

# Reproducible Workflows for Better Science and Efficient Collaboration

---

Francisco Rodriguez-Sanchez

@frod\_san

<https://frodriguezsanchez.net>

# The Reproducibility Crisis Revolution

---

NEWS | 09 December 2021

## **Half of top cancer studies fail high-profile reproducibility effort**

- Goal: Replicate 193 experiments from 53 papers

NEWS | 09 December 2021

## **Half of top cancer studies fail high-profile reproducibility effort**

- Goal: Replicate 193 experiments from 53 papers
- Finally: 50 experiments from 23 papers

NEWS | 09 December 2021

## Half of top cancer studies fail high-profile reproducibility effort

- Goal: Replicate 193 experiments from 53 papers
- Finally: 50 experiments from 23 papers
- ~Half **didn't replicate** (much smaller effect sizes)

NEWS | 09 December 2021

## Half of top cancer studies fail high-profile reproducibility effort

- Goal: Replicate 193 experiments from 53 papers
- Finally: 50 experiments from 23 papers
- ~Half **didn't replicate** (much smaller effect sizes)
- **No paper reported** all required data

NEWS | 09 December 2021

## Half of top cancer studies fail high-profile reproducibility effort

- Goal: Replicate 193 experiments from 53 papers
- Finally: 50 experiments from 23 papers
- ~Half **didn't replicate** (much smaller effect sizes)
- **No paper reported** all required data
- **Impossible to repeat** experiments w/o contacting authors

NEWS | 09 December 2021

## Half of top cancer studies fail high-profile reproducibility effort

- Goal: Replicate 193 experiments from 53 papers
- Finally: 50 experiments from 23 papers
- ~Half **didn't replicate** (much smaller effect sizes)
- **No paper reported** all required data
- **Impossible to repeat** experiments w/o contacting authors
- 1/3 authors **didn't respond or help**





Sylvain Deville ❄️ 🧑

@DevilleSy

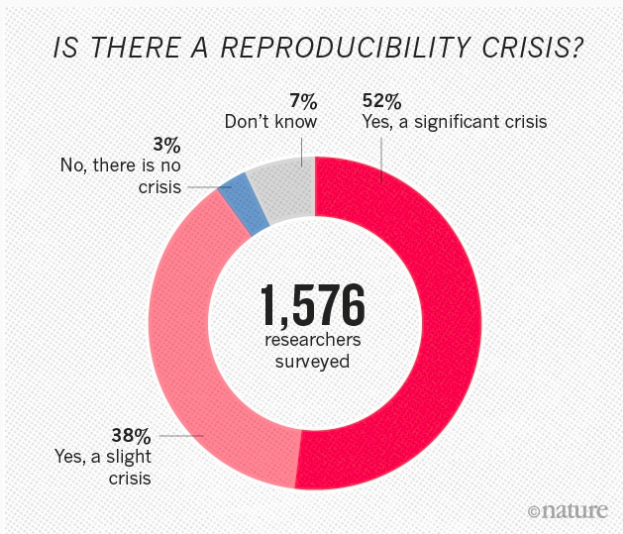


Trying to reproduce the results of a paper using only what's in the Methods section



Most scientific articles

**are NOT reproducible**



Reproducibility

~~CRISIS~~

**REVOLUTION**

# Reproducibility vs Replicability

		Data	
		Same	Different
Analysis	Same	Reproducible	Replicable
	Different	Robust	Generalisable

The Turing Way

We can't guarantee that  
our results are **replicable**.

But at least  
they should be **reproducible**.

Most scientific articles

**are NOT reproducible**

## The prevalence of statistical reporting errors in psychology (1985–2013)

Michèle B. Nuijten<sup>1</sup> · Chris H. J. Hartgerink<sup>1</sup> · Marcel A. L. M. van Assen<sup>1</sup> · Sacha Epskamp<sup>2</sup> · Jelte M. Wicherts<sup>1</sup>

### WHAT STATCHECK LOOKS FOR

This computer algorithm scans papers for statistical tests, uses reported results to recompute the  $P$  value and flags up inconsistencies.

#### Type of test

The  $t$ -test assesses differences between two groups.

#### Test statistic

Compares observed values with those expected under the null hypothesis.

$$t(37) = 4.93, P < 0.01$$

#### Degrees of freedom

Accounts for size of sample.

#### $P$ value

The likelihood of observing differences as extreme, or more so, if the null hypothesis is true.



## **The prevalence of statistical reporting errors in psychology (1985–2013)**

Michèle B. Nuijten<sup>1</sup> · Chris H. J. Hartgerink<sup>1</sup> · Marcel A. L. M. van Assen<sup>1</sup> ·  
Sacha Epskamp<sup>2</sup> · Jelte M. Wicherts<sup>1</sup>

1/2 articles: **inconsistencies** in p-values

1/8 articles: **grossly inconsistent** p-values

(affecting conclusions -> significance)
































In ecology

< 20% articles are reproducible

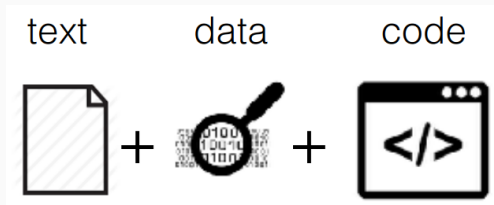
Culina et al 2020

# We can't even reproduce our own work

## Data/Code lost or unusable

 qualitative_data.csv	04/07/2016 15:50
 cleandata.xlsx	25/06/2015 01:14
 cleandata_YC.xlsx	30/06/2015 16:22
 COORDENADAS PACO_20-05-2016 CON REVIEWS.xlsx	20/05/2016 16:23
 COORDENADAS PACO_20-05-2016 CON REVIEWS_FRS.xlsx	27/05/2016 19:41
 COORDENADAS_paper195(Girella_elevata).xlsx	08/06/2016 13:09
 coordenadas_raw_2016-06-08.xlsx	09/06/2016 15:53
 coordenadas_raw_2016-06-08_eld.xlsx	08/06/2016 16:00
 coordenadas_raw_2016-06-21.xlsx	21/06/2016 16:12
 coords_2015-09-09_modif.xlsx	05/11/2015 15:23
 coords_2015-10-11_modif_YC.xlsx	17/11/2015 13:37
 coords_2015-10-11_modif_YC_PACO.xlsx	17/11/2015 17:06
 coords_2015-10-18_modif_YC.xlsx	18/11/2015 17:24
 coords_2015-12-26_modif_YC.xlsx	30/03/2016 19:38
 coords_2016-04-02.xlsx	06/04/2016 17:46
 coords_2016-04-02_YC.xlsx	06/04/2016 18:03
 coords_2016-04-08_YC.xlsx	11/04/2016 13:51
 dataset_y_coords_09_09_15.xlsx	23/09/2015 17:38
 Datos metaanálisis_18-04-2016.xlsx	19/04/2016 16:24
 FINAL METAANALISIS_14-6-2016_WITH REVIEWS.xlsx	21/06/2016 16:15
 FINAL METAANALISIS_16-6-2016_WITH REVIEWS.xlsx	21/06/2016 16:13
 FINAL METAANALISIS_2016-04-27_WITH REVIEWS.xlsx	25/05/2016 18:05
 FINAL METAANALISIS_2016-04-27_WITH REVIEWS_FRS.xlsx	27/05/2016 18:44
 FINAL METAANALISIS_2016-04-29_EXCLUDING REVIEWS.xlsx	08/06/2016 13:06
 FINAL VOTECOUNTING_1-7-2016.xlsx	04/07/2016 15:46
 fitnessdata_2016-06-22.xlsx	22/06/2016 21:00
 IFs for Bastien_19-3-2016_YC.xlsx	28/03/2016 19:26
 Metaanalysis final_01-05-2015 with coordinates.xlsx	18/05/2015 19:20
 Metaanalysis final_22-05-2015 coords.xlsx	24/06/2015 15:50
 Metaanalysis final_25-06-2015.xlsx	30/06/2015 16:55
 Metaanalysis y coords revisadas_06-08-2015_AH_JE.xlsx	23/09/2015 12:57

## What's a reproducible article?



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

A scientific article is **advertising**, not scholarship.

The actual scholarship is the **full software environment, code and data**, that produced the result.

Claerbout & Karrenback 1992

Are we sharing the data?

PERSPECTIVE

## Public Data Archiving in Ecology and Evolution: How Well Are We Doing?

Dominique G. Roche<sup>1,2\*</sup>, Loeske E. B. Kruuk<sup>1,3</sup>, Robert Lanfear<sup>1,4</sup>, Sandra A. Binning<sup>1,2</sup>

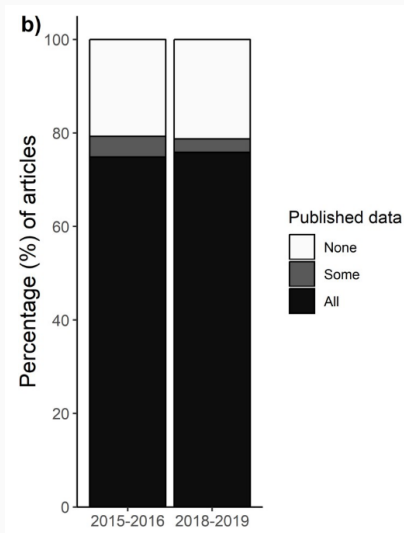
**1** Division of Evolution, Ecology and Genetics, Research School of Biology, The Australian National University, Canberra, Australian Capital Territory, Australia, **2** Éco-Éthologie, Institut de Biologie, Université de Neuchâtel, Neuchâtel, Switzerland, **3** Institute of Evolutionary Biology, School of Biological Sciences, University of Edinburgh, Edinburgh, United Kingdom, **4** Department of Biological Sciences, Macquarie University, Sydney, Australia

\* [dominique.roche@mail.mcgill.ca](mailto:dominique.roche@mail.mcgill.ca)

### Abstract

Policies that mandate public data archiving (PDA) successfully increase accessibility to data underlying scientific publications. However, is the data quality sufficient to allow reuse and reanalysis? We surveyed 100 datasets associated with nonmolecular studies in journals that commonly publish ecological and evolutionary research and have a strong PDA policy. Out of these datasets, 56% were incomplete, and 64% were archived in a way that partially or entirely prevented reuse. We suggest that cultural shifts facilitating clearer benefits to authors are necessary to achieve high-quality PDA and highlight key guidelines to help authors increase their data's reuse potential and compliance with journal data policies.

# Are we sharing data?





Quickly getting better

## **Scientific Life**

Early Career  
Researchers Embrace  
Data Sharing

Hamish A. Campbell,<sup>1,\*</sup>  
Mariana A. Micheli-Campbell,<sup>1</sup>  
and Vinay Udyawer<sup>2</sup>

Are we sharing the code?

# Code exists but rarely shared

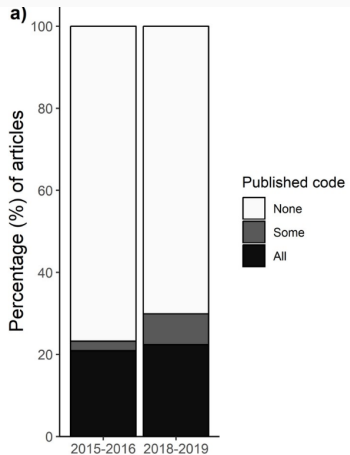
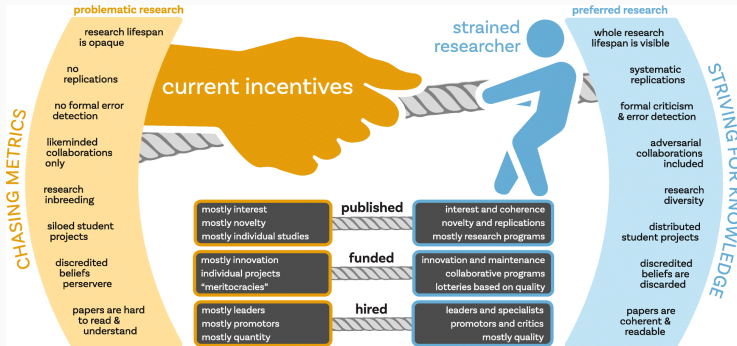


Fig 1. Code-sharing is at its infancy in ecology, when

WHY?

