

Christopher Gandrud

Reproducible Research with R and RStudio

2

Getting Started with Reproducible Research

Researchers often start thinking about making their work reproducible near the end of the research process when they write up their results or maybe even later when a journal requires their data and code be made available for publication. Or maybe even later when another researcher asks if they can use the data from a published article to reproduce the findings. By then there may be numerous versions of the data set and records of the analyses stored across multiple folders on the researcher's computers. It can be difficult and time consuming to sift through these files to create an accurate account of how the results were reached. Waiting until near the end of the research process to start thinking about reproducibility can lead to incomplete documentation that does not give an accurate account of how findings were made. Focusing on reproducibility from the beginning of the process and continuing to follow a few simple guidelines throughout your research can help you avoid these problems. Remember “reproducibility is not an afterthought—it is something that must be built into the project from the beginning” (Donohue, 2010, 386).

This chapter first gives you a brief overview of the reproducible research process: a workflow for reproducible research. Then it covers some of the key guidelines that can help make your research more reproducible.

2.1 The Big Picture: A Workflow for Reproducible Research

The three basic stages of a typical computational empirical research project are:

- data gathering,
- data analysis,
- results presentation.

Each stage is part of the reproducible research workflow covered in this book. Tools for reproducibly gathering data are covered in Part II. Part III teaches tools for tying the data we gathered to our statistical analyses and presenting

the results with tables and figures. Part IV discusses how to tie these findings into a variety of documents you can use to advertise your findings.

Instead of starting to use the individual tools of reproducible research as soon as you learn them I recommend briefly stepping back and considering how the stages of reproducible research *tie* together overall. This will make your workflow more coherent from the beginning and save you a lot of backtracking later on. Figure 2.1 illustrates the workflow. Notice that most of the arrows connecting the workflow’s parts point in both directions, indicating that you should always be thinking how to make it easier to go backwards through your research, i.e. reproduce it, as well as forwards.

