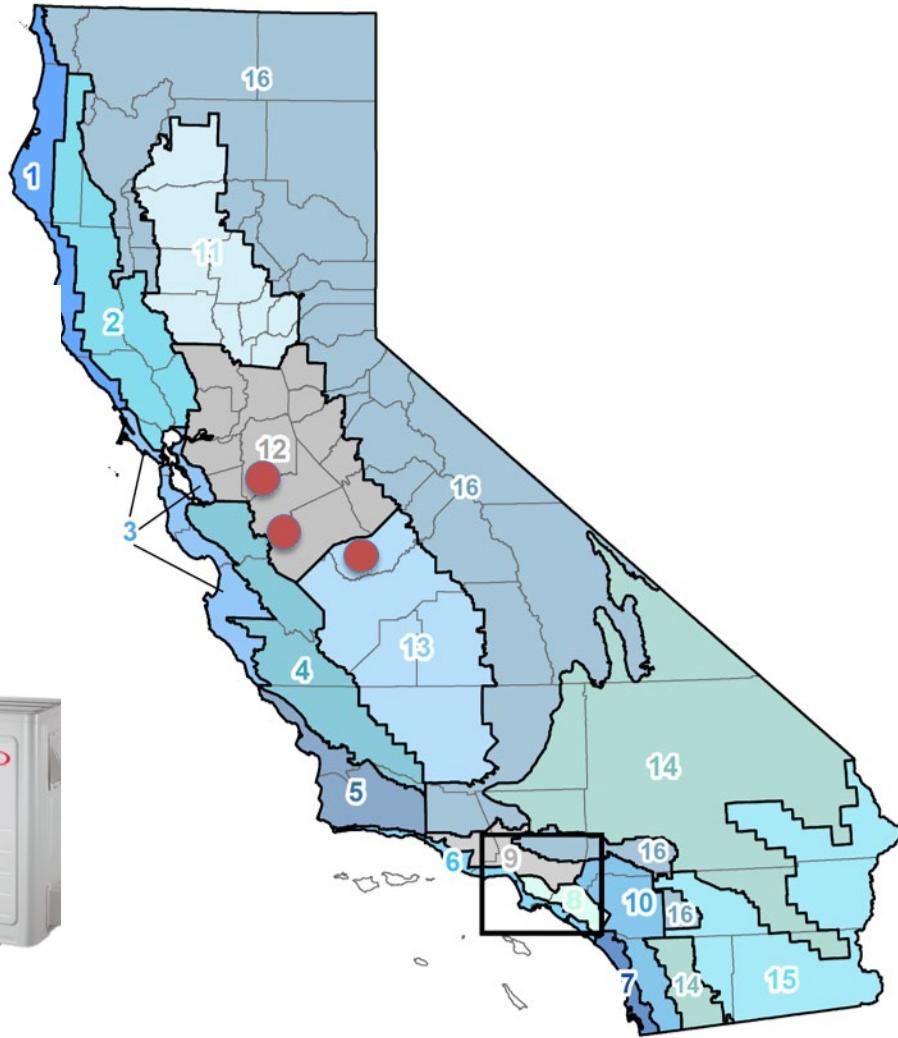


Towards more reproducible + open research in building science: why, what, and some tools for how

Dana Miller
September 27, 2018
Arch 241

Current research – field study

In existing buildings,
can high-efficiency ceiling fans
with automatic controls
reduce energy consumption
while maintaining comfort?



Goals: after this class, able to answer:

- Why is reproducible science important?
- What does reproducible mean anyway?
- Why is this hard?
- What tools can we start using now?
- Hands-on practice
 - How to create a project in RStudio?
 - (If time) What is R Markdown?

US edition

The Guardian

Attempt to replicate major social scientific findings of past decade fails

Scientists and the design of experiments under scrutiny after a major project fails to reproduce results of high profile studies



<https://www.theguardian.com/science/2018/aug/27/attempt-to-replicate-major-social-scientific-findings-of-past-decade-fails>

Mother Jones

This Cornell Food Researcher Has Had 13 Papers Retracted. How Were They Published in the First Place?

The social sciences have a problem.

KIERA BUTLER SEPTEMBER 25, 2018 6:00 AM

<https://www.motherjones.com/food/2018/09/cornell-food-researcher-brian-wansink-13-papers-retracted-how-were-they-published/>

Is Economics Research Replicable? Sixty Published Papers from Thirteen Journals Say “Usually Not”

<https://www.federalreserve.gov/econresdata/feds/2015/files/2015083pap.pdf>

Hi Ozge,

Glad you had a chance to take an initial look at the data.

I don't think I've ever done an interesting study where the data "came out" the first time I looked at it. The interesting stories come from seeing when things -- like the 1/2 price buffet -- works and when it doesn't.

I would like you to really dig into this to find a number of situations or people for which this relationship does hold -- that is where the 1/2 price buffet did result in a difference.

Here's some things to do.

Males

Females

Lunch goers

Dinner goers

People sitting alone

People eating with groups of 2

People eating in groups of 2+

People who order alcohol

People who order soft drinks

People who sit close to buffet

People who sit far away

and so on . . .

Second, think of all the different ways you can cut the data and analyze subsets of it to see when this relationship holds. For instance, if it works on men but not women, we have a moderator. Here are some groups you'll want to break out separately:

Is the evidence for austerity based on an Excel spreadsheet error?



“Growth in a time of debt”
Reinhart and Rogoff 2011

Thomas Herdon, PhD student
UMass Amherst

<https://www.washingtonpost.com/news/wonk/wp/2013/04/16/is-the-best-evidence-for-austerity-based-on-an-excel-spreadsheet-error>

<https://www.bbc.com/news/magazine-22223190>

<http://theconversation.com/the-reinhart-rogoff-error-or-how-not-to-excel-at-economics-13646>

Z	B	C	I	J	K	L	M
2			Real GDP growth Debt/GDP				
3							
4	Country	Coverage	30 or less	30 to 60	60 to 90	90 or above	30 or less
26			3.7	3.0	3.5	1.7	5.5
27	Minimum		1.6	0.3	1.3	-1.8	0.8
28	Maximum		5.4	4.9	10.2	3.6	13.3
29							
30	US	1946-2009	n.a.	3.4	3.3	-2.0	n.a.
31	UK	1946-2009	n.a.	2.4	2.5	2.4	n.a.
32	Sweden	1946-2009	3.6	2.9	2.7	n.a.	6.3
33	Spain	1946-2009	1.5	3.4	4.2	n.a.	9.9
34	Portugal	1952-2009	4.8	2.5	0.3	n.a.	7.9
35	New Zealand	1948-2009	2.5	2.9	3.9	-7.9	2.6
36	Netherlands	1956-2009	4.1	2.7	1.1	n.a.	6.4
37	Norway	1947-2009	3.4	5.1	n.a.	n.a.	5.4
38	Japan	1946-2009	7.0	4.0	1.0	0.7	7.0
39	Italy	1951-2009	5.4	2.1	1.8	1.0	5.6
40	Ireland	1948-2009	4.4	4.5	4.0	2.4	2.9
41	Greece	1970-2009	4.0	0.3	2.7	2.9	13.3
42	Germany	1946-2009	3.9	0.9	n.a.	n.a.	3.2
43	France	1949-2009	4.9	2.7	3.0	n.a.	5.2
44	Finland	1946-2009	3.8	2.4	5.5	n.a.	7.0
45	Denmark	1950-2009	3.5	1.7	2.4	n.a.	5.6
46	Canada	1951-2009	1.9	3.6	4.1	n.a.	2.2
47	Belgium	1947-2009	n.a.	4.2	3.1	2.6	n.a.
48	Austria	1948-2009	5.2	3.3	-3.8	n.a.	5.7
49	Australia	1951-2009	3.2	4.9	4.0	n.a.	5.9
50							
51				4.1	2.8	2.8	=AVERAGE(L30:L44)

Attempt to replicate major social scientific findings of past decade fails

<https://www.theguardian.com/science/2018/aug/27/attempt-to-replicate-major-social-scientific-findings-of-past-decade-fails>

Letter | Published: 27 August 2018

Evaluating the replicability of social science experiments in *Nature* and *Science* between 2010 and 2015

<https://www.nature.com/articles/s41562-018-0399-z>

“The replications follow analysis plans **reviewed by the original authors** and pre-registered prior to the replications. The replications are high powered, with **sample sizes on average about five times higher** than in the original studies. We find a significant effect in the same direction as the original study for 13 (62%) studies, and the effect size of the replications is on average about 50% of the original effect size.”

Attempt to replicate major social scientific findings of past decade fails

<https://www.theguardian.com/science/2018/aug/27/attempt-to-replicate-major-social-scientific-findings-of-past-decade-fails>

Letter | Published: 27 August 2018

Evaluating the replicability of social science experiments in *Nature* and *Science* between 2010 and 2015

<https://www.nature.com/articles/s41562-018-0399-z>

“I don’t think it’s a crisis, it’s a reformation. We’re in the midst of a dramatic increase in the rigour and transparency of research in the social sciences.”

- Brian Nosek, U Virginia (corresponding author)

Why reproducible research? Robustness

PHYSICS TODAY

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.