# Poor incentives



O'Dea et al 2021

# The Costs of Reproducibility

Russell A. Poldrack<sup>1,\*</sup>

<sup>1</sup>Department of Psychology, Stanford University, Stanford, CA, USA

\*Correspondence: [poldrack@stanford.edu](mailto:poldrack@stanford.edu)

<https://doi.org/10.1016/j.neuron.2018.11.030>

PERSPECTIVE

Open science challenges, benefits and tips in  
early career and beyond

Christopher Allen<sup>1</sup>, David M. A. Mehler<sup>1,2</sup>

## Credit data generators for data reuse

To promote effective sharing, we must create an enduring link between the people who generate data and its future uses, urge **Heather H. Pierce** and colleagues.

[Pierce et al 2019](#)

# Publish your computer code: it is good enough

*Freely provided working code — whatever its quality — improves programming and enables others to engage with your research, says **Nick Barnes**.*

Barnes 2010

- Improve training
- Avoid shaming -> constructive critique
- Ugly code better than no code



## Why doing reproducible research?

---

Reproducibility: good for you,  
good for everyone

---

## Automation (good code) saves time

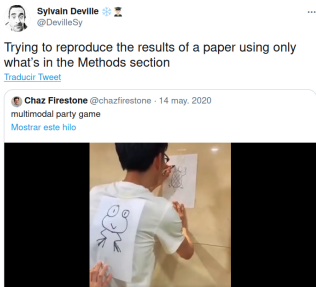


**Trevor Branch**  
@TrevorABranch



My rule of thumb: every analysis you do on a dataset will have to be redone 10–15 times before publication. Plan accordingly. [#Rstats](#)

# Code = fully traceable, reproducible analysis



## Code advantages:

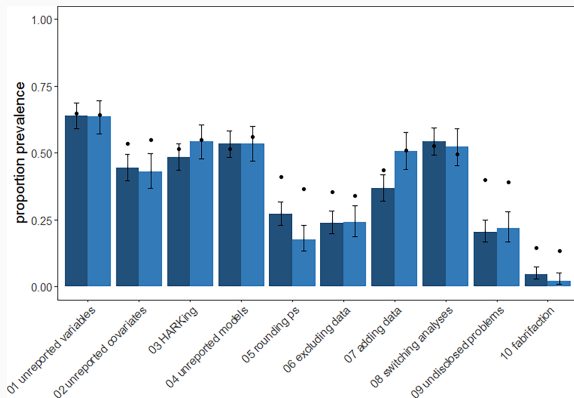
- Easier writing
- Easier, deeper review
- Reusable

# Transparency prevents bad practices

RESEARCH ARTICLE

## Questionable research practices in ecology and evolution

Hannah Fraser<sup>1\*</sup>, Tim Parker<sup>2</sup>, Shinichi Nakagawa<sup>3</sup>, Ashley Barnett<sup>1</sup>, Fiona Fidler<sup>1,4</sup>



p-hacking, HARKing, data fabrication...

DOI:10.1063/PT.6.1.20180822a

22 Aug 2018 in [Research & Technology](#)

## The war over supercooled water

How a hidden coding error fueled a seven-year dispute between two of condensed matter's top theorists.

**Ashley G. Smart**

Over the next seven years, the perplexing discrepancy would ignite a bitter conflict, with junior scientists caught in the crossfire. At stake were not only the reputations of the two groups but also a peculiar theory that sought to explain some of water's deepest and most enduring mysteries. Earlier this year, the dispute was finally settled. And as it turns out, the entire ordeal was the result of botched code.

# Transparency brings better science



**Alexey Shiklomanov**

@ashiklom711



I'm co-author on a study currently published only as a publicly available discussion paper. My code was on GitHub.

A colleague read the paper, thought the results looked weird, checked my code, found a bug and emailed me about it.

This is how science should work. [#openscience](#)

As a condition for publication in ESA journals, all underlying data and statistical code pertinent to the results presented in the publication must be made available in a permanent, publicly accessible data archive or repository, with rare exceptions (see







'Papers with exemplary **data and code archiving**  
are **more valuable** for future research and [...]   
will be given **higher priority** for publication'  
(*Molecular Ecology*)

RESEARCH ARTICLE

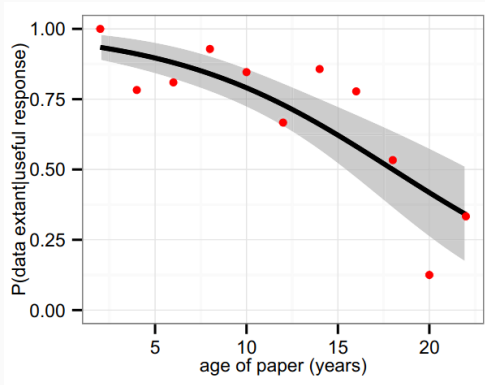
## The citation advantage of linking publications to research data

Giovanni Colavizza <sup>1,2</sup>, Iain Hrynaszkiewicz <sup>3,4</sup>, Isla Staden<sup>1,5</sup>, Kirstie Whitaker <sup>1,6</sup>,  
Barbara McGillivray<sup>1,6\*</sup>

[Colavizza et al 2020](#)

# Let's stop losing data & code

## The Availability of Research Data Declines Rapidly with Article Age

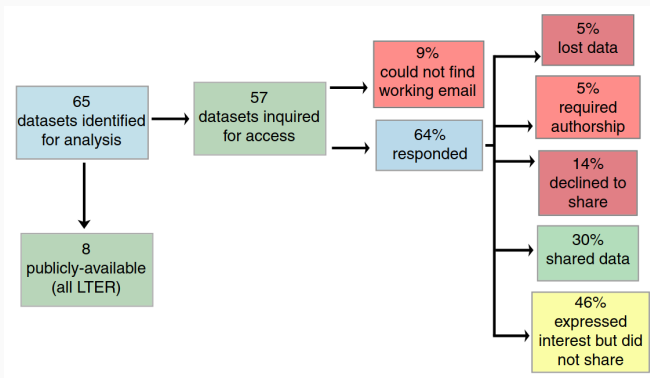


# Open data & code enable synthesis

REVIEW

## Advances in global change research require open science by individual researchers

ELIZABETH M. WOLKOVICH<sup>†</sup>, JAMES REGETZ<sup>‡</sup> and MARY I. O'CONNOR<sup>†</sup>



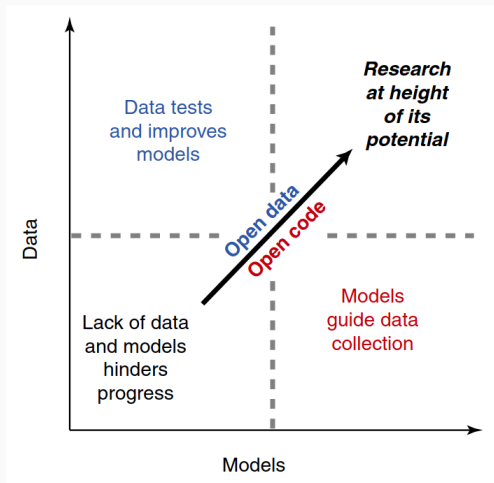
Wolkovich et al 2012

# Open data & code enable synthesis

REVIEW

Advances in global change research require open science by individual researchers

ELIZABETH M. WOLKOVICH<sup>\*,†</sup>, JAMES REGETZ<sup>‡</sup> and MARY I. O'CONNOR<sup>†</sup>



# Reproducible workflows facilitate collaboration

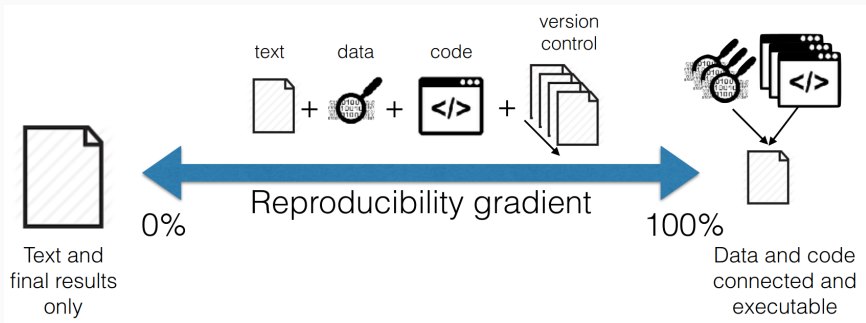
and make everyone happier



# How to do reproducible research

---

# Reproducibility is a gradient



Rodríguez-Sánchez et al. 2016 (modif. Peng 2011)



## Basic reproducibility

---

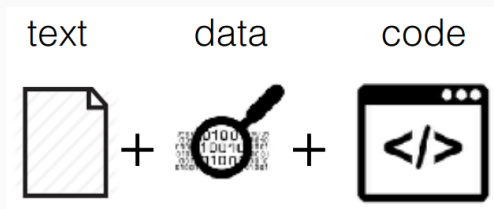
- **MANUSCRIPT** (Text + Tables + Figures)
- **DATA** in permanent archive (see [Tierney & Ram 2020](#))
- **CODE** in permanent archive (see [Eglen et al 2016](#))

*Permanent archive:*

- Zenodo, Dryad, OSF, Figshare, Data Paper...
- NOT GitHub, website...

- **Open** format (csv, txt)
- **README** (who, what, when, where, why, how)
- **Describe variables**
- **Licence** (CC0, CC-BY, ODbL)
- **Citation** (DOI)
- **Metadata** standardised (JSON, XML)

- Scripts: **plain text (.R)**
- **Permanent archive** (eg. Zenodo) with **DOI (citable)**
- [Licence](#)
- **README**

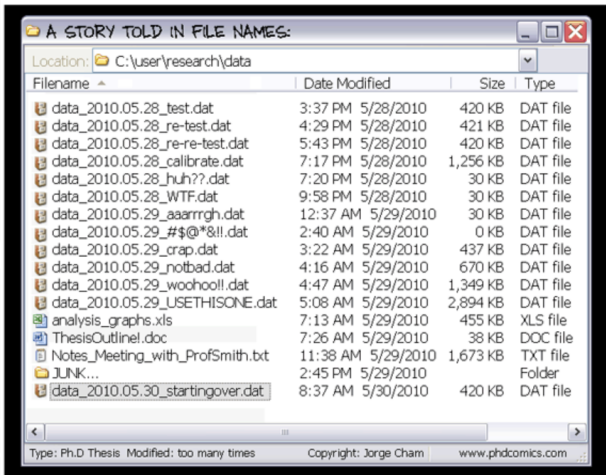


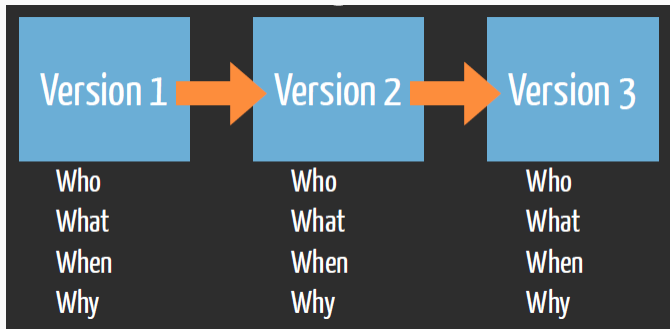
## DATA + CODE

- analysis fully traceable
- results can be regenerated

## Version control

---



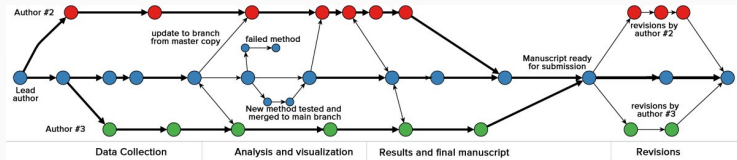




# Much to learn from software engineering

## Git can facilitate greater reproducibility and increased transparency in science

Karthik Ram





Ram 2013
















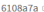





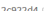





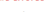










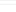

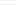

# Automatic checks with Continuous Integration

Reproducibility of computational workflows is automated using continuous analysis

