**Doing reproducible science:
from your hard-won data
to a publishable manuscript
without going mad**

Francisco Rodriguez-Sanchez (@frod_san)
February 2017

# A typical research workflow
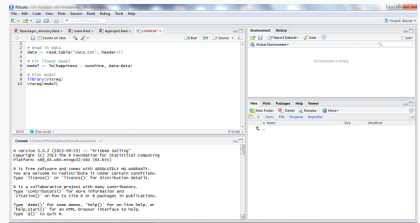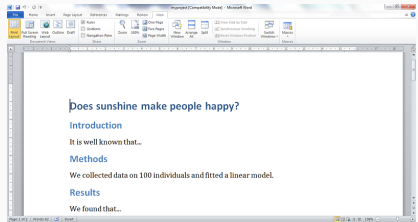
1. Prepare data (**EXCEL**)

## A typical research workflow

1. Prepare data (**EXCEL**)
2. Analyse data (**R**)

## A typical research workflow

1. Prepare data (**EXCEL**)
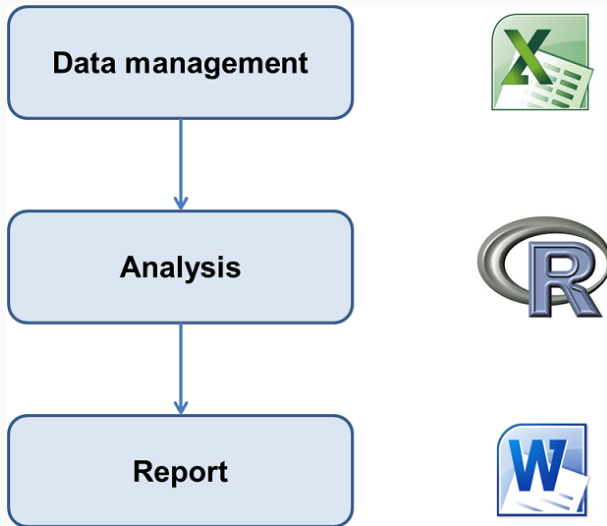2. Analyse data (**R**)
3. Write report/paper (**WORD**)

## A typical research workflow

1. Prepare data (**EXCEL**)
2. Analyse data (**R**)
3. Write report/paper (**WORD**)
4. Start the email attachments **nightmare**. . .

## Problems of a broken workflow

- How did you do this? What analysis is behind this figure? Did you account for . . . ?

## Problems of a broken workflow

- How did you do this? What analysis is behind this figure? Did you account for . . . ?

- What dataset was used? Which individuals were left out? Where is the clean dataset?

**Problems of a broken workflow**

- How did you do this? What analysis is behind this figure? Did you account for . . . ?

- What dataset was used? Which individuals were left out? Where is the clean dataset?

- Oops, there is an error in the data. Can you repeat the analysis? And update figures/tables in Word!

**Trevor A. Branch**
@TrevorABranch

My rule of thumb: every analysis you do on a dataset will have to be redone 10–15 times before publication. Plan accordingly. #Rstats
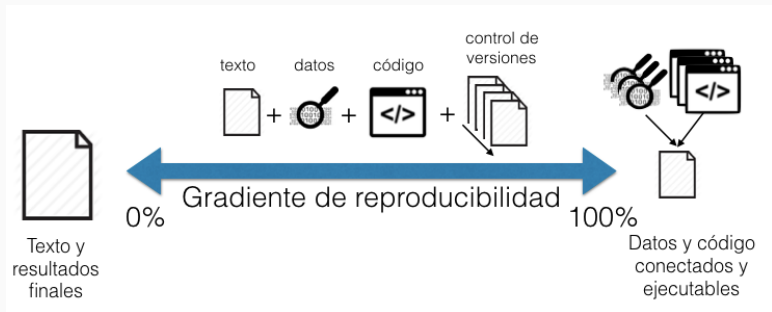
# Our everyday scary movie

https://youtu.be/s3JldKoAOzw

# WHAT is Reproducible Science?

## Reproducible Science: WHAT

A scientific article is **reproducible** if there is computer **code** that can **regenerate** all results and figures from the original data.

