

Reproducible Results and the Workflow of Data Analysis

Scott Long

Departments of Sociology and Statistics
www.indiana.edu/~jslsoc/ftp/

Workshop in Methods | January 2018

The reproducible results movement

- Open science
- Transparency in science
- Teaching integrity in research

Changing expectations for researchers

- NAS 2018 Committee on Reproducibility and Replicability in Science
- Journals require data and analysis files
- Funding agencies strengthen requirements for data access
- Haverford College requires reproducibility for undergraduates

With access comes accountability

- Retraction Watch (retractionwatch.com)
- Recent examples of flawed research...

Reproducible Results and Workflow | 1

Retraction due to coding error

RETRACTED: In Sickness and in Health? Physical Illness as a Risk Factor for Marital Dissolution in Later Life

Abstract

The health consequences of marital dissolution are well known, but little work has examined the impact of health on the risk of marital dissolution. In this study we use a sample of 2,701 marriages from the Health and Retirement Study (1992–2010) to examine the role of serious physical illness onset (i.e., cancer, heart problems, lung disease, and/or stroke) in subsequent marital dissolution due to either divorce or widowhood. We use a series of discrete-time event history models with competing risks to estimate the impact of husband's and wife's physical illness onset on risk of divorce and widowhood. We find that only wife's illness onset is associated with elevated risk of divorce, while either husband's or wife's illness onset is associated with elevated risk of widowhood. These findings suggest the importance of health as a determinant of marital dissolution in later life via both individual and gendered social pathways.

Keywords

aging, chronic disease, gender inequality, marital health

A large body of literature has identified marital status as a strong predictor of health and well-being. Not only are unmarried healthier than the married (e.g., Lillard and Willis 1995; Emmons 1992), but studies find that both divorce and widowhood are predictors of subsequent physical and mental health (e.g., Hughes and Waite 2009; Williams and Uchino 2003). This attention, however, has been paid to the health as a determinant of marital status. Work in this area has tended to focus on the positive selection of the healthier into marriage (e.g., Byrne et al. 1989; Smith and Smith 2010), but poor health may be an equally important force for selection

Booth, and Johnson 2006). Illness may initiate changes to spouses' roles—in particular, increasing caregiving responsibilities for the healthy spouse—which can tax marital relationship dynamics (Wolff and Kasper 2006). Illness may also decrease household income due to the inability of one or both spouses to work (Teachman 2010), which may increase marital strain.

Only a few studies have examined the role of poor health in subsequent divorce, and these studies are mixed in their findings, with some finding

Reproducible Results and Workflow | 2

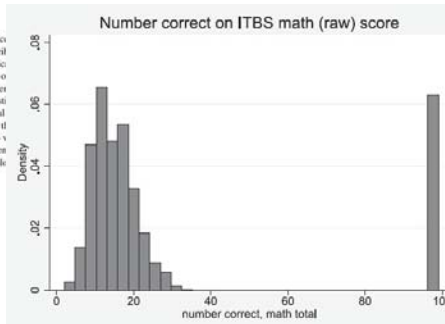
Incorrect data in published research

Distributional Analysis in Educational Evaluation: A Case Study from the New York City Voucher Program

Marianne Bitler, Thurston Domina, and Emily Penner
University of California, Irvine, Irvine, California, USA

Hilary Hoynes
University of California, Berkeley, Berkeley, California, USA

Abstract: We use quantile treatment effects estimation to examine the consequences of the New York City School Choice Scholarship Program across the distribution of achievement. Our analyses suggest that the program had negligible and statistically insignificant effects across the skill distribution. In addition to contributing to the literature on the article illustrates several ways in which distributional effects estimation can be useful in research: First, we demonstrate that moving beyond a focus on mean effects can be useful to generate and test new hypotheses about the heterogeneity of educational effects that speak to the justification for many interventions. Second, we demonstrate that effects can uncover issues even with well-studied data sets by forcing analysis to be done in new ways. Finally, such estimates highlight where in the overall national achievement scores of children exposed to particular interventions lie; this is important for the external validity of the intervention's effects.



Reproducible Results and Workflow | 3

Fragility of published results

Measurement, methods, and divergent patterns: Reassessing the effects of same-sex parents[☆]

Simon Cheng^{a,1}, Brian Powell^{b,1}

^a 344 Mansfield Rd., Department of Sociology, University of Connecticut, Storrs, CT 06269, United States

^b 744 Ballantine Hall, 1020 E. Kirkwood Ave., Department of Sociology, Indiana University, Bloomington, IN 47405-7103, United States

ARTICLE INFO

Article history:

Received 8 October 2013

Revised 24 March 2015

Accepted 8 April 2015

Available online 23 April 2015

Keywords:

Children

Family structure

Methodology

Same-sex parenting

Sexuality

ABSTRACT

Scholars have noted that survey analysis of small subsamples—for example, same-sex parent families—is sensitive to researchers' analytical decisions, and even small differences in coding can profoundly shape empirical patterns. As an illustration, we reassess the findings of a recent article by Regnerus regarding the implications of being raised by gay and lesbian parents. Taking a close look at the New Family Structures Study (NFSS), we demonstrate the potential for misclassifying a non-negligible number of respondents as having been raised by parents who had a same-sex romantic relationship. We assess the implications of these possible misclassifications, along with other methodological considerations, by reanalyzing the NFSS in seven steps. The reanalysis offers evidence that the empirical patterns showcased in the original article are fragile—so fragile that they appear largely a function of these possible misclassifications and other methodological choices. Our replication and reanalysis of the study offer a cautionary illustration of the importance of double checking and critically assessing the implications of measurement and other methodological decisions in our and others' research.

Reproducible Results and Workflow | 4

Is science broken?

Misconduct, fraud, and retractions

- Peer review was circumvented at prestigious journals
- Two journals published Maggie Simpson & Edna Krabappel's "Fuzzy, Homogeneous Configurations"
- Retraction at *Science* when data not found

Science Isn't Broken by Christie Aschwanden

"I've learned that the headline-grabbing cases of misconduct and fraud are mere distractions. The state of our science is strong, but it's plagued by a universal problem: **Science is hard – really f*ing hard.**"

"If we're going to rely on science as a means for reaching the truth - and it's still the best tool we have - **it's important that we understand and respect just how difficult it is to get a rigorous result.**"

Reproducible Results and Workflow | 5

Replication and reproduction of results

Reproducibility requires identical results with the same data.

Replicability required confirmation of results with new data.

