

Reproducible Data Science

Fundraising Analytics with GitHub and R Markdown

Paul Hively, Northwestern University



The DRIVE/ Conference

March 11-13, 2019 Baltimore, Maryland

What to expect today

- Paul Hively, Director of Analytics at Kellogg School of Management, Northwestern University
- A few bad jokes
- Plenty of real-world examples and analogies
- How the reproducible approach can **greatly increase** future efficiency and impact

Outline

- The case for reproducibility
 - What?
 - Why?
 - How?
- Tools & use cases
 - GitHub
 - R Markdown
- Demonstration

The case for reproducibility

Do any of these situations look familiar?

Sudden show-stoppers: “But it worked yesterday!”



Extract Refresh Failed

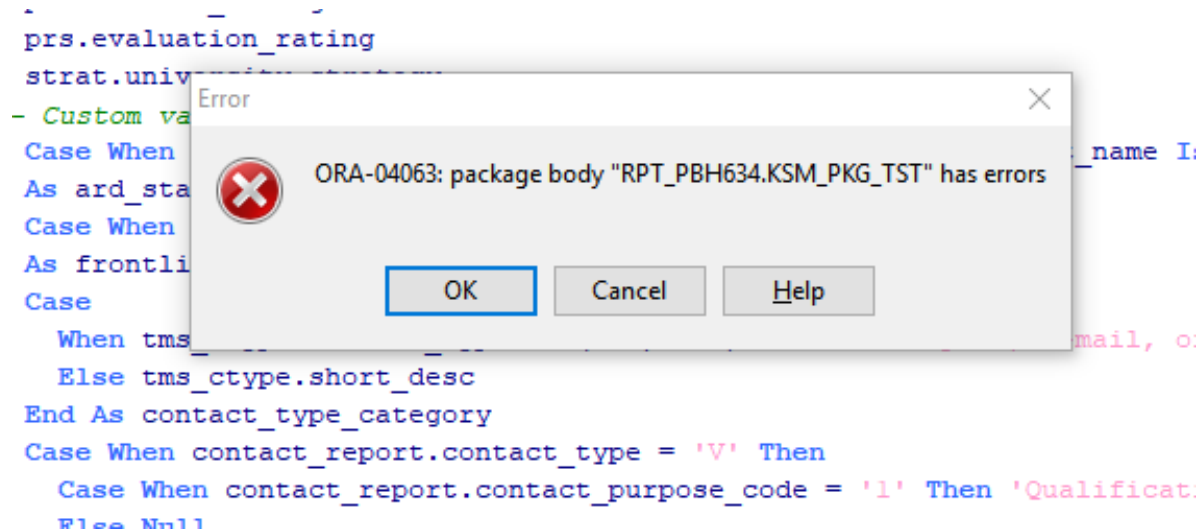
Site: Alumni Relations and Development

Workbook: KSM Fundraising Metrics

The case for reproducibility

Do any of these situations look familiar?

