

A demonstration of combined scientific data processing and publication using Literate Programming in R

Andy Clifton

2019-10-18

Contents

1	Introduction	1
1.1	Linking analysis and publication workflows	1
1.2	How Literate Programming was used to write this document	2
1.3	Please not R	3
2	Implementing a coupled analysis and publication workflow	3
2.1	Setting up the computing environment	3
2.2	Load packages	4
2.3	Loading our own routines	4
2.4	Load the data	5
2.5	Plot input data	5
2.6	Operate on the data	5
2.7	Plot the results	5
2.8	Connect processing with publication	6
2.9	Save the processed data	7
2.10	Saving packages	7
2.11	Applying Journal formating	7
2.12	Thoughts on audit trails and confidentiality	7
2.13	Packaging for storage	8
3	Conclusions	8
	Referencing this document	8
	Acknowledgements	8

1 Introduction

Something something reproducible research, mumble, grumble, get off my lawn, grumble.

1.1 Linking analysis and publication workflows

Anecdotally, the separation of the publication from the analysis process has been a barrier to reproducible research, as it is impossible to ensure the link between source data and final publication.

This document demonstrates the application of Literate Programming to reproducible research. Literate programming means that the program documentation is complete and contained within

the program itself [Knuth, 1984]¹. It is important to note that the documentation is effectively a publication, and thus it is possible to combine data analysis with the creation of a publication in the same file. The use of literate programming therefore mitigates this barrier to reproducible research.

Furthermore, this project has been structured so that the data required for this publication are in a subdirectory of the project. This means that all of the files required to reproduce the analysis results can be included in a repository.

1.2 How Literate Programming was used to write this document

In this example, an output PDF document and results are generated from a file called *main.rmd*. *main.rmd* is an **R markdown file**. R markdown is a flavor of markdown that can be processed by the R programming language [R Core Team, 2017] to run code (i.e, do analysis) and create documentation from the same document. This is done using a package called *knitr*. Instructions for how to run *knitr* are included in the *howto.md* file in this repository. A far more detailed guide to writing using R markdown can be found in Xie et al. [2019].

The markdown document contains a mixture of documentation – written in markdown or LaTeX – and so-called “code chunks”, which here are written in R. The output is a PDF.

- The .rmd document is written in Pandoc markdown, which looks like normal text.
- The document contains code chunks that look like this:

```
```{r, echo=TRUE}
y = 40 + 2
print(y)
```
```

.. which evaluates to

```
y = 40 + 2
print(y)
```

```
## [1] 42
```

- The code chunks can be configured so that their outputs are echoed to the document (or not), which in turn allows the output PDF to show only those parts of the data processing that are relevant. You can thus completely hide the data operations in your output PDF and just display the results that are relevant.
- It is possible to use other programming languages by replacing the ‘{r,’ in the code chunk with the name of another language (see “Please not R”)
- There are a lot of different possible output formats, including PDF, HTML, Notebooks, other markdown formats, and many others. Corporate formatting can usually be applied without modifying their content. The details of this are out of scope for this paper; instead, see [Xie et al., 2019] or <https://bookdown.org/yihui/rmarkdown/documents.html> for more information.

I suggest reading this PDF together with the R markdown file (*main.rmd*) and possibly the *knitr* instructions². This will greatly help in understanding what is done in the processing and what makes it to the publication.

¹Yes, that’s the same ‘Knuth’ who invented LaTeX

²See <https://yihui.name/knitr/>

1.3 Please not R

If you can't handle learning yet another new language, this next statement might interest you:

“A less well-known fact about R Markdown is that many other languages are also supported, such as Python, Julia, C++, and SQL. The support comes from the knitr package, which has provided a large number of language engines.”

— [Xie et al. \[2019\]](#)

The currently available language engines are:

```
require(knitr)
names(knitr::knit_engines$get())
```

| | | | | | |
|---------|--------------|--------------|-----------|-------------|---------------|
| ## [1] | "awk" | "bash" | "coffee" | "gawk" | "groovy" |
| ## [6] | "haskell" | "lein" | "mysql" | "node" | "octave" |
| ## [11] | "perl" | "psql" | "Rscript" | "ruby" | "sas" |
| ## [16] | "scala" | "sed" | "sh" | "stata" | "zsh" |
| ## [21] | "highlight" | "Rcpp" | "tikz" | "dot" | "c" |
| ## [26] | "fortran" | "fortran95" | "asy" | "cat" | "asis" |
| ## [31] | "stan" | "block" | "block2" | "js" | "css" |
| ## [36] | "sql" | "go" | "python" | "julia" | "sass" |
| ## [41] | "scss" | "theorem" | "lemma" | "corollary" | "proposition" |
| ## [46] | "conjecture" | "definition" | "example" | "exercise" | "proof" |
| ## [51] | "remark" | "solution" | | | |

So, you have no excuse. You can write your code in any of those 52 languages, and off you go. A possible source of confusion might be that the rest of the document is R-flavoured markdown, but honestly there's not really that much R about it.