		Reproducibility	
		High	Low
Replicability	High	<i>Scientific Ideal</i>	<i>Careless Research</i>
	Low	<i>Fragile Findings</i>	<i>Unscientific Work</i>

Reproducible Results and Workflow | 6

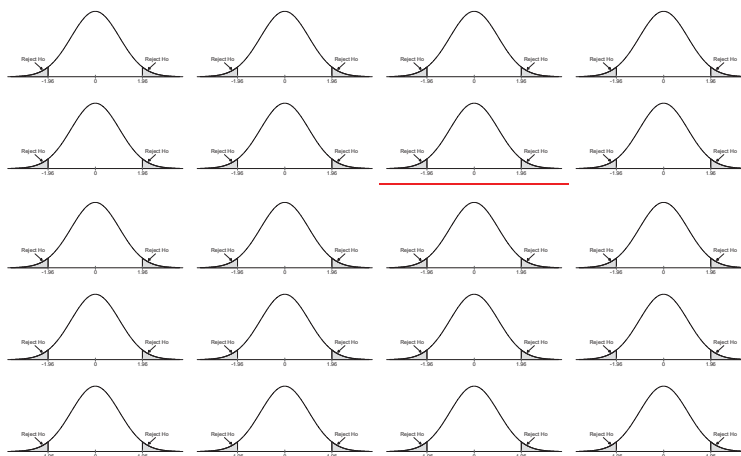
Challenges to replicability

Sample driven analyses

1. Decisions based on unique characteristics of the sample invalidate statistical tests.
2. Examples
 - o Data mining portrayed as theory testing
 - o Post analysis hypothesis construction
 - o Undocumented specification searches and p-hacking
 - o “Cherry picking” the sample
3. Consider the effect on the sampling distribution of a test statistics...

Reproducible Results and Workflow | 7

If $\alpha=0$, twenty tests of $H_0: \alpha=0$ at the 5% level



Reproducible Results and Workflow | 8

Example of using sample to select a model

1. Randomly select six sub-samples.
2. Use stepwise logit to select a model predicting diabetes.
3. Seven different models were selected.

Variable	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6
bmi	1.067***	1.066***	1.004	1.074***	1.101***	0.971
white	0.518***	0.547***	0.521***	0.543***	0.505***	0.562***
age	1.262***	1.351***	1.324***	1.288***	1.282***	1.341***
agesq	0.999***	0.998***	0.998***	0.998***	0.998***	0.998***
hsdegree	0.720***	0.680***	0.662***	0.749***	0.780***	0.650***
weight	1.006***	1.006***	1.016***	1.004**		1.022***
height			0.936**			0.909***
female				0.854*	0.733***	

Legend: p<.1; ** p<.05; *** p<.01

Reproducible Results and Workflow | 9

Model variability versus sampling variability

Young and Holsteen. 2015. Model Uncertainty and Robustness. SMR.

- o Estimates are sensitive to credible changes in model specification.
- o Point estimates capture just “one ad-hoc route through the thicket of possible models” (Leamer 1985:308)
- o For example, do higher income tax rates cause taxpayers to “vote with their feet” and migrate to states with lower taxes?

Reproducible Results and Workflow | 10

Effects of tax rate on migration

- o Estimate is significant in only 1.5% of 24,567 models.
- o The mean estimate is roughly zero.

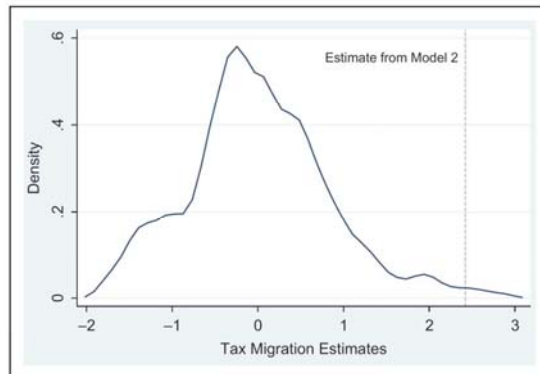


Figure 4. Modeling distribution of tax migration estimates.
Note: Kernel density graph of estimates from 24,576 models.

Reproducible Results and Workflow | 11

Reproducibility with same data

Changing expectations for reproducibility

- AJPB requires verification of results before a paper is published.
 - Only ~5 of 200 submissions succeeded.
- Many journals required that data and script files are distributed.

Challenges to reproducibility

- Reproducibility requires a systematic workflow built around the requirement of reproducibility.

My talk focuses on the workflow for reproducibility

Reproducible Results and Workflow | 12

What is a workflow for data analysis?

A workflow is a set of coordinated procedures for all aspects of data management, analysis, and presentation.

- Planning research
- Organizing and documenting
- Importing and cleaning data
- Analyzing data
- Presenting and publishing results
- Revising results
- Preserving files



Reproducible Results and Workflow | 13

You have a workflow

1. Your WF might be:

- **Planned**
- **Ad hoc**
- **Planned in an ad hoc way**

2. You can improve your WF with a modest investment of time.

- The less experience you have, the easier it is.
- It takes time, but saves more time.
- It prevents errors.
- It makes you a better data analyst.
- It is critical for reproducibility.

Reproducible Results and Workflow | 14

Origins of the workflow project

1. Incorrect results with clever explanations
2. Dissertation delayed 18 months to determine provenance
3. Unreproducible results from a 743 line do-file
4. Analyzing the wrong data set:
"The datasets are exactly the same except for the married variable."
5. The wrong variable when writing a report for the NAS
6. Mislabeled gene in a study of alcoholism
7. Collaborations that multiply the ways things go wrong
8. Misleading output such as...

Reproducible Results and Workflow | 15

Definitel a problem

```
. tabulate female sdchild_v1
```

R is female?	Q15 Would let X care for children				Total
	Defintel	Probably	Probably	Defintel	
Male	41	99	155	197	492
Female	73	98	156	215	542
Total	114	197	311	412	1,034

Reproducible Results and Workflow | 16

How important is it to...

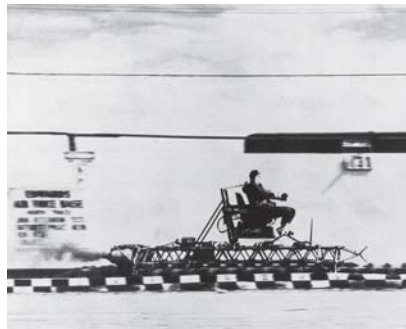
```
. codebook tc1*, compact
```

Variable	Obs	Unique	Mean	Min	Max	Label
tc1doc	1074	10	8.714153	1	10	Q46 How important is it to go to ...
tc1fam	1074	10	8.755121	1	10	Q43 How important is it to turn t...
tc1friend	1073	10	7.799627	1	10	Q44 How important is it to turn t...
tc1mhprof	1045	10	7.58756	1	10	Q48 How important is it to go to ...
tc1psy	1050	10	7.567619	1	10	Q47 How important is it to go to ...
tc1relig	1039	10	5.66025	1	10	Q45 How important is it to turn t...