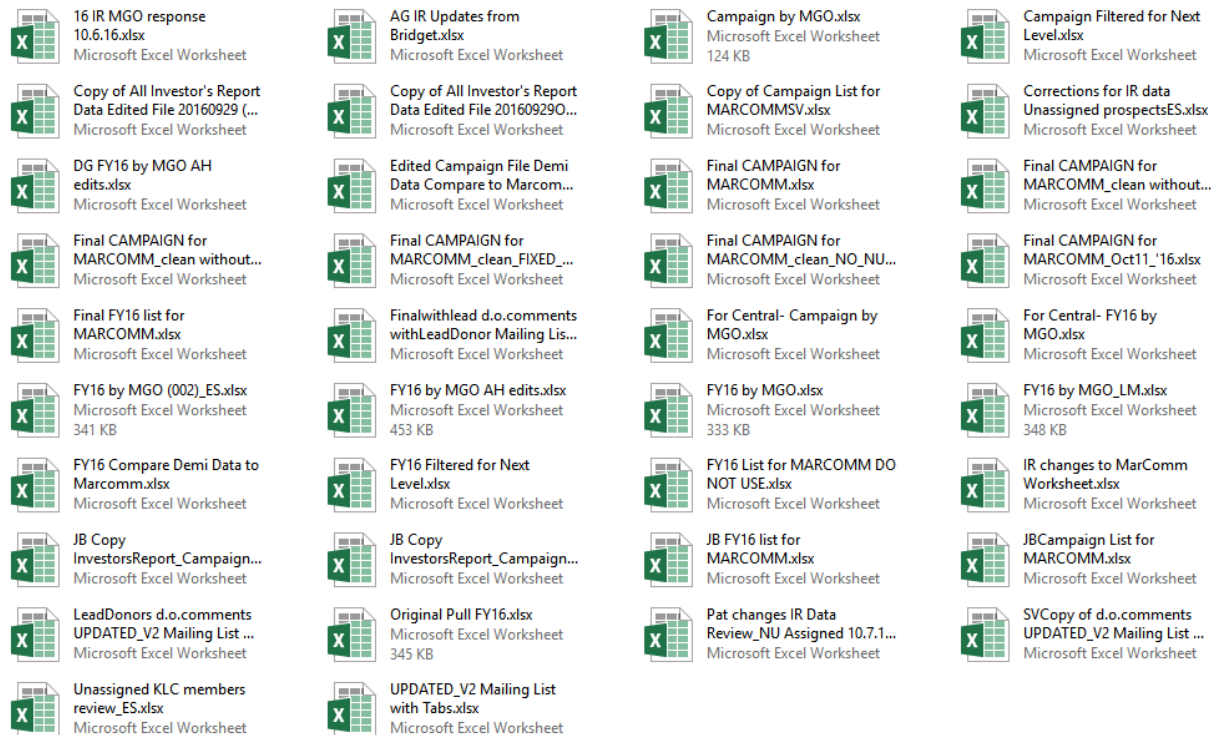
Even better: “But it worked until I hit Save!”



The case for reproducibility

Do any of these situations look familiar?

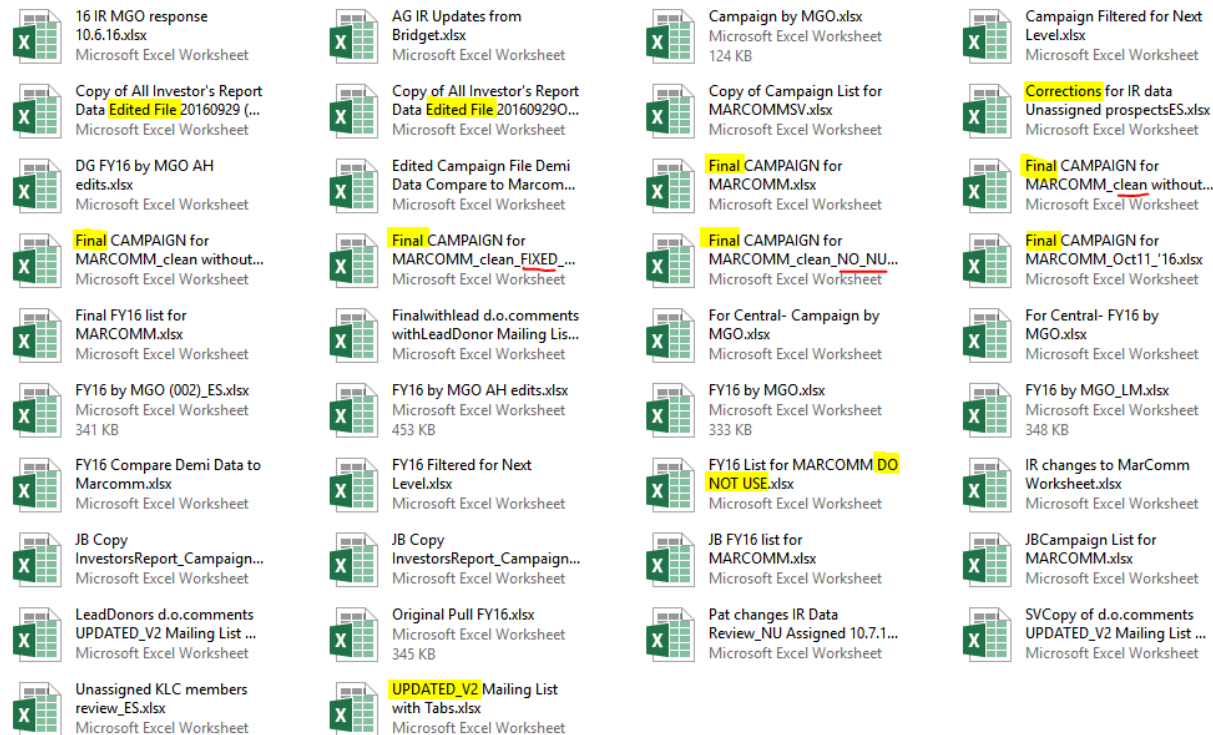
File mayhem: “Find the master list”



The case for reproducibility

Do any of these situations look familiar?

