

# Human-behavioral challenges in reproducibility & replicability

@jtleek

[jtleek.com/talks](https://jtleek.com/talks)

[@jtleek](#)

## FIXING SCIENCE

# Most science research findings are false. Here's how we can change that



POLICY &amp; ETHICS

# Is There a Reproducibility Crisis in Science?

By Nature Video on May 28, 2016





NATURE | NEWS FEATURE



[E-alert](#)

[RSS](#)

[Facebook](#)

[Twitter](#)

# 1,500 scientists lift the lid on reproducibility

Survey sheds light on the 'crisis' rocking research.

**Monya Baker**

25 May 2016 | Corrected: [28 July 2016](#)

[PDF](#)

[Rights & Permissions](#)

Is there a reproducibility crisis in science?



## Space race



### China's quest to become a space science superpower

With major spaceflight milestones behind it, China is working to build an international reputation for space science.

**GILSON**

See how Gilson makes lab life easier

All

Images

Videos

News

Shopping

More

Settings

Tools

About 176,000,000 results (0.43 seconds)

## Most Scientific Findings Are Wrong or Useless - Reason.com

[reason.com/archives/2016/08/26/most-scientific-results-are-wrong-or-useless](http://reason.com/archives/2016/08/26/most-scientific-results-are-wrong-or-useless)

Aug 26, 2016 - ScientistYanlevDreamstime Yanlev/Dreamstime"Science, the pride of modernity, our one source of objective knowledge, is in deep trouble.

## PLOS Medicine: Why Most Published Research Findings Are False

[journals.plos.org/plosmedicine/article?id=10.1371/journal.pmed.0020124](http://journals.plos.org/plosmedicine/article?id=10.1371/journal.pmed.0020124)

by JPA Ioannidis - 2005 - Cited by 4846 - Related articles

Aug 30, 2005 - Moreover, for many current **scientific** fields, claimed research findings ... Citation: Ioannidis JPA (2005) Why **Most** Published Research Findings Are False. .... what might have gone **wrong** with their data, analyses, and results.

## Is Most Published Research Wrong? - YouTube

<https://www.youtube.com/watch?v=42QuXLucH3Q>

Aug 11, 2016 - Uploaded by Veritasium

Why **Most** Published Research Findings Are False: .... The problem with the approach to **science** is that ...

## Believe It Or Not, Most Published Research Findings Are Probably ...

[bigthink.com/.../believe-it-or-not-most-published-research-findings-are-probably-fals...](http://bigthink.com/.../believe-it-or-not-most-published-research-findings-are-probably-fals...)

Ten years ago, a researcher claimed **most** published research findings are false; ... of the Internet has worked wonders for the public's access to **science**, but this ... the case, experiments are underpowered,

176,000,000!



# Replication crisis

From Wikipedia, the free encyclopedia

The **replication crisis** (or **replicability crisis**) refers to a [methodological](#) crisis in [science](#) in which scientists have found that the results of many scientific studies are difficult or impossible to [replicate](#) on subsequent investigation, either by independent researchers or by the original researchers themselves.<sup>[1]</sup> While the crisis has long-standing roots, the phrase was coined in the early 2010s as part of a growing awareness of the problem.


Since the reproducibility of experiments is an essential part of the [scientific method](#), the inability to replicate the studies of others has potentially grave consequences for many fields of science in which significant theories are grounded on unreproducible experimental work.

The replication crisis has been particularly widely discussed in the field of [psychology](#) (and in particular, [social psychology](#)) and in [medicine](#), where a number of efforts have been made to re-investigate classic results, and to attempt to determine both the validity of the results, and, if invalid, the reasons for the failure of replication.<sup>[2][3]</sup>

## Contents [\[hide\]](#)

- [General](#)
- [Medicine](#)
- [Psychology](#)
  - [Replication rates in psychology](#)
  - [A disciplinary social dilemma](#)
- [Marketing](#)

A hypothesis



I HAVE  
AN IDEA.

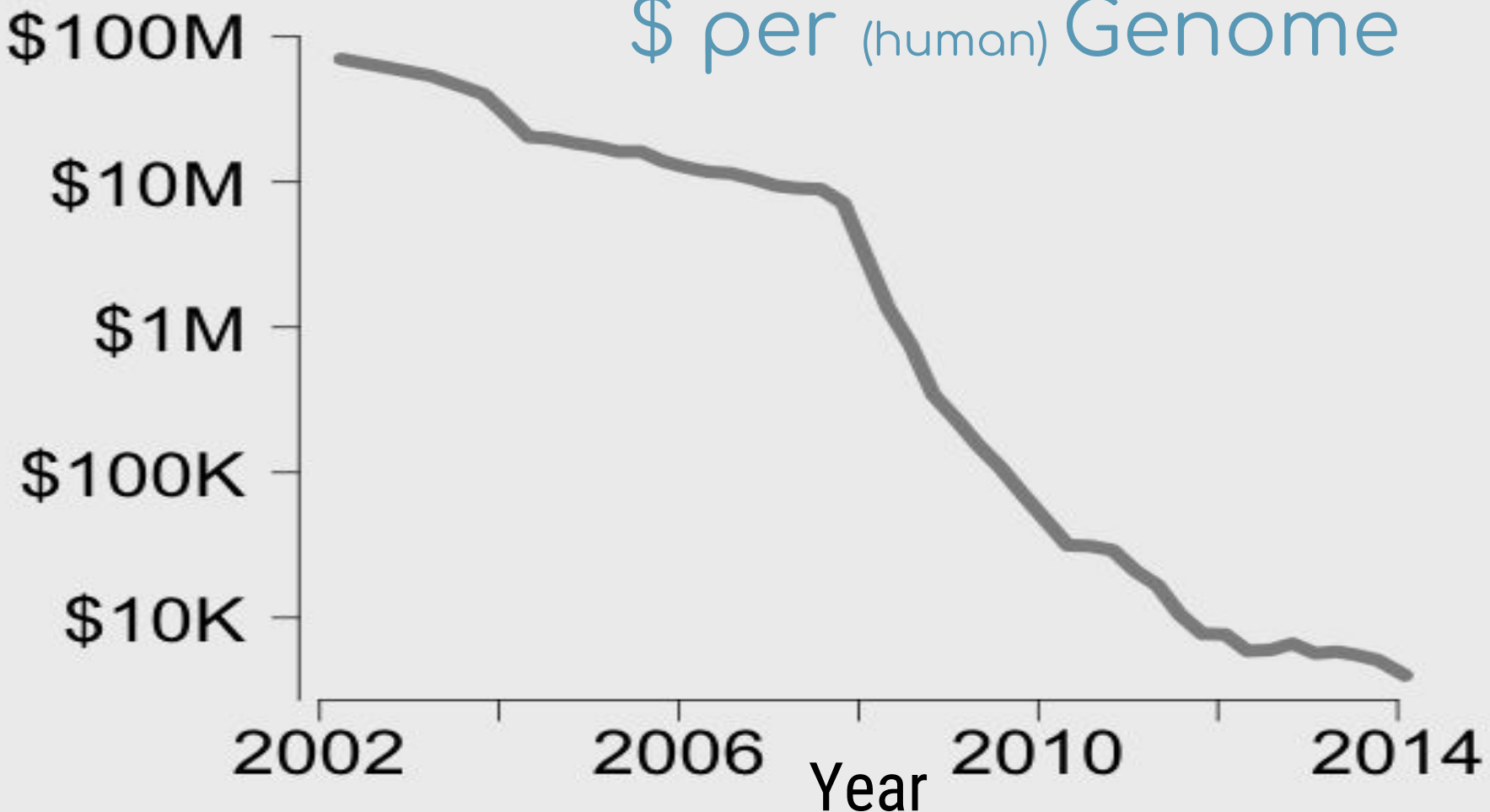
A GOOD  
ONE?

LET'S NOT  
GET AHEAD  
OF OURSELVES.

N = SAMPLE SIZE

$$N = \frac{(\text{€ YOU HAVE})}{(\text{€ PER SAMPLE})}$$

# \$ per (human) Genome





I DON'T KNOW HOW  
TO DO STATISTICS BUT  
IT DOESN'T MATTER  
BECAUSE I DIDN'T  
HAVE DATA.



# Medical school entrance requirements (U.S.)

One year of **biology**

One year of **physics**

One year of **English**

Two years of **chemistry**

# The vast majority of statistical analysis is not performed by statisticians

 Jeff Leek  2013/06/14

Whether you know it or not, everything you do produces data - from the websites you read to the rate at which your heart beats. Until pretty recently, most of the data you produced wasn't collected, it floated off unmeasured. The only data that were collected were painstakingly gathered by scientists one number at a time in small experiments with a few people. This laborious process meant that data were expensive and time-consuming to collect. Yet many of the most amazing scientific discoveries over the last two centuries were squeezed from just a few data points. But over the last two decades, the unit price of data has dramatically dropped. New technologies touching every aspect of our lives from our money, to our health, to our social interactions have made data collection cheap and easy (see e.g. [Camp Williams](#)).

To give you an idea of how steep the drop in the price of data has been, in 1967 Stanley Milgram [did an experiment](#) to determine the number of degrees of separation between two people in the U.S. In his experiment he sent 296 letters to people in Omaha, Nebraska and Wichita, Kansas. The goal was to get the letters to a specific person in Boston, Massachusetts. The trick was people had to send the letters to someone they knew, and they then sent it to someone they knew and so on. At the end of the experiment, only 64 letters made it to the individual in Boston. On average, the letters had gone through 6 people to get there. This is where the idea of “6-degrees of Kevin Bacon” comes from. Based on 64 data points. [A 2007 study](#) updated that number to “7

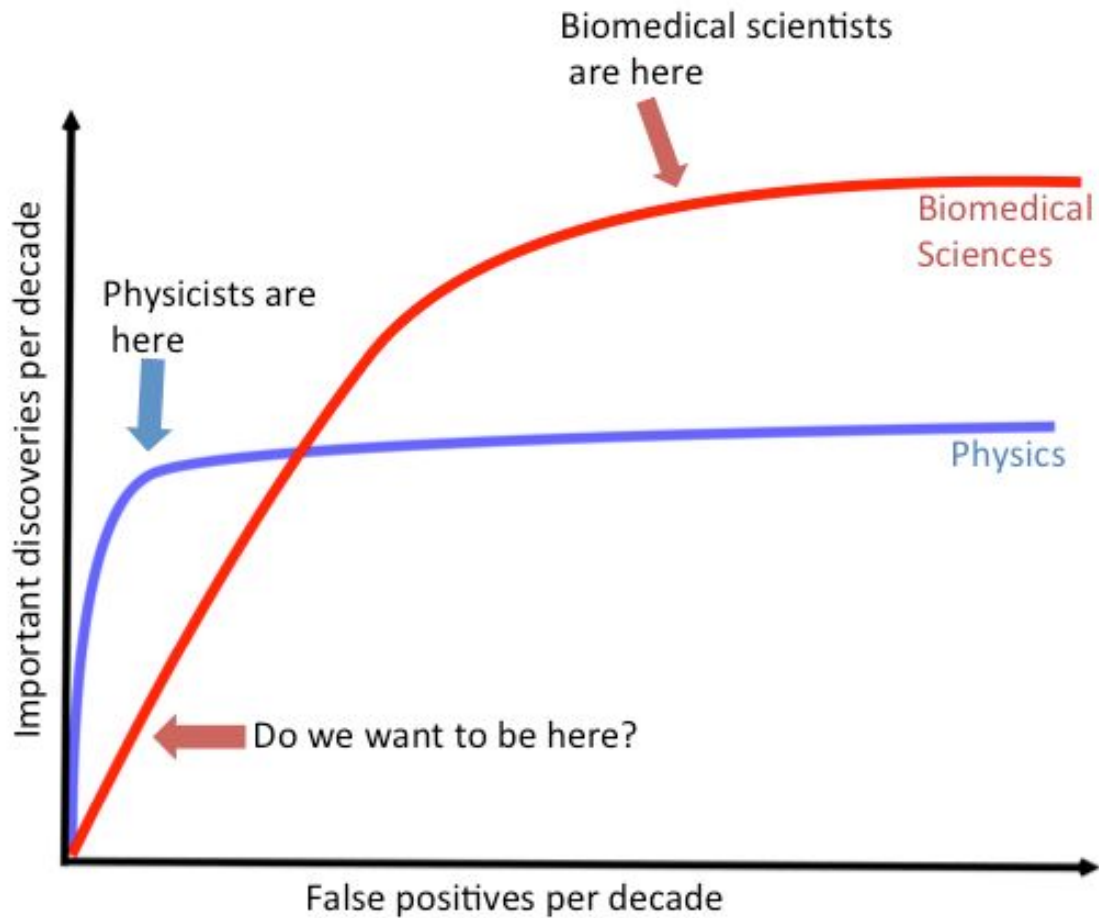
<https://simplystatistics.org/2013/06/14/the-vast-majority-of-statistical-analysis-is-not-performed-by-statisticians/>

3

We need to say what we mean

The tools already exist

Humans are the problem



3

We need to say what we mean

The tools already exist

Humans are the problem



Population



Question



Hypothesis



Experimental Design



Experimenter



Data



Analysis Plan



Analyst



Code



Estimate



Claim

**SHARE**

PERSPECTIVE | SCIENTIFIC INTEGRITY



0



0

# What does research reproducibility mean?

Steven N. Goodman\*, Daniele Fanelli and John P. A. Ioannidis

+ See all authors and affiliations

*Science Translational Medicine* 01 Jun 2016:  
Vol. 8, Issue 341, pp. 341ps12  
DOI: 10.1126/scitranslmed.aaf5027

**Article**

Figures &amp; Data

Info &amp; Metrics

eLetters

PDF

You are currently viewing the abstract.