“Limmer maintains that he and his mentor weren’t trying to hide anything. “I had and was very willing to share the code,” he says. What he didn’t have, he says, was the time or personnel to prepare the code in a form that could be useful to an outsider”

“Suddenly it made sense that the Berkeley researchers hadn’t seen a second liquid phase; they were effectively running their simulations at temperatures well above the critical point.”

<https://physicstoday.scitation.org/do/10.1063/PT.6.1.20180822a/full/>

Why reproducible research? Records

Email from grad student to Berkeley professor in 2017:

“...The citation on the slide is [your report on a state government website].

I downloaded that PDF and do not see [important statistic] anywhere in it.
Table 3 seems to have [different value] , Table 15 has [related value] in California...

...The closest match that I've found is [other statistic] from Table 15...

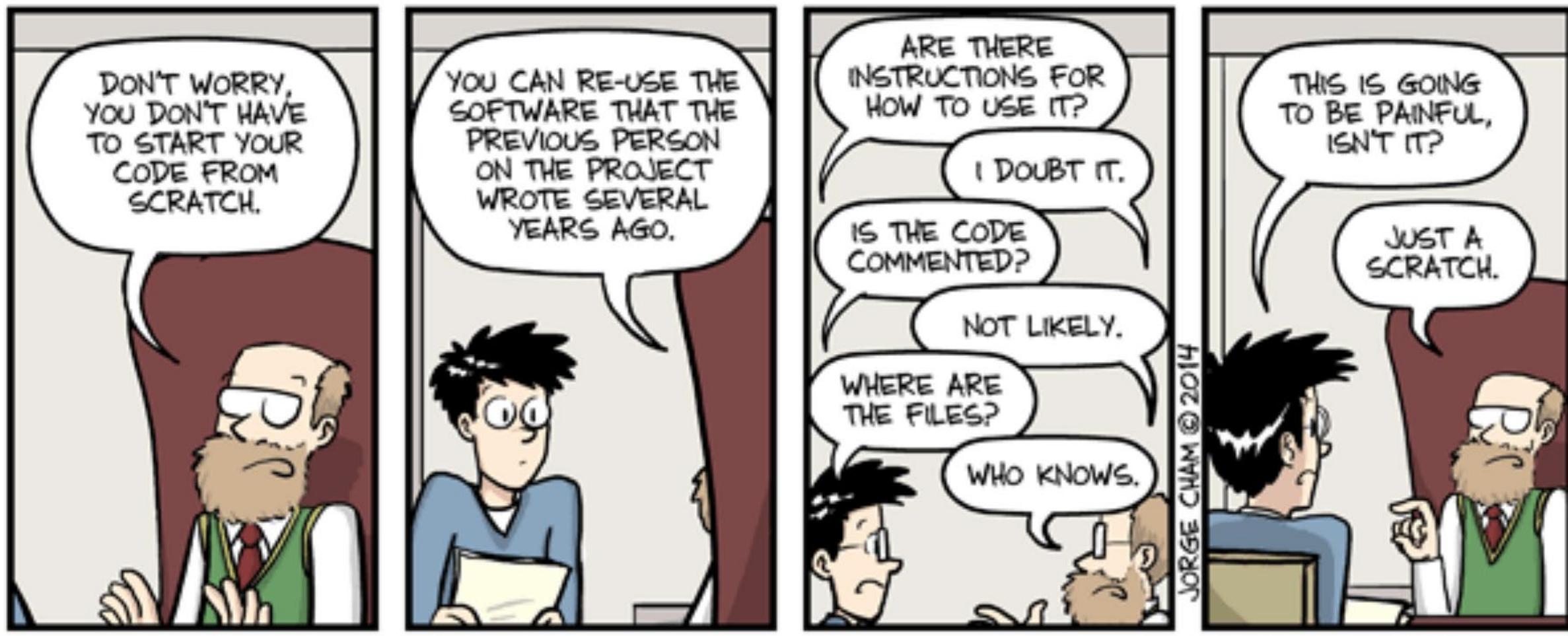
Is it possible there was an error in transcription? If not, can you help me figure out what's going on here? What does [important statistic] on [slide] represent?”

Response from Berkeley professor:

“I don't remember the details. [Former student], do you? Thanks in advance.”

[Former student] never responded

Why reproducible research? Reusability



Why is this important?

- Show evidence of the correctness of our results
- Enable others to make use of our methods and results (including future you)
- Increasingly encouraged or required by funders and journals

<http://ropensci.github.io/reproducibility-guide/sections/introduction/>

<https://www.stat.berkeley.edu/~stark/Seminars/reproNE16.htm>



<https://cos.io/our-services/open-science-badges>

“reproducibility”

What does reproducibility mean? Three ideas to distinguish

1. *Reproducibility* (sometimes called repeatability or even ‘preproducibility’) [6]

- *Data, and/or statistical methods, code, software environment are available so data analysis can be repeated and get a similar result [1, 2, 3]*

2. *Replicability* (sometimes called reproducibility)

- *Another researcher is able to conduct the same experiment and get similar results*

3. *Correctness*

- *Is your analysis appropriate? Is your hypothesis correct?
→ Reproducible research can still be wrong! [4]*

1 - *Peng Science* 2011, DOI: 10.1126/science.1213847

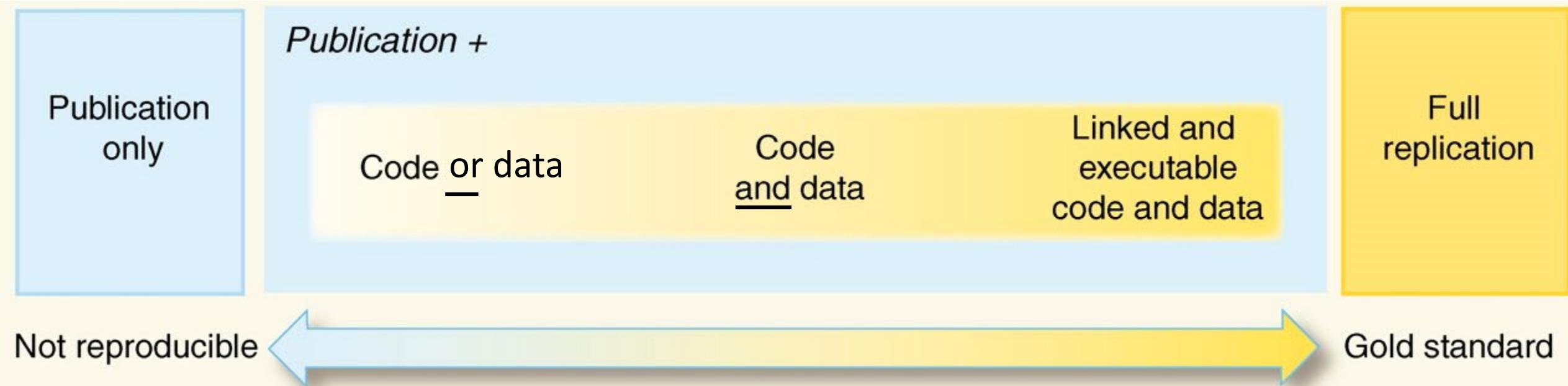
2 - <http://ropensci.github.io/reproducibility-guide/sections/introduction/>

3 - <https://simplystatistics.org/2014/06/06/the-real-reason-reproducible-research-is-important/>, Peng 2014

4 - Reproducible research can still be wrong, Leek, Peng
PNAS 2015, DOI: 10.1073/pnas.1421412111

5- <https://www.nature.com/news/muddled-meanings-hamper-efforts-to-fix-reproducibility-crisis-1.20076/#/b1>, citing <https://www.ascb.org/wp-content/uploads/2015/11/How-can-scientist-enhance-rigor.pdf> citing Schmidt, Stefan. (2009). Shall We Really Do It Again? The Powerful Concept of Replication Is Neglected in the Social Sciences. *Review of General Psychology*. 13. 90-100. 10.1037/a0015108.
6 – Barba 2018 <https://arxiv.org/pdf/1802.03311.pdf>

A spectrum of practice (depends on project)



Adopted from:

Reproducible Research in Computational Science

By Roger D. Peng

Science 02 Dec 2011 : 1226-1227

<http://science.sciencemag.org/content/334/6060/1226.full>

Why might reproducibility be hard?

Stodden (2010) Survey of NIPS:

Code	Complaint/Excuse	Data
77%	Time to document and clean up	54%
52%	Dealing with questions from users	34%
44%	Not receiving attribution	42%
40%	Possibility of patents	N/A
34%	Legal Barriers (ie. copyright)	41%
N/A	Time to verify release with admin	38%
30%	Potential loss of future publications	35%
30%	Competitors may get an advantage	33%
20%	Web/disk space limitations	29%

Philip Stark,
[https://www.stat.berkeley.edu/~stark/
Seminars/reproNE16.htm](https://www.stat.berkeley.edu/~stark/Seminars/reproNE16.htm)

I.e., fear, greed, ignorance, & sloth.

What's possible – interactive publications!