Reproducible Results and Workflow | 17

The foundation of WF is **ironical optimism**

The **universal aptitude for ineptitude** makes any human accomplishment an incredible miracle. – *Dr. John Paul Stapp*



Reproducible Results and Workflow | 21

40G's: From 0 to 995mph and back in 3 seconds...



"I was fine, only blind for a few days."

Reproducible Results and Workflow | 22

Why are results hard to reproduce?

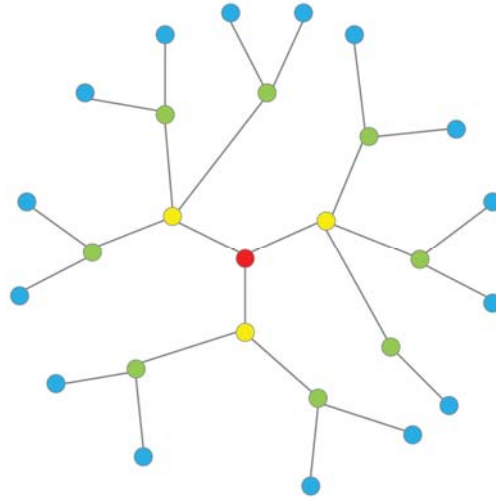
1. **The curse of dimensionality:** Research involves many decisions.

- Where to truncate a variable?
- What seed for the RN generator?
- How to scale with partially missing data?
- Which cases to keep for analysis?
- How to code education?
- What values to assign to income greater than \$200,000?
- And so on...

With only 10 such decisions, there are 1,024 combinations.

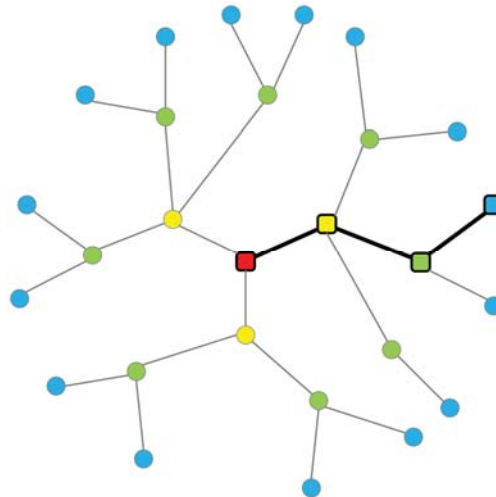
Reproducible Results and Workflow | 23

Decisions in the path to analysis: the choices that could be made



Reproducible Results and Workflow | 24

Decisions in the path to analysis: the choices made



Reproducible Results and Workflow | 25

Why are results hard to reproduce? (continued)

2. With missing documentation, you might not find the right path.
3. Changes in software can lead to different results.
 - A colleague spent weeks to reproduce results because he forgot **version 7** in a do-file.
4. Lost or changed files make reproducibility impossible.
 - Retraction in *Science* because of lost data
 - \$2,000 to retrieve a file that was "backed up"
 - Virtual servers might have 30 day rolling backups

Reproducible Results and Workflow | 26

Criteria for choosing your workflow

Accuracy

- Given reproducibility, you want the correct result

Efficiency

- Completing work quickly
- Working quickly competes with accuracy
- Requires investing time to save time

Scalability

- Adapts to projects of different sizes
- Works with individuals and teams

Reproducible Results and Workflow | 27

Standardization

- Uniform decisions for how to do things
- Increases efficiency and accuracy

Automation

- Saves times and prevents errors
- Time learning automation saves time executing

Usability

- If you won't do it, it is not a good workflow

Transferability

- Can someone else continue your work?