[View Full Text](#)

## Abstract

The language and conceptual framework of “research reproducibility” are nonstandard and unsettled across the sciences. In this Perspective, we review an array of explicit and implicit definitions of reproducibility and related terminology, and discuss how to avoid potential misunderstandings when these terms are used as a surrogate for “truth.”

Copyright © 2016, American Association for the Advancement of Science

[View Full Text](#)

## Science Translational Medicine

Vol 8, Issue 341  
01 June 2016[Table of Contents](#)**ARTICLE TOOLS**

Email

Print

Alerts

Citation tools

Download Powerpoint

Save to my folders

Request Permissions

Share

Advertisement





**SHARE****PERSPECTIVE**

0



0

# Reproducible Research in Computational Science

**Roger D. Peng**[+ See all authors and affiliations](#)

*Science* 02 Dec 2011:  
Vol. 334, Issue 6060, pp. 1226-1227  
DOI: 10.1126/science.1213847

**Article****Figures & Data****Info & Metrics****eLetters** **PDF**

## Abstract

Computational science has led to exciting new developments, but the nature of the work has exposed limitations in our ability to evaluate published findings. Reproducibility has the potential to serve as a minimum standard for judging scientific claims when full independent replication of a study is not possible.

...

The rise of computational science has led to exciting and fast-moving developments in many scientific areas. New technologies, increased computing power, and methodological advances have dramatically improved our ability to collect complex high-dimensional data (1, 2). Large data sets have led to scientists doing more

**Science**

Vol 334, Issue 6060  
02 December 2011

[Table of Contents](#)  
[Print Table of Contents](#)  
[Advertising \(PDF\)](#)  
[Classified \(PDF\)](#)  
[Masthead \(PDF\)](#)

**ARTICLE TOOLS**

- Email
- Download Powerpoint
- Print
- Save to my folders
- Alerts
- Request Permissions
- Citation tools
- Share

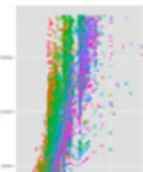
Advertisement

**DataCamp**

basic plotting



ggplot2



`...I'm actually not sure we all agree on meaning of "replication" / "replicability" / "replicate", nor is this term necessarily used consistently in the literature'

Original

Reproduction



01100  
10110  
11110

01100  
10110  
11110



# Reproduce



Original



Unobserved



Different



Incorrect

Original



01100  
10110  
11110



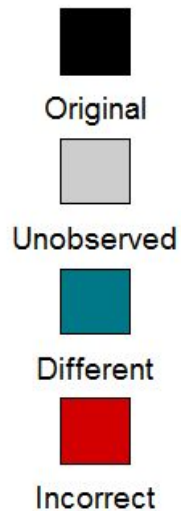
Replication



01100  
10110  
11110



# Replicate





01100  
10110  
11110



# False discovery



Original



Unobserved



Different



Incorrect

## SHARE

### RESEARCH ARTICLE

replicability

# Estimating the ~~reproducibility~~ of psychological science



0



19

Open Science Collaboration<sup>\*†</sup>

+ See all authors and affiliations

Science 28 Aug 2015:  
Vol. 349, Issue 6251, aac4716  
DOI: 10.1126/science.aac4716

[Article](#)

[Figures &  
Data](#)

[Info &  
Metrics](#)

[eLetters](#)

[PDF](#)



## Science

Vol 349, Issue 6251  
28 August 2015

- [Table of Contents](#)
- [Print Table of Contents](#)
- [Advertising \(PDF\)](#)
- [Classified \(PDF\)](#)
- [Masthead \(PDF\)](#)

### ARTICLE TOOLS



[Email](#)



[Print](#)



[Alerts](#)



[Citation tools](#)



[Download Powerpoint](#)



[Save to my folders](#)



[Request Permissions](#)



[Share](#)



Speaking of Science

# Many scientific studies can't be replicated. That's a problem.



By **Joel Achenbach** August 27, 2015





Over the course of four years, 270 researchers attempted to ~~reproduce~~ replicate the results of 100 experiments that had been published in three prestigious psychology journals. It was awfully hard. They ultimately concluded that they'd succeeded just 39 times.



Payne et. al.

Vianello (OSF)



01100  
10110  
11110

01100  
10110  
11110



Original



Unobserved



Different



Incorrect

3

We need to say what we mean

The tools already exist

Humans are the problem



Overview Repositories 81 Stars 7 Followers 4.5k Following 6

Popular repositories

Customize your pinned repositories

**datasharing**

The Leek group guide to data sharing

★ 4k 🍴 175k

**dataanalysis**

The lecture slides for Coursera's Data Analysis class

JavaScript ★ 636 🍴 631

**rpackages**

R package development - the Leek group way!

★ 337 🍴 247

**genomicspapers**

The Leek group guide to genomics papers

★ 245 🍴 114

**reviews**

Writing reviews of academic papers

★ 206 🍴 61

**capitalIn21stCenturyinR**

Piketty in R

HTML ★ 197 🍴 125

**Jeff L.**  
jtleek

[Add a bio](#)

Developer Program Member

Baltimore, MD

<http://biostat.jhsph.edu/~jleek/>

# R Markdown

from  Studio

[Get Started](#)

[Gallery](#)

[Formats](#)

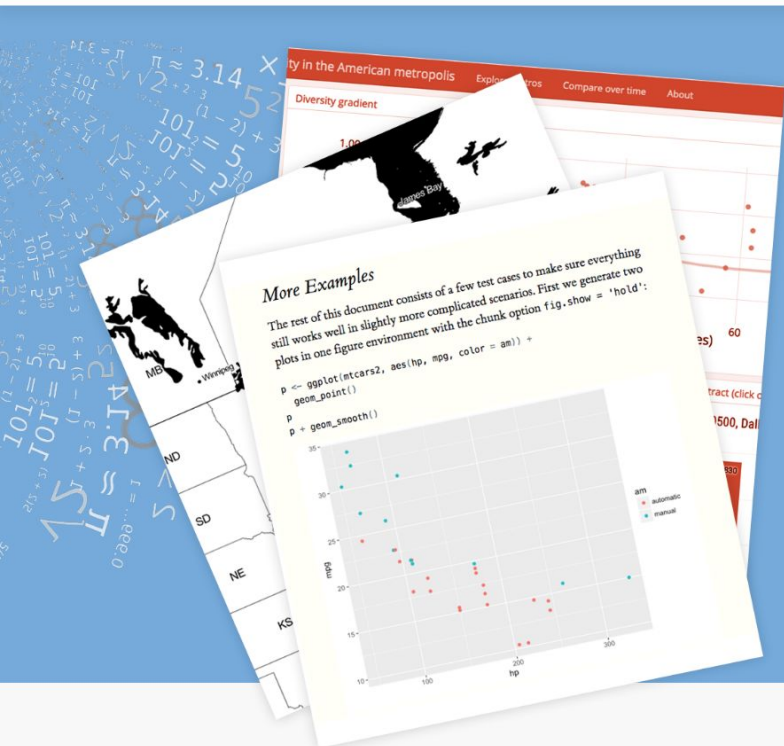
[Articles](#)



Analyze. Share. Reproduce.

Your data tells a story. Tell it with R Markdown.

Turn your analyses into high quality documents, reports, presentations and dashboards.





store, share, discover **research**

get more citations for all of the outputs of your academic research  
over 5000 citations of figshare content to date

ALSO FOR **INSTITUTIONS** & **PUBLISHERS**

*"figshare wants to open up scientific data to the world"*

**WIRED**



Help support open science today.

[Donate Now](#)

Preregistration makes your science better.



Increase the credibility of your research with preregistration. Preregistered research that is published by **December 31, 2018** can be



## About *Bioconductor*

Bioconductor provides tools for the analysis and comprehension of high-throughput genomic data. Bioconductor uses the R statistical programming language, and is open source and open development. It has two releases each year, [1560 software packages](#), and an active user community. Bioconductor is also available as an [AMI](#) (Amazon Machine Image) and a series of [Docker](#) images.

## News

- [BioC 2018](#) conference web site.
- **Core team job opportunities for [scientific programmer / analyst](#) and [senior programmer / analyst](#)!**
- Bioconductor [3.7](#) is available.
- Bioconductor [F1000 Research Channel](#) available.
- Orchestrating high-throughput genomic analysis with *Bioconductor* ([abstract](#)) and other [recent literature](#).
- View recent [course material](#).
- Use the [support site](#) to get help installing, learning and using Bioconductor.

### Install »

Get started with *Bioconductor*

- [Install Bioconductor](#)
- [Explore packages](#)
- [Get support](#)
- [Latest newsletter](#)
- [Follow us on twitter](#)
- [Install R](#)

### Learn »

Master *Bioconductor* tools

- [Courses](#)
- [Support site](#)
- [Package vignettes](#)
- [Literature citations](#)
- [Common work flows](#)
- [FAQ](#)
- [Community resources](#)
- [Videos](#)

### Use »

Create bioinformatic solutions with *Bioconductor*

- [Software](#), [Annotation](#), and [Experiment](#) packages
- [Amazon Machine Image](#)
- [Latest release announcement](#)
- [Support site](#)

### Develop »

Contribute to *Bioconductor*

- [Developer resources](#)
- Use [BIOC 'dev'](#)
- 'Dev!' [Software](#), [Annotation](#) and [Experiment](#) packages
- [Package guidelines](#)
- [New package submission](#)
- [Git source control](#)
- [Build reports](#)

3

We need to say what we mean

The tools already exist

Humans are the problem



reproducibility

## Evolution of Reporting P Values in the Biomedical Literature, 1990-2015.

[Chavalarias D](#)<sup>1</sup>, [Wallach JD](#)<sup>2</sup>, [Li AH](#)<sup>3</sup>, [Ioannidis JP](#)<sup>4</sup>.

### Author information

### Abstract

**IMPORTANCE:** The use and misuse of P values has generated extensive debates.

**OBJECTIVE:** To evaluate in large scale the P values reported in the abstracts and full text of biomedical research articles over the past 25 years and determine how frequently statistical information is presented in ways other than P values.

**DESIGN:** Automated text-mining analysis was performed to extract data on P values reported in 12,821,790 MEDLINE abstracts and in 843,884 abstracts and full-text articles in PubMed Central (PMC) from 1990 to 2015. Reporting of P values in 151 English-language core clinical journals and specific article types as classified by PubMed also was evaluated. A random sample of 1000 MEDLINE abstracts was manually assessed for reporting of P values and other types of statistical information; of those abstracts reporting empirical data, 100 articles were also assessed in full text.


**MAIN OUTCOMES AND MEASURES:** P values reported.

**RESULTS:** Text mining identified 4,572,043 P values in 1,608,736 MEDLINE abstracts and 3,438,299 P values in 385,393 PMC full-text articles. Reporting of P values in abstracts increased from 7.3% in 1990 to 15.6% in 2014. In 2014, P values were reported in 33.0% of abstracts from the 151 core clinical journals (n = 29,725 abstracts), 35.7% of meta-analyses (n = 5620), 38.9% of clinical trials (n = 4624), 54.8% of randomized controlled trials (n = 13,544), and 2.4% of reviews (n = 71,529). The distribution of reported P values in abstracts and in full text showed strong clustering at P values of .05 and of .001 or smaller. Over time, the "best" (most statistically significant) reported P values were modestly smaller and the "worst" (least statistically significant) reported P values became modestly less significant. Among the MEDLINE abstracts and PMC full-text articles with P

### Full text links



### Save items

 Add to Favorites ▾

### Similar articles

[Bias due to selective inclusion and reporting of outcomes](#) [Cochrane Database Syst Rev. 2014]

[Review Aquatic exercise training for fibromyalgia](#) [Cochrane Database Syst Rev. 2014]

[Review Screening for prostate cancer.](#) [Cochrane Database Syst Rev. 2013]

[Community-based care for the management of type 2 diabetes](#) [Ont Health Technol Assess Ser...]

[Review Resistance exercise training for fibromyalgia](#) [Cochrane Database Syst Rev. 2013]

[See reviews...](#)

[See all...](#)

Hi John

I read with interest your recent paper in JAMA on p-values:

<http://jama.jamanetwork.com/article.aspx?articleid=2503172#>

But could not find the data or code. Would you mind letting me know where they are?

Thanks!

Dear Jeff,

I still have to publish the code (I managed it on a private git). I plan to do it early june since I am quite busy until then. I just want to properly explain how it works when I release it. I hope this won't be too long

As for  
I will  
the M

Reg:

David

**“So if I have time I will make a website with an API to retrieve data on requests.”**

time,  
t's

Hi,











The dataset is now online on dataverse <http://dx.doi.org/10.7910/DVN/6FMTT3>

After import of the sql you should have

- 1,985,670 rows for the table `medline\_full\_txt\_list`
- 12,436,631 rows for the table `medline\_full\_txt\_pv`
- 16,116,061 rows for the table `medline\_pt`
- 9,088,701 rows for the table `medline\_pvalues`

Tell me if there is any issue. Source code will follow.