Brett K Beaulieu-Jones<sup>1</sup> & Casey S Greene<sup>2</sup>

Pakillo / Carex.bipolar  build passing

Current   Branches   Build History   Pull Requests   More options 

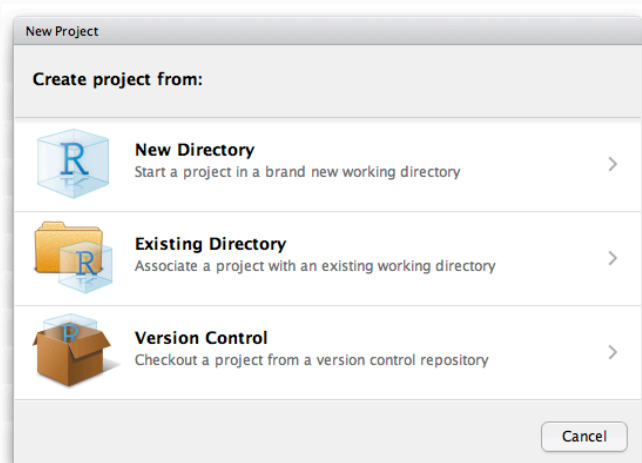
 master  Pakillo	add two more articles to pkgdown	 #7 passed  1c006ff <a href="#">↗</a>	 3 min 22 sec  a day ago
 master  Pakillo	added leaflet occurrence maps to appear as	 #6 passed  57f5374 <a href="#">↗</a>	 5 min 23 sec  a day ago
 master  Pakillo	build site with pkgdown	 #5 passed  6108a7a <a href="#">↗</a>	 17 min 35 sec  a day ago
 master  Pakillo	still trying to fix error with sf in travis (via rnat	 #4 failed  2c922d4 <a href="#">↗</a>	 16 min 58 sec  2 days ago
 master  Pakillo	adding more sf dependencies to travis	 #3 errored  5a60b49 <a href="#">↗</a>	 13 min 59 sec  2 days ago
 master  Pakillo	trying to fix error with rgdal on travis	 #2 errored  076af29 <a href="#">↗</a>	 14 min 15 sec  2 days ago
 master  Pakillo	add travis	 #1 errored  4bce6e8 <a href="#">↗</a>	 18 min 54 sec  3 days ago

## Structuring projects

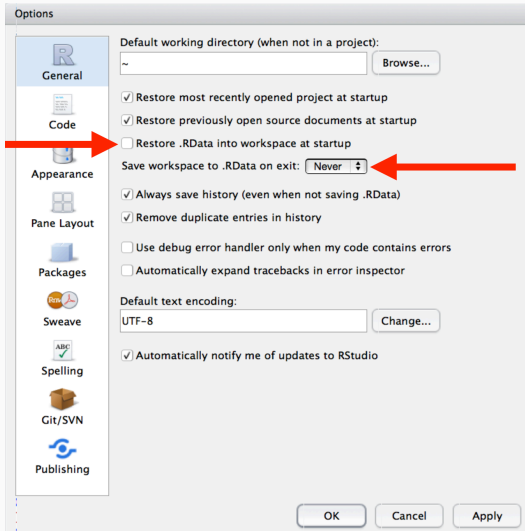
---

```
myproject
|
|- data
|
|- code
|
|- output (figures etc)
|
|- manuscript
```

- Self-contained
- Easy to navigate (file paths)
- Easy to share



# Avoid saving workspace



<https://rstats.wtf>

# Use here for file paths



```
setwd('C:/Users/PACO/myproject')
```

```
mydata <- read.csv('data/mydata.csv')
```



```
library('here')
```

```
mydata <- here('data', 'mydata.csv')
```





## fertile package: real-time feedback on reproducibility

```
library('fertile')  
  
setwd("C:/Users/FRS")
```

*Error: setwd() is likely to break reproducibility. Use here::here() instead.*

<https://github.com/baumer-lab/fertile>

## Structuring projects: guidelines

---

- All files in same directory

# Guidelines for structuring projects

- All files in **same directory**
- **Raw data** separate from **clean data**

Noble 2009, Rodriguez-Sanchez et al 2016, Wilson et al 2017

# Guidelines for structuring projects

- All files in **same directory**
- **Raw data** separate from **clean data**
- **Modular code** (functions)

Noble 2009, Rodriguez-Sanchez et al 2016, Wilson et al 2017

# Guidelines for structuring projects

- All files in **same directory**
- **Raw data** separate from **clean data**
- **Modular code** (functions)
- **Output disposable & separate** from code

## Guidelines for structuring projects

- All files in **same directory**
- **Raw data** separate from **clean data**
- **Modular code** (functions)
- **Output disposable & separate** from code
- `makefile` runs analyses in **appropriate order**

## Guidelines for structuring projects

- All files in **same directory**
- **Raw data** separate from **clean data**
- **Modular code** (functions)
- **Output disposable & separate** from code
- `makefile` runs analyses in **appropriate order**
- **Software dependencies** under control



## Guidelines for structuring projects

- All files in **same directory**
- **Raw data** separate from **clean data**
- **Modular code** (functions)
- **Output disposable & separate** from code
- `makefile` runs analyses in **appropriate order**
- **Software dependencies** under control
- README

## Guidelines for structuring projects

- All files in **same directory**
- **Raw data** separate from **clean data**
- **Modular code** (functions)
- **Output disposable & separate** from code
- `makefile` runs analyses in **appropriate order**
- **Software dependencies** under control
- README
- License

## Project organisation example

myproject

- data
  - data-raw
  - data-clean
- code
- output (figures etc)
- manuscript
- README
- License
- Makefile

- What
- Who
- How
- Licence
- Citation
- etc

README.md

## pandanuisotopes



This repository contains the data and code for our paper:

Florin, A. et al. (2020). *Palaeoprecipitation data from Madjedbebe, northern Australia: A novel proxy from ancient pandanus.*

### How to cite

Please cite this compendium as:

Marwick, B., A. Florin et al., (2020). *Compendium of R code and data for Palaeoprecipitation data from Madjedbebe, northern Australia: A novel proxy from ancient pandanus.* Accessed 16 Oct 2020. Online at <https://doi.org/xxx/xxx>

### How to download

You can download the compendium as a zip from from this URL: <https://github.com/benmarwick/pandanuisotopes/archive/master.zip>

### Licenses

**Text and figures** : [CC-BY-4.0](#)

**Code** : See the [DESCRIPTION](#) file

**Data** : [CC-0](#) attribution requested in reuse

<https://docs.ropensci.org/dataspice/>

```
library("dataspice")
create_spice() # create CSV templates for metadata

edit_creators() # open Shiny apps to edit the CSVs
prep_access()
edit_access()
prep_attributes()
edit_attributes()
edit_biblio()

write_spice() # write machine-readable metadata

build_site() # build human-readable metadata report
```

Break up scripts

```
prepare_data.R
```

```
run_analysis.R
```

```
make_figures.R
```

(and `makefile` will run them in the right order)

## makefile runs code in appropriate order

makefile.R

```
source("prepare_data.R")
```

```
source("run_analysis.R")
```

```
source("make_figures.R")
```

## Don't Repeat Yourself (DRY)

```
dataset %>%  
  filter(species == "Laurus nobilis") %>%  
  ggplot() +  
  geom_point(aes(x, y))
```

```
dataset %>%  
  filter(species == "Laurus azorica") %>%  
  ggplot() +  
  geom_point(aes(x, y))
```



## Don't Repeat Yourself

Write functions (documented + tested)

```
plot_species <- function(sp, data) {  
  data %>%  
    filter(species == sp) %>%  
    ggplot() +  
    geom_point(aes(x, y))  
}
```

Use functions

```
plot_species(sp = "Laurus nobilis", dataset)
```

```
plot_species(sp = "Laurus azorica", dataset)
```

Use for loops

```
for (i in species) {  
  plot_species(sp = i, dataset)  
}
```

Good ol' `lapply`

```
lapply(species, plot_species, data = dataset)
```

## Don't Repeat Yourself

```
library("purrr")  
  
map(species, plot_species, data = dataset)
```

Why rather than What

```
## Response is not linear, so fit gam rather than lm  
model.height <- gam(height ~ s(diameter), data = trees)
```

## Use meaningful names for objects

```
m1 <- lm(height ~ diameter, data = trees)
m2 <- gam(height ~ s(diameter), data = trees)
```

## Use meaningful names for objects

```
m1 <- lm(height ~ diameter, data = trees)
m2 <- gam(height ~ s(diameter), data = trees)
```

```
model.linear <- lm(height ~ diameter, data = trees)
model.gam <- gam(height ~ s(diameter), data = trees)
```










## Project templates

---

## Automatic project creation with template

```
library('template')
```

```
new_project("mynewproj",  
            package = FALSE)
```

 analyses data data-raw manuscript R .Rproj.user makefile.R mynewproj.Rproj README.Rmd .gitignore

# template: New projects also on GitHub

```
new_project("mynewproj",  
            package = FALSE,  
            github = TRUE)
```

Pakillo / mynewproj Private

<> Code Issues Pull requests Actions Projects Wiki Security Insights Settings

master 1 branch 0 tags

Go to file Add file Code

<b>Pakillo</b> Initial commit	654e46f	2 minutes ago	1 commit
.gitignore	Initial commit	2 minutes ago	
README.Rmd	Initial commit	2 minutes ago	
makefile.R	Initial commit	2 minutes ago	
mynewproj.Rproj	Initial commit	2 minutes ago	



**workflow:**  
reproducible projects with  
website

---

## wflow\_start creates project scaffolding

```
library('workflowr')  
  
wflow_start("newproject")
```

 analysis


 code


 data


 docs


 output

 .git


 newproject.Rproj

 README.md

 \_workflowr.yml

 .gitattributes

 .gitignore

 .Rprofile

## wflow\_open starts new analysis

```
wflow_open("analysis/first-analysis.Rmd")
```

```
---  
title: "first-analysis"  
author: "Pakillo"  
date: "2021-06-15"  
output: workflowr::wflow_html  
editor_options:  
  chunk_output_type: console  
---  
|  
## Introduction  
  
```${r}``  
data(iris)  
plot(iris)  
```\n
```

## wflow\_build()

newproject Home About License

Introduction

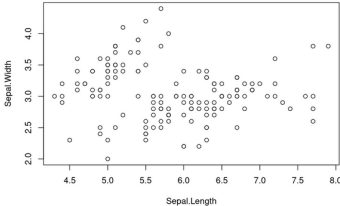
### first-analysis

Pakillo  
2021-06-15

workflow ✓

#### Introduction

```
data(iris)  
plot(iris[, 1:2])
```



Past versions of unnamed-chunk-1-1.png

Session information

The image shows a screenshot of a web application interface. At the top, there is a navigation bar with 'newproject', 'Home', 'About', and 'License'. Below this, a blue button labeled 'Introduction' is visible. The main content area features a title 'first-analysis' with the author 'Pakillo' and the date '2021-06-15'. A 'workflow' icon with a green checkmark is present. The 'Introduction' section contains a code block with the R commands 'data(iris)' and 'plot(iris[, 1:2])'. Below the code is a scatter plot with 'Sepal.Length' on the x-axis (ranging from 4.5 to 8.0) and 'Sepal.Width' on the y-axis (ranging from 2.0 to 4.0). The plot shows a positive correlation between the two variables. At the bottom, there are buttons for 'Past versions of unnamed-chunk-1-1.png' and 'Session information'.