New paper (open access)

The screenshot shows the full text of the journal article. It includes the title, authors, abstract, keywords, and the main text. The abstract discusses the concept of noise as a nuisance in ecological studies and how it can generate novel phenomena and reveal information. The text is presented in a standard academic format with headings and footnotes.

<https://onlinelibrary.wiley.com/doi/epdf/10.1111/ele.13085>

Corresponding “compendium” (with text, code, data)

The screenshot shows a GitHub repository page for 'noise-phenomena'. It displays the repository's statistics (134 commits, 2 branches, 4 releases, 1 contributor) and a list of files. Below this is a preview of the 'noise-phenomena compendium' page, which contains a summary of the repository's purpose, links to various files, and a preview of the R code for the Gillespie algorithm.

<https://github.com/cboettig/noise-phenomena>

Interactive R notebook to explore code (with all required files and software ready to run in your browser with no installation, thanks to Binder)

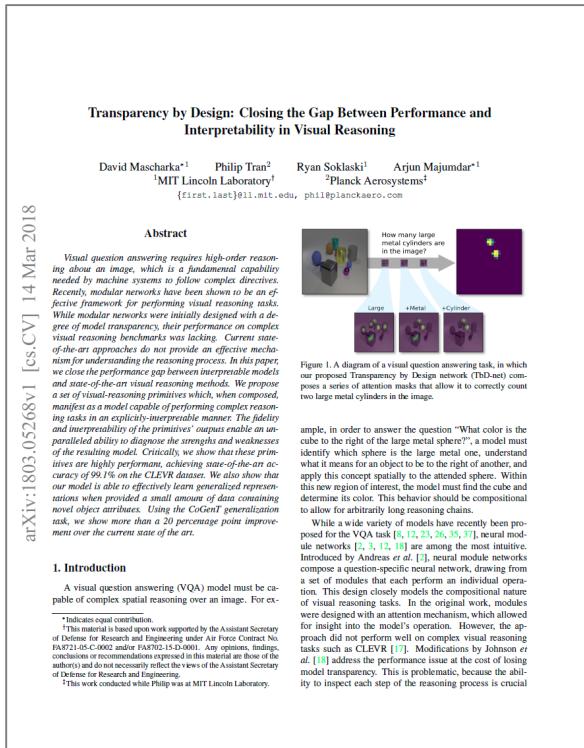
The screenshot shows an RStudio session running on Binder. The top navigation bar includes 'File', 'Edit', 'Code', 'View', 'Plots', 'Session', 'Build', 'Debug', 'Profile', 'Tools', and 'Help'. The main area shows an R notebook with code for the Gillespie algorithm. The code is as follows:

```
dn = cn (1 - n/N) - en
      birth rate, b(n)   death rate, d(n)
dn / dt = cn (1 - n/N) - en
      birth rate, b(n)   death rate, d(n)
```

The notebook also includes a mathematical equation for the derivative of the state variable n with respect to time t. The bottom pane shows the file tree and a list of files, including 'AppendixA.Rmd', 'AppendixA.pdf', 'Dai-Figure.R', 'elsarticle.cls', 'gillespie.csv', 'inflation.csv', 'noisy_switch.csv', 'quasicycles.csv', 'silent_tipping.csv', and 'tipping.csv'.

Interactive publications – in python too!

New paper
(in this case, a preprint)



<https://arxiv.org/abs/1803.05268>

Website explaining paper

(has code to replicate experiments + plots and explicit software versions)

The image shows a GitHub repository page for "Transparency-by-Design networks (TbD-nets)". It includes a "Launch Binder" button, Python 3.5.3.6, pytorch 0.2.0.3 0.41, and 80% completion. The repository description is "Transparency by Design: Closing the Gap Between Performance and Interpretability in Visual Reasoning". It lists authors: David Mascharka, Philip Tran, Ryan Soklaski, Arjun Majumdar. The repository contains code for replicating experiments and visualizations from the paper. A diagram shows a visual reasoning task with attention masks for "Large", "+Metal", and "+Cylinder". A note says "if you find this code useful in your research, please cite". The code section contains the paper's citation information.

<https://github.com/davidmascharka/tbd-nets>

+ Interactive Jupyter notebook with code to apply their method to your own images

(with all required files and software ready to run in your browser with no installation required, thanks to Binder)

The image shows a Jupyter notebook titled "jupyter full-vqa-example (unsaved changes)". It has tabs for File, Edit, View, Insert, Cell, Kernel, Widgets, Help, and a "Not Trusted" warning. The notebook content describes the VQA pipeline using TbD-Net. It loads models and generates programs. A code cell shows the VQA task: "How many large metal cylinders are in the image?". Another cell shows the generated program: "from tbdd_nets import Path from tbdd_nets import load_tbdd_net from tbdd_nets import load_tbdd_program_generator, generate_single_program from tbdd_nets import download_pretrained_models, download from tbdd_nets import extract_features, load_feature_extractor, extract_image_feats". The notebook concludes with a note about loading all models if not present.

<https://mybinder.org/v2/gh/davidmascharka/tbd-nets/binder?filepath=full-vqa-example.ipynb>

Additional reference: [Link](#) to tweet by first author publicizing this work

Demo – reproduce a figure from paper on the previous slide

See the paper: <https://onlinelibrary.wiley.com/doi/epdf/10.1111/ele.13085>

- Note the PDF includes a link to a GitHub repo

2. Inspect the code

- Open the associated GitHub repository containing the code and data used to produce the figures for the paper
 - Repo: <https://github.com/cboettig/noise-phenomena>
 - Code for the paper: <https://github.com/cboettig/noise-phenomena/blob/master/paper/paper.Rmd>

3. Reproduce a figure without installing anything on your local computer!

- Click the  button on the main page of the repo

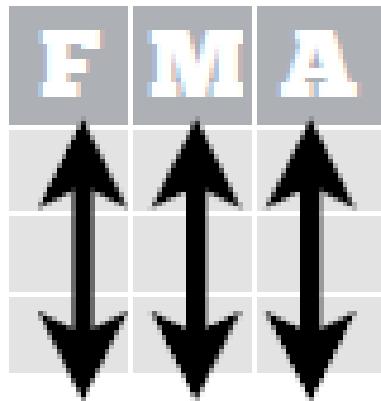
What are tools and practices
we can start using now?

Data collection

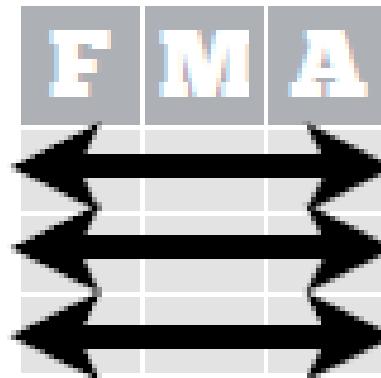
1. Back up your data regularly



2. “Tidy” spreadsheet formats



&



Each **variable** is saved
in its own **column**

Each **observation** is
saved in its own **row**

Is this a tidy dataset?

Subject	Thermal sensation votes
Person1	0.5, 0.7
Person2	-0.3
Person3	1.2, 1.5, 2



No, has more than one observation per row

Ok, is it a tidy dataset now?

Subject	TS_vote_1	TS_vote_2	TS_vote_3
Person1	0.5	0.7	
Person2	-0.3		
Person3	1.2	1.5	2



No, has more than one observation per row,
and more than one column per type of information

A tidy dataset

Subject	ThermalSensationVote
Person1	0.5
Person1	0.7
Person2	-0.3
Person3	1.2
Person3	1.5
Person3	2



Is this weather data tidy?

id	year	month	element	d1	d2	d3	d4	d5	d6	d7	d8
MX17004	2010	1	tmax	—	—	—	—	—	—	—	—
MX17004	2010	1	tmin	—	—	—	—	—	—	—	—
MX17004	2010	2	tmax	—	27.3	24.1	—	—	—	—	—
MX17004	2010	2	tmin	—	14.4	14.4	—	—	—	—	—
MX17004	2010	3	tmax	—	—	—	—	—	32.1	—	—
MX17004	2010	3	tmin	—	—	—	—	—	14.2	—	—
MX17004	2010	4	tmax	—	—	—	—	—	—	—	—
MX17004	2010	4	tmin	—	—	—	—	—	—	—	—
MX17004	2010	5	tmax	—	—	—	—	—	—	—	—
MX17004	2010	5	tmin	—	—	—	—	—	—	—	—



No, has more than one observation per row

What improved?

Wickham, H. (2014). Tidy Data. *Journal of Statistical Software*, 59(10), 1 - 23. doi:<http://dx.doi.org/10.18637/jss.v059.i10>

id	date	tmax	tmin
MX17004	2010-01-30	27.8	14.5
MX17004	2010-02-02	27.3	14.4
MX17004	2010-02-03	24.1	14.4
MX17004	2010-02-11	29.7	13.4
MX17004	2010-02-23	29.9	10.7
MX17004	2010-03-05	32.1	14.2
MX17004	2010-03-10	34.5	16.8
MX17004	2010-03-16	31.1	17.6
MX17004	2010-04-27	36.3	16.7
MX17004	2010-05-27	33.2	18.2

(b) Tidy data



What might be a problem in these spreadsheets?