- Transparent
- Traceable
- Comprehensive
- Useful

# Most science is not reproducible



Even **you** will struggle to reproduce **your own results** from a few weeks/months ago.
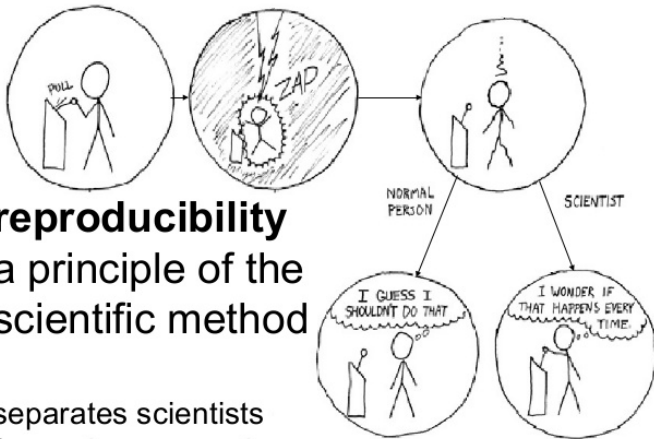
*You can't reproduce if you don't understand where a number came from.*

*You can't reproduce what you don't remember. And trust me: you won't.*

*You can't reproduce what you've lost. What if you need access to a file as it existed 1, 10, 100, or 1000 days ago?*

*Ben Bond-Lamberty*

# WHY Reproducible Science?

Noam Ross
@noamross

Follow

Gelman: "Reproducible research is even better when you're wrong" #stancon2017

16

- Fundamental pillar of **scientific method**

## Reproducible Science: WHY

- Fundamental pillar of **scientific method**

- Much less prone to **errors**

## Reproducible Science: WHY

- Fundamental pillar of **scientific method**

- Much less prone to **errors**

- Regenerate results **automatically**

## Reproducible Science: WHY

- Fundamental pillar of **scientific method**

- Much less prone to **errors**

- Regenerate results **automatically**

- **Code reuse** & sharing accelerates scientific progress

## Reproducible Science: WHY

- Fundamental pillar of **scientific method**

- Much less prone to **errors**

- Regenerate results **automatically**

- **Code reuse** & sharing accelerates scientific progress

- Increasingly required by **journals**

## Reproducible Science: WHY

- Fundamental pillar of **scientific method**

- Much less prone to **errors**

- Regenerate results **automatically**

- **Code reuse** & sharing accelerates scientific progress

- Increasingly required by **journals**

- Higher publication **impact** (citations, future collaborations, etc)

# HOW TO DO Reproducible Science?

## Reproducible Science: HOW

1. File **organisation**.
2. **Data management**. Spreadsheet good practices.
3. **Code-based** data analysis. **Rmarkdown**
4. Software **dependencies**.
5. **Version control** & collaborative writing.

## File organisation

- All files in **same directory** (Rstudio project).

- All files in **same directory** (Rstudio project).

- **Raw data untouched** in independent folder.

## File organisation

- All files in **same directory** (Rstudio project).

- **Raw data untouched** in independent folder.

- Derived, **clean data** in another folder.

## File organisation

- All files in **same directory** (Rstudio project).

- **Raw data untouched** in independent folder.

- Derived, **clean data** in another folder.

- Figures, code, etc also have their own folder.

## File organisation example

```
myproject

|- README      # general info about the project

|- analysis.R  # master script that executes everything

|- data-raw/   # original raw data

|- data/       # clean data (produced w/ script)

|- R/          # functions definitions

|- doc/        # manuscript files

|- figs/       # final figures
```

# Data management

## Editorial expression of concern

IN THE 3 June issue, *Science* published the Report "Environmentally relevant concentrations of microplastic particles influence larval fish ecology" by Oona M. Lönnstedt and Peter Eklöv (*1*). The authors have notified *Science* of the theft of the computer on which the raw data for the paper were stored. These data were not backed up on any other device nor deposited in an appropriate repository. *Science* is publishing this Editorial Expression of Concern to alert our readers to the fact that no further data can be made available, beyond those already presented in the paper and its supplement, to enable readers to understand, assess, reproduce, or extend the conclusions of the paper.