## wflow\_publish commits changes & updates everything

```
wflow_publish(c("analysis/first-analysis.Rmd",  
               "analysis/index.Rmd",  
               "analysis/about.Rmd",  
               "analysis/license.Rmd"),  
             message = "Publish initial analyses")
```



## Connect with GitHub/GitLab and deploy website

```
wflow_use_github("Pakillo")
```

```
wflow_git_push()
```

## Research compendia: projects as packages

---

- Standard structure

Rodríguez-Sánchez et al. 2016, Marwick et al 2018, but see McBain 2020

# Projects as packages

- Standard structure
- Promotes modular code, documented and tested

Rodríguez-Sánchez et al. 2016, Marwick et al 2018, but see McBain 2020

# Projects as packages

- Standard structure
- Promotes modular code, documented and tested
- Easy to share and run

Rodríguez-Sánchez et al. 2016, Marwick et al 2018, but see McBain 2020

# Projects as packages

- Standard structure
- Promotes modular code, documented and tested
- Easy to share and run
- Automatic checks (Continuous Integration)

Rodríguez-Sánchez et al. 2016, Marwick et al 2018, but see McBain 2020

# Projects as packages

- Standard structure
- Promotes modular code, documented and tested
- Easy to share and run
- Automatic checks (Continuous Integration)
- Automatic code review (**goodpractice**)

Rodríguez-Sánchez et al. 2016, Marwick et al 2018, but see McBain 2020

# Projects as packages

- Standard structure
- Promotes modular code, documented and tested
- Easy to share and run
- Automatic checks (Continuous Integration)
- Automatic code review (**goodpractice**)
- Easily create website with **pkgdown**

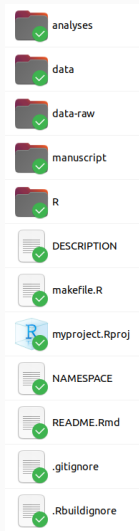
[Rodríguez-Sánchez et al. 2016](#), [Marwick et al 2018](#), but see [McBain 2020](#)



# Creating package structure with template

```
library('template')
```

```
new_project("myproject",  
           package = TRUE)
```



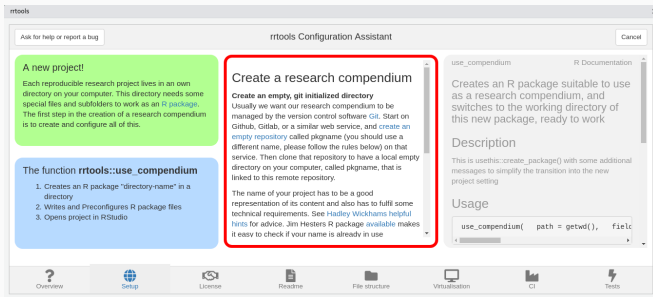
rrtools

---

# rrtools creates research compendia

```
library("rrtools")  
  
use_compendium("~/myproject/")
```

Rstudio addin: <https://github.com/nevrome/rrtools.addin>



## rrtools: project structure

```
- README
- LICENSE
- DESCRIPTION
- travis.yml
- Dockerfile
- analysis/
  |
  |- paper/
    |- paper.Rmd
    |- references.bib
  |
  |- figures/
  |
  |- data/
    |- raw_data/
    |- derived_data/
```



rcompedium

---

## rcompendium creates new project with all scaffolding


```
library('rcompendium')
```


```
new_compendium()
```



















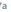





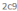






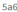







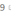






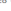


- R package structure
- GitHub repository
- Automatic testing & website update

# Continuous Integration (GitHub Actions, GitLab CI...)

Automatic testing with every commit!

Pakillo / Carex.bipolar  **build passing**

Current Branches Build History Pull Requests More options 

 <b>✓ master</b>  Pakillo	add two more articles to pkgdown	 <b>#7 passed</b>  1c006ff 	 3 min 22 sec  a day ago
 <b>✓ master</b>  Pakillo	added leaflet occurrence maps to appear as a	 <b>#6 passed</b>  57f5374 	 5 min 23 sec  a day ago
 <b>✓ master</b>  Pakillo	build site with pkgdown	 <b>#5 passed</b>  6108a7a 	 17 min 35 sec  a day ago
 <b>✗ master</b>  Pakillo	still trying to fix error with sf in travis (via mat	 <b>#4 failed</b>  2c922d4 	 16 min 58 sec  2 days ago
 <b>! master</b>  Pakillo	adding more sf dependencies to travis	 <b>#3 errored</b>  5a60b49 	 13 min 59 sec  2 days ago
 <b>! master</b>  Pakillo	trying to fix error with rgdal on travis	 <b>#2 errored</b>  076af29 	 14 min 15 sec  2 days ago
 <b>! master</b>  Pakillo	add travis	 <b>#1 errored</b>  4bce6e8 	 18 min 54 sec  3 days ago

<https://github.com/cboettig/compendium>

- DESCRIPTION (dependencies)
- Manuscript (Rmd)
- GitHub Actions



## Data management

---

See <https://dataoneorg.github.io/Education/bestpractices/>

1. Planification (e.g. DMPTool)
2. Collection
3. Metadata description (dataspice, EML, Data Packages, DataPackageR)
4. Quality control (e.g. assertr, validate, pointblank)
5. Storage

<https://docs.ropensci.org/dataspice/>

```
library("dataspice")
create_spice() # create CSV templates for metadata

edit_creators() # open Shiny apps to edit the CSVs
prep_access()
edit_access()
prep_attributes()
edit_attributes()
edit_biblio()

write_spice() # write machine-readable metadata

build_site() # build human-readable metadata report
```

## Check data before analysis

```
library("assertr")

dataset %>%
  assert(within_bounds(0, 0.20), fruit.weight) %>%
  assert(in_set("black", "red"), colour)
```

Check out also [pointblank](#)

## *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 (*J*). 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

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

- [Open Science Framework](#)
- GitHub
- Dropbox
- Figshare, Zenodo, etc
- See all data repositories in [www.re3data.org](http://www.re3data.org)

## Tidy data

---

# Tidy data

country	year	cases	population
Afghanistan	1999	745	15007071
Afghanistan	2000	2666	20095360
Brazil	1999	37737	172006362
Brazil	2000	80488	174004898
China	1999	212258	127200272
China	2000	213766	128000583

variables

country	year	cases	population
Afghanistan	1999	745	15007071
Afghanistan	2000	2666	20095360
Brazil	1999	37737	172006362
Brazil	2000	80488	174004898
China	1999	212258	127200272
China	2000	213766	128000583

observations

country	year	cases	population
Afghanistan	1999	745	15007071
Afghanistan	2000	2666	20095360
Brazil	1999	37737	172006362
Brazil	2000	80488	174004898
China	1999	212258	127200272
China	2000	213766	128000583

values

country	year	cases
Afghanistan	1999	745
Afghanistan	2000	2666
Brazil	1999	37737
Brazil	2000	80488
China	1999	212258
China	2000	213766

country	1999	2000
Afghanistan	745	2666
Brazil	37737	80488
China	212258	213766

table4



## COMMENT

## Open Access



# Gene name errors are widespread in the scientific literature

Mark Ziemann<sup>1</sup>, Yotam Eren<sup>1,2</sup> and Assam El-Osta<sup>1,3\*</sup>

### Abstract

The spreadsheet software Microsoft Excel, when used with default settings, is known to convert gene names to dates and floating-point numbers. A programmatic scan of leading genomics journals reveals that approximately one-fifth of papers with supplementary Excel gene lists contain erroneous gene name conversions.

frequently reused. Our aim here is to raise awareness of the problem.

We downloaded and screened supplementary files from 18 journals published between 2005 and 2015 using a suite of shell scripts. Excel files (.xls and .xlsx suffixes) were converted to tabular separated files (tsv) with `ssconvert` (v1.12.9). Each sheet within the Excel file was converted to a separate tsv file. Each column of data in the tsv file was screened for the presence of gene sym-

# 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 (or GForms) 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 (or GForms) 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 (or GForms) 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 through 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 (or GForms) 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 through 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 (or GForms) 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 through 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 (or GForms) 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 through code.
- Export data as plain text (txt, csv).
- <http://www.datacarpentry.org/spreadsheet-ecology-lesson/>
- <http://kbroman.org/dataorg/>

# 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 (or GForms) 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 through code.
- Export data as plain text (txt, csv).
- <http://www.datacarpentry.org/spreadsheet-ecology-lesson/>
- <http://kbroman.org/dataorg/>
- Broman & Woo: [Data organization in spreadsheets](#)

## Common spreadsheet errors

---

## More than one variable per column

Date collected	Plot	Species-Sex	Weight
1/9/78	1	DM-M	40
1/9/78	1	DM-F	36
1/9/78	1	DS-F	135
1/20/78	1	DM-F	39
1/20/78	2	DM-M	43
1/20/78	2	DS-F	144
3/13/78	2	DM-F	51
3/13/78	2	DM-F	44
3/13/78	2	DS-F	146

Date collected	Plot	Species	Sex	Weight
1/9/78	1	DM	M	40
1/9/78	1	DM	F	36
1/9/78	1	DS	F	135
1/20/78	1	DM	F	39
1/20/78	2	DM	M	43
1/20/78	2	DS	F	144
3/13/78	2	DM	F	51
3/13/78	2	DM	F	44
3/13/78	2	DS	F	146

Source: Data Carpentry



# Multiple tables

1	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG			
2	lake site May 29 2012				29-May				lake site Jun 12 2012				12-Jun				lake site Jun 19 2012				19-Jun				Lake site Jun 26 2012				26-Jun							
3		plot	bug1	bug2	gen	avr	SEM		plot	bug	bug2	gen	avr	SEM		plot	bug1	bug2	gen	avr	SEM		plot	bug1	bug2	gen	avr	SEM		plot	bug1	bug2	gen	avr	SEM	
4	1	T1	1	1	2	T1	2.6	0.51	1	T1	6	85	91	T1	30.4	15.47126	1	T1	17	80	97				1	T1	52	191	243							
5	2	T1	1	2	3	T2	0.2	0.2	2	T1	8	13	21	T2	0.2	0.2	2	T1	44	136	180	T1	77.8	30.384865	2	T1	50	270	320				T1	141.6	50.313	
6	3	T1	1	3	4	control	0.2	0.2	3	T1	11	0	11	control	0.6	0.6	3	T1	18	0	18	T2	1.8	0.5620499	3	T1	6	0	6			T2	0.2	0.2		
7	4	T1	1	0	1				4	T1	0	6	6			4	T1	0	14	14	control	0.4	0.244949	4	T1	0	39	39	control	0	0			0		
8	5	T1	0	3	3				5	T1	3	20	23			5	T1	10	70	80				5	T1	4	96	100								
9	6	T2	1	0	1				6	T2	0	0	0			6	T2	1	7	8				6	T2	0	1	1								
10	7	T2	0	0	0				7	T2	0	0	0			7	T2	0	1	1				7	T2	0	0	0								
11	8	T2	0	0	0				8	T2	1	0	1			8	T2	0	0	0				8	T2	0	0	0								
12	9	T2	0	0	0				9	T2	0	0	0			9	T2	0	0	0				9	T2	0	0	0								
13	10	T2	0	0	0				10	T2	0	0	0			10	T2	0	0	0				10	T2	0	0	0								
14	11	control	0	0	0				11	control	0	0	0			11	control	0	0	0				11	control	0	0	0								
15	12	control	0	0	0				12	control	0	0	0			12	control	0	0	0				12	control	0	0	0								
16	13	control	0	0	0				13	control	0	0	0			13	control	0	0	0				13	control	0	0	0								
17	14	control	0	0	0				14	control	0	0	0			14	control	0	1	1				14	control	0	0	0								
18	15	control	1	0	1				15	control	3	0	3			15	control	0	1	1				15	control	0	0	0								
19																																				
20																																				
21	Barn site May 29 2012				29-May				Barn site Jun 12 2012				12-Jun				Barn site Jun 19 2012				19-Jun				Barn Site Jun 26 2012				26-Jun							
22		plot	bug1	bug2	gen	avr	SEM		plot	bug	bug2	gen	avr	SEM		plot	bug1	bug2	gen	avr	SEM		plot	bug1	bug2	gen	avr	SEM		plot	bug1	bug2	gen	avr	SEM	
23	1	T1	3	3	6				1	T1	21	0	21			1	T1	5	0	5				1	T1	0	0	0								
24	2	T1	1	4	5				2	T1	36	74	110			2	T1	65	502	562				2	T1	44	2057	2101	T1	431.8	417.33					
25	3	T1	0	0	0	T1	2.4	1.288	3	T1	13	0	13	T1	30.6	10.10124	3	T1	10	7	17	T1	119.4	11.92882	3	T1	12	20	32	T2	0.4	0.4				
26	4	T1	0	0	0	T2	0.4	0.245	4	T1	7	0	7	T2	1	0.774597	4	T1	0	16	6	T2	5	1.908902	4	T1	0	16	16	control	1.2	0.5831				
27	5	T1	0	1	1	control	1	0.316	5	T1	2	0	2	control	2.2	1.714643	5	T1	0	2	2	control	2.8	0.969536	5	T1	0	10	10							
28	6	T2	0	0	0				6	T2	1	0	1			6	T2	0	8	8				6	T2	0	0	0								
29	7	T2	0	0	0				7	T2	0	4	4			7	T2	0	12	12				7	T2	0	0	0								
30	8	T2	0	1	1				8	T2	0	0	0			8	T2	0	0	0				8	T2	0	0	0								
31	9	T2	0	1	1				9	T2	0	0	0			9	T2	3	0	3				9	T2	0	0	0								
32	10	T2	0	0	0				10	T2	0	0	0			10	T2	2	0	2				10	T2	0	1	1								
33	11	control	0	0	0				11	control	1	0	1			11	control	0	5	5				11	control	0	2	2								
34	12	control	0	1	1				12	control	0	0	0			12	control	1	1	2				12	control	1	0	1								
35	13	control	0	1	1				13	control	0	0	0			13	control	0	0	0				13	control	0	0	0								
36	14	control	0	1	1				14	control	8	1	9			14	control	0	5	5				14	control	0	3	3								
37	15	control	0	2	2				15	control	0	1	1			15	control	0	2	2				15	control	1	0	0								
38																																				
39																																				

Could you avoid new tab by adding a column to original spreadsheet?