A

	A	B	C
1	id	date	glucose
2	101	2015-06-14	149.3
3	102		95.3
4	103	2015-06-18	97.5
5	104		117.0
6	105		108.0
7	106	2015-06-20	149.0
8	107		169.4

B

	A	B	C	D	E
1	strain	genotype	min	replicate	response
2	A	normal	1	1	147
3	A	normal	1	2	139
4	B	normal	1	1	246
5	B	normal	1	2	240
6	A	mutant	1	1	166
7	A	mutant	1	2	179
8	B	mutant	1	1	178
9	B	mutant	1	2	172
10	A	normal	5	1	334
11	A	normal	5	2	354
12	B	normal	5	1	514
13	B	normal	5	2	611
14	A	mutant	5	1	451
15	A	mutant	5	2	474
16	B	mutant	5	1	412
17	B	mutant	5	2	447

How would you enter this data into R?

Karl W. Broman & Kara H. Woo (2018) Data Organization in Spreadsheets, *The American Statistician*, 72:1, 2-10, DOI: 10.1080/00031305.2017.1375989

	A	B	C	D	E	F	G
1							
2	Date	11/3/14					
3	Days on diet	126					
4	Mouse #	43					
5	sex	f					
6	experiment		values		mean	SD	
7	control		0.186	0.191	1.081	0.49	0.52
8	treatment A		7.414	1.468	2.254	3.71	3.23
9	treatment B		9.811	9.259	11.296	10.12	1.05
10							
11	fold change		values		mean	SD	
12	treatment A		15.26	3.02	4.64	7.64	6.65
13	treatment B		20.19	19.05	23.24	20.83	2.17



How would you enter this data into R?

A

	A	B	C
1	id	date	glucose
2	101	2015-06-14	149.3
3	102	2015-06-14	95.3
4	103	2015-06-18	97.5
5	104	2015-06-18	1.1
6	105	2015-06-18	108.0
7	106	2015-06-20	149.0
8	107	2015-06-20	169.4



B

	A	B	C	D
1	id	date	glucose	outlier
2	101	2015-06-14	149.3	FALSE
3	102	2015-06-14	95.3	FALSE
4	103	2015-06-18	97.5	FALSE
5	104	2015-06-18	1.1	TRUE
6	105	2015-06-18	108.0	FALSE
7	106	2015-06-20	149.0	FALSE
8	107	2015-06-20	169.4	FALSE



More on tidy data and spreadsheet organization

Karl W. Broman & Kara H. Woo (2018) Data Organization in Spreadsheets, *The American Statistician*, 72:1, 2-10, DOI: [10.1080/00031305.2017.1375989](https://doi.org/10.1080/00031305.2017.1375989)

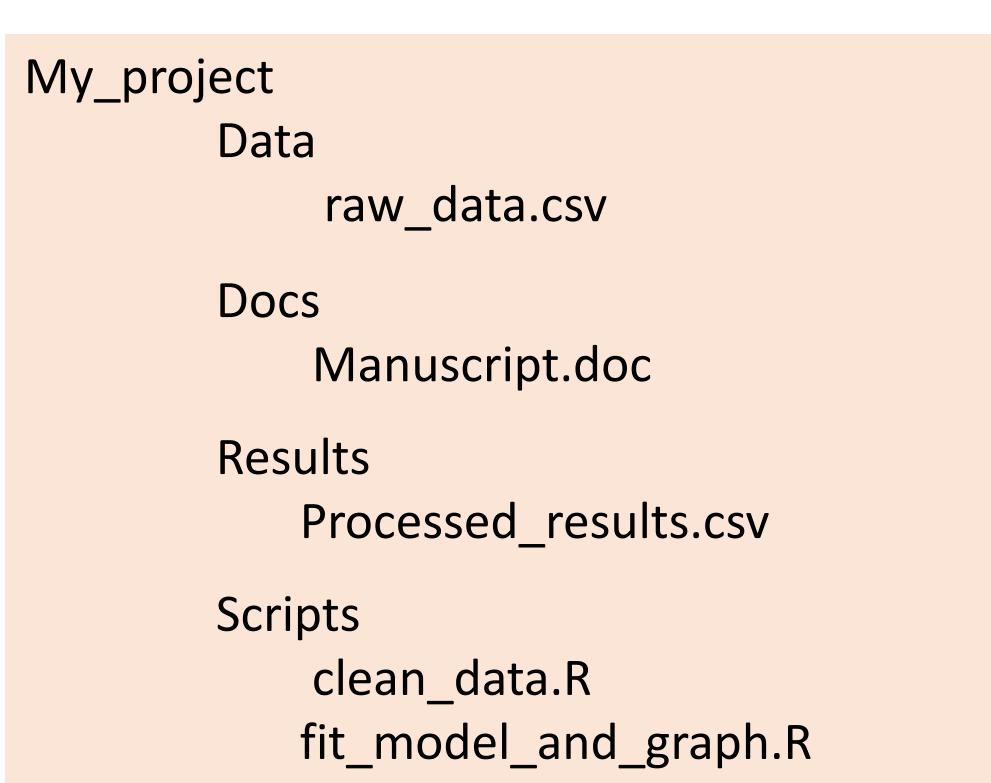
Wickham, H. (2014). Tidy Data. *Journal of Statistical Software*, 59(10), 1 - 23.
doi:<http://dx.doi.org/10.18637/jss.v059.i10>

- Be consistent
- Write dates as YYYY-MM-DD
- No empty cells
- Put just one thing in a cell
- Make it a rectangle
- Create a data dictionary
- No calculations in the raw data files
- Font color or highlighting isn't data
- Make backups
- Use data validation to avoid errors
- Save the data in plain text files

Data analysis

1. File organization

- Suggest one folder per project, eg:



Resources:

- [A Quick Guide to Organizing Computational Biology Projects](http://happygitwithr.com/)
- [Good enough practices in scientific computing](#)

2. Version control (for code)

- Like MS Word's Previous Versions → Restore, or unlimited undo for your R scripts

Subject	Author
HEAD -> refs/heads/master (origin/master) Include plos PDF	Dana <d>
Reformat and tidy setup notes	Dana <d>
Update readme	Dana <d>
added a line for demo	Dana <d>
Add background to readme	Dana <d>
Add .Renviron file	Dana <d>

Resources:

- <http://happygitwithr.com/>
- [A quick introduction to version control with Git and GitHub](#)

Sharing + publication of results

1. Data repositories

- UC system – [dash](#)
- Example: [ASHRAE Global Thermal Comfort Database II](#)

Resources:

- [Berkeley Research Data Management librarians!](#)
- Other repositories include Zenodo, Open Science Framework, DataDryad, Figshare

2. Code repositories

The screenshot shows a GitHub organization page for 'CenterForTheBuiltEnvironment'. At the top, there's a header with a search bar, a 'Pull requests' button, an 'Issues' button, a 'Marketplace' button, and an 'Explore' button. Below the header is the organization's logo, which is a stylized blue and white graphic next to the letters 'CBE'. To the right of the logo, the organization's name 'CenterForTheBuiltEnvironment' is written in a sans-serif font. Underneath the name, it says 'Berkeley, CA' and provides a link 'http://cbe.berkeley.edu'. Below this section, there are four tabs: 'Repositories 25', 'People 39', 'Teams 6', and 'Projects 0'. Further down, there are search and filter options: 'Search repositories...', 'Type: All', and 'Language: All'. The main content area displays three repository cards. The first card is for 'radiant' (Private), described as a 'Web tool for the design and operation of high thermal mass radiant systems'. It was last updated 7 minutes ago and is written in JavaScript. The second card is for 'epic-fans-energy' (Private), described as 'M&V of field study energy data'. It was last updated 8 days ago. The third card is for 'epic_fans_field_studies' (Private), described as 'R'. It was last updated 8 days ago. To the right of the repository cards, there's a 'Top languages' section showing JavaScript, Python, R, HTML, and TeX, each represented by a colored circle. Below that is a 'People' section showing a grid of 39 user profiles.