Some of my favorite features:



The case for reproducibility

Do any of these situations look familiar?

Excel Olympics: “Just follow these 60 quick steps!”

Pipeline Steps

- Open macro [Macro 20170719](#)
- Go to the BI, and click CARD Reporting-Proposal Pipeline Report
- Under “Select a Program,” select Kellogg
- Under “Select Fiscal Year,” select the current Fiscal Year, click OK
- After report appears on the screen, click the HTML icon on the upper right hand corner of the screen
- Click View in Excel Options, click View in Excel 2007 Data to export the file
- Enable macros in the exported data file
- Run Macro by clicking [CopySheet04](#)
- LABELLING SPLIT PROPOSALS:
 - Scroll to column AB called “Programs” and filter out any rows that are Kellogg only or “Kellogg, Financial Aid”. You will be left with a group of split proposals.
 - Click on the top cell in column AC and type one Asterisk in the field
 - Autofill the rest of the column so that the split proposals are labelled with the asterisk
- HIGHLIGHTING SPLIT PROPOSAL ERRORS (to prompt PAs to fill program specific ask and anticipated amounts):
 - Compare column Y (Program Level Anticipated) and Column V (KSM Anticipated Commitment). If there is a zero in any cells in column Y, highlight the cell to the RIGHT of the zero.
 - Repeat the exercise with columns O (Program Level Ask) and P (KSM Ask Amount). If there are any zeros in column O, highlight cell to the RIGHT of the zero (see below).

N222									
500000									
1	Ask Amount	Program Level Ask Amt	Kellogg Ask Amt	Closed	Probability				
14	\$ 400,000	250000	250,000	2/19/2018	90% Probability (High)	\$			
20	\$ 100,000	50000	50,000	12/19/2018	75% Probability (Medium)	\$			
24	\$ 100,000	50000	50,000	12/22/2018	90% Probability (High)	\$			
31	\$ 250,000	250000	250,000	12/28/2018	90% Probability (High)	\$			
32	\$ 500,000	250000	250,000	12/31/2018	75% Probability (Medium)	\$			
41	\$ 500,000	250,000	250,000	8/2/2019	75% Probability (Medium)	\$			
51	\$ 250,000	125000	125,000	8/2/2019	75% Probability (Medium)	\$			
61	\$ 2,000,000	1,000,000	1,000,000	4/1/2019	75% Probability (Medium)	\$			
106	\$ 1,000,000	500,000	500,000	8/1/2018	50% Probability	\$			
136	\$ 125,000	62500	62,500	1/30/2018	50% Probability	\$			
189	\$ 500,000	250,000	250,000	4/29/2018	90% Probability (High)	\$			
190	\$ 100,000	50,000	50,000	3/31/2016	75% Probability (Medium)	\$			
191	\$ 500,000	250,000	250,000	6/9/2018	90% Probability (High)	\$			
198	\$ 250,000	125,000	125,000	8/19/2018	75% Probability (Medium)	\$			
201	\$ 500,000	250,000	250,000	6/1/2018	75% Probability (Medium)	\$			
241	\$ 100,000	50,000	50,000	10/31/2018	50% Probability	\$			
171	\$ 250,000	125,000	125,000	12/31/2018	75% Probability	\$			
196	\$ 100,000	50,000	50,000	12/30/2018	50% Probability	\$			
199	\$ 100,000	50,000	50,000	11/1/2018	75% Probability (Medium)	\$			
193	\$ 100,000	50,000	50,000	8/2/2018	50% Probability	\$			