**Jeremy Berg**
Editor in Chief

http://science.sciencemag.org/content/354/6317/1242.1

23

## Storage

Use the **cloud**: safe, persistent, easy to share

- Dropbox
- OSF
- Figshare, etc
- See all data repositories in www.re3data.org

# Tidy data

## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

- Each **observation** in one **row** (e.g. individuals).

## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

- Each **observation** in one **row** (e.g. individuals).

- Avoid spaces, numbers, and special characters in column names.

## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

- Each **observation** in one **row** (e.g. individuals).

- Avoid spaces, numbers, and special characters in column names.

- Always write zero values, to distinguish from blank/missing data.

## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

- Each **observation** in one **row** (e.g. individuals).

- Avoid spaces, numbers, and special characters in column names.

- Always write zero values, to distinguish from blank/missing data.

- Use blank/empty cells, or NA, for missing data.
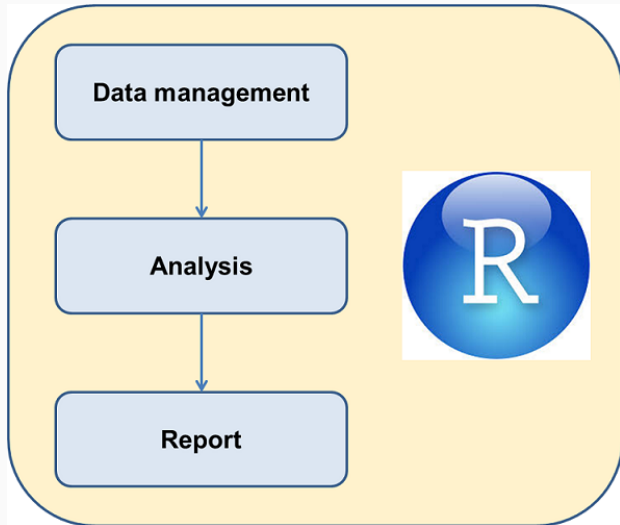
## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

- Each **observation** in one **row** (e.g. individuals).

- Avoid spaces, numbers, and special characters in column names.

- Always write zero values, to distinguish from blank/missing data.

- Use blank/empty cells, or NA, for missing data.

- Input dates as year, month, day in separate columns. Or YYYY-MM-DD as text.

## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

- Each **observation** in one **row** (e.g. individuals).

- Avoid spaces, numbers, and special characters in column names.

- Always write zero values, to distinguish from blank/missing data.

- Use blank/empty cells, or NA, for missing data.

- Input dates as year, month, day in separate columns. Or YYYY-MM-DD as text.

- Use 'Data validation' in Excel to constrain data entry to accepted values.

## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

- Each **observation** in one **row** (e.g. individuals).

- Avoid spaces, numbers, and special characters in column names.

- Always write zero values, to distinguish from blank/missing data.

- Use blank/empty cells, or NA, for missing data.

- Input dates as year, month, day in separate columns. Or YYYY-MM-DD as text.

- Use 'Data validation' in Excel to constrain data entry to accepted values.

- Don't combine multiple pieces of information in one cell.

## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

- Each **observation** in one **row** (e.g. individuals).

- Avoid spaces, numbers, and special characters in column names.

- Always write zero values, to distinguish from blank/missing data.

- Use blank/empty cells, or NA, for missing data.

- Input dates as year, month, day in separate columns. Or YYYY-MM-DD as text.

- Use 'Data validation' in Excel to constrain data entry to accepted values.

- Don't combine multiple pieces of information in one cell.

- Don't touch raw data. Do all data manipulation with R code.

## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

- Each **observation** in one **row** (e.g. individuals).

- Avoid spaces, numbers, and special characters in column names.

- Always write zero values, to distinguish from blank/missing data.

- Use blank/empty cells, or NA, for missing data.

- Input dates as year, month, day in separate columns. Or YYYY-MM-DD as text.

- Use 'Data validation' in Excel to constrain data entry to accepted values.

- Don't combine multiple pieces of information in one cell.

- Don't touch raw data. Do all data manipulation with R code.

- Export data as plain text (txt, csv)

## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

- Each **observation** in one **row** (e.g. individuals).

- Avoid spaces, numbers, and special characters in column names.