## Using formatting, comments, etc to convey information

Plot: 2			
Date collect	Species	Sex	Weight
1/8/14	NA		
1/8/14	DM	M	44
1/8/14	DM	M	38
1/8/14	OL		
1/8/14	PE	M	22
1/8/14	DM	M	38
1/8/14	DM	M	48
1/8/14	DM	M	43
1/8/14	DM	F	35
1/8/14	DM	M	43
1/8/14	DM	F	37
1/8/14	PF	F	7
1/8/14	DM	M	45
1/8/14	OT		
1/8/14	DS	M	157
1/8/14	OX		
2/18/14	NA	M	218
2/18/14	PF	F	7
2/18/14	DM	M	52

measurement device not calibrated

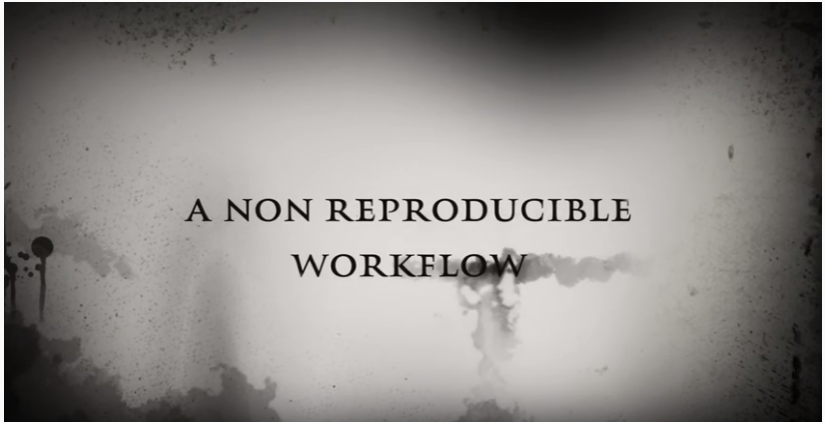
Date collect	Species	Sex	Weight	Calibrated
1/8/14	NA			
1/8/14	DM	M	44	Y
1/8/14	DM	M	38	Y
1/8/14	OL			
1/8/14	PE	M	22	Y
1/8/14	DM	M	38	Y
1/8/14	DM	M	48	Y
1/8/14	DM	M	43	Y
1/8/14	DM	F	35	Y
1/8/14	DM	M	43	Y
1/8/14	DM	F	37	Y
1/8/14	PF	F	7	Y
1/8/14	DM	M	45	Y
1/8/14	OT			
1/8/14	DS	M	157	N
1/8/14	OX			
2/18/14	NA	M	218	N
2/18/14	PF	F	7	Y
2/18/14	DM	M	52	Y

Your turn: tidy up this messy dataset

<https://ndownloader.figshare.com/files/2252083>

## Reproducible dynamic documents with Rmarkdown

---



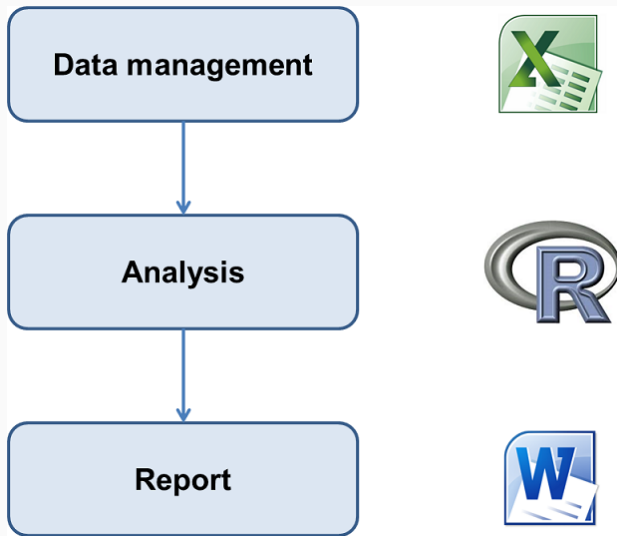
<https://youtu.be/s3JldKoA0zw>

## A typical research workflow

1. Prepare data (**s**preadsheet)
2. Analyse data (**R**)
3. Write report/paper (**W**ord)
4. Start the email attachments nightmare...



# This workflow is broken





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

- 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?

- 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!

## Manual copy-paste is tedious & problematic

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-0.0651657	0.4264970	-0.153	0.879
sunshine	0.0100228	0.0004232	23.683	<2e-16

'Transcribing numbers from stats software by hand was the largest source of errors'

(Eubank 2016)



**Trevor A. Branch**

@TrevorABranch



Follow

My rule of thumb: every analysis you do on a dataset will have to be redone 10–15 times before publication. Plan accordingly. [#Rstats](#)

Your **closest collaborator** is you 6 months ago,  
and you don't respond to emails.

(P. Wilson)

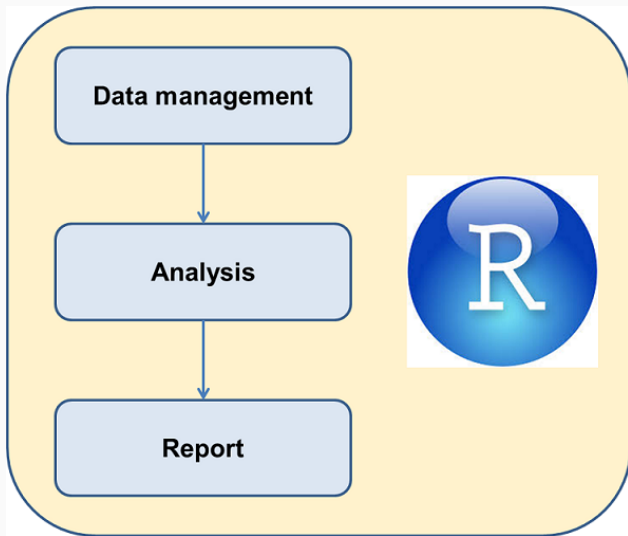
Even **you** will struggle to reproduce  
**your own results** from a few weeks/months ago.

Writing reproducible manuscripts is hard

Revising non-reproducible manuscripts is even harder

.

**Also, please note that because rev#1  
asked to re-calculate effect sizes (...)  
we need to change every single  
number in the main text.**



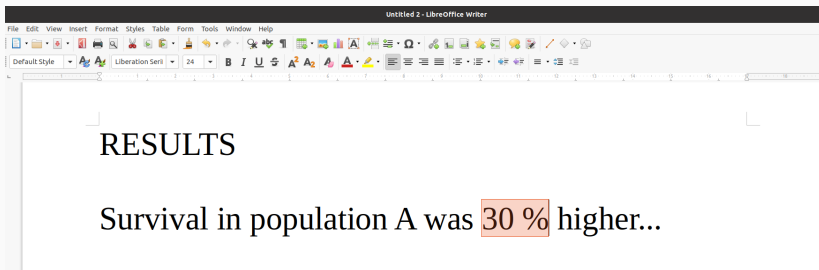


# Rmarkdown documents

- Fully reproducible (trace all results inc. tables and plots)
- Dynamic (regenerate with 1 click)
- Multiple outputs:
  - documents (HTML, Word, PDF)
  - presentations (HTML, PDF, PowerPoint)
  - books
  - websites...



# Where does this value come from?



*Rmarkdown:*

Survival in population A was ``r surv.diff`` % higher

*Output:*

Survival in population A was **30** % higher

```
mydata <- read.csv("data.txt")
```

*Rmarkdown:*

We measured `nrow(mydata)` individuals

*Output:*

We measured **100** individuals

*Much better than copy-paste!*

# Rmarkdown: code (R, Python, etc) + text (Markdown)

```
---  
title: "Does sunshine make people happy?"  
author: "FRS"  
output: word_document  
---
```

## ## Introduction

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

## ## Methods

```
```{r echo=FALSE}  
## Read data  
data <- read.table("data.txt", header = TRUE)  
  
# Fit linear model  
model <- lm(happiness ~ sunshine, data = data)  
```
```

We collected data on `nrow(data)` individuals and fitted a linear model.

Metadata  
(YAML)

Text  
(Markdown)

Code  
(R, Python...)

```
```{r echo=FALSE, eval=TRUE, cache=TRUE, fig.height=3}  
plot(iris)  
```
```

<https://yihui.org/knitr/options/>

```
```{r}  
#| echo = FALSE  
#| eval = TRUE  
#| fig.cap = "My figure caption"  
plot(iris)  
```
```

# Naming chunks helps debugging

```
processing file: test.Rmd
|.....| 14%
ordinary text without R code

|.....| 29%
label: setup (with options)
List of 1
$ include: logi FALSE

|.....| 43%
ordinary text without R code

|.....| 57%
label: read.data

|.....| 71%
ordinary text without R code

|.....| 86%
label: plot (with options)
List of 1
$ echo: logi FALSE

Quitting from lines 28-29 (test.Rmd)
Error in eval(predvars, data, env) : object 'specie' not found
Calls: <Anonymous> ... plot.formula -> eval -> eval -> <Anonymous> -> eval -> eval
Execution halted
```



## Naming chunks helps navigating long docs

```
1 ---
2 title: "My Analysis"
3 author: "FRS"
4 output: html_document
5 ---
6
7 ```{r setup, include=FALSE}
8 knitr::opts_chunk$set(echo = TRUE)
9 ```
10
11 This is an R Markdown document. Markdown is a simple
12 for authoring HTML, PDF, and MS Word
13 for details on using R Markdown see
14 .rstudio.com>.
```

**My Analysis**  
Chunk 1: setup  
Chunk 2: read.data  
Chunk 3: plot

11:60 (Top Level) ↕ R Markdown ↕

## Naming chunks: figure files take chunk name



unnamed-chunk-1-1.png



unnamed-chunk-1-2.png



unnamed-chunk-1-3.png



unnamed-chunk-1-4.png

knitr engines:

|      |             |             |            |         |          |           |
|------|-------------|-------------|------------|---------|----------|-----------|
| [1]  | "asis"      | "asy"       | "awk"      | "bash"  | "block"  | "block2"  |
| [7]  | "bslib"     | "c"         | "cat"      | "cc"    | "coffee" | "comment" |
| [13] | "css"       | "ditaa"     | "dot"      | "embed" | "eviews" | "exec"    |
| [19] | "fortran"   | "fortran95" | "gawk"     | "go"    | "groovy" | "haskell" |
| [25] | "highlight" | "js"        | "julia"    | "lein"  | "mysql"  | "node"    |
| [31] | "octave"    | "perl"      | "php"      | "psql"  | "python" | "R"       |
| [37] | "Rcpp"      | "Rscript"   | "ruby"     | "sas"   | "sass"   | "scala"   |
| [43] | "scss"      | "sed"       | "sh"       | "sql"   | "stan"   | "stata"   |
| [49] | "targets"   | "tikz"      | "verbatim" | "zsh"   |          |           |

# Markdown: easy text formatting

# Header

## Subheader

*\*italic\**

**\*\*bold\*\***

[a link](https://example.com)

.

Handy: <https://thinkr-open.github.io/remedy/>

Or use [Visual Markdown Editor](#)

```
----
title: "Does sunshine make people happy?"
output: pdf_document
bibliography: refs.bib
----

# Introduction

Climate influences individual well-being [Rehdanz_2005].
However, ...

# Methods

```{r echo=FALSE}
# read data
data <- read.table("data.txt", header=T)
data[10,1] <- 11 # correct error

# fit linear model
model <- lm(happiness ~ sunshine, data=data)
```

We collected data on `r nrow(data)` individuals and fitted a
linear model.

# Results

We found that...