Reproducible Results and Workflow | 28

Collaboration and workflow

1. Collaboration makes it harder to have an effective workflow.
2. Why is workflow harder when you collaborate?

Reproducible Results and Workflow | 29

Coordinating multiple workflows



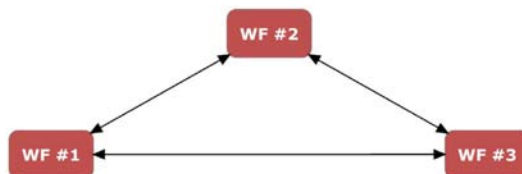
Reproducible Results and Workflow | 30

Coordinating multiple workflows



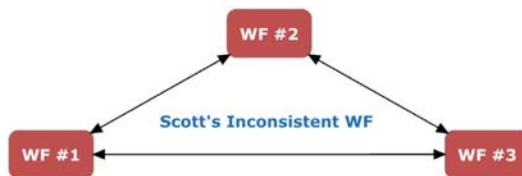
Reproducible Results and Workflow | 31

Coordinating multiple workflows



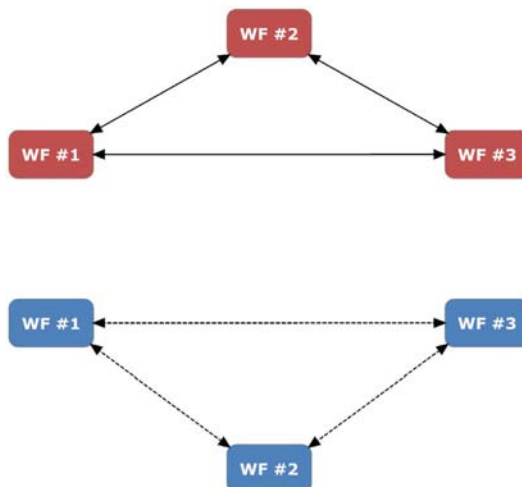
Reproducible Results and Workflow | 32

Coordinating multiple workflows



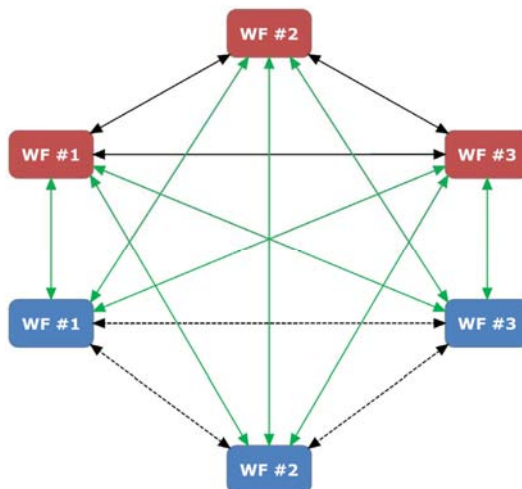
Reproducible Results and Workflow | 33

Coordinating multiple workflows starts here



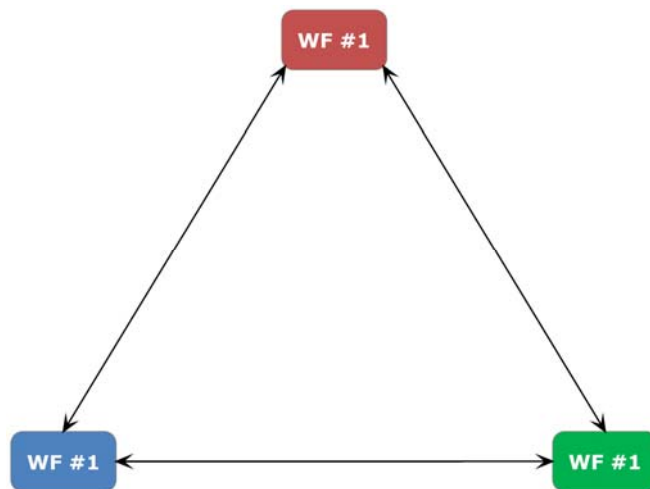
Reproducible Results and Workflow | 34

Coordinating 30 pairs of workflows



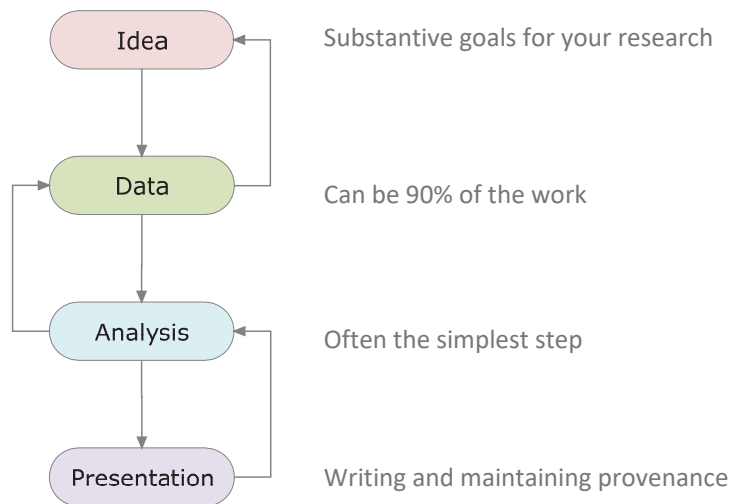
Reproducible Results and Workflow | 35

Coordinating multiple workflows: Agree on a WF



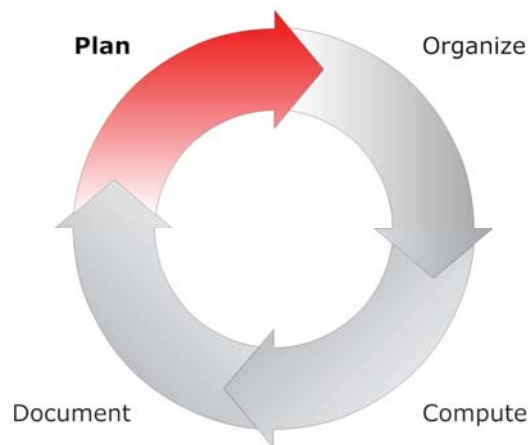
Reproducible Results and Workflow | 36

Steps in your workflow



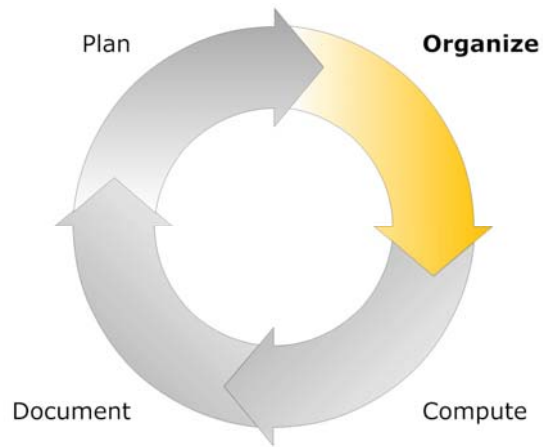
Reproducible Results and Workflow | 37

Tasks within each step



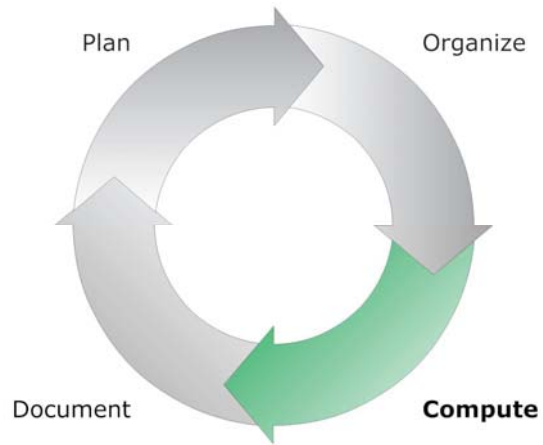
Reproducible Results and Workflow | 38

Tasks within each step



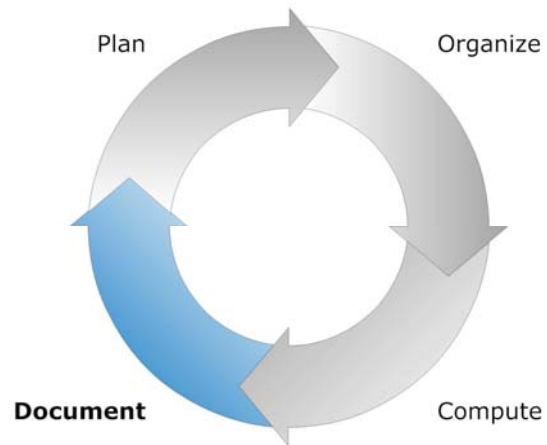
Reproducible Results and Workflow | 39

Tasks within each step



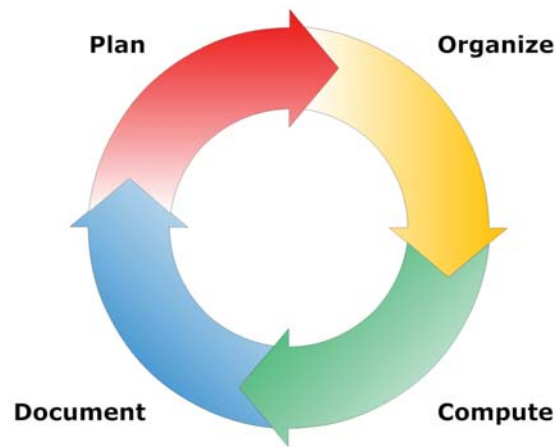
Reproducible Results and Workflow | 40

Tasks within each step



Reproducible Results and Workflow | 41

Tasks within each step



Reproducible Results and Workflow | 42

Planning

- The project timeline
- Division of labor
- Scheduling your time
- How to document and organize research
- Variable names, labels and metadata
- Procedures for missing data
- Analyses
- Writing
- Preserving files
- And more...

