard that R was good for data analysis, and it absolutely is. And so is Python. They're just... different. Personally, I prefer to use R but I understand that other people don't.

Importantly though, never feel as though you've missed a trick by picking a particular language. Programming is not just a practice, it's a way of thinking, and experience is almost always transferable across languages.

For reference however, here are a few of the things that you can use R for:

- data analysis
- reporting and writing
- web apps
- text analysis

If you're interested in any of these, then you're in the right place.

2.1.1.2 Using R

R is very simple. There is a console where you type commands and get responses. Like the classic command-line interfaces you see when the stereotypical nerd has to hack into the FBI database, you type commands, one at a time, into the console, R processes it, and then produces a response if appropriate. For example, if you type `2 + 2` into the R console and hit enter, you'll get 4.

Writing commands out one at a time can be quite time consuming if you want to make changes however. So we use scripts to store multiple lines of code that can then be run altogether. When you execute a script, each line gets passed one by one to the console and executed. For example, I might make a script with this code:

```
variable1 <- 2 + 2
variable1 / 10
```

```
## [1] 0.4
```


So when I run the script, it will run the first line, then the second without me having to type anything else in.

2.1.2 RStudio

RStudio is separate from R. R is a programming language and RStudio is an integrated development environment or *IDE*. This means that RStudio doesn't actually run any code, it just passes it to R for you, meaning that you'll need R to really use RStudio.

RStudio is a massive part of how you interact with R however. For example, with the exception of a few days when I was waiting for RStudio to be installed, I can't ever remember using R without RStudio.

In the previous chapter, we talked very briefly of the R console and scripts. RStudio helps with this workflow. It makes it easier to create scripts, providing extra tools to help write code quicker, and then acts as a window to R when you want to execute the script.

2.1.2.1 What is an IDE?

At its simplest definition, an IDE helps you get work done in your programming language of choice. It can help you save blocks of code, organise projects, save plots and everything in between. R comes with a basic user interface when you install it, but RStudio provides lot more functionality to help you interact with the R console.

2.1.2.2 RStudio Panes

TO DO

2.2 Operators

2.2.1 Arithmetic operators

At the base of lots of programming languages are the arithmetic operators. These are your symbols that perform things like addition, subtraction, multiplication, etc. Because these operations are so ubiquitous however, the symbols that are used are often very similar across languages, so if you've used Excel or Python or SPSS or anything similar before, then these should be fairly straightforward.

Here are the main operators in use:

```
2 + 2 # addition
## [1] 4
10 - 5 # subtraction
## [1] 5
5 * 4 # multiplication
## [1] 20
100 / 25 # division
## [1] 4
```

2.2.2 Logical operators

Logical operators are slightly different to arithmetic operators - they are used to evaluate a particular criteria. For example, are two values equal. Or, are two values equal *and* two other values different.

To compare whether two things are equal, we use two equal signs (==) in R:

```
1 == 1 # equal
## [1] TRUE
```

Why two I hear you say? Well, a bit later on we'll see that we use a single equals sign for something else.

To compare whether two things are different (not equal), we use !=:

```
1 != 2 # not equal
## [1] TRUE
```

The ! sign is also used in other types of criteria, so the best way to think about it is that it inverts the criteria you're testing. So in this case, it's inverting the "equals" criteria, making it "not equal".

Testing whether a value is smaller or larger than another is done with the < and > operators:

```
2 > 1 # greater than
## [1] TRUE
2 < 4 # less than
## [1] TRUE
```

Applying our logic with the ! sign, we can also test whether something is *not* smaller or *not* larger:

```
1 >! 2 # not greater than
```

```
## [1] TRUE
```

```
2 <! 4 # not less than
```

```
## [1] FALSE
```

Why is the `!` sign before the equals sign in the “not equal” to code, but after the “less than/greater than” sign? No idea. It’d probably make more sense if they were the same, but I suppose worse things happen at sea.

There are three more logical operators, and they are the “and”, “or”, and “xor” operators. These are used to test whether at least one or more than one or only one of the logical comparisons are true or false:

```
1 == 1 | 2 == 3 # or (i.e. are either of these TRUE)
```

```
## [1] TRUE
```

```
1 == 1 & 2 == 3 # and (i.e. are these both TRUE)
```

```
## [1] FALSE
```

The xor operator is a bit different:

```
xor(1 == 1, 2 == 3) # TRUE because only 1 is
```

```
## [1] TRUE
```

```
xor(1 == 1, 2 == 2) # FALSE because both are
```

```
## [1] FALSE
```

For `xor()`, you need to provide your criteria in brackets, but this will make much more sense once we look at functions.