2 Implementing a coupled analysis and publication workflow

An analysis and publication workflow usually follows a similar path:

1. Set up the computing environment
2. Load some external packages
3. Load our own data processing routines
4. Import data
5. Plot it
6. Do some operations
7. Plot some more
8. Write
9. Format for a journal
10. Iterate around items 1-10 for a while
11. Submit

Fortunately, all of this can be captured in a markdown document.

2.1 Setting up the computing environment

Like most scripts, *main.rmd* includes a few variables that the user must set to run the analysis.

- The *project.root* variable defines the location of the files required for this analysis.
- The *made.by* variable forms part of a label that will be added to the plots.

An advantage of *knitr* is that we can simply execute the code and show the code and results inline:

```
# Where can files be found?
project.root <- file.path('/Users/andyc/Documents/public/GitHub/LiterateSciencePublishingDemo')
project.root

# Who ran this script
made.by = "A. Clifton"
made.by
```

```
## [1] "/Users/andyc/Documents/public/GitHub/LiterateSciencePublishingDemo"
## [1] "A. Clifton"
```

We can also show the value of those variables in the documentation using a relatively simple format, e.g. `r _project.root_`. This lets us then record in the documentation that

- *project.root* is /Users/andyc/Documents/public/GitHub/LiterateSciencePublishingDemo
- *made.by* is A. Clifton.

We've already set up several important subdirectories in *project.root*:

- **/code** contains functions required for the analysis
- **/data** contains the data files to be analyzed.

Let's tell the code where these are. We can also change our the working directory (*working.dir*) to the root directory of the project.

We'll also create a new directory for the results of the analysis.

Looking at your file system, you'll see there is now a new directory called **/Users/andyc/Documents/public/GitHub/LiterateSciencePublishingDemo**.

2.2 Load packages

2.3 Loading our own routines

Every data processing workflow requires its own scripts or functions to run. In this example, they are included in the *codes* directory and sourced during the preparation of this document. I have included output below to show these codes being called.

```
# source these functions
code.files = dir(code.dir, pattern = "\\R$")
for (file in code.files){
  source(file = file.path(code.dir,file))
  print(paste0("Sourcing ", file, "."))
}
```

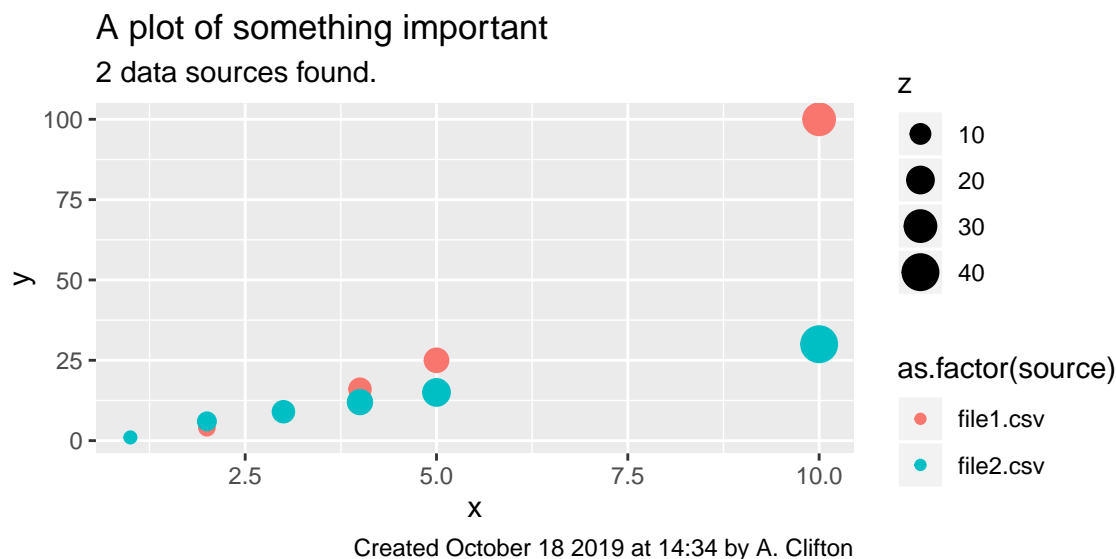
```
## [1] "Sourcing cleanPlot.R."
## [1] "Sourcing plotInfoLabel.R."
## [1] "Sourcing plotSomething.R."
## [1] "Sourcing theme_Literate.R."
```