Reproducible Results and Workflow | 43

Blau and Duncan's *The American Occupational Structure*

- Analyses were specified 9 months before output was received.
- Book was written from a single set of output.
- Later books with full access to the data were not better.

Michael Faraday's advice

Work. Finish. Publish.

A plan is a reminder to stay on track, finish the project, and publish results.

Reproducible Results and Workflow | 44

Organizing

1. Organization is motivated by two goals

- Finding things
- Avoiding duplication

2. Organization

- Lets you work faster
- Rewards consistency and uniformity
- Is contagious (and so is disorganization)
- Requires regular maintenance to overcome entropy

Reproducible Results and Workflow | 45

Signs of poor organization

1. Can't find a file and you think you deleted it.

2. Multiple versions of a file and you don't know which is which.

- You and a co-author edit different versions of a paper. You have incompatible, incomplete drafts.
- You need the file for draft submitted for review, but you have two (or 16) files with "final" in the name.

This: final report v16.docx

Or this: NSF report 2010-10-21.docx

3. Finally, after this talk a student showed me a text:

- Urgent: don't analyze **final.dta**, use **lastversion.dta** for presentation tomorrow."
- Surely this is a rooky mistake....

Reproducible Results and Workflow | 46

The final paper

The screenshot shows a journal article page for *Political Analysis*. At the top, there is a banner that reads "#1 JOURNAL IN Political Science". Below the banner, the article title is "Erratum for Keele, Linn, and Webb (2016)". The authors listed are Luke Keele, Suzanna Linn, and Clayton McLaughlin Webb. The article is from *Political Analysis*, Spring 2016, issue 24 (2), pages 389-396. The article is available for free full-text download. The page also includes a search bar, a table of contents, and a current issue section.

Reproducible Results and Workflow | 47

Organization should be like a Model T



Reproducible Results and Workflow | 48

Too often it is more like this



Reproducible Results and Workflow | 49

With predictable consequences



Reproducible Results and Workflow | 50

Digital assets and the curse of cheap storage

1. It is easier to create a file than to find a file.
2. It is easier to find a file than to know what is in a file.
3. It is easy to create lots of files.

- 115,000 files on a research center's LAN
- 2,000,000 files accumulated in 10 years

Files are scatter across multiple, overlapping locations

1. Office computer
2. Home computer
3. Laptop
4. LAN
5. Dropbox
6. Box
7. USB sticks
8. Old laptop
9. External drives
10. Mom's computer

Reproducible Results and Workflow | 51

Operating systems organize files for entertainment

Win

Desktop
Music
Pictures
Videos
Documents

Mac

Desktop
Music
Pictures
Movies
Documents

Reproducible Results and Workflow | 52

Digital asset management (DAM)

How important is this?

- How much time do you waste dealing with files?
- How many PDFs do you have of the same article?

How to manage files

1. Name files carefully and systematically.
2. Use a planned directory structure.
 - Every file has one place it belongs.
 - A file's location documents the file.

For example...

Reproducible Results and Workflow | 53

A planned set of primary directories

\- To shelve	Files to put in the correct directory
\Active	Active projects
\Admin	Administration and service
\Bookshelf	Books, articles, reprints, etc.
\Inactive	Projects that are on hold
\Shared	Files shared with others on the cloud
\Teaching	Teaching materials
\Templates	Files used as templates
\Vault	Completed work that will <i>never</i> change

Reproducible Results and Workflow | 54

A structure for projects in \Active, \Inactive and \Teaching

\Group Differences

- \- Hold then delete
- \- To shelve
- \Admin
- \Posted
- \Resources
- \Work
- \Write

Reproducible Results and Workflow | 55

File naming

Writing

- groups 2017-11-07.docx
- groups 2018-01-17.docx

PDF files

- Long 1978 ASR productivity position.PDF

Datasets

- groups-hrs1.dta
- groups-hrs2.dta

Script files

- groups-data03-recoding.do
- groups-data04-scales.do

Reproducible Results and Workflow | 56

Organization: uniform formats for robust script files

```
capture log close
log using wftalk01-example, replace text
version 15.1
clear all
set linesize 80

// project: wf talk
local pgm wftalk01
local dte 2018-01-18
local who Scott Long
local tag "`pgm'.do `who' `dte'"
di "`tag'" // for provenance

// #1 describe task

// #2 describe task

log close
exit
```

Reproducible Results and Workflow | 57

Documentation

1. Without documentation,
 - Reproduction is much more difficult.
 - Mistakes are more likely.
 - Work takes longer.
2. Long's Law: It is faster to document it today than tomorrow.
 - Nobody likes to write documentation.
 - Nobody regrets having documentation.
3. More codified fields demand documentation.
 - *The Research Log* (American Chemical Society)

Reproducible Results and Workflow | 58

Suggestions for writing documentation

1. Do it today.
2. Check it next week even if it makes sense today.
3. Review it at key stages of your work, like finishing a draft.
4. Include full dates and names.
5. Use reinforcing, non-redundant forms of documentation.
6. Start with a research diary for each project.

Reproducible Results and Workflow | 59

A simple research diary

First complete set of analysis for FLIM measures paper

f2alt01a.do - 24May2002

Descriptive information on all rhs, lhs, and flim measures

f2alt01b.do - 25May2002

Compute bic' for each of four outcomes and all flim measures.

```
** Outcome: Can Work          global lhs "qcanwrk95"  
** Outcome: Work in three categories global lhs "wdlthwrk95"  
** Outcome: bath trouble     global lhs "bathdif95"  
** Outcome: adlsum95 - sum of adls global lhs "adlsum95"
```

f2alt01c.do - 25May2002

Compute bic' for each of four outcomes and with only these restricted flim measures.

```
* 1. ln(x+.5) and ln(x+1)  
* 2. 9 counts: >=5&5 >=7&7 (50% and 75%)  
* 3. 8 counts: >=4&4 >=6&6 (50% and 75%)  
* 4. 18 counts: >=9&9 >=14&14 (50% and 75%)  
* 5. probability splits at .5; these don't work well in prior tests
```

f2alt01d.do - 25May2002

bic' for all four outcomes in models that include all raw flim measures (fla*p5; fill*p5); pairs of u/l measures; groups of LCA measures