Resources:

- [Software Carpentry's Intro to Version Control for Novices](#)
- [A quick introduction to version control with Git and GitHub](#)

Demo – make a new project in RStudio

1. Open RStudio (this is a good time for questions)

2. File → New project → Choose directory → New project

(Detailed instructions: RStudio ["Using Projects"](#))

3. What are benefits of Rstudio Projects?

- All files for one project in one place
- Easily switch between projects
- Starts you in a fresh R session each time
- Integration with git (version control)

4. If time allows:

- R Markdown files
- Look at example of paper being written in Markdown
- Using git inside RStudio

R Markdown documents

R Markdown (.Rmd)
document mixing
code and text

Code can be executed in-line
(above) or in console (below)

See file explorer, help documentation, plots, packages

Convert file to .md, .doc, .pdf, html, html notebook, slides,
book, blog, dashboards, Shiny app, etc

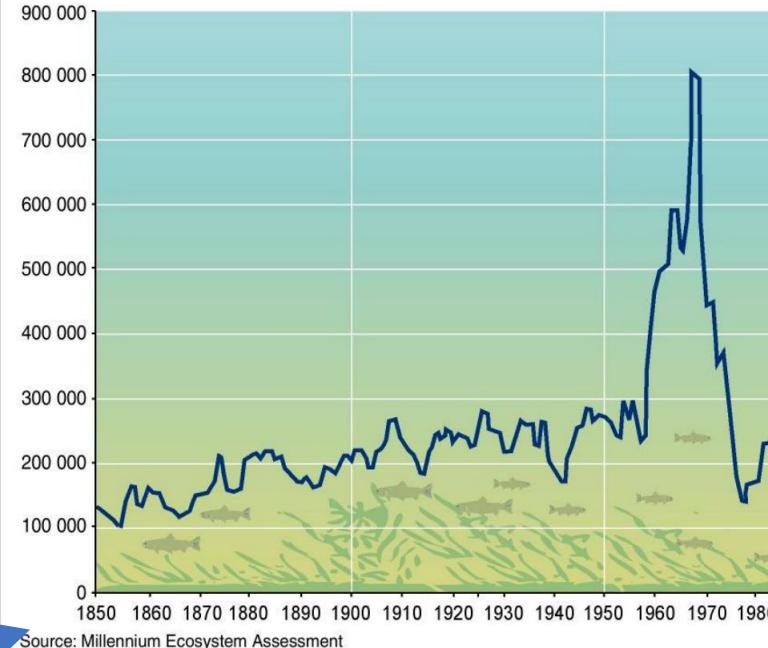
See R Markdown [output formats](#)

The screenshot shows the RStudio interface. On the left, the code editor displays an R Markdown (.Rmd) file with both R code and explanatory text. In the middle, the terminal window shows the execution of the R code, displaying the results of data joins and the resulting data frame. On the right, the file explorer shows the project structure, including files like 'fish-assignment.Rmd' and 'fish-assignment.html'. A blue arrow points from the explanatory text 'Code can be executed in-line (above) or in console (below)' to the terminal window.

Exercise 1: Investigating the North-Atlantic Cod

First, we seek to replicate the following figure from the Millennium Ecosystem Assessment Project using the RAM data.

Fish landings in tons



Task 1: Joining the necessary data

To replicate this plot, we need a table with the following columns: "country", "ssb_unit", "catch_landings_unit", "scientificname", "commonname", "year", "ssb", and "TC".

Using the `select()` and `join()` functions you were introduced to in in Module 1, build a tidy table with the desired columns.

Suggestion: separate blocks within each category added first, e.g. timeseriesunits <- select(ts=TC....)

Code to look at categories of all the data tables

Version control
tracker
See [happygitwithr](#)

Literate programming with R Markdown

Starting with one
.Rmd file with both
R code and text

This button or
command:



executes all R
code
(output shown in
console)



A screenshot of the RStudio interface. The top menu bar shows 'File Edit Code View Plots Session Build Debug Profile Tools Help'. Below the menu is a toolbar with icons for file operations, search, and knit. The main workspace shows an R Markdown file named 'fish-assignment.Rmd'. The code includes a title, author information, and a section about fisheries collapse. The 'Console' tab at the bottom shows the execution of the code, including library imports like tidyverse, readxl, and scales, and the loading of the RAM Legacy Stock Assessment Database. The 'Terminal' tab is also visible.

This button or command:



renders formatted document
with executed code and text
(.md, .pdf, .doc etc)

Overfishing Exercise

Chelsea Andreozzi and Dana Miller

Unit 3: Fisheries Collapse Module

This module will focus on understanding and replicating fisheries stock assessment data and fisheries collapse.

The Database

We will use data from the RAM Legacy Stock Assessment Database

First, load in the necessary libraries. Note that this time we need a package we haven't used before `readxl`. This package is useful for reading in .xls or .xlsx files. As always if you want more info on a package run `?readxl` after loading it.

```
library("tidyverse")
library("readxl")
library("scales") # For y-axis labels not in scientific notation - is there a better way to do this since 2012?
```

Reading in the tables

```
download.file("https://depts.washington.edu/ramlegac.wordpress/databaseVersions/RLSADB_v3.0_(assessment_data_only).zip")
path <- unzip("ramlegacy.zip") # unzip the .xls files
sheets <- readxl::excel_sheets(path) # use the readxl package to identify sheet names

# purrr::map is the tidyverse version of lapply
ram <- lapply(sheets, readxl::read_excel, path = path) # read the data from all 3 sheets into a list
names(ram) <- sheets # give the list of datatables their assigned sheet names

## check the names
names(ram)
```

What next?

UCB Prof Phillip Stark on reproducibility:

“Science is *show me*, not *trust me*”

“Perfection is impossible but improvement is easy”

There are many communities who can help!

On campus: BIDS [Data Analysis Tool Series](#), [D-Lab](#), [Librarians](#), [Research Computing](#), [Software Carpentry Workshops...](#)

Beyond: R for Data Analysis (R4DS), RBloggers, PyData (also has R videos), RLadies, #rstats...

And we can learn a lot from each other ☺

Some more R resources

R for Data Science:
<https://r4ds.had.co.nz>

From Excel to R:
<http://rpubs.com/acolumbus/how-to-use-r-with-excel>
(includes list of common
Excel functions in R)

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

Free resources aimed at beginners:

1. [introduction to programming in R](#)
2. [Reproducible scientific analysis using RStudio and R](#)

Useful but not fully free:

Datacamp - Each of these video tutorial series + hands on exercises is free to start and then unfortunately requires a paid subscription to continue. If you are in a class with more than 10 students, [the instructor can sign up and then all the students can get access for the whole semester for](#) free

- Introduction to the [RStudio IDE](#)
- Introduction to the tidyverse: <https://www.datacamp.com/courses/introduction-to-the-tidyverse>
- Generating reports with R markdown: <https://www.datacamp.com/courses/reporting-with-r-markdown>

Bonus slides!

(for reference)

Where do we go from here?

ROpenSci's perspective:

- **Train students** by putting homework, assignments & dissertations on the reproducible research spectrum
- **Publish examples** of reproducible research in our field
- **Request code & data** when reviewing
- **Submit to & review for journals** that support reproducible research
- **Critically review & audit** data management plans in grant proposals
- Consider reproducibility wherever possible in **hiring, promotion & reference letters.**

<http://ropensci.github.io/reproducibility-guide/sections/introduction/> - Leveque et al

From: [Brian Wansink](#)
To: [David Just](#)
Cc: [Collin Payne](#); [Sandra Cuellar](#)
Subject: Can Branding Improve School Lunches?
Date: Saturday, January 7, 2012 7:17:42 AM
Attachments: [Elmo Icon-AJPH - 1-7-12.doc](#)
[ATT00001.htm](#)

Hi David,

Here's the Elmo study we are going to spin off and submit.
I think we start with the AJPH as a Brief (80 word abstract and 800 word paper),
and go from there. I'll give Sandra a list of the journals and the priority order we
should consider. Let's consider these two first:

Brief -- American Journal of Public Health

Research Letter – Archives of Pediatric and Adolescent Medicine

One sticking point is that although the stickers increase apple selection by 71%, for some reason this is a p value of .06. It seems to me it should be lower. Do you want to take a look at it and see what you think. If you can get the data, and it needs some tweeking, it would be good to get that one value below .05.



Contents lists available at ScienceDirect

Journal of Experimental Social Psychology

journal homepage: www.elsevier.com/locate/jesp



Breaking the prejudice habit: Mechanisms, timecourse, and longevity^{☆,☆☆,★}



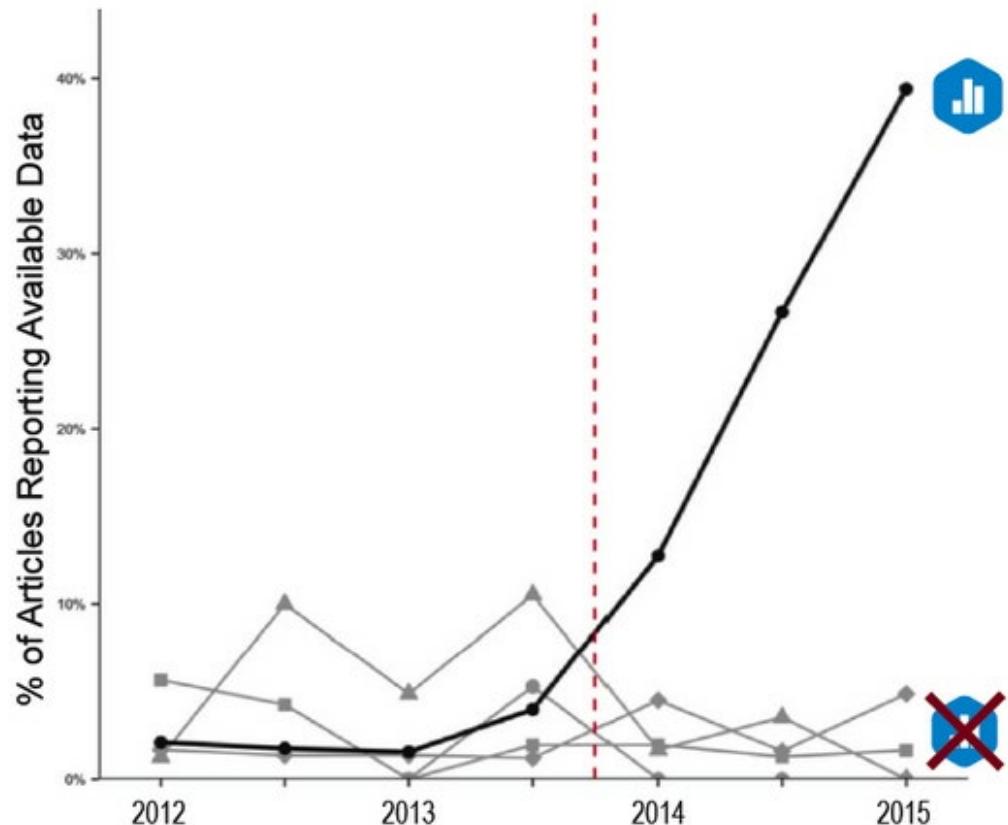
Patrick S. Forscher*, Chelsea Mitamura, Emily L. Dix, William T.L. Cox, Patricia G. Devine*

Department of Psychology, University of Wisconsin, Madison, United States

“Badges seem silly. Do they work?