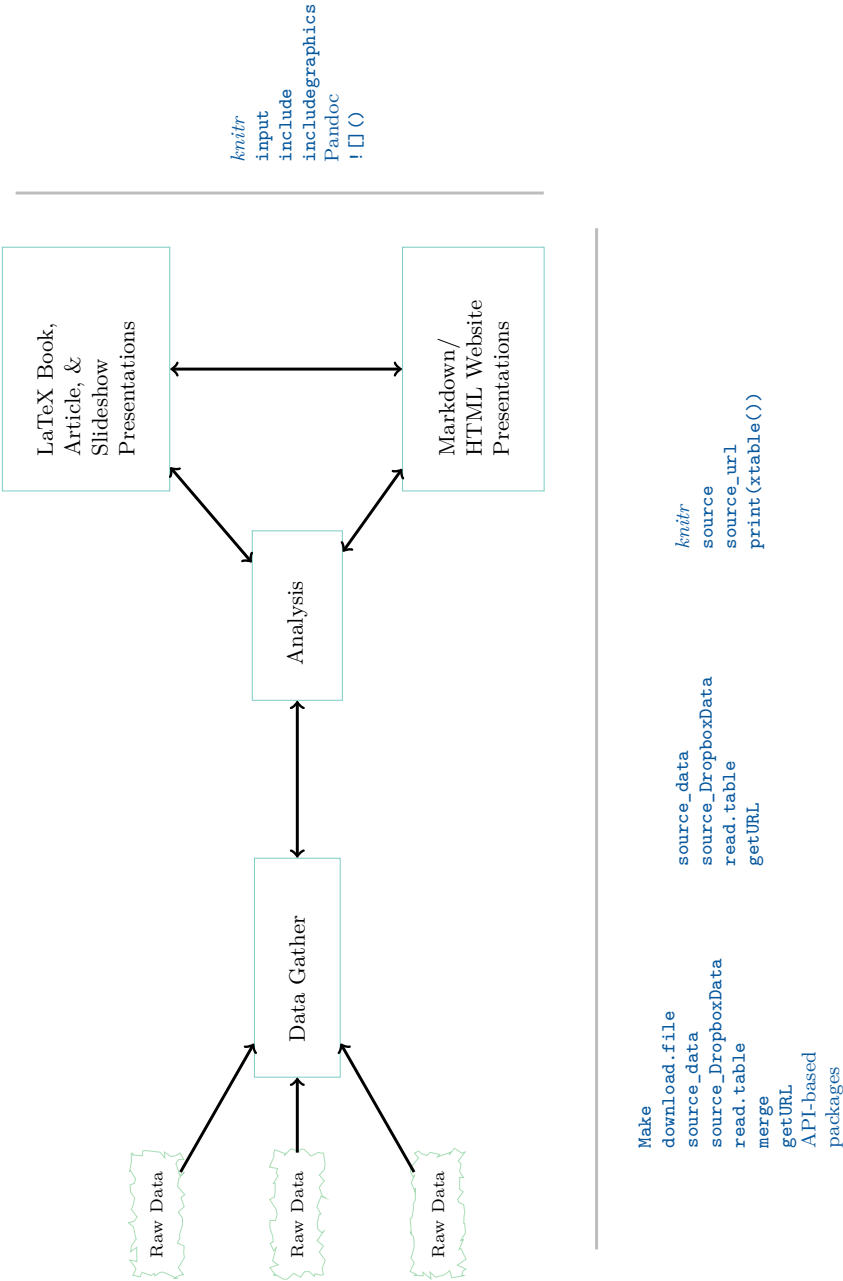
Around the edges of the figure are some of the commands you will learn to make it easier to go forwards and backwards through the process. These commands tie your research together. For example, you can use API-based R packages to gather data from the internet. You can use R’s `merge` command to combine data gathered from different sources into one data set. The `getURL` from R’s *RCurl* package (Temple Lang, 2013a) and the `read.table` commands can be used to bring this data set into your statistical analyses. The *knitr* package then ties your analyses into your presentation documents. This includes the code you used, the figures you created, and, with the help of tools such as the *xtable* package, tables of results. You can even tie multiple presentation documents together. For example, you can access the same figure for use in a LaTeX article and a Markdown created website with the `includegraphics` and `` commands, respectively. This helps you maintain a consistent presentation of results across multiple document types. We’ll cover these commands in detail throughout the book. See Table 2.1 for a brief but more complete overview of the main *tie commands*.

2.1.1 Reproducible theory

An important part of the research process that I do not discuss in this book is the theoretical stage. Ideally, if you are using a deductive research design, the bulk of this work will precede and guide the data gathering and analysis stages. Just because I don’t cover this stage of the research process doesn’t mean that theory building can’t and shouldn’t be reproducible. It can in fact be “the easiest part to make reproducible” (Vandewalle et al., 2007, 1254). Quotes and paraphrases from previous works in the literature obviously need to be fully cited so that others can verify that they accurately reflect the source material. For mathematically based theory, clear and complete descriptions of the proofs should be given.

Though I don’t actively cover theory replication in depth in this book, I do touch on some of the ways to incorporate proofs and citations into your presentation documents. These tools are covered in Part IV.

FIGURE 2.1
Example Workflow & A Selection of Commands to Tie it Together



2.2 Practical Tips for Reproducible Research

Before we start learning the details of the reproducible research workflow with R and RStudio, it's useful to cover a few broad tips that will help you organize your research process and put these skills in perspective. The tips are:

1. Document everything!,
2. Everything is a (text) file,
3. All files should be human readable,
4. Explicitly tie your files together,
5. Have a plan to organize, store, and make your files available.

Using these tips will help make your computational research really reproducible.

2.2.1 Document everything!

In order to reproduce your research others must be able to know what you did. You have to tell them what you did by documenting as much of your research process as possible. Ideally, you should tell your readers how you gathered your data, analyzed it, and presented the results. Documenting everything is the key to reproducible research and lies behind all of the other tips in this chapter and tools you will learn throughout the book.