f2alt01e.do - all LCA probabilities - 25May2002

:::

f2alt01j.do - use three probability measures from LCA - 29May2002

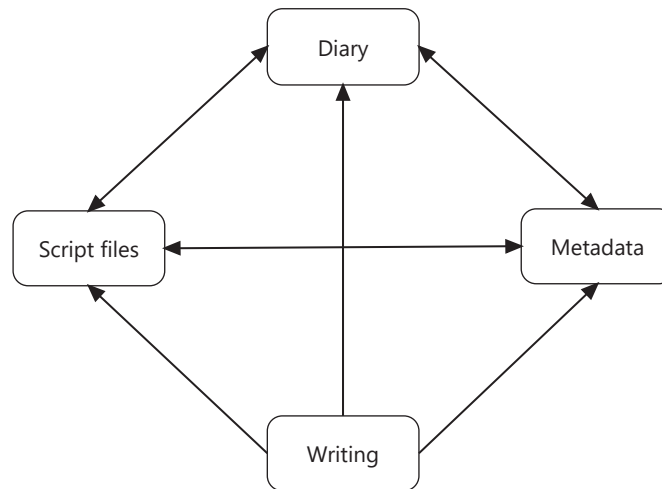
:::

f2alt02c.do - 29May2002

use three binary variables, not just LC class numbers.
: dummies work better than the class number;
: effects of LCA and group variables are highly correlated

Results and Workflow | 60

Reinforcing forms of documentation



Reproducible Results and Workflow | 61

Execution and computing

1. Execution involves carrying out tasks within each step.

2. Effective execution requires mastering tools.

- Software

- File manager

- Macro program

- Text editor

- Statistical software

- Hardware

3. Planning is more important than computing power.

- Consider the changes in computing...

Reproducible Results and Workflow | 62

Cornell 1975: the entire computing infrastructure



IBM 370 with 240K memory



Winchester drives with 3MB storage

- **Cost of computing \$1,000,000.**
- **Mean time to degree 7.6 years.**

Reproducible Results and Workflow | 63

Laptop 2009



Asus 1000HE with 2GB memory
10,000 times more



Free Agent with 1TB storage
350,000 times more

- **Cost of computing \$400.**
- **Mean time to degree 7.6 years.**

Reproducible Results and Workflow | 64

A thought experiment on planning and computing

1. Divide yourselves into two groups.

Computers compute any time they want to.

Planners compute only 12 hours a week.

2. Who will finish their dissertation first?

Reproducible Results and Workflow | 65

Principles for a computing workflow

1. Legible and robust script files
2. Posting files
3. Dual workflow for data management and analysis
4. Run order naming of scripts

Reproducible Results and Workflow | 66

Robust and legible script files

1. Programs must run on another computer without *any* changes.
 - Self-contained
 - Version control
 - No hard coded directory information
 - Explicit seeds for random numbers
 - Archived user written programs
2. Careful internal documentation of what the script does.
3. Formatting to improve legibility.

Reproducible Results and Workflow | 67

The *essential* posting principle

1. Posting is defined by two simple rules.

The share rule

Only share results after the files are posted.

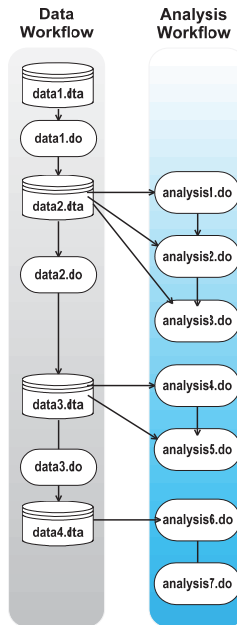
The no change rule

Once a file is posted, *never* change it.

2. Without posting, you cannot reproduce your results.
 - If you don't have the dataset, how can you confirm the results?
 - If you don't have the scripts, how were the results produced?

Reproducible Results and Workflow | 68

Dual workflow and run order naming



Reproducible Results and Workflow | 69

Data cleaning, including names and labels

Planning labels

Bad labels

```
. codebook tc1*, compact
```

Variable	Obs	Unique	Mean	Min	Max	Label
tc1doc	1074	10	8.714153	1	10	Q46 How important is it to go to ...
tc1fam	1074	10	8.755121	1	10	Q43 How important is it to turn t...
tc1friend	1073	10	7.799627	1	10	Q44 How important is it to turn t...
tc1mhprof	1045	10	7.58756	1	10	Q48 How important is it to go to ...
tc1psy	1050	10	7.567619	1	10	Q47 How important is it to go to ...
tc1relig	1039	10	5.66025	1	10	Q45 How important is it to turn t...

Better labels

```
. codebook tc2*, compact
```

Variable	Obs	Unique	Mean	Min	Max	Label
tc2doc	1074	10	8.714153	1	10	Q46 How Impt: Go to a gen med doc...
tc2fam	1074	10	8.755121	1	10	Q43 How Impt: Turn to family for ...
tc2friend	1073	10	7.799627	1	10	Q44 How Impt: Turn to friends for...
tc2mhprof	1045	10	7.58756	1	10	Q48 How Impt: Go to a mental heal...
tc2psy	1050	10	7.567619	1	10	Q47 How Impt: Go to a psych for Help
tc2relig	1039	10	5.66025	1	10	Q45 How Impt: Turn to a religious...

Reproducible Results and Workflow | 70

Even better labels

```
. codebook tc3*, compact
```

Variable	Obs	Unique	Mean	Min	Max	Label
tc3doc	1074	10	8.714153	1	10	Q46 Med doctor help important
tc3fam	1074	10	8.755121	1	10	Q43 Family help important
tc3friend	1073	10	7.799627	1	10	Q44 Friends help important
tc3mhprof	1045	10	7.58756	1	10	Q48 MH prof help important
tc3psy	1050	10	7.567619	1	10	Q47 Psychiatric help important
tc3relig	1039	10	5.66025	1	10	Q45 Relig leader help important

Planning labels

	A	B	C	D
1	Number	Name	Value label	Variable labels
2	1	id_iu		Respondent Number
3	2	cntry_iu	cntry_iu	IU Country Number
4	3	vignum	vignum	Vignette
5	4	serious		Q1 How serious would you consider Xs situation to be?
6	5	opfam	Ldummy	Q2_1 What X should do:Talk to family
7	6	opfriend	Ldummy	Q2_2 What X should do:Talk to friends
8	7	tospi	Ldummy	Q2_7 What X should do:Go to spiritual or traditional healer
9	8	tonpm	Ldummy	Q2_8 What X should do:Take non-prescription medication
10	9	oppme	Ldummy	Q2_9 What X should do:Take prescription medication