- Remove filters from data file, then add filter to “KSM Ask Amount” >= \$100K+

- Create three additional new tabs in the excel workbook. Label them
 - Closed
 - Open FY (Insert current fiscal year)
 - Open FY (Next fiscal year) and beyond
- Go to the tab with your data. Copy and paste all rows that say “Closed Proposals: 2016” (column c), in the new tab that you just created called “Closed”
- Next, copy and paste all rows that say “Open Proposals: [current fiscal year]” (column c), in the new tab that you just created called “Open FY (insert current fiscal year)”
- Next, copy and paste all rows that say “Open Proposals: [next fiscal year] [Beyond]” (column c), in the new tab that you just created called “Open FY (insert next fiscal year) and Beyond”
- Go to your newly created “Closed” tab and sort the proposals by “KSM Granted Amount”
- Go to your newly created “Open FY (insert current fiscal year)” tab and sort the tab by KSM Ask Amount first, and then status (Approved by Donor, then Submitted, then Anticipated)
- Go to your newly created “Open FY (insert next fiscal year) and Beyond” tab and sort the tab by KSM Ask Amount first, and then status (Approved by Donor, then Submitted, then Anticipated)
- Open proposal pipeline template saved [here](#)
- Paste the closed proposals in the “Closed Tab” in the template. Highlight proposals with granted amount of \$100K+ in green
- Check closed tab for data integrity
 - Highlight close dates and ask dates in the future
 - Highlight statuses other than Declined, Withdrawn, Funded
 - Highlight blank proposal manager fields
- Go back to your data sheet and paste the “Open Proposals: [current fiscal year]” proposals in the “Open” tab on the template. In template, add another header under that section that says “Open Proposals: FY (Next FY) and Beyond.”
- Paste the “Open Proposals: [next fiscal year] [Beyond]” in the “Open” tab under that section
- Check the “Open” tab for data integrity
 - Highlight statuses other than Approved By donor, Anticipated, and Submitted
 - Highlight blank proposal manager fields
 - Highlight past due Ask Dates for proposals with statuses “Approved by Donor” and “Submitted”
 - Highlight Ask Dates that have passed for “Anticipated” proposals
 - Highlight all close dates that are past due
- Make sure each page of the “Open” tab in the template has the appropriate header.
- Save the new report in the following location:
 - G:\Ext\rel\ADVANCEMENT\REPORTING\Damli Giannaras\01 Proposal Pipeline

The case for reproducibility

Do any of these situations look familiar?

Excel Olympics: Sorting SNAFU

The diagram illustrates a common Excel sorting mistake. On the left, a spreadsheet shows columns A (A to Z) and B (1 to 26). A context menu is open over column B, with the 'Sort' option selected. The 'Sort' submenu shows 'Sort Smallest to Largest' and 'Sort Largest to Smallest' options. A blue arrow points from the 'Sort' menu to the right-hand spreadsheet, which shows the result of sorting column B by the values in column A, resulting in a reversed order (Z to A). The word 'oops' is written below the arrow.

A	B
A to Z	1 to 26
A	1
B	2
C	3
D	4
E	5
F	6
G	7
H	8
I	9
J	10
K	11
L	12
M	13
N	14
O	15
P	16
Q	17
R	18
S	19
T	20
U	21
V	22
W	23
X	24
Y	25
Z	26

oops

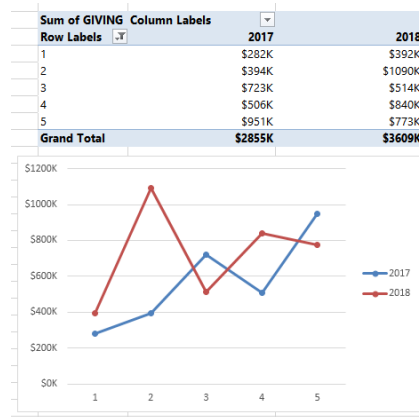
A	B
A to Z	1 to 26
A	26
B	25
C	24
D	23
E	22
F	21
G	20
H	19
I	18
J	17
K	16
L	15
M	14
N	13
O	12
P	11
Q	10
R	9
S	8
T	7
U	6
V	5
W	4
X	3
Y	2
Z	1

The case for reproducibility

Do any of these situations look familiar?

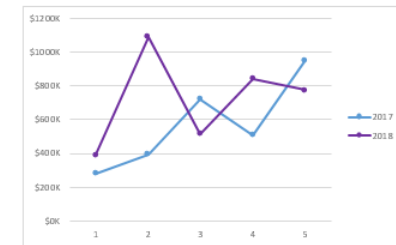
Copy-Paste Purgatory

	A	B	C	D	E	F
1	FISCAL_YEAR	DATE_OF_RECORD	GIVING	MO	DAY	
2	2017	2/8/2017	33805	2	8	
3	2018	2/1/2018	8700	2	1	
4	2018	12/31/2017	647935	12	31	
5	2018	5/1/2018	97601.23	5	1	
6	2017	4/13/2017	17853.87	4	13	
7	2017	10/17/2016	21432.31	10	17	
8	2018	11/3/2017	11574.5	11	3	
9	2018	2/6/2018	19098.18	2	6	
10	2018	9/29/2017	83676.02	9	29	
11	2017	4/26/2017	5334.84	4	26	
12	2017	8/1/2017	79601.5	8	1	
13	2017	11/15/2016	33209.08	11	15	
14	2017	6/30/2017	34598.49	6	30	
15	2017	10/18/2016	16944.34	10	18	
16	2017	11/17/2016	7366.87	11	17	