Document your R session info

Before discussing the other tips it's important to learn a key part of documenting with R. You should *record your session info*. Many things in R have stayed the same since it was introduced in the early 1990s. This makes it easy for future researchers to recreate what was done in the past. However, things can change from one version of R to another and from one version of an R package to another. Also, the way R functions and especially how R packages are handled can vary across different operating systems, so it's important to note what system you used. Finally, you may have R set to load packages by default (see Section 3.1.8 for information about packages). These packages might be necessary to run your code, but other people might not know what packages and what versions of the packages were loaded from just looking at your source code. The `sessionInfo` command in R prints a record of all of these things. The information from the session I used to create this book is:

```
# Print R session info
sessionInfo()
```

```
## R version 3.0.1 (2013-05-16)
## Platform: x86_64-apple-darwin10.8.0 (64-bit)
##
## locale:
## [1] en_GB.UTF-8/en_GB.UTF-8/en_GB.UTF-8/C/en_GB.UTF-8/en_GB.UTF-8
##
## attached base packages:
## [1] tools      grid      splines    stats      graphics  grDevices  utils
## [8] datasets  methods   base
##
## other attached packages:
## [1] ZeligBayesian_0.1      MCMCpack_1.3-3      coda_0.16-1
## [4] Zelig_4.1-3            sandwich_2.2-10     boot_1.3-9
## [7] xtable_1.7-1          WDI_2.2             RAuth_0.9.3
## [10] treebase_0.0-6        ape_3.0-8           texreg_1.25
## [13] tables_0.7            sqldf_0.4-6.4       RSQLite.extfuns_0.0.1
## [16] RSQLite_0.11.3        chron_2.3-43        gsubfn_0.6-5
## [19] proto_0.3-10          DBI_0.2-7           slidify_0.3.3
## [22] stargazer_3.0.1       RJSONIO_1.0-3       RCurl_1.95-4.1
## [25] bitops_1.0-5          reshape2_1.2.2      repmis_0.2.5
## [28] quantmod_0.4-0        TTR_0.22-0          xts_0.9-3
## [31] zoo_1.7-9             Defaults_1.1-1      plyr_1.8
## [34] openair_0.8-5         memisc_0.96-4       MASS_7.3-26
## [37] lattice_0.20-15       markdown_0.5.4      knitr_0.4-4
## [40] bibtex_0.3-5          Hmisc_3.10-1.1     survival_2.37-4
## [43] httr_0.2              googleVis_0.4.2     ggplot2_0.9.3.1
## [46] gdata_2.12.0.2        formatR_0.7         extrafont_0.14
## [49] estout_1.2            digest_0.6.3        devtools_1.2
## [52] data.table_1.8.8      countrycode_0.13    brew_1.0-6
## [55] apsrtable_0.8-8       animation_2.2       knitr_1.2
##
## loaded via a namespace (and not attached):
## [1] car_2.0-17            cluster_1.14.4      colorspace_1.2-2
## [4] dichromat_2.0-0       evaluate_0.4.3      gtable_0.1.2
## [7] gtools_2.7.1          labeling_0.1         latticeExtra_0.6-24
## [10] Matrix_1.0-12         memoise_0.1         mgcv_1.7-22
## [13] munsell_0.4           nlme_3.1-109        parallel_3.0.1
## [16] RColorBrewer_1.0-5    rjson_0.2.12        Rttf2pt1_1.2
## [19] scales_0.2.3          stringr_0.6.2       tcltk_3.0.1
## [22] twitter_1.1.6         whisker_0.3-2       XML_3.95-0.2
## [25] yaml_2.1.7
```

Chapter 4 gives specific details about how to create files with dynamically included session information.

If you use non-R tools such as Pandoc, you should also record what versions you used.

2.2.2 Everything is a (text) file

Your documentation is stored in files that include data, analysis code, the write up of results, and explanations of these files (e.g. data set codebooks, session

info files, and so on). Ideally, you should use the simplest file format possible to store this information. Usually the simplest file format is the humble, but versatile, text file.¹

Text files are extremely nimble. They can hold your data in, for example, comma-separated values (`.csv`) format. They can contain your analysis code in `.R` files. And they can be the basis for your presentations as markup documents like `.tex` or `.md`, for LaTeX and Markdown files, respectively. All of these files can be opened by any program that can read text files.

