

# Writing Reproducible Code

Literate, linear programming and modularity

Daniela Palleschi

Fri Aug 23, 2024

## Table of contents

<b>Learning objectives</b>	<b>1</b>
<b>Modular analyses</b>	<b>2</b>
<b>Reproducible code</b>	<b>2</b>
Writing linear code . . . . .	3
Literate programming . . . . .	3
Example R script . . . . .	3
<b>Dynamic reports</b>	<b>4</b>
Structure your reports . . . . .	5
Session Information . . . . .	5
Printing session info . . . . .	6
<b>Tips and tricks</b>	<b>6</b>
<b>Session Info</b>	<b>6</b>
<b>Hands-on: working with Quarto</b>	<b>7</b>

## Learning objectives

- learn about literate programming and modular analyses
- create and render a dynamic report with Quarto
- load data
- include a table and figure

## Modular analyses

- recall our `scripts` folder (which you might've named `analysis` or something else)
- ideally, this would also contain subfolders, one for each stage of your analysis
  - or at least, multiple scripts
- this is the concept of *modularity* (Bowers & Voors, 2016; Nagler, 1995)
  - separating data cleaning, pre-processing, recoding, merging, analyses, etc. into files/scripts

### Modularity

The concept of modularity relating to data analysis refers to saving different stages of data manipulation and analysis in separate scripts (Nagler, 1995). An important step is to also ensure the scripts are appropriately named so that the structure of your analysis/scripts folder “can be a map itself” (Bowers & Voors, 2016, p. 838).

A suggestion I have is to add numerical prefixes before scripts to explicitly state the order in which they must be run. For example, the script that reads in the raw data and cleans/preprocesses it (e.g., renaming or removing irrelevant/redundant variables) and exports a new data file that will be used in subsequent scripts could be called `00_cleaning.qmd` or something of the sort. My current preferred analysis script structure for a given data set is as follows:

```
scripts/  
  00_cleaning.qmd  
  01_preprocessing-and-eda.qmd  
  02_analyses.qmd  
  03_figures.qmd  
  04_tables.qmd
```

Where EDA refers to Exploratory Data Analysis, where I look at things like distributions of variables and demographic information. You'll note the `03_figures.qmd` and `04_tables.qmd` scripts follow the `02_analyses.qmd` script. This is because my figures and tables include model results or predictions, and so need to follow the analyses.

## Reproducible code

- how you write your code is the first step in making it reproducible
- the first principle is that your code must be *linear*

- this means code must be written in a linear fashion
- this is because we typically run a script from top-to-bottom

```
read_csv(here("data", "my_data.csv"))  
  
library(readr)  
library(here)
```

## Writing linear code

- you need to load a package *before* you call a function from it
  - if we're just working in an R session, *before* means temporally prior
  - with linear code, *before* means higher up in the script
- such pre-requisite code must
  - a. be present in the script
  - b. appear above the first line of code that uses a function from this package
- missing pre-requisite code might not throw an error message
  - but might produce output we aren't expecting
  - e.g., forgetting to filter out certain observations
  - or forgetting *that* some observations have been filtered out

## Literate programming

- introduced in Knuth (1984)
- refers to writing and documenting our code so that humans can understand it
  - important for us: we are (generally) not professional programmers, nor are our peers
- we need to not only know what our code is doing when we look back at it in the future/share it
- the easiest way: informative comments
  - the length and frequency of these comments is your choice

## Example R script

```

# Analysis script for phoneme paper
# author: Joe DiMaggio
# date: Feb. 29, 2024
# purpose: analyse cleaned dataset

# Set-up ###

# load required packages
library(dplyr)
library(readr)
library(ggplot2)
library(lme4)
library(broom.mixed) # tidy model summaries
library(ggeffects) # model predictions
library(here) # project-relative file path

# load-in data
df_phon <- read_csv(here("data", "phoneme_tidy_data.csv"))

# Explore data ###

```

- begins with some meta-information about the document, including its purpose
  - aids in knowing which scripts to run in which sequence
- there are three hashtags after some headings (###)
  - this is helpful because it structures the outline of the document in RStudio
- the purpose of chunks of code are written above
  - description of specific lines of code are also given

## Dynamic reports

- R scripts are useful, but don't show the code output
  - and commenting can get clunky
- dynamic reports combine prose, code, and code output
  - R markdown (.Rmd file extension) and Quarto (.qmd) are extensions of markdown
    - \* can embed R code 'chunks' in a script, thus producing 'dynamic' reports

- produce a variety of output files which contain text, R code chunks, and the code chunk outputs all in one

#### **i** Task: New Quarto document

1. Navigate to **File > New file > Quarto document**
2. Write some title, your name (Author), make sure ‘Visual markdown Editor’ is *unchecked*
3. Click ‘Create’
4. A new tab will open in R Studio. Press the ‘Render’ button above the top of the document, you will be prompted to save the document. Store it in a folder called **scripts** and save it as **01-literate-programming.qmd**.
5. What happens?