## 7 Files

		 Download
<input type="checkbox"/>	<div data-bbox="170 303 260 423">  </div> <p data-bbox="293 285 575 312"><a href="#">medline_full_txt_list.sql</a></p> <p data-bbox="293 321 821 345">Unknown - 228.3 MB - Apr 12, 2016 - 4 Downloads</p> <p data-bbox="293 352 757 376">MD5: 5d78f42859d4044660cf636675281384</p> <p data-bbox="293 384 637 408">List of all PMC papers processed</p> <div data-bbox="293 423 479 447"> <span>Data</span> <span>P-values</span> </div>	 Download
<input type="checkbox"/>	<div data-bbox="170 500 260 620">  </div> <p data-bbox="293 481 575 508"><a href="#">medline_full_txt_pv.sql</a></p> <p data-bbox="293 517 795 541">Unknown - 1.0 GB - Apr 12, 2016 - 9 Downloads</p> <p data-bbox="293 549 757 573">MD5: 9e3ee983cf8bee7d75153124abfae6e2</p> <p data-bbox="293 580 830 604">sql table of all P-values extracted from PMC full text</p> <div data-bbox="293 620 479 644"> <span>Data</span> <span>P-values</span> </div>	 Download
<input type="checkbox"/>	<div data-bbox="170 707 260 827">  </div> <p data-bbox="293 678 575 705"><a href="#">medline_full_txt_pv.tab</a></p> <p data-bbox="293 714 869 738">Tabular Data - 174.3 MB - Apr 12, 2016 - 14 Downloads</p> <p data-bbox="293 745 672 769">7 Variables, 3438298 Observations -</p> <p data-bbox="293 777 703 801">UNF:6:tRQMiS7wwbyq7H4scJ6DAQ==</p> <p data-bbox="293 809 792 833">CSV of all P-values extracted from PMC full text</p> <div data-bbox="293 848 479 873"> <span>Data</span> <span>P-values</span> </div>	<div data-bbox="1387 740 1568 805">  Explore         </div> <div data-bbox="1588 740 1804 805">  Download ▾         </div>
<input type="checkbox"/>	<div data-bbox="170 915 260 1035">  </div> <p data-bbox="293 896 479 923"><a href="#">medline_pt.sql</a></p> <p data-bbox="293 932 784 956">Unknown - 1.5 GB - Apr 12, 2016 - 1 Download</p> <p data-bbox="293 964 768 988">MD5: 68ebdedabc670c188d9b49629db931d3</p> <p data-bbox="293 995 823 1019">sql of all P-values extracted from Medline abstracts</p> <div data-bbox="293 1035 479 1059"> <span>Data</span> <span>P-values</span> </div>	 Download

```
> library(readr)
> dat = read_csv("~/data/medicine/medline_full_txt_pv.csv")
```

```
Parsed with column specification:
```

```
cols(
  `7669595` = col_integer(),
  `0370635` = col_character(),
  `=` = col_character(),
  `0.14` = col_double(),
  `1995` = col_integer(),
  plain = col_character(),
  `1` = col_integer()
)
```

```
=====| 100% 174 MB
=====| 64% 112 MB
```











```
> head(dat)
```

```
# A tibble: 6 x 7
```

	`7669595` <int>	`0370635` <chr>	`=` <chr>	`0.14` <dbl>	`1995` <int>	plain <chr>	`1` <int>
1	7669596	0370635	=	0.001	1995	plain	1
2	8611396	0370635	<	0.010	1996	plain	1
3	8611396	0370635	<	0.010	1996	plain	1
4	8611396	0370635	<	0.010	1996	plain	1
5	8611397	0370635	<	0.010	1996	plain	1
6	8611398	0370635	<	0.010	1996	plain	1

```
> |
```

## 7 Files

		 Download
<input type="checkbox"/>	<div data-bbox="170 303 260 423">  </div> <p data-bbox="293 285 575 312"><a href="#">medline_full_txt_list.sql</a></p> <p data-bbox="293 321 821 345">Unknown - 228.3 MB - Apr 12, 2016 - 4 Downloads</p> <p data-bbox="293 352 757 376">MD5: 5d78f42859d4044660cf636675281384</p> <p data-bbox="293 384 637 408">List of all PMC papers processed</p> <div data-bbox="293 423 479 447"> <span>Data</span> <span>P-values</span> </div>	 Download
<input type="checkbox"/>	<div data-bbox="170 500 260 620">  </div> <p data-bbox="293 481 569 508"><a href="#">medline_full_txt_pv.sql</a></p> <p data-bbox="293 517 795 541">Unknown - 1.0 GB - Apr 12, 2016 - 9 Downloads</p> <p data-bbox="293 549 757 573">MD5: 9e3ee983cf8bee7d75153124abfae6e2</p> <p data-bbox="293 580 830 604">sql table of all P-values extracted from PMC full text</p> <div data-bbox="293 620 479 644"> <span>Data</span> <span>P-values</span> </div>	 Download
<input type="checkbox"/>	<div data-bbox="170 707 260 827">  </div> <p data-bbox="293 678 575 705"><a href="#">medline_full_txt_pv.tab</a></p> <p data-bbox="293 714 869 738">Tabular Data - 174.3 MB - Apr 12, 2016 - 14 Downloads</p> <p data-bbox="293 745 672 769">7 Variables, 3438298 Observations -</p> <p data-bbox="293 777 705 801">UNF:6:tRQMiS7wwbyq7H4scJ6DAQ==</p> <p data-bbox="293 809 792 833">CSV of all P-values extracted from PMC full text</p> <div data-bbox="293 848 479 873"> <span>Data</span> <span>P-values</span> </div>	<div data-bbox="1387 740 1568 805">  Explore         </div> <div data-bbox="1568 740 1804 805">  Download ▾         </div>
<input type="checkbox"/>	<div data-bbox="170 915 260 1035">  </div> <p data-bbox="293 896 473 923"><a href="#">medline_pt.sql</a></p> <p data-bbox="293 932 782 956">Unknown - 1.5 GB - Apr 12, 2016 - 1 Download</p> <p data-bbox="293 964 768 988">MD5: 68ebdedabc670c188d9b49629db931d3</p> <p data-bbox="293 995 821 1019">sql of all P-values extracted from Medline abstracts</p> <div data-bbox="293 1035 479 1059"> <span>Data</span> <span>P-values</span> </div>	 Download



# P-values from Chavalarias et al. 2016 for the tidypvals package

**Jeff Leek**

**26 July 2017**

## Contents

---

### [1 Set up](#)

#### [1.1 Load packages](#)

#### [1.2 Load data](#)

### [2 Tidy p-values](#)

#### [2.1 Format p-values](#)

#### [2.2 Select the appropriate columns and clean](#)

### [3 Save data](#)

### [4 Session information](#)

These p-values come from the paper: [Evolution of Reporting P Values in the Biomedical Literature](#). The csv file for the p-values from medline did not have column names, so to ensure we had the right data we downloaded the MySQL dump from the Dataverse <https://dataverse.harvard.edu/file.xhtml;jsessionid=94274f10cbdbecaaaf6da71ca209?fileId=2801917&version=RELEASED&version=.0> on 2017-07-24. We re-loaded it into a MySQL database and that is where the code starts.

## 1 Set up

---

### 1.1 Load packages

# How to share data for collaboration

Shannon E. Ellis

Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health  
and

Jeffrey T. Leek \*

Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health

September 1, 2017

replicability

**This Issue**Citations **1,078**

PDF



More ▾



Cite



Permissions

**Original Contribution****FREE**

March 8, 2006

# Fine Particulate Air Pollution and Hospital Admission for Cardiovascular and Respiratory Diseases

Francesca Dominici, PhD; Roger D. Peng, PhD; Michelle L. Bell, PhD; [et al](#)[Article Information](#)

“In the first stage, single lag and distributed lag overdispersed Poisson regression models were used for estimating county-specific RRs of hospital admissions associated with ambient levels of PM<sub>2.5</sub>”

$\log(E[\text{Hospital Admissions} \mid X]) =$

PM2.5

+  $\log(\# \text{ at risk})$  + factor(day of week) + ns(day of year, 8) +  
ns(temperature, 6)

+ ns(dew point temperature, 3) + ns(lagged temperature, 6)

+ ns(lagged dew point temperature, 3)

Copyrighted Material

Monographs  
on Statistics and  
Applied Probability 37

# Generalized Linear Models

SECOND EDITION

P. McCullagh and  
J.A. Nelder FRS



Springer-Science+Business  
Media, B.V.

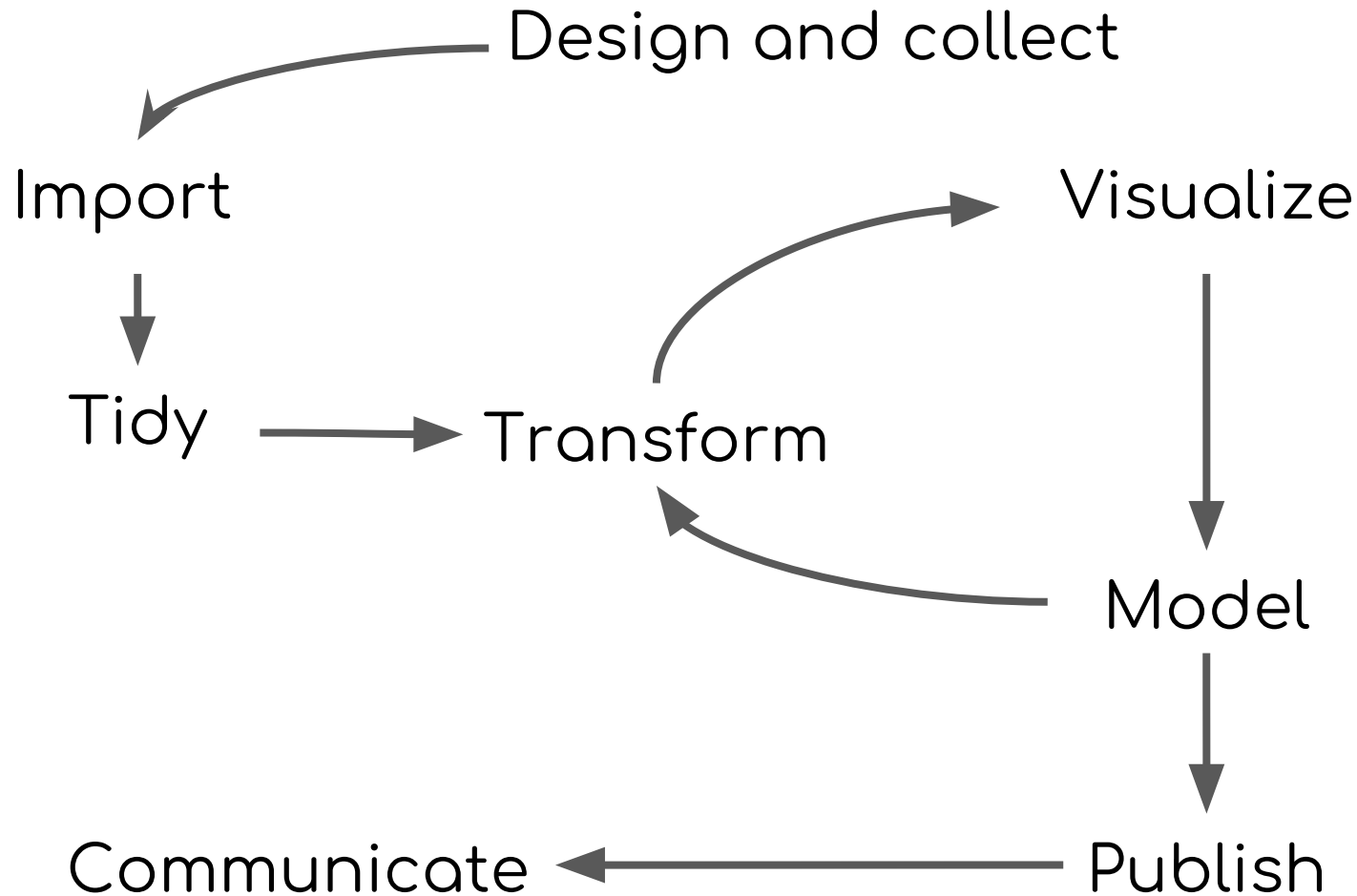
Eight outcomes were considered based on the ICD-9 codes for 5 cardiovascular outcomes (heart failure [428], heart rhythm disturbances [426-427], cerebrovascular events [430-438], ischemic heart disease [410-414, 429], peripheral vascular disease [440-448]), 2 respiratory outcomes (chronic obstructive pulmonary disease [COPD; 490-492], respiratory tract infections [464-466, 480-487]), and hospitalizations caused by injuries and other external causes (800-849). The county-wide daily hospitalization rates for each outcome for 1999-2002 appear in Table 1.