Reproducible Results and Workflow | 71

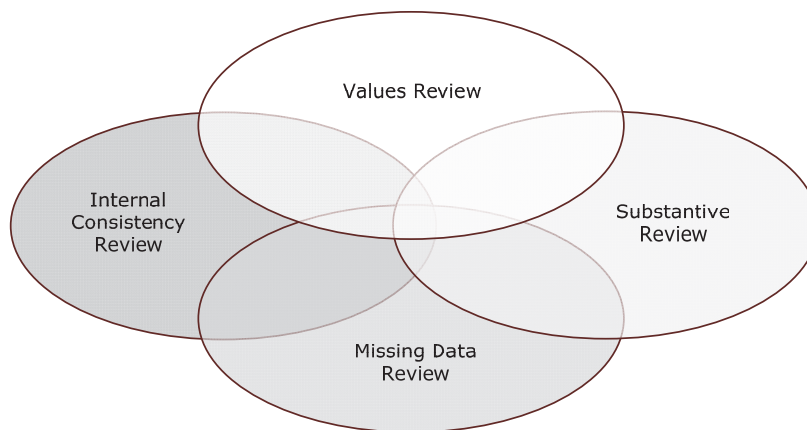
Cost of not planning variables names

1. Confusion between `ownsex` and `ownsexu` caused weeks of delay.
2. Do you want `R003189` or `R001389`?
3. Is `timetophd` elapsed time or enrolled time?

Reproducible Results and Workflow | 72

Data cleaning and preventing retractions

Statistical analysis assumes the variables are clean.



Reproducible Results and Workflow | 73

A two-way table would have detected the problem

RETRACTED: In Sickness and in Health? Physical Illness as a Risk Factor for Marital Dissolution in Later Life

Abstract

The health consequences of marital dissolution are well known, but little work has examined the impact of health on the risk of marital dissolution. In this study we use a representative sample of 2,701 marriages from the Health and Retirement Study (1992–2010) to examine the role of various physical illness onset (i.e., cancer, heart problems, lung disease, and/or stroke) in subsequent marital dissolution due to either divorce or widowhood. We use a series of discrete-time event history models with competing risks to estimate the impact of husband's and wife's physical illness onset on risk of divorce and widowhood. We find that only wife's illness onset is associated with elevated risk of divorce, while either husband's or wife's illness onset is associated with elevated risk of widowhood. These findings suggest the importance of health as a determinant of marital dissolution in later life via both individual and gendered social pathways.

Keywords

aging, chronic disease, gender inequality, marital health

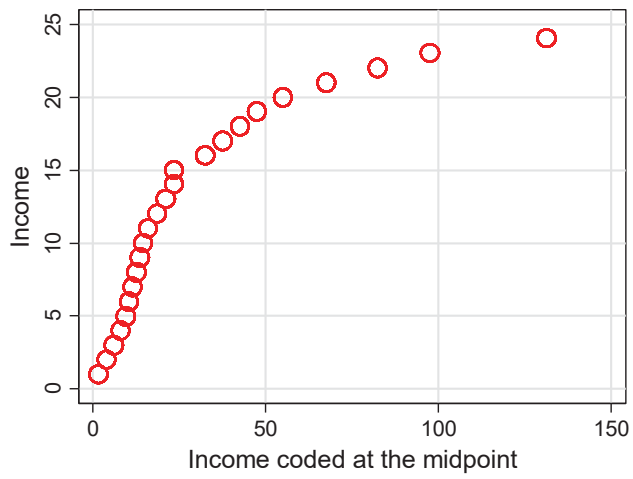
A large body of literature has identified marital status as a strong predictor of health and well-being. Not only are unmarried healthier than the married (e.g., Lillard and Willis 1995; Emmons 1992), but studies find that both divorce and widowhood are pathways to increased physical and mental health (e.g., Hughes and Waite 2009; Williams and Uchino 2003). Little attention, however, has been paid to how health may be a determinant of marital status. Work in this area has tended to focus on the positive selection of the healthier into marriage (e.g., Byrne et al. 1989; Smith and Smith 2010), but poor health may be an equally important force for selection.

Booth, and Johnson 2006). Illness may initiate changes to spouses' roles—in particular, increasing caregiving responsibilities for the healthy spouse—which can tax marital relationship dynamics (Wolff and Kasper 2006). Illness may also decrease household income due to the inability of one or both spouses to work (Teachman 2010), which may increase marital strain.

Only a few studies have examined the role of poor health in subsequent divorce, and these studies are mixed in their findings, with some finding

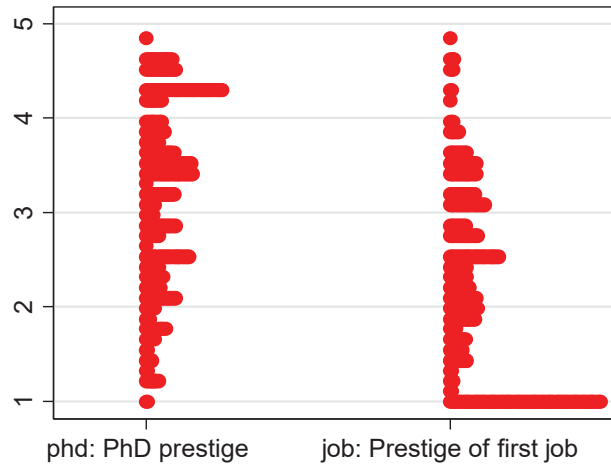
Reproducible Results and Workflow | 74

Use graphs to find errors



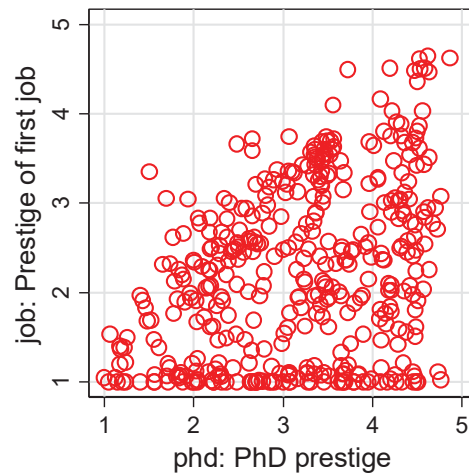
Reproducible Results and Workflow | 75

Graphs highlight forgotten coding decisions



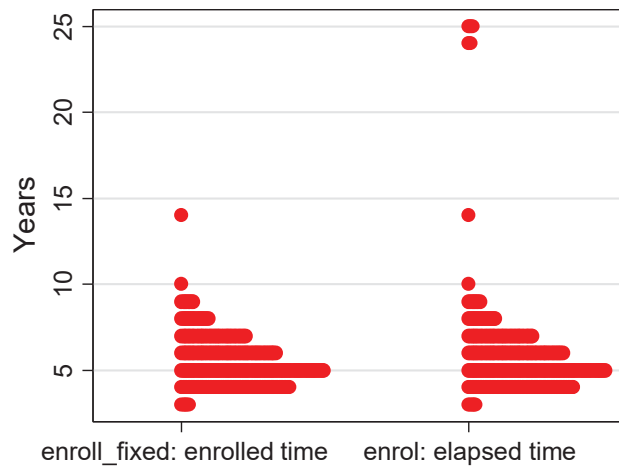
Reproducible Results and Workflow | 76

Locating outliers and gaining substantive insights



Reproducible Results and Workflow | 77

Avoiding expensive mistakes from misread documentation



Reproducible Results and Workflow | 78

Statistical analysis

This can be the simplest part of the project.