Reporting Requirements

This is a super cool report with a bunch of text. Many astute observations were made, and they are supported by colorful graphics such as the below chart. However, this document is not quite what was requested.



As you can see, 2016 does not appear here, which turns out to be an oversight. When 2016 is requested, since this was created as a manual process that began with manual data aggregation in Excel, followed by creation of a chart, which was then pasted into this document and manually re-colored and formatted, all that work will have to be repeated when 2016 data is added in the initial step.

A better workflow would have been to use R Markdown instead. In the words of the authors of development environment RStudio, "R Markdown documents are fully reproducible. Use a productive

Manual data pull → Manual aggregation → Manual formatting and content

Need to include more data

“Here we go again...”

The case for reproducibility

Reproducibility can help:

- Roll back and compare previous versions
- Effectively build upon past work
- Create unambiguous and automated procedures
- Save time, sanity, hairline

What is reproducibility?

“An analysis that can be passed from one person to another and, using the same data, generate the same results in an unambiguous manner.”

—Bray, Çetinkaya-Rundel, and Stangl

“Same data + Same script = Same results”

—Daniel Marcelino

What is reproducibility?



CartoonStock.com

<https://deevybee.blogspot.com/2018/02/improving-reproducibility-future-is.html>

Whose idea was reproducibility?

Replication crisis

Essay

Why Most Published Research Findings Are False

John P. A. Ioannidis

Bloomberg

OPINION | ECONOMICS

Why 'Statistical Significance' Is Often Insignificant

Researchers who want professorships are sometimes driven to publish suspect findings.

By Noah Smith

19 November 2, 2017, 7:00 AM CDT

The New York Times

Many Psychology Findings Not as Strong as Claimed, Study Says

General Article

False-Positive Psychology: Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant

Joseph P. Simmons¹, Leif D. Nelson², and Uri Simonsohn¹

¹The Wharton School, University of Pennsylvania, and ²Haas School of Business, University of California, Berkeley

Archive > Volume 497 > Issue 7447 > News > Article

NATURE | NEWS



Disputed results a fresh blow for social psychology

Failure to replicate intelligence-priming effects ignites row in research community.

Alison Abbott

30 April 2013 | Clarified: 17 May 2013

By Daniel Engber

TWEET

SHARE

Daryl Bem Proved ESP Is Real

Which means science is broken.

MAY 17, 2017 • COVER STORY

Whose idea was reproducibility?

Replication crisis – what about ESP?!

- Psychologist Daryl Bem used large sample sizes and accepted experimental/statistical methods to demonstrate ESP exists
- Is ESP real or do the accepted methods have issues?



Lisa Larson-Walker for Slate magazine

Whose idea was reproducibility?

Replication crisis

- The results of many published scientific studies do not hold up to further scrutiny
- Methodology, statistical, and analysis issues
- Reproducibility can help: it requires a detailed methodology and provides evidence that correct results were produced

Whose idea was reproducibility?



“But we’re trying to fundraise, not publish research!”

Whose idea was reproducibility?

Collaboration

- Reproducibility makes it easier to catch mistakes
- Reproducibility enables others to catch mistakes
- Reproducibility is a teaching tool
- Reproducibility spreads the current state of the art

Whose idea was reproducibility?