2.3 Variable assignment

Do you ever tell a story to a friend, and then someone else walks in once you’ve finished and so you have to tell the whole thing again?

Well, imagine after the second friend walks in, another friend comes in, and you have to start the story over again, and then another friend comes in and so on and so forth. What would be the best way to save you repeating yourself? As weird as it would look, if you wrote the story down then anyone who came in could just read it, rather than you having to go through the effort of explaining the whole thing each time.

This is essentially what we can do in R. Sometimes you’ll use the same value again and again in your script. For example, say you’re looking at total expenditure

over a year, the value for the amount spent would probably come up quite a lot. Now, you could just type that value in every time you need it, but what happens if the value changed? You'd then have to go through and change it every time it appears.

Instead, you could store the value in a variable, and then reference the variable every time you need it. This way, if you ever have to change the value, you only need to change it once.

Creating variables in R is really easy. All you need to do is provide a valid name, use the `<-` symbol, and then provide a value to assign:

```
hello_im_a_variable <- 100
hello_im_a_variable
```

```
## [1] 100
```

Now, whenever you want to use your variable, you just need to provide the variable name in place of the value:

```
hello_im_a_variable / 10
```

```
## [1] 10
```

You can even use your variable to create new variables:

```
hello_im_another_variable <- hello_im_a_variable / 20
hello_im_another_variable
```

```
## [1] 5
```

When you come across other people's work, you may see that they use `=` instead of `<-` when they create their variables. Even though it's not the end of the world if you do do that, I would recommend getting into the habit of using `<-`. `<-` is purely used for assignment, whereas `=` is actually also used when we call functions, and so it can get a bit confusing if you use them interchangeably.

As a side note, you'll see that the value of the variable isn't outputted when we assign it. If we want to see the value, we need just the name.

Variables are very flexible. You can overwrite a previously defined variable just be reassigning a new value to the same name. R will also give the variable the correct *type* based on the value you assign. So for example, if you assign 20 to a variable, then that variable will be stored as a number. If you assign something in quotation marks like "hello", then R will store it as text. Let's look in a bit more detail at the different data types.

2.4 Data types

Data can be stored in lots of different forms. For example, "TRUE" and TRUE are stored as two different types, even though they look very similar to us.

The main different data types are:

- logical
- TRUE
- FALSE
- double
- 12.5
- 19
- 99999
- integer
- 2L
- 34L
- character
- "hello"
- "my name is"
- factors
- dates
- 2019-06-01
- datetime (POSIXct)
- 2019-06-01 12:00:00

2.4.1 Logical

A logical variable can only have two *real* values, TRUE or FALSE. I say two *real* values, because you can also have things like NA, but that's true of any data type.

Logical variables are used a lot in response questionnaires, where the answer to the question is either "Yes" or "No" (TRUE or FALSE). I would recommend converting any character strings like "Yes" or "No" or "TRUE" or "FALSE" to a logical variable rather than leaving them as characters, because it'll make your analysis less verbose (use fewer lines of code), even if it doesn't change the underlying logic.

To test whether something is stored as logical, we use the `is.logical()` function:

```
is.logical(TRUE)
```

```
## [1] TRUE
```

```
is.logical("TRUE")
```

```
## [1] FALSE
```

To convert a value to logical, use the `as.logical()` function:

```
as.logical(1)

## [1] TRUE
as.logical(0)

## [1] FALSE
as.logical("TRUE")

## [1] TRUE
as.logical("FALSE")

## [1] FALSE
```

Be careful though, just because a conversion seems obvious to you, doesn't mean you'll get the expected result! For instance, what do you think `as.logical(2)` should return? See for yourself.

2.4.2 Double

The best way to think of a `double` value is as a number. It can be a whole number (but see `Integers`) or a decimal. R will often take care of any implicit number conversion that needs to be done under the hood, so the only thing you really need to keep in mind is that when you assign a number, be it a whole number or decimal, it will be stored as `double` by default.

As an aside, it's called `double` because it's stored using double precision.