### Structure your reports

- describe the function/purpose at the beginning
- document your train of thought and findings throughout the script
  - e.g., why are you producing this plot, what does it tell you?
- give an overview of the findings/end result at the end
- it’s wise to avoid very long, multi-purpose scripts
  - rule of thumb: one script per product or purpose
  - e.g., data cleaning, exploration, analysis, publication figures, etc.

### Session Information

- R and R package versions are both open source, and are frequently updated
  - you might’ve run your code using **dplyr** version 1.1.0 or later, which introduced the **.by** per-operation grouping argument
  - what happens when somebody who has an older version of **dplyr** tries to run your code?
    - \* They won’t be able to!
  - the reverse of this situation is more common:
    - \* a newer version of a package no longer supports a deprecated function or argument

## Printing session info

- so, print your session info at the end of every script:

```
sessionInfo()
```

## Tips and tricks

- when you start a new script make sure you *always* start with a clean R environment:  
Session > Restart R or Cmd/Ctrl+Shift+O
  - this means no packages, data, functions, or any other dependencies are loaded
- at the top of your script, always load packages required below
  - you can always add more packages to the list as you add to your script
- **Render often:** when you make changes to your script make sure you re-render your document
  - checks you haven't introduced any errors
  - easier to troubleshoot if smaller changes have been made
- if you can run your script manually from source but it won't render, restart your R session and see if you can still run it from source
  - often the problem is some dependency in your environment that is not linearly introduced in the script

## Session Info

My session info.

```
sessionInfo()
```

R version 4.4.0 (2024-04-24)

Platform: aarch64-apple-darwin20

Running under: macOS Ventura 13.2.1

Matrix products: default

BLAS: /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRblas.0.dylib

LAPACK: /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRlapack.dylib;

```

locale:
[1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8

time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics  grDevices datasets  utils      methods    base

loaded via a namespace (and not attached):
[1] compiler_4.4.0    fastmap_1.2.0     cli_3.6.3         htmltools_0.5.8.1
[5] tools_4.4.0       rstudioapi_0.16.0 yaml_2.3.10       rmarkdown_2.28
[9] knitr_1.48        jsonlite_1.8.8    xfun_0.47         digest_0.6.36
[13] rlang_1.1.4       renv_1.0.7        evaluate_0.24.0

```

## Hands-on: working with Quarto

If you've never worked with Rmarkdown or Quarto before, try the following task. If you have, then try looking at a current or past analysis project you have, and check to see if it has the following:

- a designated folder containing all required files
- an `.RProj` file in this folder
- a transparent and organised folder structure
- a consistent and clear folder and file naming convention
- a README file in the project folder, and any other large folders
- code that runs if you try it

### Task: Editing your Quarto document

1. In the YAML (meta document info between `---`):

- Change the title
- Add `date: "08/21/2024"` (MM/DD/YYYY format)
- Change `format: html` to the following:

```

format:
  html:
    toc: true

```

Render the document. If it works, continue. If not, check that you aren't missing any

punctuation (e.g., quotation marks).

## 2. Adding a heading and text

- Remove all text below the YAML.
- Add a heading `# Packages`
- Navigate two lines below your heading, and write `Install the tidyverse and here packages. This will not run when we render the script.`

## 3. Adding a code chunk

- Navigate two lines below the text, and insert a new code chunk (either `Code > Insert Chunk` or the keyboard shortcut `Cmd+Option+I` on Mac or `Ctrl+Alt+I` on Windows)
- in the very first line of the code chunk, add `#| eval: false`
- below this line, write `# install packages`
- below this line, write `install.packages("tidyverse")` and hit `Cmd/Ctrl+Enter` and wait for the package to install (this may take some minutes)
- below this line, write `install.packages("here")` and hit `Cmd/Ctrl+Enter`

## 4. Loading packages

- Add a new line of text: `Print the project directory filepath using the ``here`` package.`
- Insert a new code chunk two lines below the first code chunk
- below this line, write `# load packages`
- below this line, write `library(tidyverse)`
- below this line, write `library(here)`
- hit `Cmd/Ctrl+Shift+Enter`, or click the little green triangle at the top right of the chunk

## 5. Printing our project directory path

- In a new code chunk, add the comment `# print project directory` and the code `here()`

## 6. Render the document

### Reveal Quarto syntax

```
---
title: "Literate programming"
author: "Daniela Palleschi"
date: "08/22/2024"
```



```

format:
  html:
    toc: true
---

# Packages

Install the tidyverse and here packages. This will not run when we render the script.

```{r}
#| eval: false
# install packages
install.packages("tidyverse")
install.packages("here")
```

Load required packages.

```{r}
# load packages
library(tidyverse)
library(here)
```

Print the project directory filepath using the `here` package.

```{r}
# print project directory
here()
```

```

- Bowers, J., & Voors, M. (2016). How to improve your relationship with your future self. *Revista de Ciencia Política (Santiago)*, 36(3), 829–848. <https://doi.org/10.4067/S0718-090X2016000300011>
- Knuth, D. (1984). Literate programming. *The Computer Journal*, 27(2), 97–111.
- Nagler, J. (1995). Coding Style and Good Computing Practices. *PS: Political Science & Politics*, 28(3), 488–492. <https://doi.org/10.2307/420315>