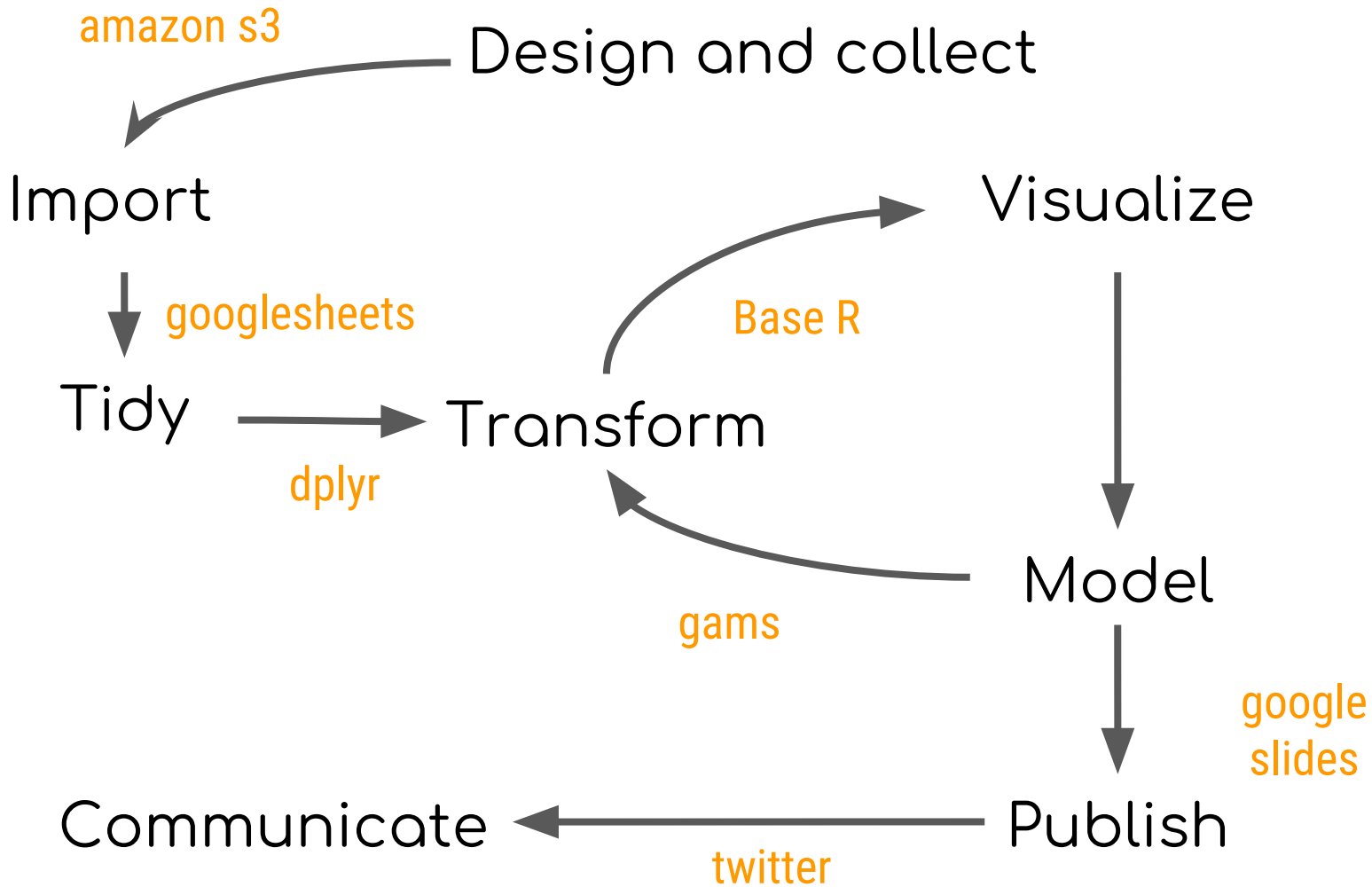
The study population includes 11.5 million Medicare enrollees residing an average of 5.9 miles from a PM<sub>2.5</sub> monitor. The analysis was restricted to the 204 US counties with populations larger than 200 000. Of these 204 counties, 90 had daily PM<sub>2.5</sub> data across the study period and the remaining counties had PM<sub>2.5</sub> data collected once every 3 days for at least 1 full year. The locations of the 204 counties appear in Figure 1. The counties were clustered into 7 geographic regions by applying the K-means clustering algorithm to longitude and latitude for the counties.<sup>10,11</sup>

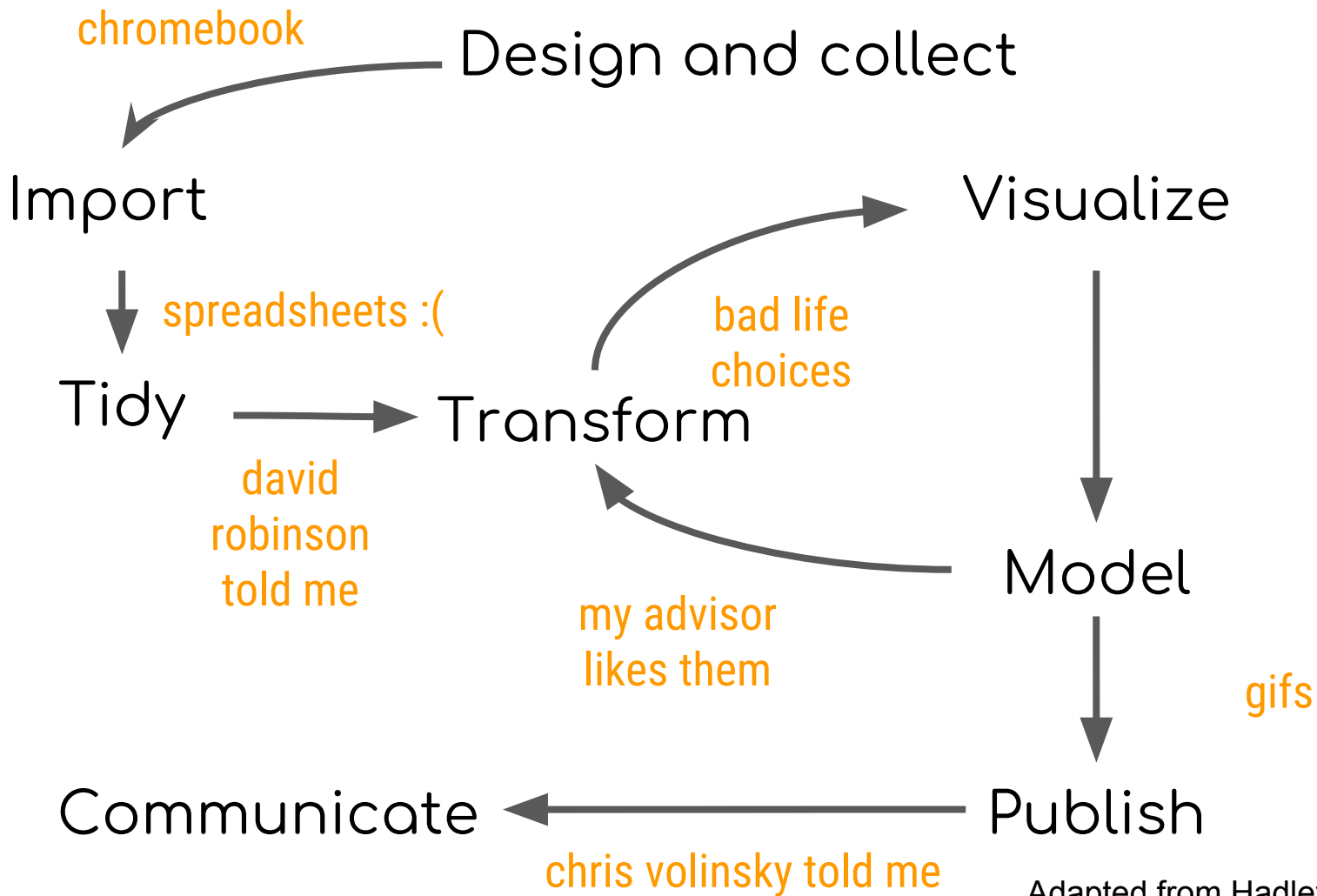


**Eight outcomes** were considered based on the ICD-9 codes for **5 cardiovascular outcomes** (heart failure [428], heart rhythm disturbances [426-427], cerebrovascular events [430-438], ischemic heart disease [410-414, 429], peripheral vascular disease [440-448]), **2 respiratory outcomes** (chronic obstructive pulmonary disease [COPD; 490-492], respiratory tract infections [464-466, 480-487]), and hospitalizations caused by injuries and other external causes (800-849). The county-wide daily hospitalization rates for each outcome for 1999-2002 appear in Table 1.

The study population includes 11.5 million Medicare enrollees residing an average of 5.9 miles from a PM2.5 monitor. The **analysis was restricted to the 204 US counties with populations larger than 200 000**. Of these 204 counties, 90 had daily PM2.5 data across the study period and the remaining counties had PM2.5 data collected once every 3 days for at least 1 full year. The locations of the 204 counties appear in Figure 1. **The counties were clustered into 7 geographic regions** by applying the K-means clustering algorithm to longitude and latitude for the counties.<sup>10,11</sup>







Suppose we have a sample with 30 lawyers and 70 engineers.

What is the probability Dick is a lawyer?

Dick is a 30 year old man. He is married with no children. A man of high ability and high motivation, he promises to be quite successful in his field. He is well liked by his colleagues.

What is the probability Dick is a lawyer?

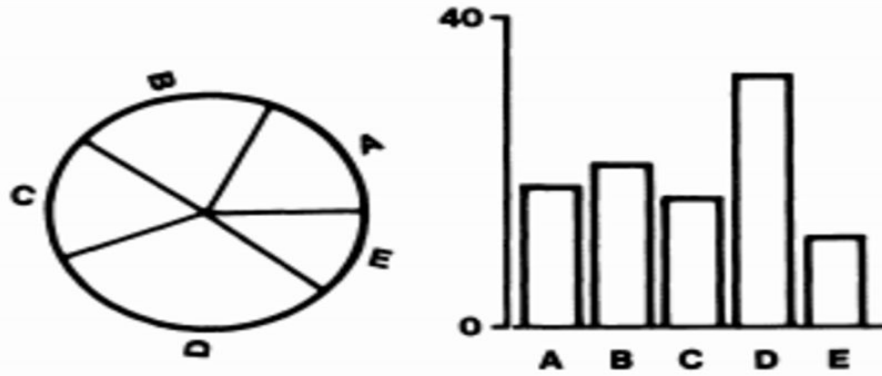


Figure 3. Graphs from position-angle experiment.

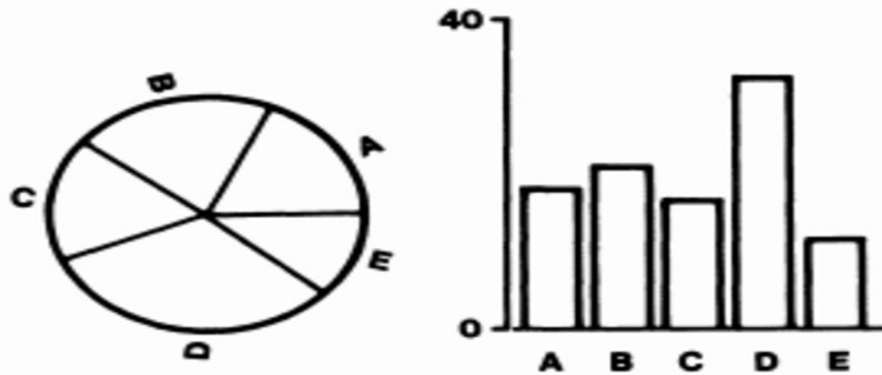
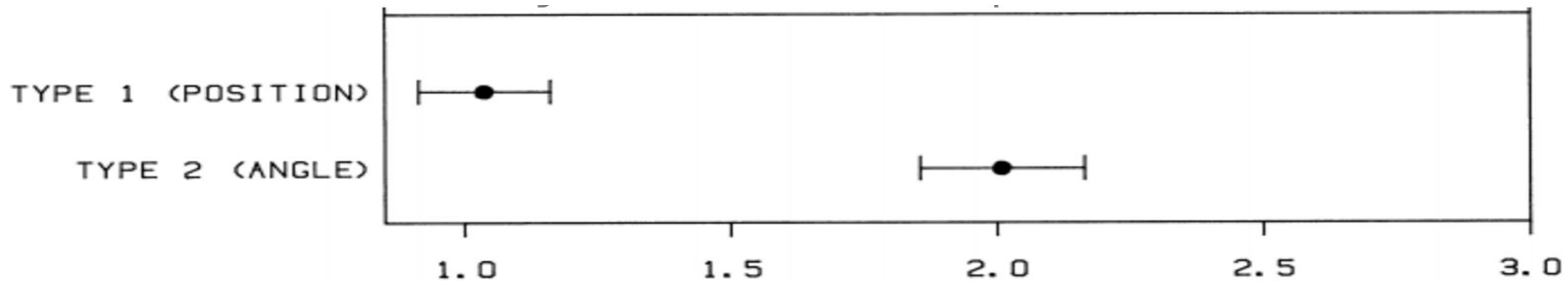


Figure 3. Graphs from position-angle experiment.