To check whether a value is stored as `double` (or more generally `numeric`), use the `is.double()` and `is.numeric()` functions:

```
is.double(2)

## [1] TRUE
is.numeric("not numeric")

## [1] FALSE
is.double(2L) # see the next section for why this returns FALSE

## [1] FALSE
```

To convert a value to a `double`, use the `as.double()` or `as.numeric()` functions:

```
as.double("5")

## [1] 5
```

```
as.numeric("10")

## [1] 10
as.double("im going to cause an error")

## Warning: NAs introduced by coercion
## [1] NA
```

2.4.3 Integer

Whilst also storing numeric data (like double), integers are specific to whole numbers. Also, by default, even when you assign a whole number, like this: `5`, R will store that value as double rather than as an integer. To store something explicitly as an integer, suffix the value with an `L`, like this: `5L`. Attempting to store something that isn't an integer as an integer will result in a warning:

```
1.5L

## [1] 1.5
```

For the most part, I let R take care of how it stores numbers, unless I explicitly need it to be of a certain type. This is pretty rare though.

To check if something is an integer, use the `is.integer()` function:

```
is.integer(2)

## [1] FALSE
is.integer(2L)

## [1] TRUE
```

To convert to an integer, use the `L` suffix or the `as.integer()` function:

```
1L

## [1] 1
```

2.4.4 Character

Sometimes called characters, or character strings, or just strings, characters store text. If you assign a value within quotation marks, regardless of what's inside the quotation marks, it will be stored as character. For example, `"5"` stores a character string with the text "5", not the number 5. This is particularly important when you want to start combining variables. For example, `{r, error = TRUE} "5" + 5` doesn't work, because you're trying to add text to a number, which doesn't make sense.

To check whether something is stored as a character, use the `is.character()` function:

```
is.character("hello")
```

```
## [1] TRUE
```

```
is.character(5)
```

```
## [1] FALSE
```

```
is.character(TRUE)
```

```
## [1] FALSE
```

To convert something to a character, use the `as.character()` function:

```
as.character(5)
```

```
## [1] "5"
```

```
as.character(TRUE)
```

```
## [1] "TRUE"
```

2.4.5 Factors

Factors are a unique but useful data type in R. Essentially, factors store different levels that represent some sort of grouping. For example, say you were collecting some information on people from different countries, the column that holds which country the respondent is from could be stored as a factor, with the levels England, Spain, France, etc.

A factor level is made up of two things. A label and a number that represents that group. For example, in my countries example, our factor would have the labels “England”, “Spain”, “France” and the values 1, 2, 3. This means that internally, a factor is essentially a collection of integers representing the level position and character strings representing the level label.

To create a factor, we just use the `factor()` function:

```
factor(c("England", "France", "Spain"))
```

```
## [1] England France Spain
```

```
## Levels: England France Spain
```

It’s also worth remembering that you can have levels that don’t appear in the data you have. For example, in a questionnaire, you may provide the options “None”, “Some”, “All”. But in your responses, you may see that no one chose the “None” option. In that case, you would still create a factor with three levels, even though only two of them appear.

You can also specify whether a factor is *ordered*. You would use an ordered factor when the levels have meaningful order. For instance, in the above example, it would make sense that “Some” is better than “None”, and “All” is better than “Some”. To create an ordered factor, just specify `ordered = TRUE` in your function. By default, the factor will be ordered in the order the values appear, unless you specify levels (see below).

To convert something to a factor, use the `factor()` function if you want to specify levels and labels, or `as.factor()` to do it for you:

```
factor(c("Some", "All"), levels = c("None", "Some", "All"))

## [1] Some All
## Levels: None Some All

factor(c("Some", "All"), levels = c("None", "Some", "All"), ordered = TRUE)

## [1] Some All
## Levels: None < Some < All

as.factor(c("Some", "All"))

## [1] Some All
## Levels: All Some
```

Notice the difference in the output of those three lines. The first allows us to specify the levels (i.e. the values that were possible). The second does the same but we also specify the ordering of the levels, and the third just converts the provided values and generates the levels based on that data.

Note: An important change in R version 4.0.0 is that R will no longer automatically convert strings (characters) to factors when you import data using `data.frame()` or `read.table()`. Prior to 4.0.0, it would automatically convert strings to characters unless otherwise specified.

2.4.5.1 Converting from factors

Sometimes you’ll need to convert data from a factor to something else, usually a character. This is fairly straightforward using the tools we’ve already seen:

```
as.character(factor(c("Some", "None", "All")))

## [1] "Some" "None" "All"
```

2.4.6 Dates

Dates in any language are tricky. Different countries store dates in different formats and different bits of software stores dates in different ways (looking at you Excel). This can make storing values as dates tough.

