

REPRODUCIBLE RESEARCH

Ryan Womack

Data Librarian, Rutgers University, rwomack@rutgers.edu

Spring 2017



This work is licensed under a [Creative Commons Attribution
-NonCommercial-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-nc-sa/4.0/).

REPRODUCIBILITY: WHAT DO WE MEAN?

Reproducible
Research

Ryan Womack

Introduction

Individuals
(Everyone)

Collaboration

Team Science

Conclusion

- ▶ Credibility in Science
 - ▶ Scientific fraud is [on the rise](#)
- ▶ Duke “[starter set](#)” and [article](#)
 - ▶ Research misconduct is a problem, but so is [human error](#)
- ▶ Reputation, Prestige, and Funding are all affected
 - ▶ “[Set the default to Open](#)”
- ▶ Replication (redoing the experiment from scratch) is expensive, and may not be possible due to the passage of time. (see [Validation](#))
- ▶ Science on [Replication and Reproducibility](#)
- ▶ Victoria Stodden, Friedrich Leisch, and Roger D. Peng (eds.). *Implementing Reproducible Research*. CRC Press, 2014.

“In all research that utilizes a computer, instructions for the research are stored in software and scientific data are stored digitally. A typical publication in computational research is based foundationally on data, and the computer instructions applied to the data that generated the scientific findings. The complexity of the data generation mechanism and the computational instruction is typically very large, too large to capture in a traditional scientific publication. Hence when computers are involved in the research process, scientific publication must shift from a scientific article to the triple of scientific paper, and the software and data from which the findings were generated. This triple has been referred to as a “research compendia” and its aim is to transmit research findings that others in the field will be able to reproduce by running the software on the data. Hence, data and software that permits others to reproduce the findings must be made available.”

– Victoria Stodden - <http://blog.stodden.net/2014/09/28/my-input-for-theostp-rfi-on-reproducibility/>

Introduction

Individuals
(Everyone)

Collaboration

Team Science

Conclusion

- ▶ **ICPSR** has been in operation for over 50 years, with well-established archiving practices and data documentation via codebooks and metadata
- ▶ **IPUMS** is reformatting and making data compatible across many decades and different projects, to enable international comparisons of microdata
- ▶ Coming from the world in the social sciences where long-term is, if not always routine, at least well-established
- ▶ Disciplinary separation of practices is diminishing when similar computational techniques can be applied to physical sciences or digital humanities

- ▶ We will illustrate some practices in a few contexts
 - ▶ an individual researcher
 - ▶ a team or research group
 - ▶ ongoing, large-scale collaboration

Some basic practices:

- ▶ Keep raw data pristine and separate from any working data
- ▶ Document your variables and data collection
 - ▶ anything you yourself would forget when revisiting the project 3 years later in response to a query
 - ▶ that will be the same thing other users need too!
- ▶ Don't work in Excel [if you can] or other manual editing environment
 - ▶ you should write down all your steps if you are doing this
 - ▶ better to use code or an environment that will at least record your steps

DOI, the Digital Object identifier, is the great success story

- ▶ makes it easy to have a permanent reference and good citation practice
- ▶ usually associated with quality data repositories
- ▶ encapsulates a lot of good stuff
- ▶ moral: defined standards and centralized tools make adoption and use easy
- ▶ Treat your local data as if you were pulling it from a DOI, and you will be baking in reproducibility

We will discuss examples in R, but other programming environments support this as well (Mathematica notebooks, Python/Jupyter)

- ▶ originally implemented in L^AT_EX + Sweave
- ▶ can embed R code and run it as the document is generated
- ▶ Code that is “tangled” in with text can be extracted, and formatted documents can be “woven” from the literate program.
- ▶ always ensures that the latest data and results are actually incorporated
- ▶ helps to document and explain code in context (literate programming)
- ▶ PDF, document, and HTML formats are easy to obtain

- ▶ Markdown + [knitr](#) has become a popular, lightweight replacement
- ▶ Simple syntax and implementation
- ▶ Integrated into RStudio
- ▶ Publish documents with one click at [RPods](#)
- ▶ Can fall back on \LaTeX /Sweave for more complex document formatting

```
— title: "A short literate programming exercise" author: "Ryan Womack"
date: "October 10, 2016" output: pdf_document —
“{r setup, include=FALSE} knitr::opts_chunk$set(echo = TRUE) “
## Read in the data
Let's read in the data with the following commands:
“{r load} library(readxl)
download.file("https://ryanwomack.com/data/PharmaDemo.xls",
"mydata.xls")
mydata<-read_excel("mydata.xls")
names(mydata)
attach(mydata)
“
## Describe the Data
Then we will get some summary statistics on the Age and Weight variables:
“{r summary} summary(Age)
summary(Weight) “
Now plot the data:
“{r plot, echo=FALSE} library(ggplot2)
ggplot(mydata, aes(Weight, Age))+ geom_point() “
## Regression
“{r regression} summary(lm(Age~Weight))
ggplot(mydata, aes(Weight, Age))+ geom_point()+ stat_smooth() “
All done!
```