1. Take classes and go to talks on data analysis.
2. Find exemplars in the best journals.
3. Use automation and script files.
4. Maintain a dual workflow to prevent errors.

Presentations and provenance

1. Content and methods are disciplinary decisions.
2. Standards for presentations vary by discipline.
 - o Bad presentations transcend disciplines
3. Maintaining provenance is critical for reproducibility.

Reproducible Results and Workflow | 79

Documenting provenance

The provenance of every number must be documented.

1. The circled text contains results I may need to confirm later:

1922-1926 cohort, employed women have fewer limitations than those who are out for family reasons, (.48 and .73, respectively ($z=2.55, p<.01$)). However, this gap has disappeared for the 1943-1947 cohort and, indeed, employed women have slightly more limitations (.76 for non-

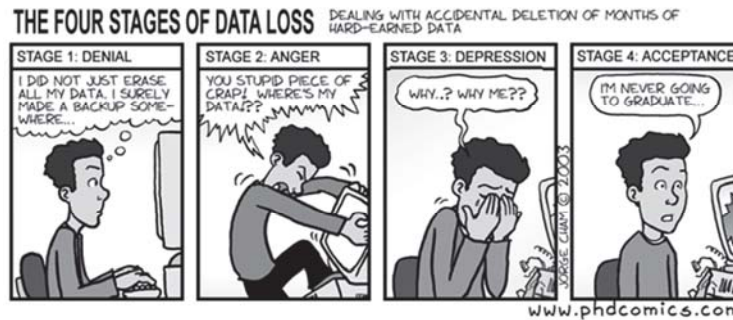
2. Turning on "show/hide ¶" reveals the provenance:

1922-1926 cohort, employed women have fewer limitations than those who are out for family reasons, (.48 and .73, respectively ($z=2.55, p<.01$ {cwhrr-fig03c-hrmemp4.do #4 jsl 17May06})). However, this gap has disappeared for the 1943-1947 cohort and, indeed, employed women have

Reproducible Results and Workflow | 80

Preserving your files

Expect things to go wrong, expect to delete the wrong file at the worst possible time, and expect a hose to be left on in the room above your computer. If you expect the worst, you might prevent it.



Reproducible Results and Workflow | 81

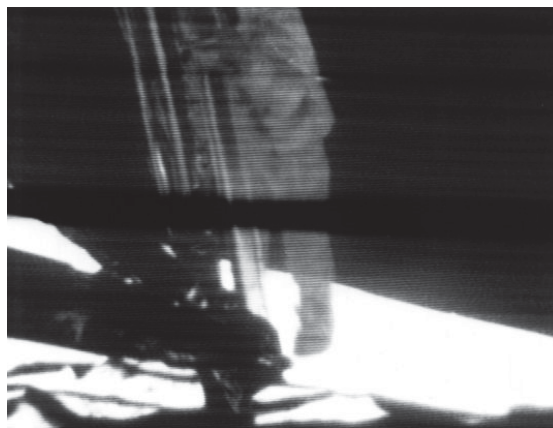
Examples of data loss

1. Water leaked above ICPSR server room.
2. 508K volumes are in obsolete formats at British Museum.
3. Data for Wolfgang's *Delinquency in a Birth Cohort* burned.
4. NASA lost 1000s of moon tapes.

Reproducible Results and Workflow | 82

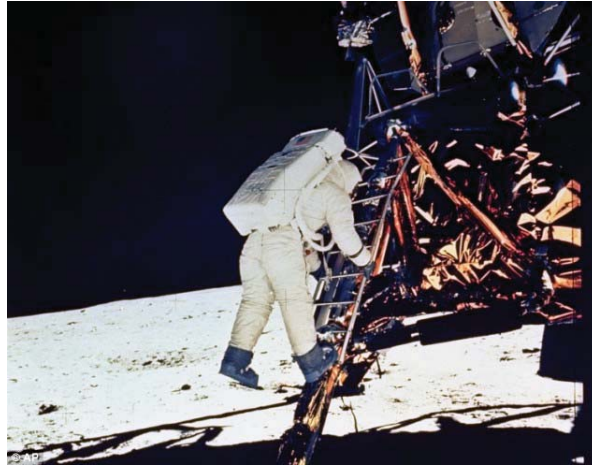
Neil Armstrong's walk on the moon

In 1969, America saw "a fuzzy gray blob wading through an inkwell".



Reproducible Results and Workflow | 83

NASA had video that was too good to show on TV, but lost the tapes.



Reproducible Results and Workflow | 84

The produce for Pink Floyd rock video archived two moon tapes!



Dark Side of the Moon

Reproducible Results and Workflow | 85

Preserving files does not preserve content

“These files were saved six years ago as Gauss FMT files. We need to revise a paper and need the data in these files, but I can’t open them. We have an old version of Gauss that doesn’t run anymore. Any ideas?”

Reproducible Results and Workflow | 86

Conclusions

Replicability and reproducibility

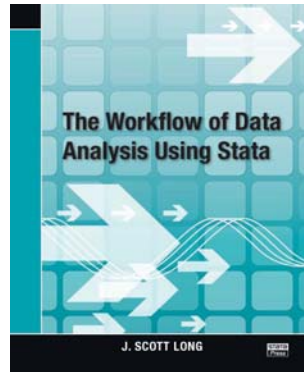
- Expectations for replicability and reproducibility are growing.
- These *positive* developments increase demands for researchers.
- An effective workflow is essential.

Changing your workflow

- Slowly, systematically, thoughtfully.
- Finish the last 5% of each change.
- Do not do it under deadline.

Whose workflow

- There are *many* viable workflows, but it is nice to have a place to start.



Reproducible Results and Workflow | 87

Thank you!

Reproducible Results and Workflow | 88