```{r echo=FALSE, results='asis'}
# make table with model output
print(xtable::xtable(model), comment = FALSE)
```

```{r echo=FALSE, fig.height=3, fig.width=3, fig.align='center'}
visreg::visreg(model) # plot
```

# Discussion

Our results confirm that happiness is related to
sunshine (slope = `r coef(model)[2]`).

# References
```

**a**

## Does sunshine make people happy?

**b**

### Introduction

Climate influences individual well-being (Rehdanz and Maddison 2005). However, ...

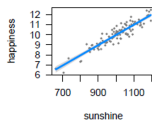
### Methods

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

### Results

We found that...

|             | Estimate | Std. Error | t value | P(> t ) |
|-------------|----------|------------|---------|---------|
| (Intercept) | -0.0986  | 0.4271     | -0.23   | 0.8180  |
| sunshine    | 0.0101   | 0.0004     | 23.75   | 0.0000  |



### Discussion

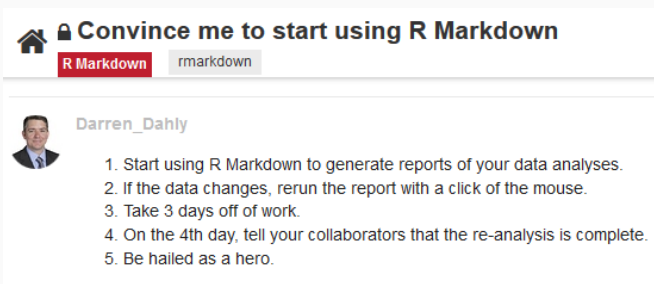
Our results confirm that happiness is related to sunshine (slope = 0.0100652).

### References

Rehdanz, Katrin, and David Maddison. 2005. "Climate and Happiness." *Ecological Economics* 52 (1). Elsevier BV: 111-25. doi:10.1016/j.ecolecon.2004.06.015.

## Spotted error in the data? No problem!


- Make changes in Rmarkdown document
- Click **Knit** in Rstudio
- Report will **update automatically!**



The screenshot shows a forum post with a home icon, a lock icon, and the title "Convince me to start using R Markdown". Below the title are two tags: "R Markdown" in a red box and "rmarkdown" in a grey box. The author's profile picture and name "Darren\_Dahly" are on the left. The post content is a numbered list of five points.

## Convince me to start using R Markdown

R Markdown rmarkdown

 Darren\_Dahly

1. Start using R Markdown to generate reports of your data analyses.
2. If the data changes, rerun the report with a click of the mouse.
3. Take 3 days off of work.
4. On the 4th day, tell your collaborators that the re-analysis is complete.
5. Be hailed as a hero.

<https://community.rstudio.com/t/convince-me-to-start-using-r-markdown/1636/12>

Your turn

---



## Create, edit and share Rmarkdown document

File > New File > Rmarkdown

Write text

Insert code chunks

Change chunk options (echo, eval, etc)

HTML/Word/PDF output

PDF generation requires LaTeX

```
library('tinytex')  
  
install_tinytex()
```

## Rmarkdown bells and whistles

---

# 'Visual Rmarkdown': Rmd as in word processor

The editor toolbar includes buttons for the most commonly used formatting commands:



Additional commands are available on the **Format**, **Insert**, and **Table** menus:

| Format                                      | Insert                          | Table                           |
|---|---------------------------------|---------------------------------|
| <b>B</b> Bold <span>⌘B</span>               | Rmd Chunk <span>⌘I</span>       | Insert Table... <span>⌘T</span> |
| <i>I</i> Italic <span>⌘I</span>             | Image... <span>⇧⌘I</span>       | ✓ Table Header                  |
| <code>&lt;/&gt;</code> Code <span>⌘D</span> | Link... <span>⌘K</span>         | Table Caption                   |
| Text ▶                                      | Horizontal Rule <span>⌘_</span> | Align Column ▶                  |
| Bullets & Numbering ▶                       | Definition ▶                    | Insert Row Above                |
| Blockquote                                  | Inline Math                     | Insert Row Below                |
| Line Block                                  | Display Math                    | Insert Column Left              |
| Div Block...                                | Footnote <span>⇧⌘F7</span>      | Insert Column Right             |
| Code Block...                               | Citation...                     | Delete Row                      |
| Raw ▶                                       | Div Block...                    | Delete Column                   |
| Clear Formatting <span>⌘\</span>            | Code Block...                   | Delete Table                    |
| Edit Attributes... <span>F4</span>          | YAML Block                      |                                 |
|   | Comment <span>⇧⌘C</span>        |                                 |

<https://rstudio.github.io/visual-markdown-editing>

## Automatic table generation

```
model <- lm(happiness ~ sunshine, data = mydata)
xtable(model)
```

|             | Estimate | Std. Error | t value | Pr(> t ) |
|-------------|----------|------------|---------|----------|
| (Intercept) | -0.0652  | 0.4265     | -0.15   | 0.8789   |
| sunshine    | 0.0100   | 0.0004     | 23.68   | 0.0000   |

Many alternatives: `gtsummary`, `modelsummary`, `huxtable`, etc

We fitted a linear model:

```
library('equatiomatic')  
model <- lm(happiness ~ sunshine, data = mydata)  
extract_eq(model)
```

$$\text{happiness} = \alpha + \beta_1(\text{sunshine}) + \epsilon \quad (1)$$

# Models that describe themselves!

```
library("report")
model <- lm(happiness ~ sunshine, data = mydata)
report(model)
```

We fitted a linear model (estimated using OLS) to predict happiness with sunshine (formula: happiness ~ sunshine). The model explains a statistically significant and substantial proportion of variance ( $R^2 = 0.85$ ,  $F(1, 98) = 560.90$ ,  $p < .001$ , adj.  $R^2 = 0.85$ ). The model's intercept, corresponding to sunshine = 0, is at -0.07 (95% CI [-0.91, 0.78],  $t(98) = -0.15$ ,  $p = 0.879$ ). Within this model:

- The effect of sunshine is statistically significant and positive ( $\beta = 0.01$ , 95% CI [9.18e-03, 0.01],  $t(98) = 23.68$ ,  $p < .001$ ; Std.  $\beta = 0.92$ , 95% CI [0.85, 1.00])

Standardized parameters were obtained by fitting the model on a standardized version of the dataset. 95% Confidence Intervals (CIs) and p-values were computed using a Wald t-distribution approximation.

Using LaTeX:

\$\$

```
y \sim N(\mu, \sigma^2)
```

\$\$

$$y \sim N(\mu, \sigma^2)$$

- Mathpix: <https://github.com/jonocarroll/mathpix>

# Citing bibliography

Insert Citation

My Sources

- Bibliography
- Zotero
- My Library
- From DOI
- Crossref
- DataCite
- PubMed

Search for citation

|  |  |   |
|--|--|---|
| @baghizadehfini2020  | Baghizadeh Fini, M 2020                            | + |
| What dentists need to know about COVID-19  |  |   |
| @bostanciklioglu2020   | Bostanciklioglu, M 2020                            | + |
| Severe Acute Respiratory Syndrome Coronavirus 2 is Penetrating to Dementia Re...   |  |   |
| @fran  | Elliott, C, and Hudak, P 1997                      | + |
| Functional reactive animation  |  |   |
| @guo2020   | Guo, Y, Cao, Q, Hong, Z, Tan, Y, Chen, et al. 2020 | + |
| The origin, transmission and clinical therapies on coronavirus disease 2019 (CO... |  |   |
| @hu2020  | Hu, B, Huang, S, and Yin, L 2020                   | + |
| The cytokine storm and COVID-19  |  |   |
| @malik2020   | Malik, Y, Kumar, N, Sircar, S et al. 2020          | + |
| Coronavirus Disease Pandemic (COVID-19): Challenges and a Global Perspective       |  |   |
| @R-base  | R Core Team 2017                                   | + |
| R: A language and environment for statistical computing                            |  |   |

Selected Citation Keys

Add to bibliography: book.bib

Insert Cancel

<https://rstudio.github.io/visual-markdown-editing/#/citations>



```
---  
title: "My awesome Rmd"  
output: html_document  
bibliography: references.bib  
---
```

## Format bibliography for any journal

```
---  
title: "Does sunshine make people happy?"  
author: "FRS"  
output: word_document  
bibliography: myrefs.bib  
csl: ecology-letters.csl  
---
```

Thousands of Citation Styles:

<https://www.zotero.org/styles>

<https://github.com/citation-style-language/styles>

# Rmarkdown templates

- rarticles
- papaja
- rrttools
- pinp
- rmdTemplates
- pagedreport
- GitHub!

## My cool paper written in Rmarkdown

F. Rodriguez Sanchez<sup>1,2</sup> and And Fritander<sup>3</sup>

<sup>1</sup>Max Institute of Technology, Department of Street, City, State, Zip; <sup>2</sup>Academic/University Department, Street, City, State, Zip

This manuscript was completed on September 10, 2010.

Please provide an abstract of no more than 200 words in a single paragraph. Abstracts should explain to the general reader the major contributions of the article. No references in the abstract may be cited or full words in the abstract may be used and the like in the text.

one | two | optional | optional | optional

This PNAS journal template is provided to help you write your work in the current journal format. Instructions for use are provided below.

Note: please start your introduction without including the word "Introduction" as a section heading (except for each article in the Physical Science section); this heading is implied in the first paragraph.

### Guide to using this template

Please note that while this template provides a preview of the typeset manuscript for submission, in this preparation, it will not necessarily be the final publication layout. For more detailed information please see the PNAS Information for Authors.

**Author Affiliations.** Include department, institution, and complete address, with the ZIP/postal code, for each author. Use lower case letters to match authors with institutions, as shown in the example. Authors with an ORCID ID may supply this information at submission.

**Submitting Manuscripts.** All authors must submit their article as PNAScentral. If you are using Overleaf to write your article, you can use the "Feedback to PNAS" option in the top bar of the editor window.

**Format.** Many authors find it useful to organize their manuscripts with the following order of sections: Title, Author Affiliations, Keywords, Abstract, Significance Statement, Results, Discussion, Materials and Methods, Acknowledgments, and References. Other orders and headings are permitted.

**Manuscript Length.** PNAS generally uses a two-column format averaging 67 characters, including spaces, per line. The maximum length of a Direct Submission research article is six pages and a PNAS PLUS research article is ten pages including all text, space, and the number of characters displayed by figures, tables, and equations. When submitting tables, figures, and/or equations in addition to text, keep the text for your manuscript under 10,000 characters (including spaces) for Direct Submissions and 72,000 characters (including spaces) for PNAS PLUS.

**References.** References should be cited in numerical order as they appear in text; this will be done automatically via bibTeX, e.g. (1) and (2, 3). All references, including for the SI, should be included in the main manuscript file. References appearing in both sections should not be duplicated. All references



Fig. 1. Photocopy image of a frog with a large sample caption to show publication width.

included in tables should be included with the main reference section.

**Data Archival.** PNAS must be able to archive the data essential to a published article. Where such archiving is not possible, deposition of data to public databases, such as GenBank, ArrayExpress, Protein Data Bank, UniProt, and others outlined in the Information for Authors, is acceptable.

**Language-Editing Services.** Prior to submission, authors who believe their manuscripts would benefit from professional editing are encouraged to use a language-editing service (see list at [www.pnas.org/author/language-editing-services](http://www.pnas.org/author/language-editing-services)). PNAS does not take responsibility for or endorse these services, and their use has no bearing on acceptance of a manuscript for publication.

### Significance Statement

Authors must submit a 120-word maximum statement about the significance of their research paper written at a level understandable to an unprejudiced, nonexpert outside their field of specialty. The primary goal of the Significance Statement is to explain the relevance of the work. In broad contrast to a broad narrative, the Significance Statement appears in the paper file and is required for all research papers.