- Yes. Implementing these badges dramatically increases the rate of data sharing ([Kidwell et al, 2016](#)).
- A recent systematic review identified this badging program as the only evidence-based incentive program that is associated with increased data sharing ([Rowhani-Farid et al., 2017](#)).
- View a list of journals and organizations that have adopted badges [here](#).

<https://cos.io/our-services/open-science-badges/>



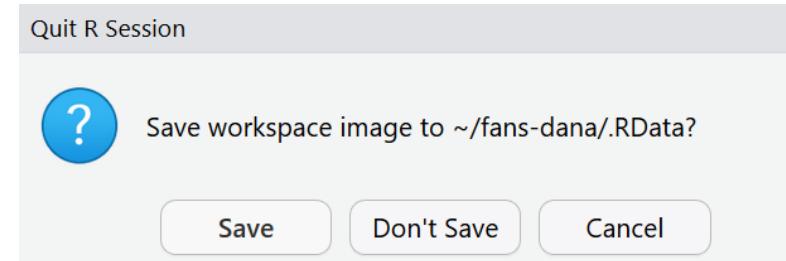
Three tips you can start using today

- here::here

```
library(here)  
db <- read_csv(here("subfolder_name", "file_name.csv"))
```

- “Save workspace image?”

- Don’t save!



- read_csv

- tidyverse version of read.csv
 - won’t coerce strings to factors, outputs a dataframe

What is markdown?

Text using Markdown syntax	Corresponding HTML produced by a Markdown processor	Text viewed in a browser
<p>Heading =====</p> <p>## Sub-heading</p> <p>Paragraphs are separated by a blank line.</p> <p>Two spaces at the end of a line produces a line break.</p> <p>Text attributes <code>_italic_</code>, <code>**bold**</code>, <code>`monospace`</code>.</p> <p>Horizontal rule:</p> <p>---</p> <p>Bullet list:</p> <ul style="list-style-type: none">* apples* oranges* pears	<pre><h1>Heading</h1> <h2>Sub-heading</h2> <p>Paragraphs are separated by a blank line.</p> <p>Two spaces at the end of a line
 produces a line break.</p> <p>Text attributes italic, bold, <code>monospace</code>. </p></pre> <p><p>Horizontal rule:</p></p> <p><hr /></p> <p><p>Bullet list:</p></p> <pre> apples oranges pears</pre>	<p>Heading [edit]</p> <p>Sub-heading [edit]</p> <p>Paragraphs are separated by a blank line.</p> <p>Two spaces at the end of a line produces a line break.</p> <p>Text attributes <i>italic</i>, bold, <code>monospace</code>.</p> <p>Horizontal rule:</p> <p>Bullet list:</p> <ul style="list-style-type: none">• apples• oranges• pears <p>Numbered list:</p> <ol style="list-style-type: none">1. wash2. rinse3. repeat

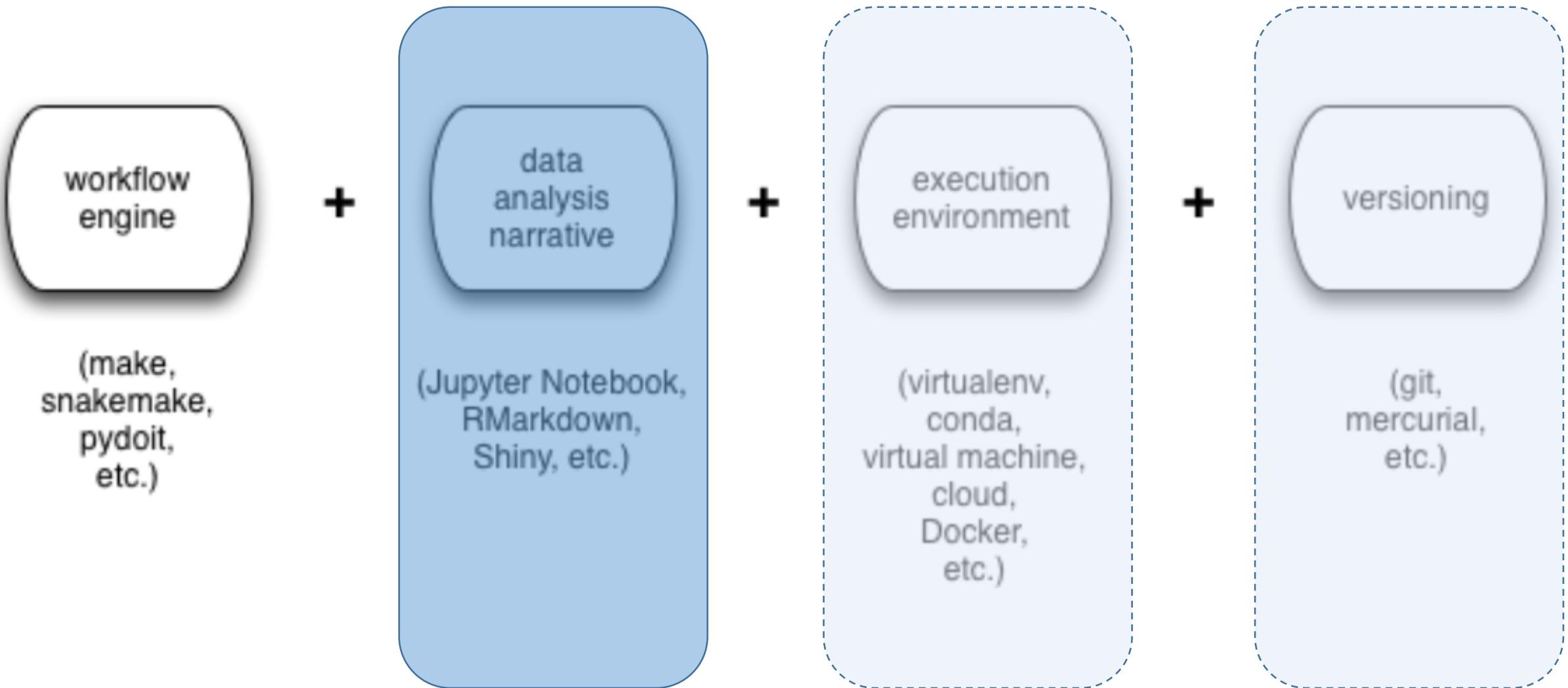
A reproducibility spectrum...example for a concept or theory

Classification created by the American Society for Cell Biology [5]:

1. **analytic replication** - reproduce results by reanalysing original data
2. **direct replication** – repeat original experiment
3. **systematic replication** - repeat with different experimental conditions
(eg different cell line)
4. **conceptual replication**– repeat validity of a concept
5. (eg in different organisms)

5- <https://www.nature.com/news/muddled-meanings-hamper-efforts-to-fix-reproducibility-crisis-1.20076#/b1>, citing <https://www.ascb.org/wp-content/uploads/2015/11/How-can-scientist-enhance-rigor.pdf> citing Schmidt, Stefan. (2009). Shall We Really Do It Again? The Powerful Concept of Replication Is Neglected in the Social Sciences. *Review of General Psychology*. 13. 90-100. 10.1037/a0015108.