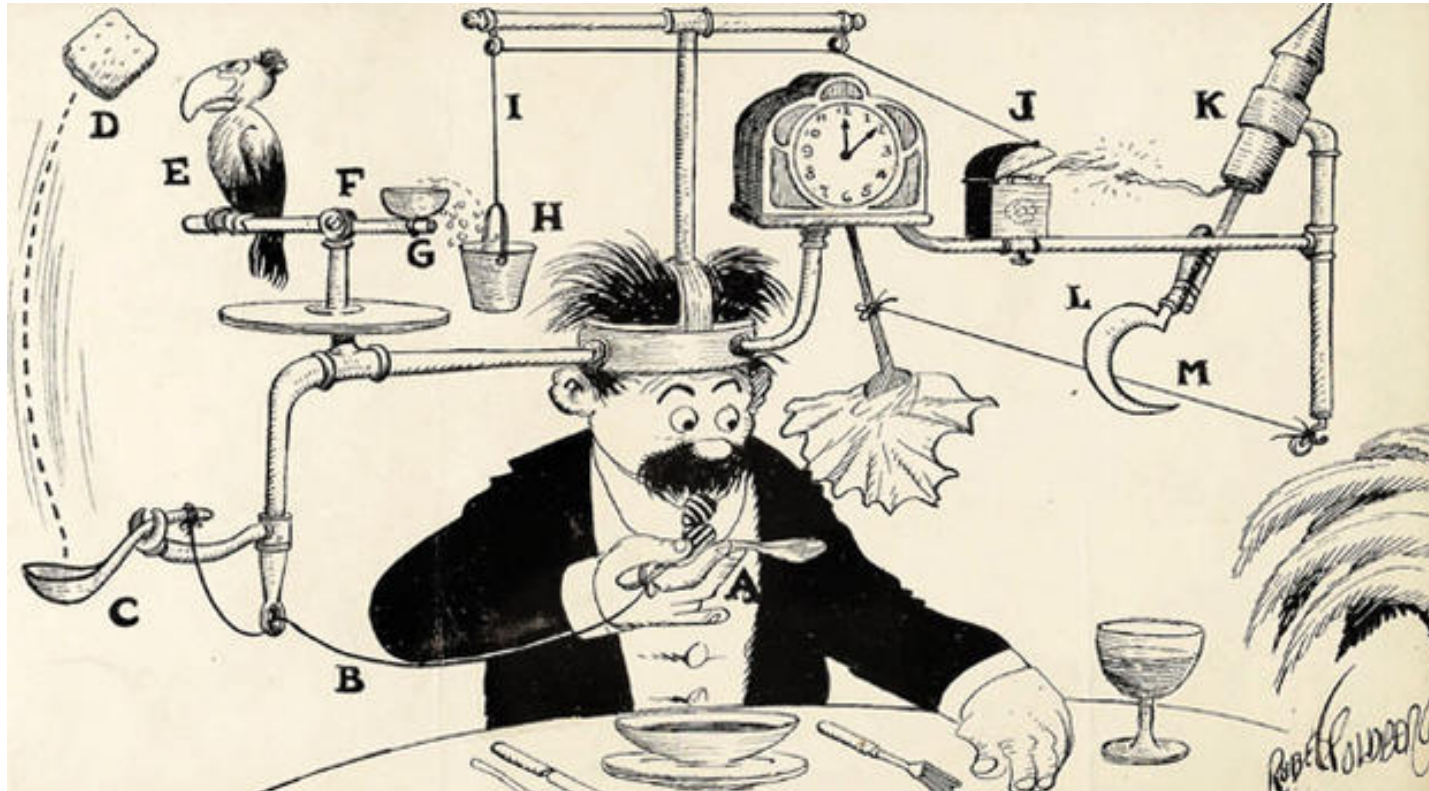
Collaboration

Recent prospect-dmm thread:

“I wanted to get some thoughts about how often others update thresholds/indicators of donors/prospects based on their models.... I would love insight on how y’all approach this part of the modeling workflow.”

Would you rather have thoughts or an example?

What does reproducibility require?



"Self-operating napkin" (Rube Goldberg, 1931)

What does reproducibility require?