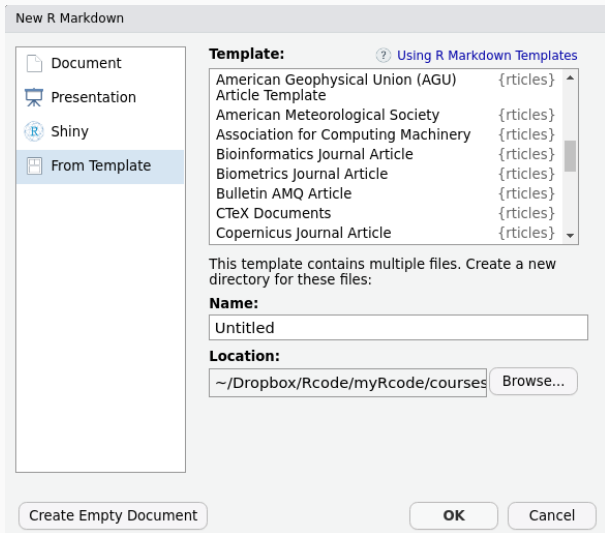
This box provides a table of contents.

This box shows the content of a research article.

[www.pnas.org/cgi/doi/10.1073/pnas.1009.2000.2000](http://www.pnas.org/cgi/doi/10.1073/pnas.1009.2000.2000)

PNAS | September 10, 2010 | vol. 87 | no. 37 | 1-4

# Accessing Rmd templates



# Revise writing style: gramr

**Ignore**

- Passive Voice
- Duplicate words (the the)
- 'So' at start of sentence
- 'There is/are; at start of sentence
- Avoid weasel words
- Wordiness
- Problematic Adverbs
- Cliches
- Avoid 'Being' words

Next Finish

**Text to Check**

So the cat was stolen. This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://markdown.rstudio.com>.

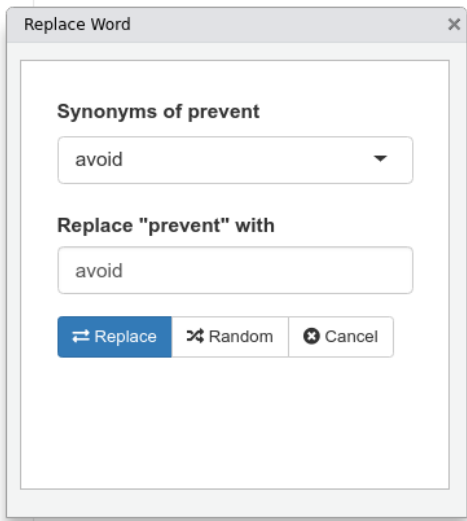
"was stolen" may be passive voice

<https://github.com/ropenscilabs/gramr>

<https://github.com/nevrome/wellspell.addin>

# Find synonyms

prevent



<https://github.com/gadenbuie/synamyn>

| Method          | koRpus      | stringi       |
|-----------------|-------------|---------------|
| Word count      | 107         | 104           |
| Character count | 604         | 603           |
| Sentence count  | 10          | Not available |
| Reading time    | 0.5 minutes | 0.5 minutes   |

<https://github.com/benmarwick/wordcountaddin>



## BOOKDOWN

### Write HTML, PDF, ePub, and Kindle books with R Markdown

The `bookdown` package is an [open-source R package](https://bookdown.org/) that facilitates writing books and long-form articles/reports with R Markdown. Features include:

- Generate printer-ready books and ebooks from R Markdown documents.
- A markup language easier to learn than LaTeX, and to write elements such as section headers, lists, quotes, figures, tables, and citations.
- Multiple choices of output formats: PDF, LaTeX, HTML, EPUB, and Word.
- Possibility of including dynamic graphics and interactive applications (HTML widgets and Shiny apps).
- Support a wide range of languages: R, C/C++, Python, Fortran, Julia, Shell scripts, and SQL, etc.
- LaTeX equations, theorems, and proofs work for all output formats.
- Can be published to GitHub, bookdown.org, and any web servers.
- Integrated with the RStudio IDE.
- One-click publishing to <https://bookdown.org>.



<https://bookdown.org/>

# Presentation Ninja



with xaringan

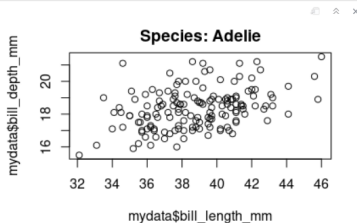
Yihui Xie

RStudio, PBC

<https://slides.yihui.org/xaringan/>

# Parameterised reports

```
---  
title: "My template report"  
output: html_document  
params:  
  sp: Adelie  
---  
|  
````{r}  
library(palmerpenguins)  
data("penguins")  
  
mydata <- subset(penguins, species == params$sp)  
  
plot(mydata$bill_length_mm, mydata$bill_depth_mm,  
      main = paste0("Species: ", params$sp))  
...  
|
```



## Render thousands of individual reports from Rmd template

```
library('rmarkdown')

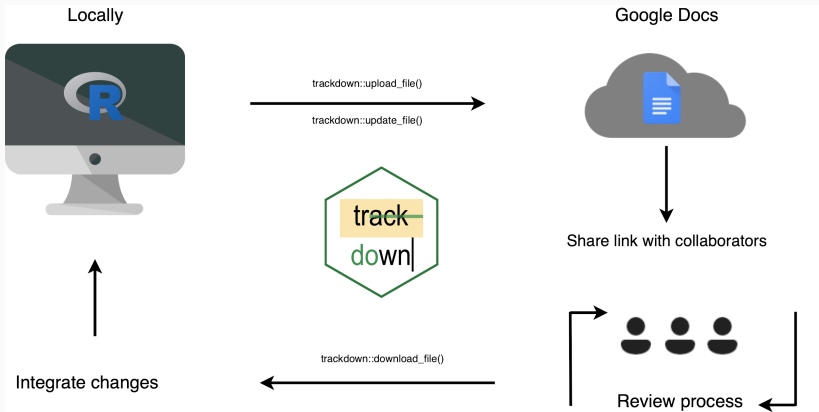
for (i in unique(penguins$species)) {

  render("template_report.Rmd",
        params = list(sp = i))

}
```

# Collaborative writing

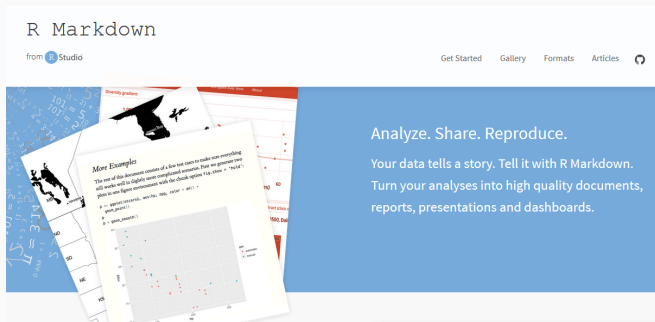
- GitHub, GitLab, etc
- Google Docs ([trackdown](#))
- [redoc](#)



## Rmarkdown resources

---

<http://rmarkdown.rstudio.com/>



The image shows a screenshot of the R Markdown website homepage. At the top left, the text "R Markdown" is displayed in a large, monospaced font, with "from RStudio" below it. To the right, there are navigation links: "Get Started", "Gallery", "Formats", "Articles", and a search icon. The main content area features a blue background with the text "Analyze. Share. Reproduce." and "Your data tells a story. Tell it with R Markdown. Turn your analyses into high quality documents, reports, presentations and dashboards." On the left side of this area, there is a collage of images including a map of the United States, a scatter plot, and a document snippet titled "More Examples" which contains R code and a plot.





# Rmarkdown reference guide



## R Markdown Reference Guide

Learn more about R Markdown at [rmarkdown.rstudio.com](http://rmarkdown.rstudio.com)  
Learn more about interactive Shiny at [shiny.rstudio.com/articles](http://shiny.rstudio.com/articles)

Contents:  
1. Markdown Syntax  
2. Knitr chunk options  
3. Pandoc options

| Syntax                                                                                                                                                                                           | Becomes                                                                                                                                                                                 |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p>Make a code chunk with three back ticks followed by an <code>r</code> in braces. End the chunk with three back ticks:</p> <pre>```{r} paste("Hello", "World!") ```</pre>                      | <p>Make a code chunk with three back ticks followed by an <code>r</code> in braces. End the chunk with three back ticks:</p> <pre>paste("Hello", "World!")  ## [1] "Hello World!"</pre> |
| <p>Place code inline with a single back tick. The first back tick must be followed by an <code>R</code>, like this: <code>`r paste("Hello", "World!")`</code>.</p>                               | <p>Place code inline with a single back tick. The first back tick must be followed by an <code>R</code>, like this: <code>Hello World!</code>.</p>                                      |
| <p>Add chunk options within braces. For example, <code>echo=FALSE</code> will prevent source code from being displayed:</p> <pre>```{r eval=TRUE, echo=FALSE} paste("Hello", "World!") ```</pre> | <p>Add chunk options within braces. For example, <code>echo=FALSE</code> will prevent source code from being displayed:</p> <pre>## [1] "Hello World!"</pre>                            |

Learn more about chunk options at <http://tinyurl.com/knitr/options>

| Chunk options            |                        |                                                                                                                                                                                                                                                                                                                         |
|--------------------------|------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| option                   | default value          | description                                                                                                                                                                                                                                                                                                             |
| <b>Code execution</b>    |                        |                                                                                                                                                                                                                                                                                                                         |
| <code>child</code>       | <code>NULL</code>      | A character vector of filenames. Knitr will knit the files and place them into the main document.                                                                                                                                                                                                                       |
| <code>code</code>        | <code>NULL</code>      | Set to <code>R</code> code. Knitr will replace the code in the chunk with the code in the <code>code</code> option.                                                                                                                                                                                                     |
| <code>engine</code>      | <code>"r"</code>       | Knitr will evaluate the chunk in the named language, e.g. <code>engine = "python"</code> . Run <code>names(knitr::knit_engines())</code> to see supported languages.                                                                                                                                                    |
| <code>eval</code>        | <code>TRUE</code>      | If <code>FALSE</code> , knitr will not run the code in the code chunk.                                                                                                                                                                                                                                                  |
| <code>include</code>     | <code>TRUE</code>      | If <code>FALSE</code> , knitr will not include the chunk but not include the chunk in the final document.                                                                                                                                                                                                               |
| <code>raw</code>         | <code>TRUE</code>      | If <code>FALSE</code> , knitr will not include the chunk when <code>raw=TRUE</code> is used to extract the source code.                                                                                                                                                                                                 |
| <b>Formatting</b>        |                        |                                                                                                                                                                                                                                                                                                                         |
| <code>collapse</code>    | <code>FALSE</code>     | If <code>TRUE</code> , knitr will collapse all the source and output blocks created by the chunk into a single block.                                                                                                                                                                                                   |
| <code>echo</code>        | <code>TRUE</code>      | If <code>FALSE</code> , knitr will not display the code in the code chunk above it's results in the final document.                                                                                                                                                                                                     |
| <code>results</code>     | <code>"markup"</code>  | If <code>"html"</code> , knitr will not display the code's results in the final document. If <code>"text"</code> , knitr will delay displaying of output pieces until the end of the chunk. If <code>"raw"</code> , knitr will pass through results without reformatting them (unless it results return use HTML, etc.) |
| <code>error</code>       | <code>TRUE</code>      | If <code>FALSE</code> , knitr will not display any error messages generated by the code.                                                                                                                                                                                                                                |
| <code>message</code>     | <code>TRUE</code>      | If <code>FALSE</code> , knitr will not display any messages generated by the code.                                                                                                                                                                                                                                      |
| <code>warning</code>     | <code>TRUE</code>      | If <code>FALSE</code> , knitr will not display any warning messages generated by the code.                                                                                                                                                                                                                              |
| <b>Code Decoration</b>   |                        |                                                                                                                                                                                                                                                                                                                         |
| <code>background</code>  | <code>"#f7f7f7"</code> | A background color for chunks in LaTeX output.                                                                                                                                                                                                                                                                          |
| <code>comment</code>     | <code>"#"</code>       | A character string. Knitr will append the string to the start of each line of results in the final document.                                                                                                                                                                                                            |
| <code>highlight</code>   | <code>TRUE</code>      | If <code>TRUE</code> , knitr will highlight the source code in the final output.                                                                                                                                                                                                                                        |
| <code>prepost</code>     | <code>FALSE</code>     | If <code>TRUE</code> , knitr will add <code>&gt;</code> to the start of each line of code displayed in the final document.                                                                                                                                                                                              |
| <code>size</code>        | <code>"normal"</code>  | Fontsize for LaTeX output.                                                                                                                                                                                                                                                                                              |
| <code>strip.white</code> | <code>TRUE</code>      | If <code>TRUE</code> , knitr will remove white spaces that appear at the beginning or end of a code chunk.                                                                                                                                                                                                              |
| <code>tidy</code>        | <code>FALSE</code>     | If <code>TRUE</code> , knitr will tidy code chunks for display with the <code>tidy_source()</code> function in the <code>formatt</code> package.                                                                                                                                                                        |

RStudio

Updated 10/30/2014

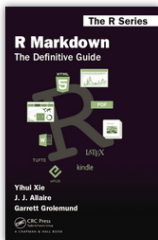
© 2014 RStudio, Inc. All rights reserved.

## R Markdown: The Definitive Guide

by Yihui Xie, J. J. Allaire, Garrett Golemund

2018-09-11

Star 239



The first official book authored by the core R Markdown developers that provides a comprehensive and accurate reference to the R Markdown ecosystem. With R Markdown, you can easily create reproducible data analysis reports, presentations, dashboards, interactive applications, books, dissertations, websites, and journal articles, while enjoying the simplicity of Markdown and the great power of R and other languages. *Read more* →

<https://bookdown.org/yihui/rmarkdown/>

<https://bookdown.org/yihui/rmarkdown-cookbook/>

# Welcome to Quarto

Quarto<sup>®</sup> is an open-source scientific and technical publishing system built on [Pandoc](#)

- Create dynamic content with [Python](#), [R](#), [Julia](#), and [Observable](#).
- Author documents as plain text markdown or [Jupyter](#) notebooks.
- Publish high-quality articles, reports, presentations, websites, blogs, and books in HTML, PDF, MS Word, ePub, and more.
- Author with scientific markdown, including equations, citations, crossrefs, figure panels, callouts, advanced layout, and more.

<https://quarto.org/>

Your turn

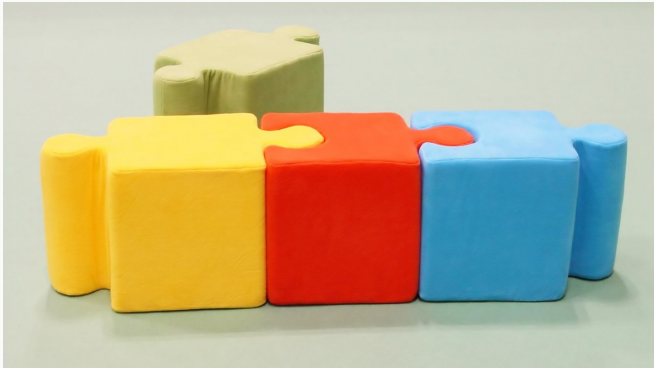
---

- Try visual markdown editor
- Add bibliography
- Try templates (rticles, rmdTemplates)
- Parameterised reports (e.g. different iris or penguin species)

# Workflow management

---

In complex projects we must **keep pieces organised**



## makefile runs all code in right order

makefile.R

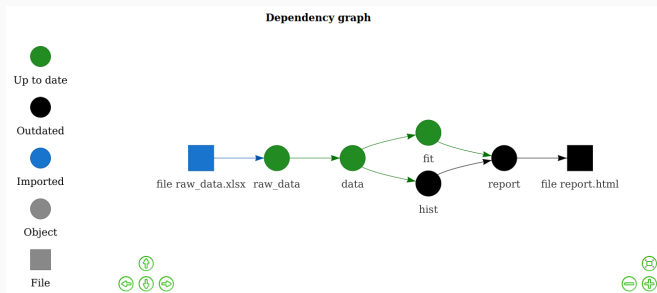
```
source("clean_data.R")
```