Key categories of tools



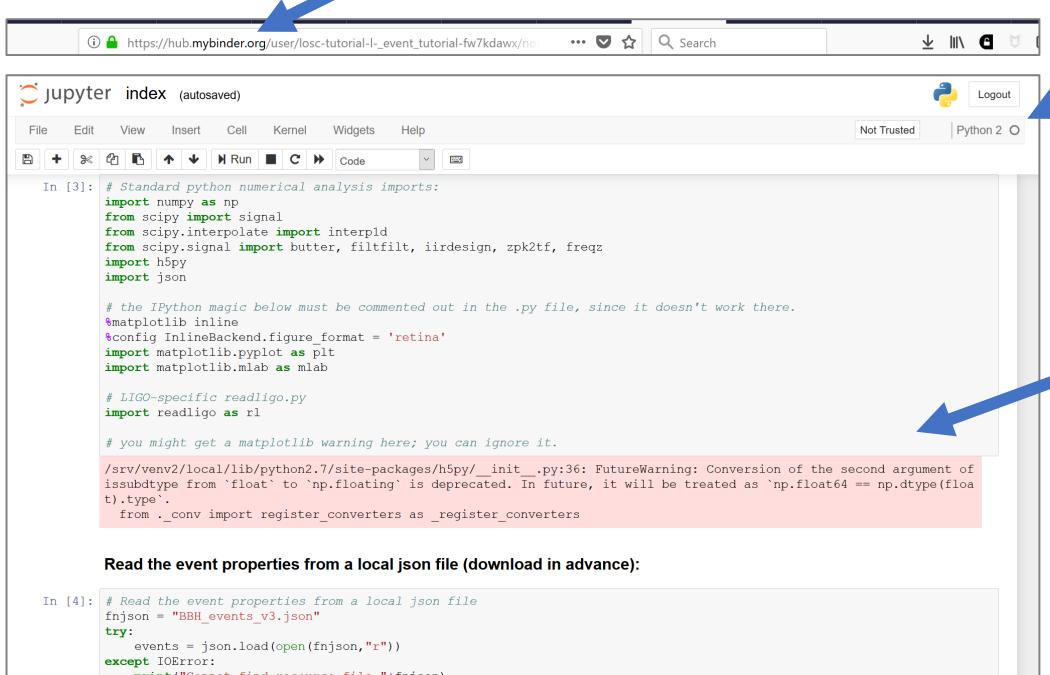
How does Binder work?

“Binder allows you to create custom computing environments that can be shared and used by many remote users”

<https://mybinder.readthedocs.io/en/latest/>

Visible in your browser,
but running on a remote machine

Your browser



```
# Standard python numerical analysis imports:
import numpy as np
from scipy import signal
from scipy.interpolate import interp1d
from scipy.signal import butter, filtfilt, iirdesign, zpk2tf, freqz
import h5py
import json

# the IPython magic below must be commented out in the .py file, since it doesn't work there.
%matplotlib inline
%config InlineBackend.figure_format = 'retina'
import matplotlib.pyplot as plt
import matplotlib.mlab as mlab

# LIGO-specific readligo.py
import readligo as rl

# you might get a matplotlib warning here; you can ignore it.

/srv/venv2/local/lib/python2.7/site-packages/h5py/_init_.py:36: FutureWarning: Conversion of the second argument of isubdtype from 'float' to 'np.floating' is deprecated. In future, it will be treated as 'np.float64 == np.dtype(float).type'.
    from .conv import register_converters as _register_converters

Read the event properties from a local json file (download in advance):
```

```
# Read the event properties from a local json file
fnjson = "BBH_events_v3.json"
try:
    events = json.load(open(fnjson, "r"))
except IOError:
    print("Cannot find resource file "+fnjson)
```

Specified version
of Python
(or other kernels eg R)

You see the code in
the .ipynb file

How does Binder work?

Code for notebook and list of dependencies and files available on GitHub (right now Binder only works with public repositories)



Your browser

https://hub.mybinder.org/user/losc-tutorial-l_-event_tutorial-fw7kdawx/notebooks/index.ipynb

jupyter index (autosaved)

In [3]:

```
# Standard python numerical analysis imports:  
import numpy as np  
from scipy import signal  
from scipy.interpolate import interp1d  
from scipy.signal import butter, filtfilt, iirdesign, zpk2tf, freqz  
import h5py  
import json  
  
# the IPython magic below must be commented out in the .py file, since it doesn't work there.  
%matplotlib inline  
config InlineBackend.figure_format = 'retina'  
import matplotlib.pyplot as plt  
import matplotlib.mlab as mlab  
  
# LIGO-specific readligo.py  
import readligo as rl  
  
# you might get a matplotlib warning here; you can ignore it.  
  
/srv/venv2/local/lib/python2.7/site-packages/h5py/_init_.py:36: FutureWarning: Conversion of the second argument of isubdtype from 'float' to 'np.floating' is deprecated. In future, it will be treated as 'np.float64 == np.dtype(float).type'.  
from _conv import register_converters as _register_converters
```

Read the event properties from a local json file (download in advance):

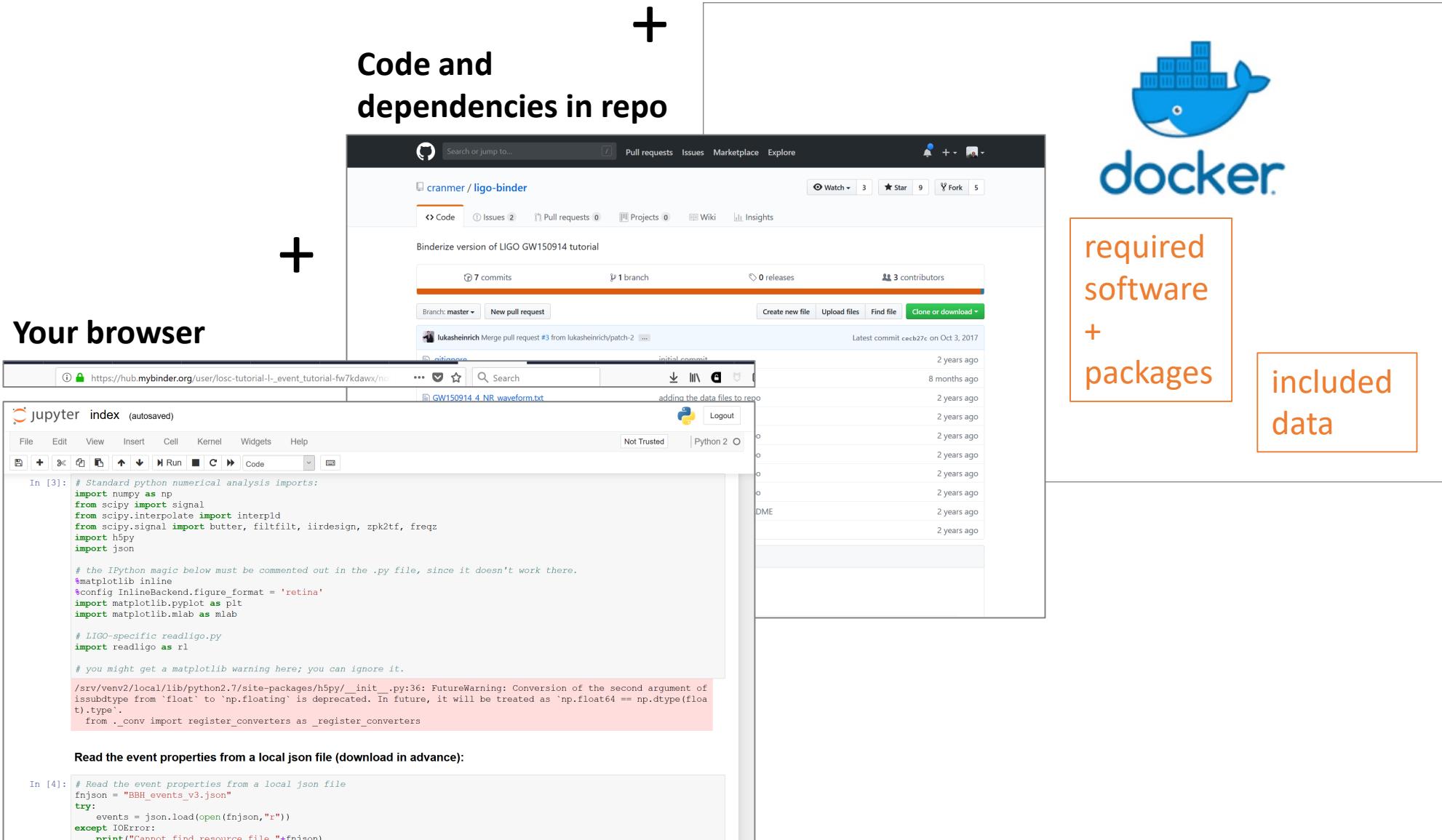
In [4]:

```
# Read the event properties from a local json file  
fnjson = "BBH_events_v3.json"  
try:  
    events = json.load(open(fnjson, "r"))  
except IOError:  
    print("Cannot find resource file "+fnjson)
```

How does Binder work?



builds Docker image based on repo and generates URL for public access



How does Binder work?