The right mindset

- Willingness to plan ahead
- Go slow to go fast
- Don't wait to document

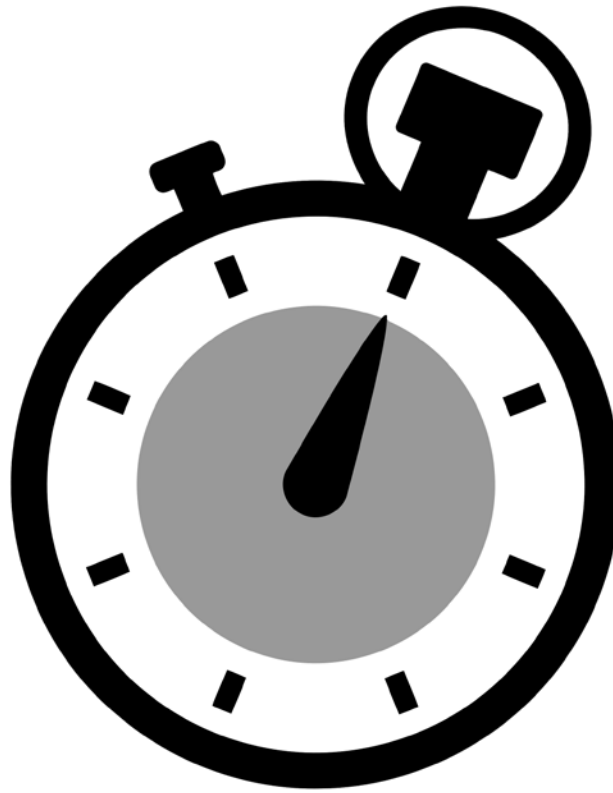
What does reproducibility require?

The right tools

- A way to save earlier revisions
- A way to automate the process

Quick pause

Questions? Comments? Funny stories? Other stories?



GitHub



Key facts

- Cloud hosting solution for Git
- Free version control platform
- Easy to save and compare changes from any two points in history
- Mascot = Octocat



GitHub

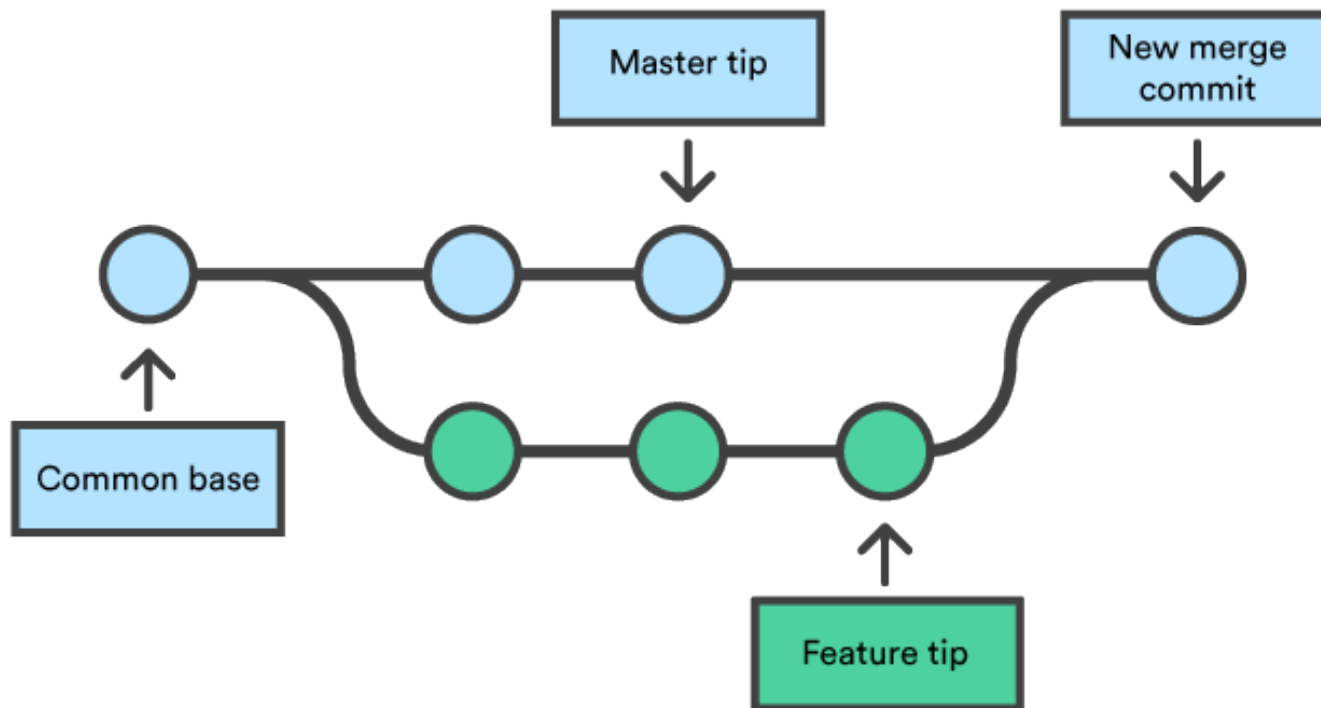


Key concepts

- Repository = project
- Branch = copy
- Merge = reconcile revisions

GitHub

Key concepts



R Markdown



Key facts

- Built into the R Studio IDE
- Formatted text, R code, plots, tables, etc. can all exist in the same document
- Supports literate programming and self-documenting code
- Create HTML, pdf, and even Microsoft Office files (with plug-ins)

R Markdown



Coding rules of thumb

- Comment everything
- Use functions
- Parameters must be easy to find
- Use descriptive naming conventions
- Format consistently (including whitespace)

“There is no ‘after’ in which to write documentation.”
—*Overheard on prospect-dmm*

Demo time!