```
source("fit_model.R")
```

```
render("report.Rmd")
```



# targets: advanced workflow management



<https://docs.ropensci.org/targets/>

Your turn

---

Write makefile.R for your project

Try `targets` minimal example

<https://github.com/wlandau/targets-minimal>

## Controlling software dependencies

---



Package changes can break your analysis

How to reproduce your analysis in a year,  
or different computer?

# sessionInfo records OS & used packages

```
sessionInfo()
```

```
R version 4.2.0 (2022-04-22)
```

```
Platform: x86_64-pc-linux-gnu (64-bit)
```

```
Running under: Ubuntu 20.04.4 LTS
```

```
Matrix products: default
```

```
BLAS: /usr/lib/x86_64-linux-gnu/blas/libblas.so.3.9.0
```

```
LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.9.0
```

```
locale:
```

```
[1] LC_CTYPE=en_GB.UTF-8      LC_NUMERIC=C
[3] LC_TIME=es_ES.UTF-8      LC_COLLATE=en_GB.UTF-8
[5] LC_MONETARY=es_ES.UTF-8  LC_MESSAGES=en_GB.UTF-8
[7] LC_PAPER=es_ES.UTF-8     LC_NAME=C
[9] LC_ADDRESS=C             LC_TELEPHONE=C
[11] LC_MEASUREMENT=es_ES.UTF-8 LC_IDENTIFICATION=C
```

```
attached base packages:
```

```
[1] stats      graphics  grDevices  utils      datasets  methods    base
```

```
other attached packages:
```

```
[1] report_0.5.5      equatiomatic_0.3.1 xtable_1.8-4      knitr_1.40
```

```
loaded via a namespace (and not attached):
```

```
[1] Rcpp_1.0.9          mvtnorm_1.1-3      lattice_0.20-45   tidyr_1.2.0
[5] zoo_1.8-10          assertthat_0.2.1   digest_0.6.29     utf8_1.2.2
[9] mime_0.12           R6_2.5.1           backports_1.4.1   evaluate_0.16
[13] coda_0.19-4         pillar_1.8.1       rlang_1.0.5       multcomp_1.4-20
[17] performance_0.9.2  rstudioapi_0.14    Matrix_1.4-1      effectsize_0.7.0.5
[21] rmarkdown_2.16     splines_4.2.0      stringr_1.4.1     shiny_1.7.2
[25] broom_1.0.1         compiler_4.2.0     httpuv_1.6.5      xfun_0.32
[29] pkgconfig_2.0.3    parameters_0.18.2  htmltools_0.5.3   insight_0.18.2
[33] tidycselect_1.1.2  tibble_3.1.8       codetools_0.2-18  fansi_1.0.3
[37] tidyr_1.2.0        MASS_7.3-55        bit_4.0.4         bit64_4.0.5
[41] R6_2.5.1           Rcpp_1.0.9         mvtnorm_1.1-3     lattice_0.20-45
[45] zoo_1.8-10          assertthat_0.2.1   digest_0.6.29     utf8_1.2.2
[49] mime_0.12           R6_2.5.1           backports_1.4.1   evaluate_0.16
[53] coda_0.19-4         pillar_1.8.1       rlang_1.0.5       multcomp_1.4-20
[57] performance_0.9.2  rstudioapi_0.14    Matrix_1.4-1      effectsize_0.7.0.5
[61] rmarkdown_2.16     splines_4.2.0      stringr_1.4.1     shiny_1.7.2
[65] broom_1.0.1         compiler_4.2.0     httpuv_1.6.5      xfun_0.32
[69] pkgconfig_2.0.3    parameters_0.18.2  htmltools_0.5.3   insight_0.18.2
[73] tidycselect_1.1.2  tibble_3.1.8       codetools_0.2-18  fansi_1.0.3
[77] tidyr_1.2.0        MASS_7.3-55        bit_4.0.4         bit64_4.0.5
```



## checkpoint reconstructs packages in given date

```
library('checkpoint')  
  
checkpoint("2019-10-08")  
  
source("analysis.R")
```

1. Detects packages used
2. Installs version from given date (only CRAN)
3. Independent install (not messing w/ main library)

# automagic records & install packages (CRAN + GitHub)

```
automagic::make_deps_file()
```

File `deps.yaml` records dependencies:

```
- Package: equatiomatic  
  Repository: CRAN  
  Version: 0.1.0  
  
- Package: report  
  GithubUsername: easystats  
  GithubRepo: report  
  GithubRef: HEAD  
  GithubSHA1: c48a4bb0a40df7116bc502aa3ce2cbbc9d70b7e2
```

To install all those dependencies:

```
automagic()
```

## renv also controls dependencies

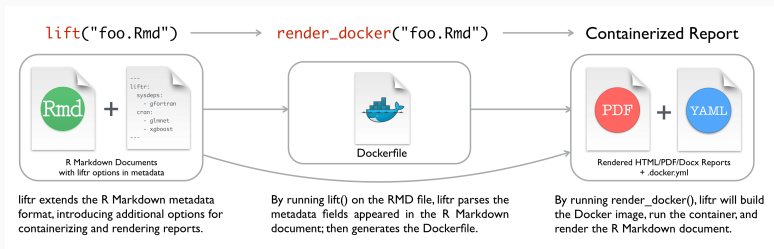
```
renv::init()  
# Create private package library for project  
  
renv::snapshot()  
# Capture dependencies in lockfile  
  
renv::restore()  
# Regenerate dependencies from lockfile
```

<https://environments.rstudio.com/>

To ensure reproducibility,  
besides R packages  
we also need to control  
**computational environment**

**Docker** recreates virtual systems  
from a **Dockerfile**

# liftr: process Rmd in Docker container



<https://liftr.me/>

## containerit creates Dockerfile

```
library("containerit")  
  
dockfile <- dockerfile(from = "mypaper.Rmd")
```

<https://o2r.info/containerit>

# holepunch: reproduce analysis in the cloud (Binder)

**BAM!**

**OMG!**

**binder**

Starting repository: karthik/friday-test/master

How to Binder? Check out the [Quickstart](#) for more information.

README.md

### Example repo for...

This repository is an example repository for the `statkit` package for the `statkit` package.

To test `holepunch`, follow these steps:

1. Click Use this template to the right.
2. Give this repo a new name and create a new repo in your account
3. Click the Clone or download button, copy the URL.
4. In RStudio Desktop, click the Project drop down on the top right, Choose **New Project > Version Control > Git**, and paste in the URL of your new GitHub repository

```
R version 3.6.3 (2019-10-25) -- "Vibrating of a Tree"
signature(22) 2020 Thu R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a distributed system with some parallelization capabilities.
Type 'parallel()' for more information and 'enable()' or 'disable()' for how to take full advantage of your hardware.

See 'help()' for more details, 'help("R")' for an online help, or
'help.search()' for an online search interface to help.
Type '?()' to see it.
```

| File                 | Size   | Modified               |
|----------------------|--------|------------------------|
| ./                   |        |                        |
| ./data               | 20 B   | Jan 24, 2018, 13:03 PM |
| ./R                  | 2.80 B | Jan 24, 2018, 13:03 PM |
| ./statkit            | 200 B  | Jan 24, 2018, 13:03 PM |
| ./statkit.Rproj      | 500 B  | Jan 24, 2018, 13:03 PM |
| ./statkit.Rproj.user | 547 B  | Jan 24, 2018, 13:03 PM |
| ./statkit.Rproj.user | 1.7 KB | Jan 24, 2018, 13:03 PM |

<https://karthik.github.io/holepunch/>



Your turn

---

- Create script/Rmd using different packages
- Call `checkpoint` on former date

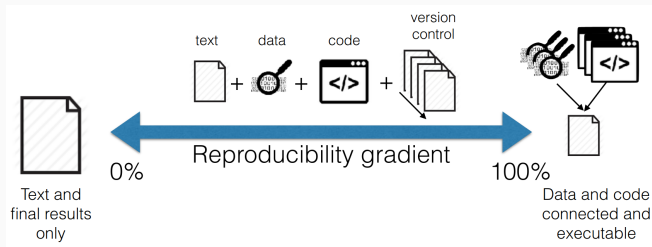
- Record dependencies:
  - `automagic::make_deps_file()`
  - `renv::snapshot`
- Recreate packages
  - `automagic()`
  - `restore()`

## How to write more reproducible code

- [BES guide to reproducible code](#)
- [Turing Way](#)
- [Good enough practices in scientific computing](#)
- [Ciencia reproducible: qué, por qué, cómo](#)
- <https://rstats.wtf>
- `fertile` package
- [CodeCheck](#)

# Reproducibility

- Good for you, good for science
- Requires systemic changes
- Reproducibility gradient: step by step



Happy collaboration!