- ▶ Open source is an important enabler of reproducibility
- ▶ Anyone can grab copies of the software to execute
- ▶ And can get older versions if necessary for compatibility
- ▶ You can also record information about your computing environment (`session.Info()` in R)
- ▶ The [checkpoint](#) package automates this process in R

- ▶ Don't save output. Where did it come from? This should be done in the code.
- ▶ Clean, well-formatted data (`tidyr`) and code (`formatR`) are a plus
- ▶ If using “readme” files approach, document everything

- ▶ **Jupyter** grew out of iPython
 - ▶ now over 40 languages supported
- ▶ Mathematica Notebooks, cloud support
- ▶ Cloud services making sharing much easier
- ▶ Becoming an expectation

- ▶ The same forces (cloud computing, shared platforms, standards) are making collaboration easier than ever
- ▶ [Github](#), [Bitbucket](#), and others enable easy collaboration on programming
 - ▶ with significant side benefits for reproducibility due to availability of code
- ▶ The [Open Science Framework](#) provides a more data-specific approach
- ▶ A key feature is that the same platform is used for private work and then public sharing
- ▶ [Psychology reproducibility study](#) uses OSF.
 - ▶ See the [Science article](#) for a start

- ▶ Collaborative projects must/should agree on:
 - ▶ data practices and sharing platform
 - ▶ coding practices and sharing platform
- ▶ This is made easier by already existing platforms and practices
 - ▶ [ropensci](#)
 - ▶ [ProjectTemplate](#)
 - ▶ Data and code can be distributed through packages, start with the `package.skeleton()` command
 - ▶ Reprozip and others at [Reproducible Science](#)

The [Yale Institute for Social and Policy Studies](#) is an example of a research group that enables reproducibility.

- ▶ Data and papers archived together onsite
- ▶ Handles (not DOIs) for data
- ▶ Code and documentation archived
- ▶ Code review for correct execution
- ▶ Good example of providing explanatory metadata for studies
- ▶ Possible because the Institute requires compliance as a condition of grants

Multi-year, multi-institutional projects that may continue beyond original PIs *require* reproducibility.

- ▶ Many people will be coming on and off of the project over time
- ▶ Many unanticipated uses are anticipated (“known unknowns” or something like that)
- ▶ Collaboration and continuity must be consciously planned for
- ▶ Decisions should be made with more consideration for future use than current convenience
- ▶ But disciplinary expertise is building in these areas
- ▶ PDB, of course
- ▶ Also, <https://www.teamsciencetoolkit.cancer.gov/>

- ▶ In big science, the main node(s) are enablers of future reuse
- ▶ They provide basic infrastructure ([OOI](#))
- ▶ But also provide a clearinghouse for other projects that link to and build on the central node
- ▶ One major future goal is to have more generic, all-purpose collaborative infrastructure
- ▶ Rutgers is developing a Virtual Data Collaboratory for this purpose
- ▶ Open infrastructure like [Zenodo](#), [Dataverse](#), [OpenICSPR](#) and the [Open Science Data Cloud](#) have been developed

- ▶ Standards such as DOI and ORCID enable the broader community to coalesce around good practice
- ▶ Data repositories are developing standards
- ▶ [Re3data.org](https://re3data.org/) is a directory of repositories
- ▶ The [Data Seal of Approval](#) is awarded to repositories using sound data practices
- ▶ [ISO 16363](#) (Trusted Digital Repositories) is a more stringent standard
- ▶ One important step is to plan for what happens when the project winds down (expectedly or unexpectedly)
- ▶ What is the equivalent standard for computing and reproducibility?

- ▶ Massive investments, massive amounts of data
- ▶ Many repositories too (NIH GDS)
- ▶ Existing repositories are useful and aggregate many software tools
- ▶ But researchers want even larger pooled databases, especially for human genome
- ▶ Technical issues are complex, but the rights and permissions involved are equally complex
- ▶ How can data be federated for maximum discoverability?

OPEN DATA -> REPRODUCIBLE RESEARCH

Reproducible
Research

Ryan Womack

Introduction

Individuals
(Everyone)

Collaboration

Team Science

Conclusion

- ▶ Increasing openness is a long-term trend
 - ▶ Internet, Government, Data, Software, Cloud Computing
- ▶ Pressure for Reproducible Research can only increase as these trends intensify
- ▶ Good news is...
 - ▶ Benefits are clear for society and for knowledge creation!
 - ▶ Tools to enable this are getting easier all the time!
 - ▶ Eventually it will be a standard we will take for granted!