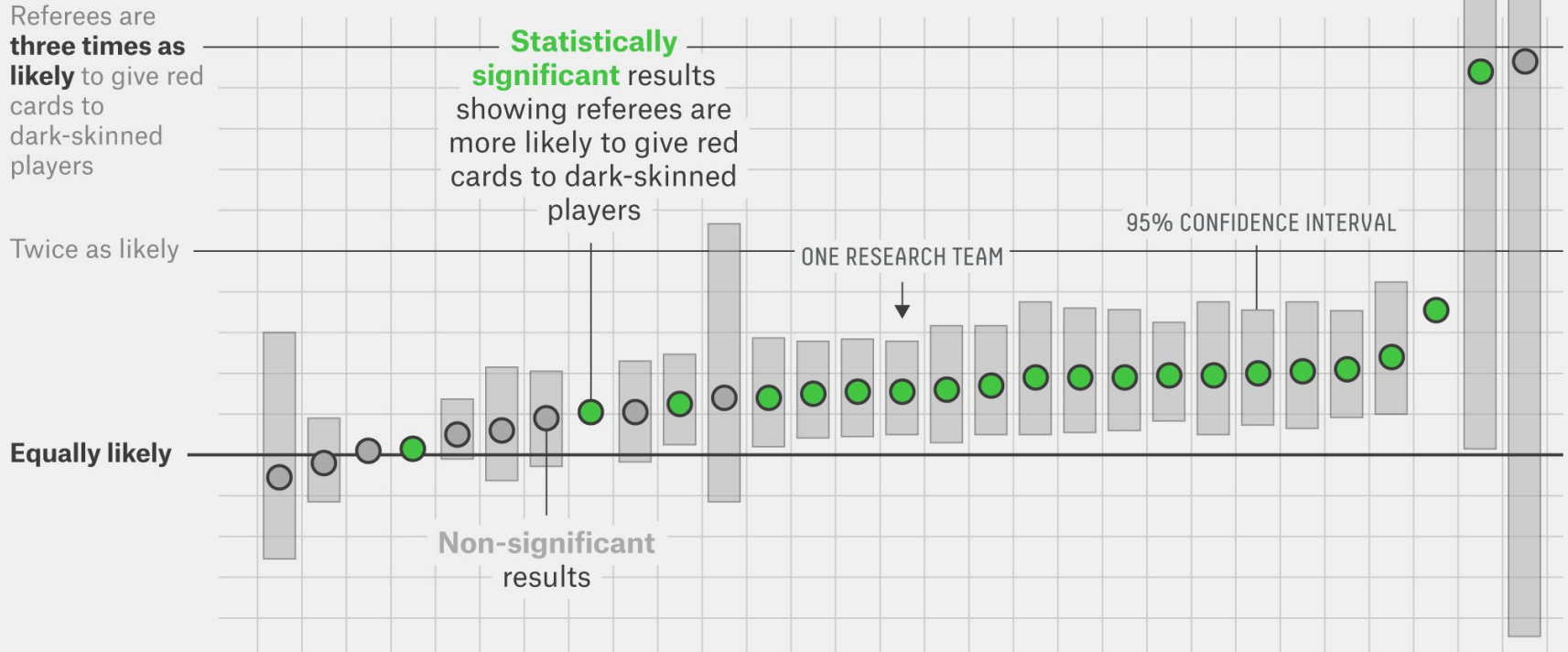


# Human data interaction

(observational epidemiology)

## Same Data, Different Conclusions

Twenty-nine research teams were given the same set of soccer data and asked to determine if referees are more likely to give red cards to dark-skinned players. Each team used a different statistical method, and each found a different relationship between skin color and red cards.





MSD in R  
6 Courses  
35K+ Enrollments



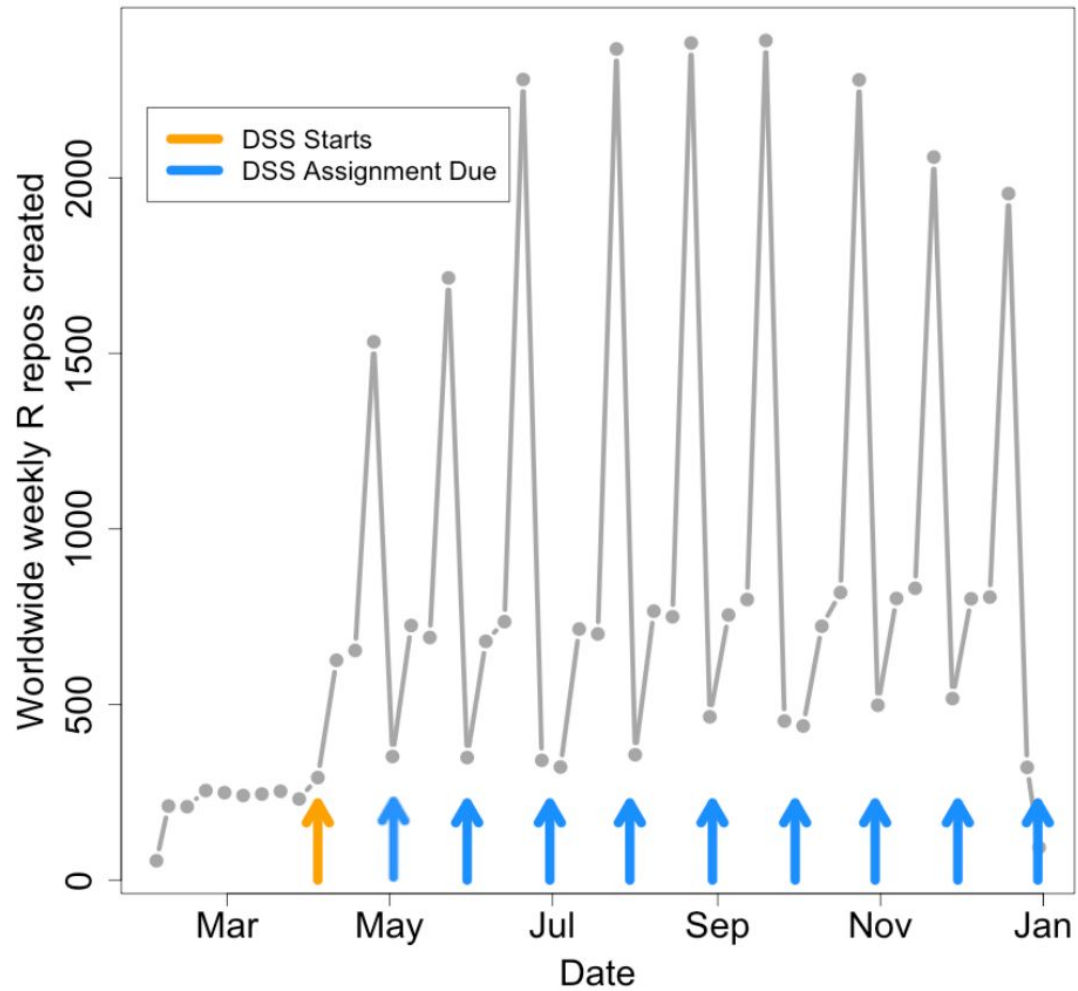
Executive Data Science  
5 Courses  
150K+ Enrollments

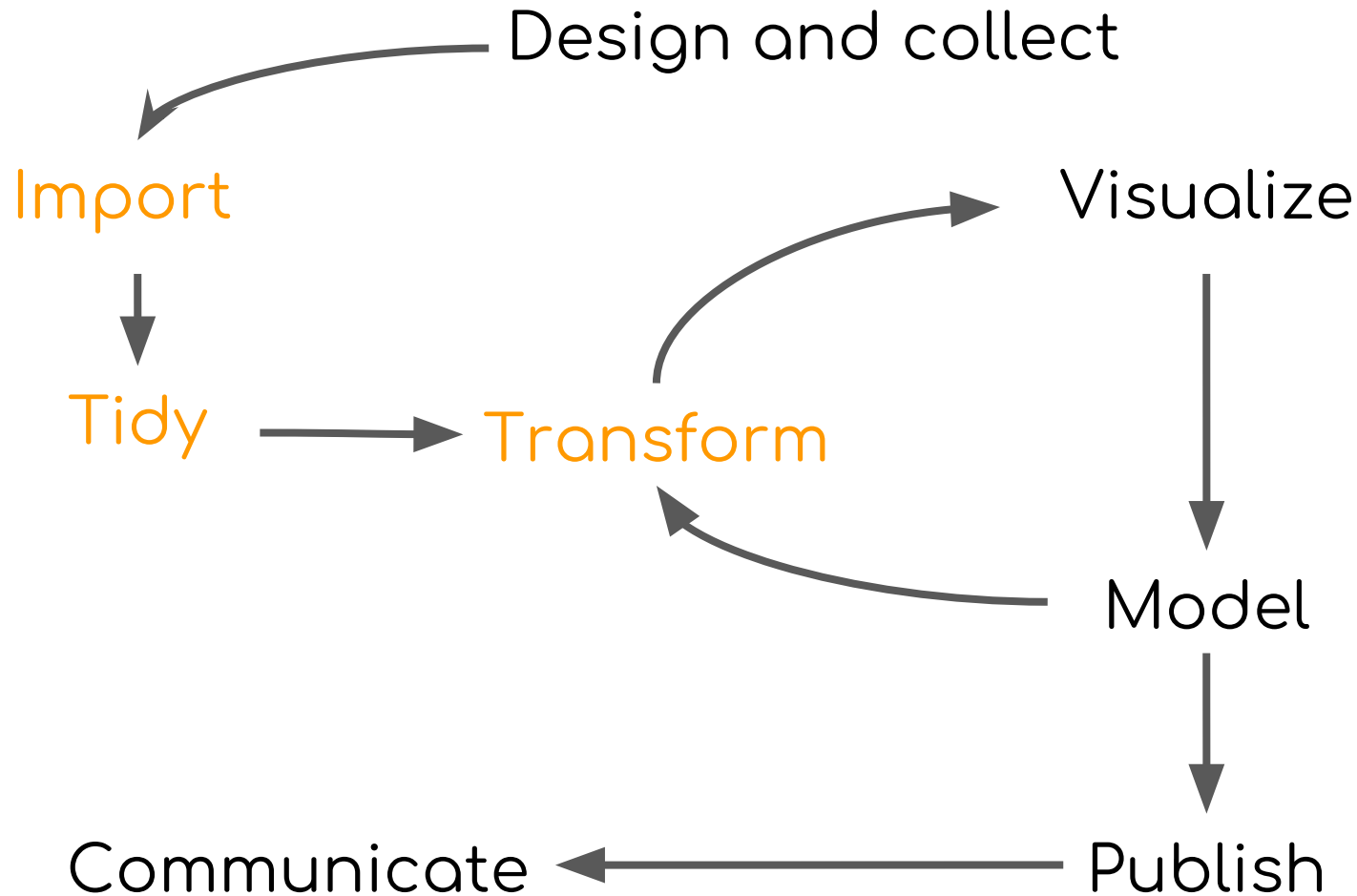


Genomic Data Science  
9 Courses  
230k+ Enrollments



Data Science  
10 Courses  
4.4M+ Enrollments





Repositories	33K
Code	2M
Commits	1M
Issues	13K
Topics	1
Wikis	2K
Users	3

Languages	
R	29,958
HTML	569
Rebol	77
Python	15
Jupyter Notebook	13
Java	7
Shell	6
CSS	5

## 33,499 repository results

Sort: Best match ▾

**benjamin-chan/GettingAndCleaningData** ● R ★ 68  
Repository for Coursera course *Getting and Cleaning Data*.

Updated on Sep 18, 2016

**Xiaodan/Coursera-Getting-and-Cleaning-Data** ● R ★ 26  
Repo for Coursera.com online course: *Getting and Cleaning Data*

Updated on Oct 11, 2015

**bgentry/course-getting-and-cleaning-data-project** ● R ★ 4  
course project for Coursera "*Getting and Cleaning Data*"

Updated on May 18, 2016

getting and cleaning data

Pull requests Issues Marketplace Explore

33,499 repository results

Sort: Best match ▾

Repositories	33,499
Code	2M
Commits	1M
Issues	13K
Topics	1
Wikis	2K
Users	3

### Languages

R	29,958
HTML	569
Rebol	77
Python	15
Jupyter Notebook	13
Java	7
Shell	6
CSS	5

### [benjamin-chan/GettingAndCleaningData](#)

● R

★ 68

Repository for Coursera course *Getting and Cleaning Data*.