- Always write zero values, to distinguish from blank/missing data.

- Use blank/empty cells, or NA, for missing data.

- Input dates as year, month, day in separate columns. Or YYYY-MM-DD as text.

- Use 'Data validation' in Excel to constrain data entry to accepted values.

- Don't combine multiple pieces of information in one cell.

- Don't touch raw data. Do all data manipulation with R code.

- Export data as plain text (txt, csv)

- http://www.datacarpentry.org/spreadsheet-ecology-lesson/

## Spreadsheet good practices

- Put **variables** in **columns** (things you are measuring: height, weight, sex)

- Each **observation** in one **row** (e.g. individuals).

- Avoid spaces, numbers, and special characters in column names.

- Always write zero values, to distinguish from blank/missing data.

- Use blank/empty cells, or NA, for missing data.

- Input dates as year, month, day in separate columns. Or YYYY-MM-DD as text.

- Use 'Data validation' in Excel to constrain data entry to accepted values.

- Don't combine multiple pieces of information in one cell.

- Don't touch raw data. Do all data manipulation with R code.

- Export data as plain text (txt, csv)

- http://www.datacarpentry.org/spreadsheet-ecology-lesson/

- http://kbroman.org/dataorg/

# Data analysis

- Reproducible
- Reusable

## Rmarkdown documents

- Fully reproducible (trace all results inc. tables and plots)
- Dynamic (regenerate with 1 click)
- Suitable for
    - documents (Word, PDF, etc)
    - presentations
    - books
    - websites
    - . . .

**Let's see Rmarkdown in action**

In Rstudio, create new Rmarkdown document and click on `Knit HTML`.

# Example: Does sunshine influence happiness?

See `myproject.Rmd` (http://bit.ly/rmdsun)

### Does sunshine make people happy?

*F. Rodriguez-Sanchez*
*Tuesday, November 25, 2014*

**Introduction**

It is well known that individual well-being can be influenced by climatic conditions. However, ...

**Methods**

We collected data on 100 individuals and fitted a linear model.

**Results**

We found that ...

|  | Estimate | Std. Error | t value | Pr(>|t|) |
|---|---|---|---|---|
| (Intercept) | -0.0651657 | 0.4264970 | -0.1527928 | 0.8788758 |
| sunshine | 0.0100228 | 0.0004232 | 23.6833264 | 0.0000000 |



**Discussion**

These results confirm that sunshine is good for happiness (slope = 0.0100228).

**Acknowledgements**

Y. Xie, J. MacFarlane, Rstudio...

**Spotted error in the data? No problem!**

Make changes in Rmarkdown document, click `knit` and report will
**update automatically!**

# Adding citations by DOI

rcrossref addin

## Adding citations from BibTeX file

citr addin

https://github.com/crsh/citr/

- rticles
- rmdTemplates

# Can write full thesis in Rmarkdown!

See thesis.Rmd.

See thesis.pdf.

http://rmarkdown.rstudio.com/index.html

https://www.rstudio.com/wp-content/uploads/2016/03/

# Managing software dependencies

## Managing package dependencies in R

- **sessionInfo** (or session_info)
- switchr
- rctrack
- **checkpoint**
- **packrat**
- docker

# Version control

## Dropbox

Dropbox keeps record of deleted/edited files for 30 days

Automatic version control, no time limit.

R. Fitzjohn
(https://github.com/richfitz/reproducibility-2014)

Ecosistemas 25(2): 83-92 [Mayo-Agosto 2016]
Doi.: 10.7818/ECOS.2016.25-2.11

Artículo publicado en Open Access bajo los términos
de Creative Commons attribution Non Comercial License 3.0.

REVISIONES

**ecosistemas**
REVISTA CIENTÍFICA DE ECOLOGÍA Y MEDIO AMBIENTE

ISSN 1697-2473 / Open access
disponible en www.revistaecosistemas.net

## Ciencia reproducible: qué, por qué, cómo

F. Rodríguez-Sánchez[1,*], A.J. Pérez-Luque[2,**], I. Bartomeus[1,**], S. Varela[3,**]

http://www.revistaecosistemas.net/index.php/ecosistemas/article/
viewFile/1178/973

# Happy writing!