<https://github.com/phively/drive19>

```
37 02 Reproducible data pipeline.Rmd View file
13 13 @@ -13,10 +13,8 @@ A fresh session will not have any packages or data loaded. I recommend Hadly Wic
14 14 ```{r}
15 15 # Run this after installing a new version of R to download the most up-to-date version of the package
16 16 - install.packages('tidyverse')
17 17 - ```
18 18 + # install.packages('tidyverse')
19 19 - ```{r}
20 20 # Run this to load the package into the current session
21 21 library(tidyverse)
22 22 ```
23 23
24 24 @@ -31,7 +29,7 @@ contribution <- read.csv(file = 'data/contribution.csv', stringsAsFactors = FALSE)
25 25
26 26
27 27
28 28
29 29
30 30 Examining the first few rows of data:
31 31
32 32 - ```{r}
33 33 + ```{r, cols.print = 12}
34 34 head(contribution)
35 35 ```
36 36
37 37
38 38 @@ -48,18 +46,28 @@ read.csv(file = 'data/contribution.csv') %>%
39 46 filter(AttendanceEvent == 1) %>%
40 47 # Define lifetime giving as the sum of 5 years of giving
41 48 mutate(LifetimeGiving = FY04Giving + FY03Giving + FY02Giving + FY01Giving + FY00Giving) %>%
42 49 - # Group the remaining data by class year, and sum our new LifetimeGiving object within each year
43 50 + # Group the remaining data by class year, and compute the statistics of interest within each year
44 51 group_by(Class.Year) %>%
45 52 summarise(ClassLifetimeGiving = sum(LifetimeGiving)) %>%
46 53 # Format as dollars
47 54 mutate(ClassLifetimeGiving = scales::dollar(ClassLifetimeGiving))
48 55 + summarise(
49 56 + Donors = sum(LifetimeGiving > 0)
50 57 + , ClassGiving = sum(LifetimeGiving)
51 58 + ) %>%
52 59 + # We can easily compute derived statistics
```

Key takeaways

- Automate everything you can
- Get help automating the things you can't
- Go slow to go fast
- Don't wait to document
- Share your work – let's grow together!

References

Reproducibility in Science

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

Disputed results a fresh blow for social psychology (Alison Abbott)

<https://www.nature.com/news/disputed-results-a-fresh-blow-for-social-psychology-1.12902>

Five Concrete Reasons Your Students Should Be Learning to Analyze Data in the Reproducible Paradigm (Bray, Çetinkaya-Rundel, and Stangl)

<http://chance.amstat.org/2014/09/reproducible-paradigm/>

Many Psychology Findings Not as Strong as Claimed, Study Says (Benedict Carey)

<https://www.nytimes.com/2015/08/28/science/many-social-science-findings-not-as-strong-as-claimed-study-says.html>

Reproducible Research (John Cook)

<https://reproducibleresearch.net>

Daryl Bem Proved ESP Is Real (Daniel Engber)

<https://slate.com/health-and-science/2017/06/daryl-bem-proved-esp-is-real-showed-science-is-broken.html>

References

Why most published research findings are false (John Ioannidis)

<http://journals.plos.org/plosmedicine/article?id=10.1371/journal.pmed.0020124>

What is reproducible research? (Daniel Marcelino)

<https://www.r-bloggers.com/what-is-reproducible-research/>

The Real Reason Reproducible Research is Important (Roger Peng)

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

What is reproducible research? (Henry Rogalin)

<https://bitesizebio.com/37187/reproducible-research/>

False-Positive Psychology (Simmons, Nelson, and Simonsohn)

<http://journals.sagepub.com/doi/abs/10.1177/0956797611417632>

Why 'statistical significance' is often insignificant (Noah Smith)

<https://www.bloomberg.com/view/articles/2017-11-02/why-statistical-significance-is-often-insignificant>