Updated on Sep 18, 2016

### [Xiaodan/Coursera-Getting-and-Cleaning-Data](#)

● R

★ 26

Repo for Coursera.com online course: *Getting and Cleaning Data*

Updated on Oct 11, 2015

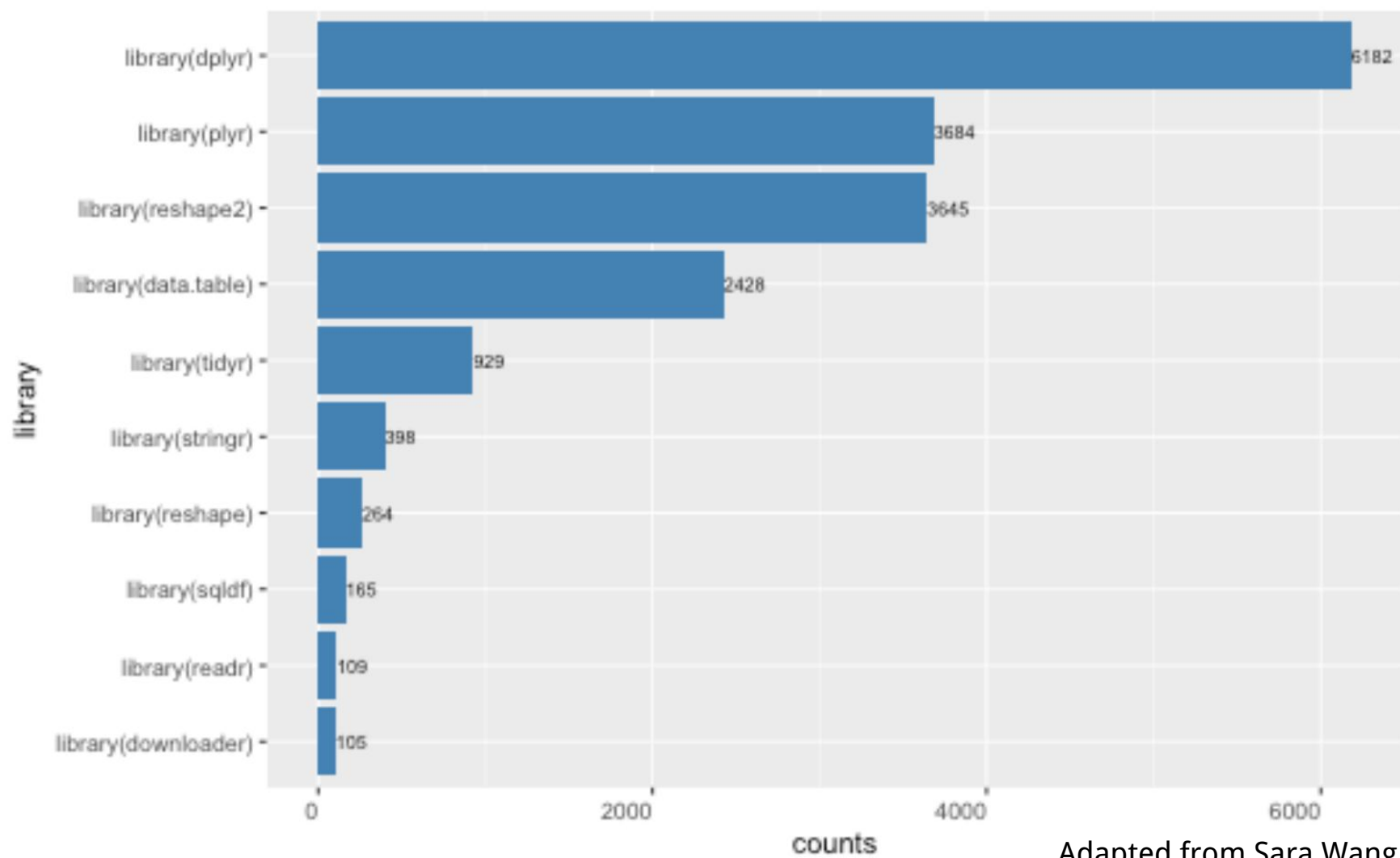
### [bgentry/coursera-getting-and-cleaning-data-project](#)

● R

★ 4

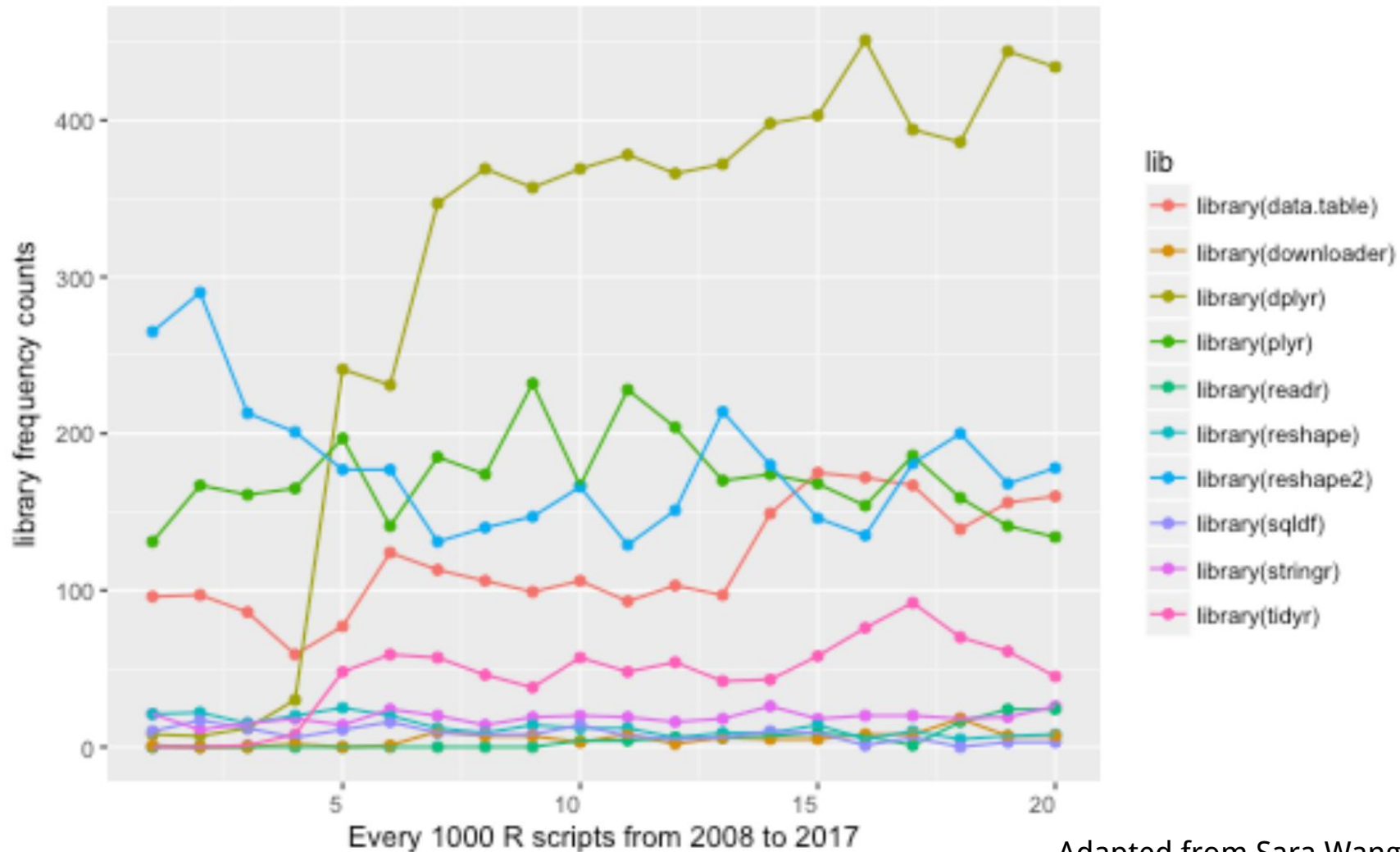
course project for Coursera "*Getting and Cleaning Data*"

Updated on May 18, 2016

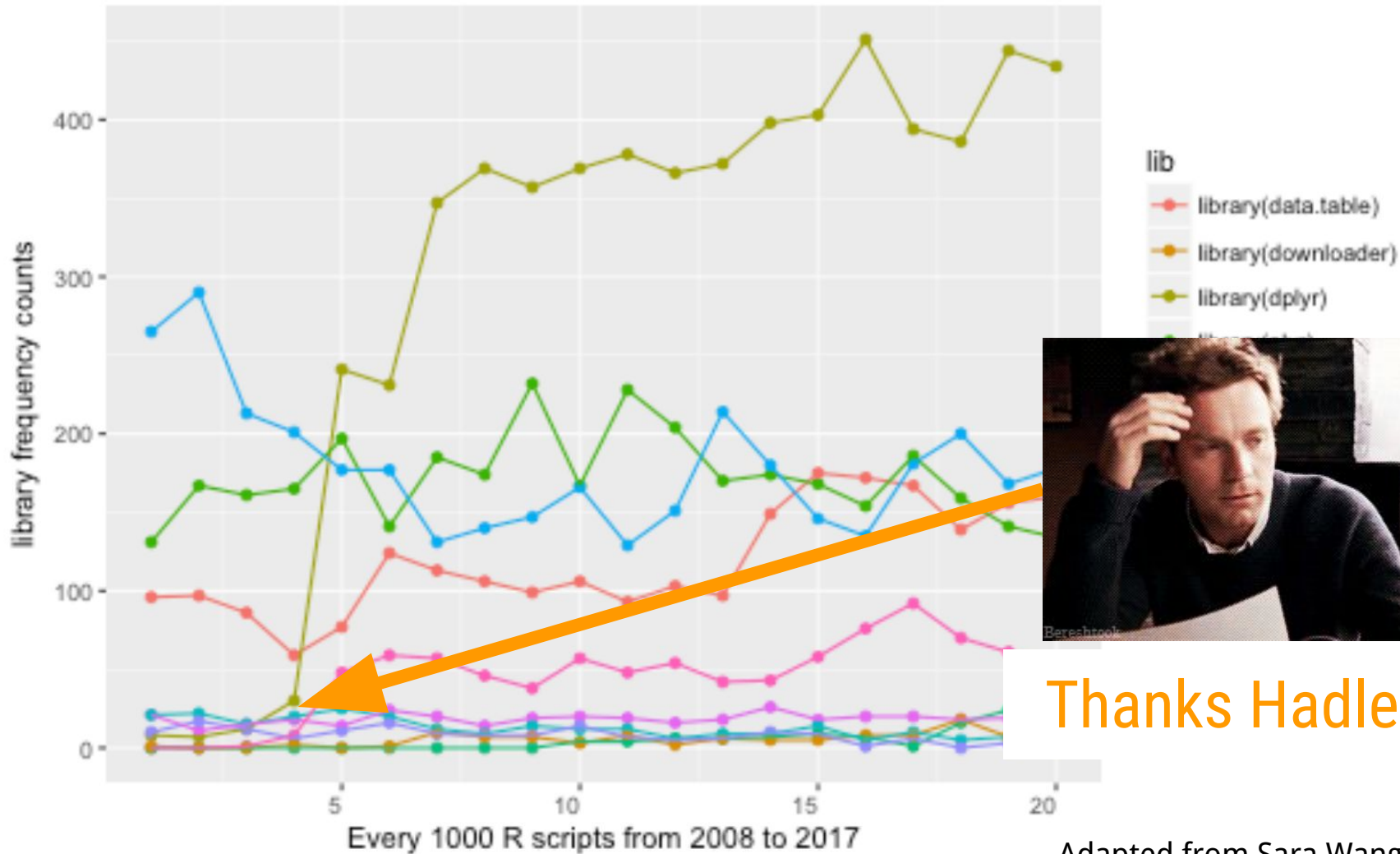


Adapted from Sara Wang





Adapted from Sara Wang



Thanks Hadley!

Adapted from Sara Wang

# Human data interaction

(randomized trials)

# (OPTIONAL) Data analysis practice with immediate feedback (NEW! 10/18/2017)

## Question 1

1 Point

VARIATION 1 You are being asked to participate in a research experiment with the purpose of better understanding how people analyze data. If you complete this quiz, you are giving your consent to participate in the study. This quiz involves a short data analysis that gives you a chance to practice the regression concepts you have learned so far.