One reason reproducible research is best stored in text files is that this helps *future proof* your research. Other file formats, like those used by Microsoft Word (`.docx`) or Excel (`.xlsx`), change regularly and may not be compatible with future versions of these programs. Text files, on the other hand, can be opened by a very wide range of currently existing programs and, more likely than not, future ones as well. Even if future researchers do not have R or a LaTeX distribution, they will still be able to open your text files and, aided by frequent comments (see below), be able to understand how we conducted your research (Bowers, 2011, 3).

Text files are also very easy to search and manipulate with a wide range of programs—such as R and RStudio—that can find and replace text characters as well as merge and separate files. Finally, text files are easy to version and changes can be tracked using programs such as Git (see Chapter 5).

2.2.3 All files should be human readable

Treat all of your research files as if someone who has not worked on the project will, in the future, try to understand them. Computer code is a way of communicating with the computer. It is ‘machine readable’ in that the computer is able to use it to understand what you want to do.² However, there is a very good chance that other people (or you six months in the future) will not understand what you were telling the computer. So, you need to make all of your files ‘human readable’. To make them human readable, you should comment on your code with the goal of communicating its design and purpose (Wilson et al., 2012). With this in mind it is a good idea to *comment frequently* (Bowers, 2011, 3) and *format your code using a style guide* (Nagler, 1995). For especially important pieces of code you should use *literate programming*—where the source code and the presentation text describing its design and purpose appear in the same document. Doing this will make it very clear to others how you accomplished a piece of research.

¹Plain text files are usually given the file extension `.txt`. Depending on the size of your data set it may not be feasible to store it as a text file. Nonetheless, text files can still be used for analysis code and presentation files.

²Of course, if the computer does not understand it will usually give an error message.

Commenting

In R everything on a line after a hash character—`#` (also known as number, pound, or sharp) is ignored by R, but is readable to people who open the file. The hash character is a comment declaration character. You can use the `#` to place comments telling other people what you are doing. Here are some examples:

```
# A complete comment line
2 + 2 # A comment after R code

## [1] 4
```

On the first line the `#` is placed at the very beginning, so the entire line is treated as a comment. On the second line the `#` is placed after the simple equation `2 + 2`. R runs the equation and finds the answer `4`, but it ignores all of the words after the hash.

Different languages have different comment declaration characters. In LaTeX everything after the `%` percent sign is treated as a comment, and in markdown/HTML comments are placed inside of `<!-- -->`. The hash character is used for comment declaration in command line shell scripts.

Nagler (1995, 491) gives some advice on when and how to use comments:

- write a comment before a block of code describing what the code does,
- comment on any line of code that is ambiguous.

In this book I follow these guidelines when displaying code.

He also suggests that all of your source code files should begin with a comment header. *At the least* the header should include:

- a description of what the file does,
- the date it was last updated,
- the name of the file's creator and any contributors.

You may also want to include other information in the header such as what files it depends on, what output files it produces, what version of the programming language you are using, or sources that may have influenced the code. Here is an example of a minimal file header for an R source code file that creates the third figure in an article titled 'My Article':


```
#####  
# R Source code file used to create Figure 3 in My 'Article'  
# Created by Christopher Gandrud  
# Updated 1 May 2013  
#####
```

Feel free to use things like the long series of hash marks above and below the header, white space, and indentations to make your comments more readable.

Style guides

In natural language writing you don't necessarily have to follow a style guide. People could probably figure out what you are trying to say, but it is a lot easier for your readers if you use consistent rules. The same is true when writing computer code. It's good to follow consistent rules for formatting your code so that it's easier for you and others to understand.

There are a number of R style guides. Most of them are similar to the Google R Style Guide.³ Hadley Wickham also has a nicely presented R style guide.⁴ You may want to use the *formatR* (Xie, 2012) package to automatically reformat your code so that it is easier to read.

Literate programming

For particularly important pieces of research code it may be useful to not only comment on the source file, but also display code in presentation text. For example, you may want to include key parts of the code you used for your main statistical models and an explanation of this code in an appendix following your article. This is commonly referred to as literate programming (Knuth, 1992).