The most common way of creating a date is to use the `as.Date()` function. To use this function, you just need to provide your date as a character string:

```
as.Date("2019/01/01")
```

```
## [1] "2019-01-01"
```

But Adam, how does R know which one is the month and which is the day? Good question, thank you for asking. By default, R expects your character string to be in the order “Year/Month/Day”. If you don’t provide it in that format, you’ll get a nonsense output:

```
as.Date("01/12/2019")
```

```
## [1] "1-12-20"
```

If your data is in a different format however, you can specify the format:

```
as.Date("01/12/2019", format = "%d/%m/%Y")
```

```
## [1] "2019-12-01"
```

Here, we’re telling R that the string is in the format “Day/Month/Year”. A list of the different codes that can be used in the format parameter can be found [here](#), or by typing “R date codes” into Google.

Because nothing in life is simple, sometimes you’ll get some data that has the date stored as a number. This is because the source of that data has the date stored as the number of days that have passed since an origin date. Because it’s a number, our `as.Date(..., format = ...)` doesn’t work. Instead, we can still use the `as.Date()` function, but we need to specify what the origin date is that the number refers to.

By default, when importing from Excel, the origin date is January 1st 1970, also known as the “epoch date”. Chances are your data source is also using this date, but always check.

Anyway, to specify your origin, we use the `origin` parameter, like this:

```
as.Date(18262, origin = "1970/01/01")
```

```
## [1] "2020-01-01"
```

Notice the format I’ve provided the origin in. It’s the same as the default that R expects, and I would recommend copying that format wherever possible. If you’re someone who just wants to watch the world burn, then you can specify a format for your origin as well...

```
as.Date(18262, origin = as.Date("01/01/1970", format = "%d/%m/%Y"))
```

```
## [1] "2020-01-01"
```

but where’s the humanity in that?

Testing whether something is a date is not as simple as the other data types unfortunately. Instead, we just use the `is()` function. If the first value returned is “Date”, then you know it’s a date:

```
is(as.Date("2020/01/01"))

## [1] "Date"      "oldClass"
```

2.4.7 Datetimes (POSIXct)

If you thought dates were annoying, datetimes are like dates’ little brother who keeps asking when his turn on the Xbox is. One of the reasons for this is that datetimes aren’t actually called datetimes. They’re called POSIXct in R. So whenever you see that dreadful word, just remember “ah, Adam told me that means datetime” and you’ll be fine.

Another thing that makes datetimes tough is that in addition to dates, datetimes (as you may have guessed) also store the time. The issue with that is that time is a more relative concept - there are lots of different timezones, so how do you know which one you’re referring to. By default, R has a locale for where you currently are and will use that location for your timezone. You override that default using the `Sys.setlocale()` function, or you can use the `tz` parameter when creating your datetime as we’ll see below.

With these annoyances aside however, creating datetimes isn’t all that different to creating dates except that we use the `as.POSIXct()` function instead. We just provide a character string (with a `format` specification if necessary), or a number with an origin. One important departure from dates though, is that now our origin is in seconds, not days, to allow us to calculate the time.

```
as.POSIXct("2020/01/01 12:00:00")

## [1] "2020-01-01 12:00:00 UTC"

as.POSIXct(1577880000, origin = "1970/01/01")

## [1] "2020-01-01 12:00:00 UTC"
```

Similar to dates, there is no `as.POSIXct()` function in base R, so we use the `is()` function instead:

```
is(as.POSIXct("2020/01/01 12:00:00"))

## [1] "POSIXct" "POSIXt"  "oldClass"
```

2.5 Data structures (and subsetting)

filler

2.6 Functions

Chapter 3

Data Analysis in R

In Module 2, we'll look more specifically at how one might do some simple data analysis in R. For a more in-depth view, I would highly recommend Hadley's R4DS

The xaringan presentation for this module can be found [here](#).

3.1 Installing packages

filler

3.2 Loading data

filler

3.3 Cleaning data

filler

3.4 Summary statistics

filler

3.5 Plots

Chapter 4

Programming in R

In this Module, we'll look at some of the programming concepts and syntax in R.

The xaringan presentation for this module can be found [here](#).

4.1 User-defined functions

filler

4.2 For loops

filler

4.3 If/else statements

filler

Chapter 5

Tidyverse

In this module, we'll take a quick look over some of the packages in the tidyverse.
The xaringan presentation for this module can be found [here](#).

5.1 tidyr

filler

5.2 dplyr

filler

5.3 stringr

filler

5.4 ggplot2

filler

Chapter 6

Modelling

6.1 Linear modelling

6.2 Clustering

6.3 Programming and modelling

Chapter 7

Exercises

Here are some interactive exercises covering the basic concepts we've looked at in the last 4 modules.

This exercises are hosted on a Shiny server [here](#)