**We anticipate that this will take about 15 minutes to complete. You will be receiving feedback on your work immediately after submission. For this reason, we ask that you do not post on the forums about this quiz to maintain the integrity of this experiment.**

Thank you for helping us learn more about data science! -Brian, Roger, Jeff

-----

Your assignment is to study how income varies across different categories of college majors. You will be using data from a study of recent college graduates. Make sure to use good practices that you have learned so far in this course and previous courses in the specialization. In particular, it is good practice to specify an analysis plan early in the process to avoid the “p-hacking” behavior of trying many analyses to find one that has desired results. If you want to learn more about “p-hacking”, you can visit

# (OPTIONAL) Data analysis practice with immediate feedback (NEW! 10/18/2017)

## Question 1

1 Point

VARIATION 1

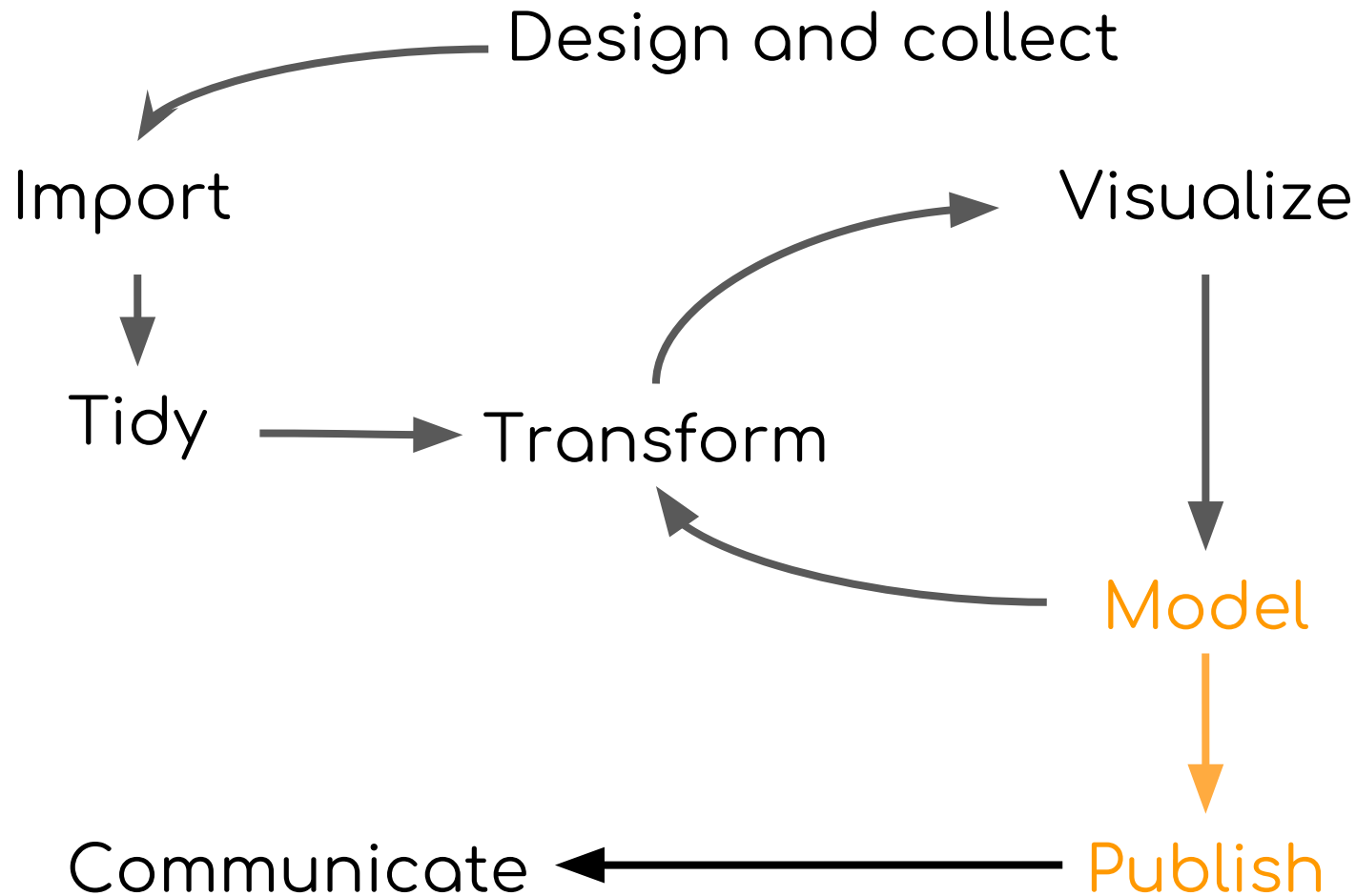
You are being asked to participate in a research experiment with the purpose of better understanding how people analyze data. If you complete this quiz, you are giving your consent to participate in the study. This quiz involves a short data analysis that gives you a chance to practice the regression concepts you have learned so far.

**We anticipate that this will take about 15 minutes to complete. You will be receiving feedback on your work immediately after submission. For this reason, we ask that you do not post on the forums about this quiz to maintain the integrity of this experiment.**

Thank you for helping us learn more about data science! -Brian, Roger, Jeff

-----

Your assignment is to study how income varies across different categories of college majors. You will be using data from a study of recent college graduates. Make sure to use good practices that you have learned so far in this course and previous courses in the specialization. In particular, it is good practice to specify an analysis plan early in the process to avoid the “p-hacking” behavior of trying many analyses to find one that has desired results. If you want to learn more about “p-hacking”, you can visit



“..Coming up with nothing is unacceptable..”

Article  
TextArticle  
infoCitation  
Tools

Share



Responses

Article  
metricsEthics  
Research

# Identifying bioethical issues in biostatistical consulting: findings from a US national pilot survey of biostatisticians

Min Qi Wang<sup>1</sup>, Alice F Yan<sup>2</sup>, Ralph V Katz<sup>3</sup>[Author affiliations +](#)

## Abstract

**Objectives** The overall purposes of this first US national pilot study were to (1) test the feasibility of online administration of the Bioethical Issues in Biostatistical Consulting (BIBC) Questionnaire to a random sample of American Statistical Association (ASA) members; (2) determine the prevalence and relative severity of a broad array of bioethical violations requests that are presented to biostatisticians by investigators seeking biostatistical consultations; and (3) establish the sample size needed for a full-size phase II study.

**Design** A descriptive survey as approved and endorsed by the ASA.



**[It has been widely noted that there is a relationship between  $x_1$  and  $Y$ .]** We are interested in studying the magnitude of this relationship and reporting results to the president of the company. Using the data available at **[data1/data2]**, use a linear regression model to investigate this relationship.

Would you be confident telling the company president that there is a meaningful relationship between  $x_1$  and  $y$ ?

- A. Yes
- B. No

**data1:  $x_2$  confounds the relationship between  $Y$  and  $x_1$**   
**data2:  $x_2$  is not a confounder but is predictive of  $Y$**

# Confounded

Primed

68% found signal

(n = 263)

Not Primed

57% found signal

(n = 252)

# Not confounded

Primed

78% found signal

(n = 263)

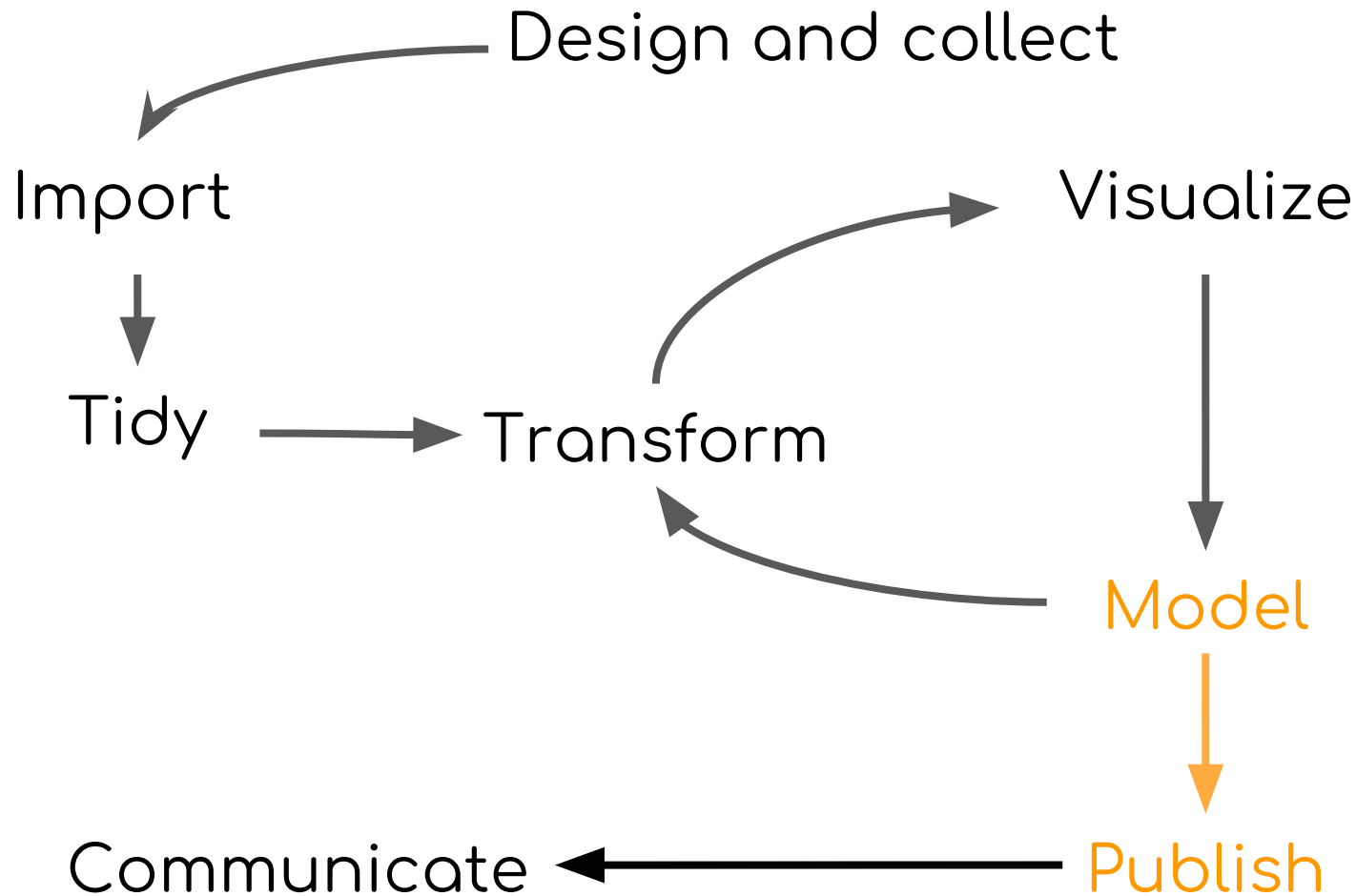
Not Primed

87% found signal


(n = 258)

# Human data interaction

(interventional randomized trials)



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

Psychological Science  
22(11) 1359–1366  
© The Author(s) 2011  
Reprints and permission:  
sagepub.com/journalsPermissions.nav  
DOI: 10.1177/0956797611417632  
<http://pss.sagepub.com>  


Joseph P. Simmons<sup>1</sup>, Leif D. Nelson<sup>2</sup>, and Uri Simonsohn<sup>1</sup>

<sup>1</sup>The Wharton School, University of Pennsylvania, and <sup>2</sup>Haas School of Business, University of California, Berkeley

## Abstract

In this article, we accomplish two things. First, we show that despite empirical psychologists' nominal endorsement of a low rate of false-positive findings ( $\leq .05$ ), flexibility in data collection, analysis, and reporting dramatically increases actual false-positive rates. In many cases, a researcher is more likely to falsely find evidence that an effect exists than to correctly find evidence that it does not. We present computer simulations and a pair of actual experiments that demonstrate how unacceptably easy it is to accumulate (and report) statistically significant evidence for a false hypothesis. Second, we suggest a simple, low-cost, and straightforwardly effective disclosure-based solution to this problem. The solution involves six concrete requirements for authors and four guidelines for reviewers, all of which impose a minimal burden on the publication process.

## Keywords

methodology, motivated reasoning, publication, disclosure

Your assignment is to study how income varies across different categories of college majors. You will be using data from a study of recent college graduates. Make sure to use good practices that you have learned so far in this course and previous courses in the specialization. [In particular, it is good practice to specify an analysis plan early in the process to avoid the “p-hacking” behavior of trying many analyses to find one that has desired results. If you want to learn more about “p-hacking”, you can visit <https://projects.fivethirtyeight.com/p-hacking/>.]

If you will proceed with the analysis, click “Yes”. Otherwise, click “No”.

	Significant	Not significant
Warned	66	90
Not Warned	59	88



$$Y = X + e$$

# The researcher degrees of freedom - recipe tradeoff in data analysis

Jeff Leek 2013/07/31

An important concept that is only recently gaining the [attention it deserves](#) is researcher degrees of freedom. From [Simmons et al.](#):

The culprit is a construct we refer to as researcher degrees of freedom. In the course of collecting and analyzing data, researchers have many decisions to make: Should more data be collected? Should some observations be excluded? Which conditions should be combined and which ones compared? Which control variables should be considered? Should specific measures be combined or transformed or both?

So far, researcher degrees of freedom has primarily been used with [negative connotations](#). This probably stems from the original definition of the idea which focused on how analysts could “manufacture” statistical significance by changing the way the data was processed without disclosing those changes. Reproducible research and distributed code would of course address these issues to some extent. But it is still relatively easy to obfuscate dubious analysis by [dressing it up in technical language](#).

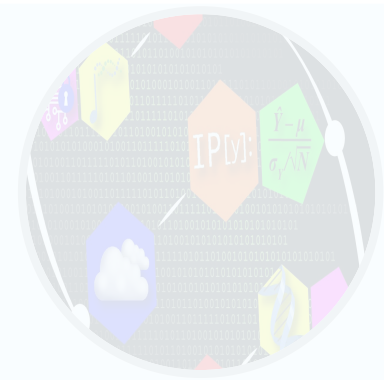
$$Y = X + e_{\text{analyst}} + e_{\text{sampling}}$$



MSD in R  
6 Courses  
35K+ Enrollments



Executive Data Science  
5 Courses  
150K+ Enrollments



Genomic Data Science  
9 Courses  
230k+ Enrollments



Data Science  
10 Courses  
4.4M+ Enrollments



MSD in R  
6 Courses  
35K+ Enrollments



Executive Data Science  
5 Courses  
150K+ Enrollments



Genomic Data Science  
9 Courses  
230k+ Enrollments



Data Science  
10 Courses  
4.4M+ Enrollments

# Data analysis subcultures

 Jeff Leek  2015/04/29

Roger and I responded to the controversy around the journal that banned p-values today [in Nature](#). A piece like this requires a lot of information packed into very little space but I thought one idea that deserved to be talked about more was the idea of data analysis subcultures. From the paper:

Data analysis is taught through an apprenticeship model, and different disciplines develop their own analysis subcultures. Decisions are based on cultural conventions in specific communities rather than on empirical evidence. For example, economists call data measured over time ‘panel data’, to which they frequently apply mixed-effects models. Biomedical scientists refer to the same type of data structure as ‘longitudinal data’, and often go at it with generalized estimating equations.