2.2.4 Explicitly tie your files together

If everything is just a text file then research projects can be thought of as individual text files that have a relationship with one another. They are tied together. A data file is used as input for an analysis file. The results of an analysis are shown and discussed in a markup file that is used to create a PDF document. Researchers often do not explicitly document the relationships between files that they used in their research. For example, the results of an analysis—a table or figure—may be copied and pasted into a presentation document. It can be very difficult for future researchers to trace the table or figure back to a particular statistical model and a particular data set without

³See: <http://google-styleguide.googlecode.com/svn/trunk/google-r-style.html>.

⁴You can find it at <https://github.com/hadley/devtools/wiki/Style>.

clear documentation. Therefore, it is important to make the links between your files explicit.

Tie commands are the most dynamic way to explicitly link your files together. These commands instruct the computer program you are using to use information from another file. In Table 2.1 I have compiled a selection of key tie commands you will learn how to use in this book. We'll discuss many more, but these are some of the most important.

2.2.5 Have a plan to organize, store, & make your files available

Finally, in order for independent researchers to reproduce your work they need to be able access the files that instruct them how to do this. Files also need to be organized so that independent researchers can figure out how they fit together. So, from the beginning of your research process you should have a plan for organizing your files and a way to make them accessible.

One rule of thumb for organizing your research in files is to limit the amount of content any one file has. Files that contain many different operations can be very difficult to navigate, even if they have detailed comments. For example, it would be very difficult to find any particular operation in a file that contained the code used to gather the data, run all of the statistical models, and create the results figures and tables. If you have a hard time finding things in a file you created, think of the difficulties independent researchers will have!

Because we have so many ways to link files together there is really no need to lump many different operations into one file. So, we can make our files modular. One source code file should be used to complete one or just a few tasks. Breaking your operations into discrete parts will also make it easier for you and others to find errors (Nagler, 1995, 490).

Chapter 4 discusses file organization in much more detail. Chapter 5 teaches you a number of ways to make your files accessible through the cloud computing services Dropbox and GitHub.

TABLE 2.1

A Selection of Commands/Packages/Programs for Tying Together Your Research Files

Command/Package/ Program	Language	Description	Chapters Discussed
<i>knitr</i>	R	R package with commands for tying analysis code into presentation documents including those written in LaTeX and Markdown.	Throughout
<code>download.file</code>	R	Downloads a file from the internet.	6
<code>read.table</code>	R	Reads a table into R. You can use this to import a plain-text file formatted data into R.	6
<code>read.csv</code>	R	Same as <code>read.table</code> with default arguments set to import <code>.csv</code> formatted data files.	6
<code>source_data</code>	R	Reads a table stored on the internet into R. You can use it to import a plain-text formatted data file into R from secure (<code>https</code>) URLs.	6
<code>source_DropboxData</code>	R	Imports a plain-text data file stored in a Dropbox non-Public folder into R.	6
API-based packages	R	Various packages use APIs to gather data from the internet.	6
<code>merge</code>	R	Merges together data frames.	7
<code>source</code>	R	Runs an R source code file.	8
<code>source_url</code>	R	From the <i>devtools</i> package. Runs an R source code file from a secure (<code>https</code>) url like those used by GitHub.	8
<code>print(xtable())</code>	R	Combining the <code>print</code> & <code>xtable</code> commands creates LaTeX & HTML tables from R objects.	9
<code>toLaTeX</code>	R	Converts R objects to LaTeX.	2
<code>input</code>	LaTeX	Includes LaTeX files inside of other LaTeX files.	12
<code>include</code>	LaTeX	Similar to <code>input</code> , but puts page breaks on either side of the included text. Usually it is used for including chapters.	12
<code>includegraphics</code>	LaTeX	Inserts a figure into a LaTeX document.	10
<code></code>	Markdown	Inserts a figure into a Markdown document.	13
Pandoc	shell	A shell program for converting files from one markup language to another. Allows you to tie presentation documents together.	12 & 13
Make	shell	A shell program for automatically building many files at the same time.	6