2.4 Load the data

We now analyse the data from the simple data set. In this case, code has been written to load all of the files in the *data.dir* directory (`/Users/andyc/Documents/public/GitHub/LiterateSciencePublishingDemo/data`). I'm also going to map the three columns in the data files to the variables *x*, *y*, and *z*.³

2.5 Plot input data

The next step is to plot the input data. In this case we plot all of the input data together in one plot, but there are many different possibilities. Figures can also be given a consistent look and feel through ggplot's themes.



For convenience, we'll also save a copy of the figure as a *.png* file to the *analysis* directory.

2.6 Operate on the data

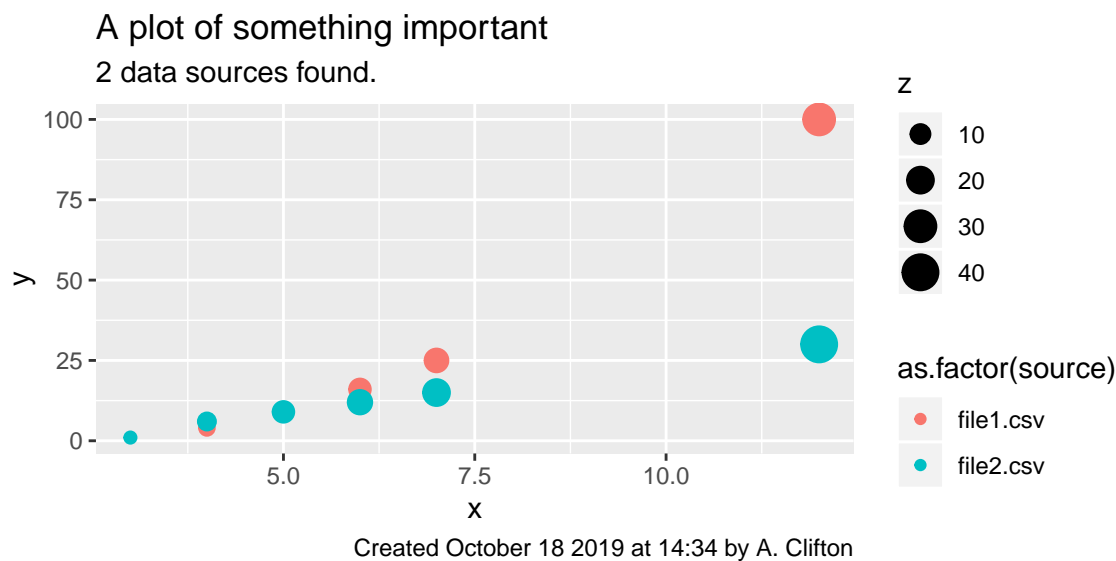
At this point we can do any number of operations on the data. For sake of demonstration, let's add 2 to all *x* values.

```
df.all <- df.in
df.all$x <- df.in$x + 2.0
```

2.7 Plot the results

Let's run that *plotSomething* routine again.

³See <https://www.calvin.edu/~rpruim/courses/s341/S17/from-class/MathinRmd.html> for more information about including maths in R markdown



And, as we can see, the data have shifted along x by a small amount.

2.8 Connect processing with publication

So far we have demonstrate that we can import and manipulate data and plot results. Another important part of a publication is the ability to generate statistics or summary information from data and include that in our text.

To demonstrate that, I can calculate that the maximum value of y in the input data sets was 100. This can be confirmed by checking the input data files. I could also include more complex logic in these statements, for example to say if one statistic is bigger or larger than another.