Cloud hosting + management

The diagram illustrates the workflow for running a Jupyter notebook. It features three main components arranged vertically:

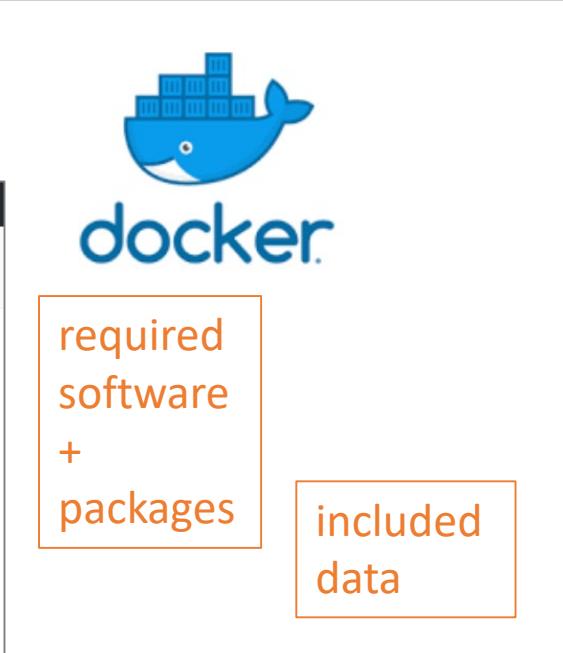
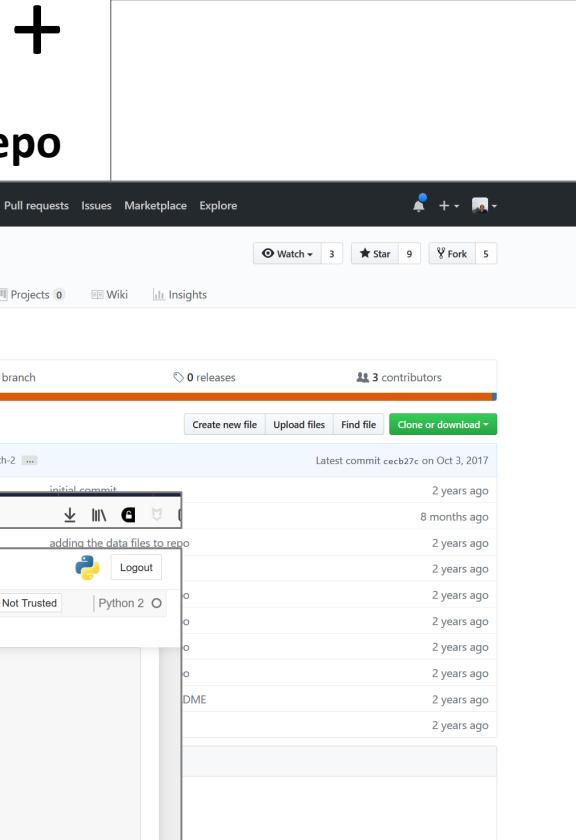
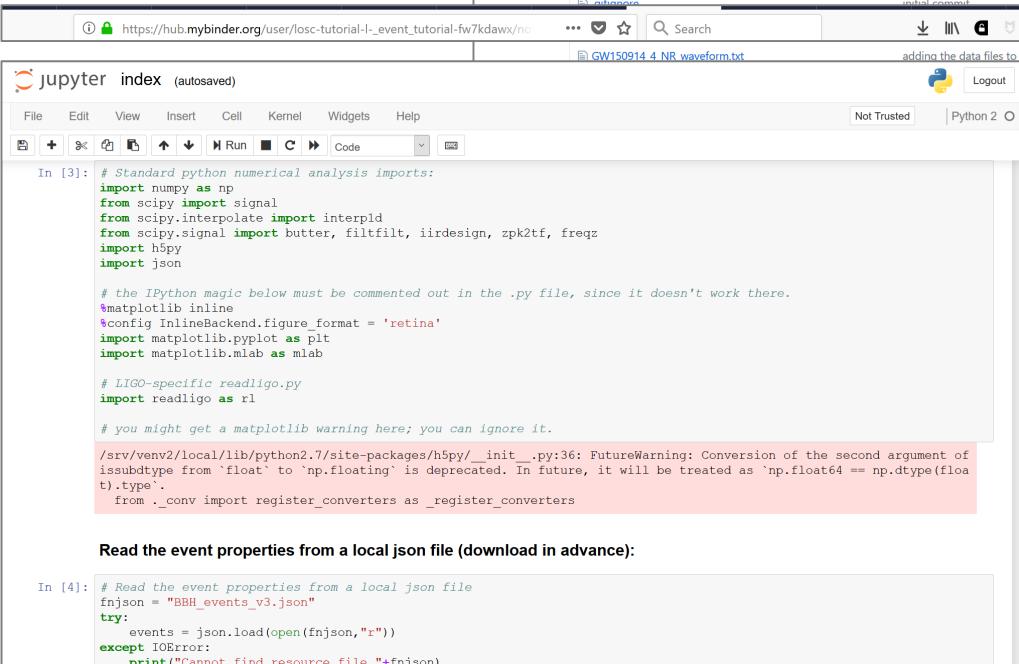
- Docker image**: At the top, a large box labeled "Docker image" contains a plus sign (+). Below it, the text "Code and dependencies in repo" is displayed.
- Your browser**: In the middle, a box labeled "Your browser" also contains a plus sign (+). It displays a screenshot of a GitHub repository page for "cranmer / ligo-binder". The repository description is "Binderize version of LIGO GW150914 tutorial". It shows 7 commits, 1 branch, 0 releases, and 3 contributors. A pull request by "lukasheinrich" is visible, along with a note about adding data files to the repository.
- Jupyter Notebook**: At the bottom, a box displays a Jupyter notebook interface. The title bar says "jupyter index (autosaved)". The code cell "In [3]" contains Python code for importing various scientific libraries like numpy, scipy, and matplotlib. It also includes IPython magic commands for configuring the backend and reading specific files. A note at the bottom of this cell states: "# you might get a matplotlib warning here; you can ignore it.". The output cell "Out [3]" shows the results of the executed code, including a list of event properties from a local JSON file.

+

Docker image

Code and dependencies in repo

Your browser



JupyterHub

+



Kubernetes

Generated binders are hosted by a hub
that deploys, scales, and manages
compute resources
using [Kubernetes](#)

(can run on any cloud platform provider)
More details: [BinderHub](#)

Is this a tidy dataset?

year	artist	track	time	date.entered	wk1	wk2	wk3
2000	2 Pac	Baby Don't Cry	4:22	2000-02-26	87	82	72
2000	2Ge+her	The Hardest Part Of ...	3:15	2000-09-02	91	87	92
2000	3 Doors Down	Kryptonite	3:53	2000-04-08	81	70	68
2000	98^0	Give Me Just One Nig...	3:24	2000-08-19	51	39	34
2000	A*Teens	Dancing Queen	3:44	2000-07-08	97	97	96
2000	Aaliyah	I Don't Wanna	4:15	2000-01-29	84	62	51
2000	Aaliyah	Try Again	4:03	2000-03-18	59	53	38
2000	Adams, Yolanda	Open My Heart	5:30	2000-08-26	76	76	74

Table 7: The first eight Billboard top hits for 2000. Other columns not shown are wk4, wk5, ..., wk75.



No, has more than one observation per row

<http://vita.had.co.nz/papers/tidy-data.html>

What improved?

year	artist	time	track	date	week	rank
2000	2 Pac	4:22	Baby Don't Cry	2000-02-26	1	87
2000	2 Pac	4:22	Baby Don't Cry	2000-03-04	2	82
2000	2 Pac	4:22	Baby Don't Cry	2000-03-11	3	72
2000	2 Pac	4:22	Baby Don't Cry	2000-03-18	4	77
2000	2 Pac	4:22	Baby Don't Cry	2000-03-25	5	87
2000	2 Pac	4:22	Baby Don't Cry	2000-04-01	6	94
2000	2 Pac	4:22	Baby Don't Cry	2000-04-08	7	99
2000	2Ge+her	3:15	The Hardest Part Of ...	2000-09-02	1	91
2000	2Ge+her	3:15	The Hardest Part Of ...	2000-09-09	2	87
2000	2Ge+her	3:15	The Hardest Part Of ...	2000-09-16	3	92
2000	3 Doors Down	3:53	Kryptonite	2000-04-08	1	81
2000	3 Doors Down	3:53	Kryptonite	2000-04-15	2	70
2000	3 Doors Down	3:53	Kryptonite	2000-04-22	3	68
2000	3 Doors Down	3:53	Kryptonite	2000-04-29	4	67
2000	3 Doors Down	3:53	Kryptonite	2000-05-06	5	66

id	artist	track	time
1	2 Pac	Baby Don't Cry	4:22
2	2Ge+her	The Hardest Part Of ...	3:15
3	3 Doors Down	Kryptonite	3:53
4	3 Doors Down	Loser	4:24
5	504 Boyz	Wobble Wobble	3:35
6	98^0	Give Me Just One Nig...	3:24
7	A*Teens	Dancing Queen	3:44
8	Aaliyah	I Don't Wanna	4:15
9	Aaliyah	Try Again	4:03
10	Adams, Yolanda	Open My Heart	5:30
11	Adkins, Trace	More	3:05
12	Aguilera, Christina	Come On Over Baby	3:38
13	Aguilera, Christina	I Turn To You	4:00
14	Aguilera, Christina	What A Girl Wants	3:18
15	Alice Deejay	Better Off Alone	6:50

id	date	rank
1	2000-02-26	87
1	2000-03-04	82
1	2000-03-11	72
1	2000-03-18	77
1	2000-03-25	87
1	2000-04-01	94
1	2000-04-08	99
2	2000-09-02	91
2	2000-09-09	87
2	2000-09-16	92
3	2000-04-08	81
3	2000-04-15	70
3	2000-04-22	68
3	2000-04-29	67
3	2000-05-06	66