We sometimes need to include formatted tables in documents. This can be done using the *kable* function (Table 1).

```
knitr::kable(df.all,
  format = "pandoc",
  # format = "markdown", # this breaks cross references
  booktabs=TRUE,
  caption = "The $df.all$ data frame.")
```

Table 1: The *df.all* data frame.

| x | y | z | source |
|----|-----|----|-----------|
| 3 | 1 | 3 | file1.csv |
| 4 | 4 | 6 | file1.csv |
| 5 | 9 | 9 | file1.csv |
| 6 | 16 | 12 | file1.csv |
| 7 | 25 | 15 | file1.csv |
| 12 | 100 | 30 | file1.csv |
| 3 | 1 | 4 | file2.csv |
| 4 | 6 | 8 | file2.csv |
| 5 | 9 | 12 | file2.csv |

| x | y | z | source |
|----|----|----|-----------|
| 6 | 12 | 16 | file2.csv |
| 7 | 15 | 20 | file2.csv |
| 12 | 30 | 40 | file2.csv |

2.9 Save the processed data

We now write our processed data to file.

```
# save the data
save(list = c("project.root",
              "made.by",
              "df.all"),
      file = file.path(output.dir, "Data.RData"),
      envir = .GlobalEnv)
```

In R it is also possible to save the whole workspace. We can do that here as well:

```
# save the workspace
save.image(file=file.path(output.dir, "workspace.RData"))
```

2.10 Saving packages

Packages are required to supplement base functions in R and many other languages. For example, this script requires the *reticulate*, *bookdown*, *ggplot2*, *grid*, *knitr*, *RColorBrewer*, *rgdal*, and *stringr* packages to run. These are called from the script using the *require()* function. This assumes that the packages are available on your system.⁴ The use of packages represents a challenge to reproducible and repeatable research as it is possible that the function and output of the packages may change over time.

2.11 Applying Journal formatting

Scientific Journals often have their own formatting requirements. These requirements can still be met using markdown. The mechanics of such a process are beyond the scope of this paper and should probably be done as the last step in the publishing process. The reader is suggested to look at the *rticles* package and to use the detailed instructions in section 13 of the R Markdown Guide [Xie et al., 2019].

2.12 Thoughts on audit trails and confidentiality

Audit trails and confidentiality have implications for data (which might be commercially sensitive) and algorithms (which represent intellectual property).

A first step would be to avoid saving the workspace during the processing. Similarly, care should be taken to not commit any temporary files to a repository.

This document assumed that the data were available in text files that are stored in the same repository. However, it's equally possible that the raw data would have been called from a remote

⁴For details of how to install packages, see the RStudio help.

database at run time (or cached) and should be kept confidential. This can be done by not saving them to file. It would be desirable instead to store data IDs that would allow traceability. This could be combined with saving the results of the data processing, rather than the raw data.

Algorithms could be called directly from third party services or accessed via APIs. This shifts the onus to those third parties to provide the tracking required for auditing, but does preserve intellectual property.

All of these considerations can be dealt with in the data processing routines that are included in the source code of this document. There are doubtless many other issues and anyone concerned about this is encouraged to consult an expert.

2.13 Packaging for storage

A simple solution to repeatability and reproducibility may be to have a generic “data processing” image of a computer system, e.g. as a Docker image, that is started for each new project and used *only* for that project. This clean system is then used for one data processing task which is managed through a file such as the one you are reading. When the project is complete, the image is simply stored for as long as required. This would also avoid problems associated with changes introduced by package and system updates. A drawback to this approach would be the need to migrate the data every 5 years or so to a new system, which would be required to avoid data being stranded on old software.

3 Conclusions

Literate programming allows the creation of a single document that captures all of the process of preparing and analysing data, and creating a publication to describe that data. This is a fundamental requirement of reproducible research.

Referencing this document

This document has been assigned the Digital Object Identifier [10.5281/zenodo.3497450](https://doi.org/10.5281/zenodo.3497450). Citations in a range of formats can be obtained through Zenodo.

DOI [10.5281/zenodo.3497450](https://doi.org/10.5281/zenodo.3497450)

The source code for this document is available through github.com/AndyClifton/LiterateDemo.

Acknowledgements

Many thanks to Nikola Vasiljevic at DTU for prompting me to get this done.

Bibliography

Donald E. Knuth. Literate programming. *The Computer Journal*, 27(2):97–111, 1984.

R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2017. URL <https://www.R-project.org/>.

Yihui Xie, J. J. Allaire, and Garrett Golemund. *R Markdown: The Definitive Guide*. 2019. URL <https://bookdown.org/yihui/rmarkdown/>.