I think this is one of the least appreciated components of modern data analysis. Data analysis is almost entirely taught through an apprenticeship culture with completely different behaviors taught in different disciplines. All of these disciplines agree about the mathematical optimality of specific methods under very specific conditions. That is why you see [methods](#) like [randomized trials](#) [Roger and I responded to the controversy around the journal that banned p-values today [in Nature](#). A piece like this requires a lot of information packed into very little space but I thought one idea that deserved to be talked about more was the idea of data analysis subcultures. From the paper:

$$Y = X + e_{\text{population}} + e_{\text{analyst}} + e_{\text{sampling}}$$

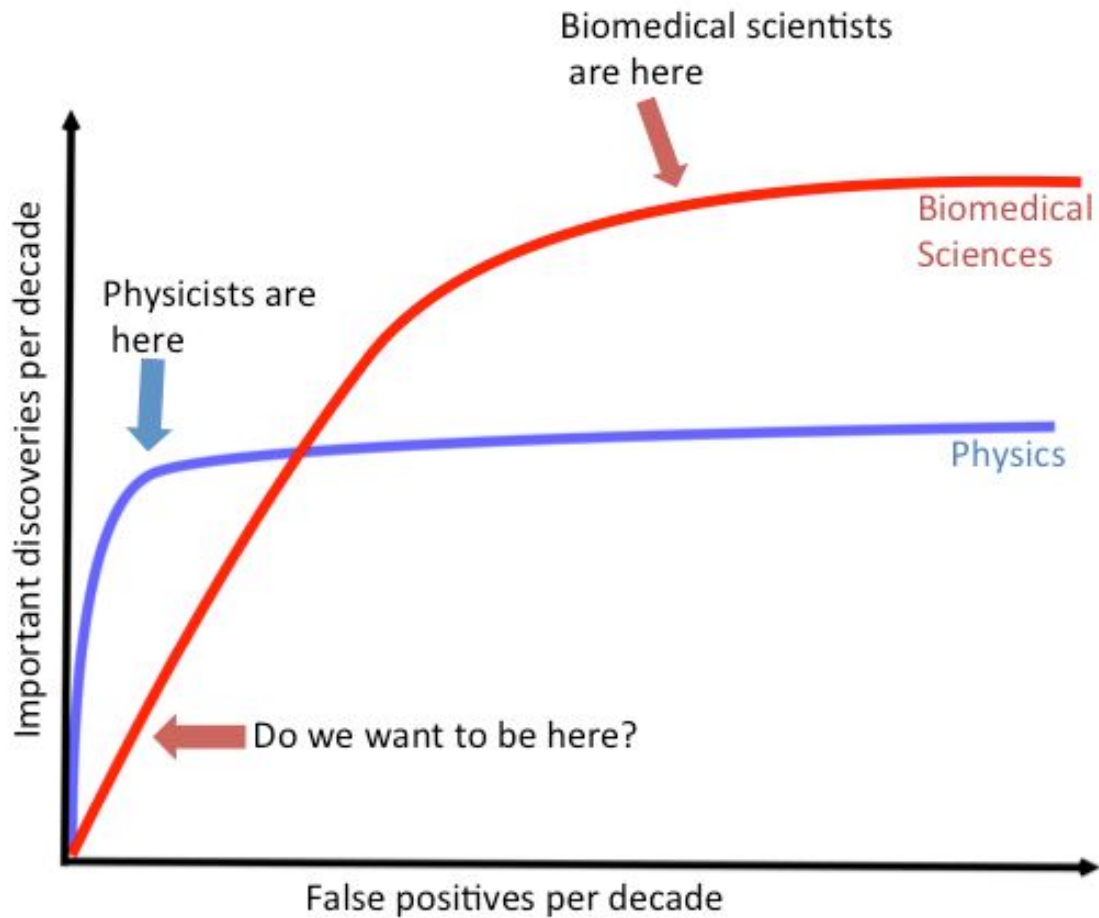
# 3

We need to say what we mean

The tools already exist

Humans are the problem





## Leek group

- Shannon Ellis
- **Aboozar Hadavand**
- **Lucy D'Agostino McGowan**
- **Leslie Myint**
- Jack Fu
- Sara Wang
- Kayode Sosina
- You?

## Alumni

- Simina Boca
- Hilary Parker
- Andrew Jaffe
- Alyssa Frazee
- Prasad Patil
- Leo Collado Torres
- Abhi Nellore
- Kai Kammers
- **Nick Carchedi**
- **Sean Kross**
- Divya Narayanan

## Collaborators

- Kasper Hansen
- Margaret Taub
- **Leah Jager**
- Ben Langmead
- **Roger Peng**
- **Aaron Fisher**
- **Brooke Anderson**
- **Ira Gooding**

```
head(college)
fit <- lm(rank ~ major_category - 1, college)
summary(fit)
```

```
str(college)
library(ggplot2)
library(dplyr)
g1 <- ggplot(college, aes(x = major_category, y = median, fill = major_category))
g1 <- g1 + geom_boxplot()
g1 <- g1 + theme(axis.text.x = element_text(angle = 45, hjust = 1), legend.position = "none")
g1
tapply(college$sample_size, college$major_category, sum)
majormean <- tapply(college$sample_size, college$major_category, mean)
max(majormean)
college$major_category <- as.factor(college$major_category)
college2 <- college
college2$major_category <- relevel(college$major_category, ref = "Business")
fit1 <- lm(median ~ major_category, data = college2)
fit2 <- update(fit1, median ~ major_category + perc_men)
fit3 <- update(fit1, median ~ major_category + perc_men + perc_employed)
fit4 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround)
fit5 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround + perc_college_jobs)
anova(fit1, fit2, fit3, fit4)
summary(fit5)
tail(hatvalues(fit5)[order(hatvalues(fit5, decreasing = T))], 5)
tail(cooks.distance(fit5)[order(cooks.distance(fit5, decreasing = T))], 5)
```

```
str(college)
library(ggplot2)
library(dplyr)
g1 <- ggplot(college, aes(x = major_category, y = median, fill = major_category))
g1 <- g1 + geom_boxplot()
g1 <- g1 + theme(axis.text.x = element_text(angle = 45, hjust = 1), legend.position = "none")
g1
tapply(college$sample_size, college$major_category, FUN = function(x) {
  majormean <- tapply(college$sample_size, college$major_category, FUN = function(x) {
    max(majormean)
  })
  college$major_category <- as.factor(college$major_category)
  college2 <- college
  college2$major_category <- relevel(college$major_category, ref = "Business")
  fit1 <- lm(median ~ major_category, data = college2)
  fit2 <- update(fit1, median ~ major_category + perc_men)
  fit3 <- update(fit1, median ~ major_category + perc_men + perc_employed)
  fit4 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround)
  fit5 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround + perc_college_jobs)
  anova(fit1, fit2, fit3, fit4)
  summary(fit5)
  tail(hatvalues(fit5)[order(hatvalues(fit5, decreasing = T))], 5)
  tail(cooks.distance(fit5)[order(cooks.distance(fit5, decreasing = T))], 5)
})
}
```

# Setup

```
str(college)
library(ggplot2)
library(dplyr)
g1 <- ggplot(college, aes(x = major_category, y = median, fill = major_category))
g1 <- g1 + geom_boxplot()
g1 <- g1 + theme(axis.text.x = element_text(angle = 45, hjust = 1), legend.position = "none")
g1
tapply(college$sample_size, college$major_category, sum)
majormean <- tapply(college$sample_size, college$major_category, mean)
max(majormean)
college$major_category <- as.factor(college$major_category)
college2 <- college
college2$major_category <- relevel(college$major_category, ref = "Business")
fit1 <- lm(median ~ major_category, data = college2)
fit2 <- update(fit1, median ~ major_category + perc_men)
fit3 <- update(fit1, median ~ major_category + perc_men + perc_employed)
fit4 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround)
fit5 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround + perc_college_jobs)
anova(fit1, fit2, fit3, fit4)
summary(fit5)
tail(hatvalues(fit5)[order(hatvalues(fit5, decreasing = T))], 5)
tail(cooks.distance(fit5)[order(cooks.distance(fit5, decreasing = T))], 5)
```

# EDA

```
str(college)
library(ggplot2)
library(dplyr)
g1 <- ggplot(college, aes(x = major_category, y = median, fill = major_category))
g1 <- g1 + geom_boxplot()
g1 <- g1 + theme(axis.text.x = element_text(angle = 45, hjust = 1), legend.position = "none")
g1
tapply(college$sample_size, college$major_category, sum)
majormean <- tapply(college$sample_size, college$major_category, mean)
max(majormean)
college$major_category <- as.factor(college$major_category)
college2 <- college
college2$major_category <- relevel(college$major_category, ref = "Business")
fit1 <- lm(median ~ major_category, data = college2)
fit2 <- update(fit1, median ~ major_category + perc_men)
fit3 <- update(fit1, median ~ major_category + perc_men + perc_employed)
fit4 <- update(fit1, median ~ major_category + perc_r
perc_employed_fulltime_yearround)
fit5 <- update(fit1, median ~ major_category + perc_r
perc_employed_fulltime_yearround + perc_college_jobs)
anova(fit1, fit2, fit3, fit4)
summary(fit5)
tail(hatvalues(fit5)[order(hatvalues(fit5, decreasing = T))], 5)
tail(cooks.distance(fit5)[order(cooks.distance(fit5, decreasing = T))], 5)
```

# Munging

```
str(college)
library(ggplot2)
library(dplyr)
g1 <- ggplot(college, aes(x = major_category, y = median, fill = major_category))
g1 <- g1 + geom_boxplot()
g1 <- g1 + theme(axis.text.x = element_text(angle = 45, hjust = 1), legend.position = "none")
g1
tapply(college$sample_size, college$major_category, :
majormean <- tapply(college$sample_size, college$maj
max(majormean)
college$major_category <- as.factor(college$major_ca
college2 <- college
college2$major_category <- relevel(college$major_category, ref = "Business")
fit1 <- lm(median ~ major_category, data = college2)
fit2 <- update(fit1, median ~ major_category + perc_men)
fit3 <- update(fit1, median ~ major_category + perc_men + perc_employed)
fit4 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround)
fit5 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround + perc_college_jobs)
anova(fit1, fit2, fit3, fit4)
summary(fit5)
tail(hatvalues(fit5)[order(hatvalues(fit5, decreasing = T))], 5)
tail(cooks.distance(fit5)[order(cooks.distance(fit5, decreasing = T))], 5)
```

# Modeling



```
str(college)
library(ggplot2)
library(dplyr)
g1 <- ggplot(college, aes(x = major_category, y = median, fill = major_category))
g1 <- g1 + geom_boxplot()
g1 <- g1 + theme(axis.text.x = element_text(angle = 45, hjust = 1), legend.position = "none")
g1
tapply(college$sample_size, college$major_category, FUN = function(x) {
  majormean <- tapply(x, college$sample_size, FUN = function(x) {
    max(x)
  })
  college$major_category <- as.factor(college$major_category)
  college2 <- college
  college2$major_category <- relevel(college2$major_category, ref = "Business")
  fit1 <- lm(median ~ major_category, data = college2)
  fit2 <- update(fit1, median ~ major_category + perc_men)
  fit3 <- update(fit1, median ~ major_category + perc_men + perc_employed)
  fit4 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround)
  fit5 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround + perc_college_jobs)
  anova(fit1, fit2, fit3, fit4)
  summary(fit5)
  tail(hatvalues(fit5)[order(hatvalues(fit5, decreasing = T))], 5)
  tail(cooks.distance(fit5)[order(cooks.distance(fit5, decreasing = T))], 5)
})
```

# Evaluation

# Setup

```
str(college)
library(ggplot2)
library(dplyr)

g1 <- ggplot(college, aes(x = major_category, y = median, fill = major_category))
g1 <- g1 + geom_boxplot()
g1 <- g1 + theme(axis.text.x = element_text(angle = 45, hjust = 1), legend.position = "none")
g1

tapply(college$sample_size, college$major_category, sum)
majormean <- tapply(college$sample_size, college$major_category,
max(majormean)

college$major_category <- as.factor(college$major_category)
college2 <- college
college2$major_category <- relevel(college$major_category, ref = "Business Administration")

fit1 <- lm(median ~ major_category, data = college2)
fit2 <- update(fit1, median ~ major_category + perc_men)
fit3 <- update(fit1, median ~ major_category + perc_men + perc_employed)
fit4 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround)
fit5 <- update(fit1, median ~ major_category + perc_men + perc_employed +
perc_employed_fulltime_yearround + perc_college_jobs)
anova(fit1, fit2, fit3, fit4)

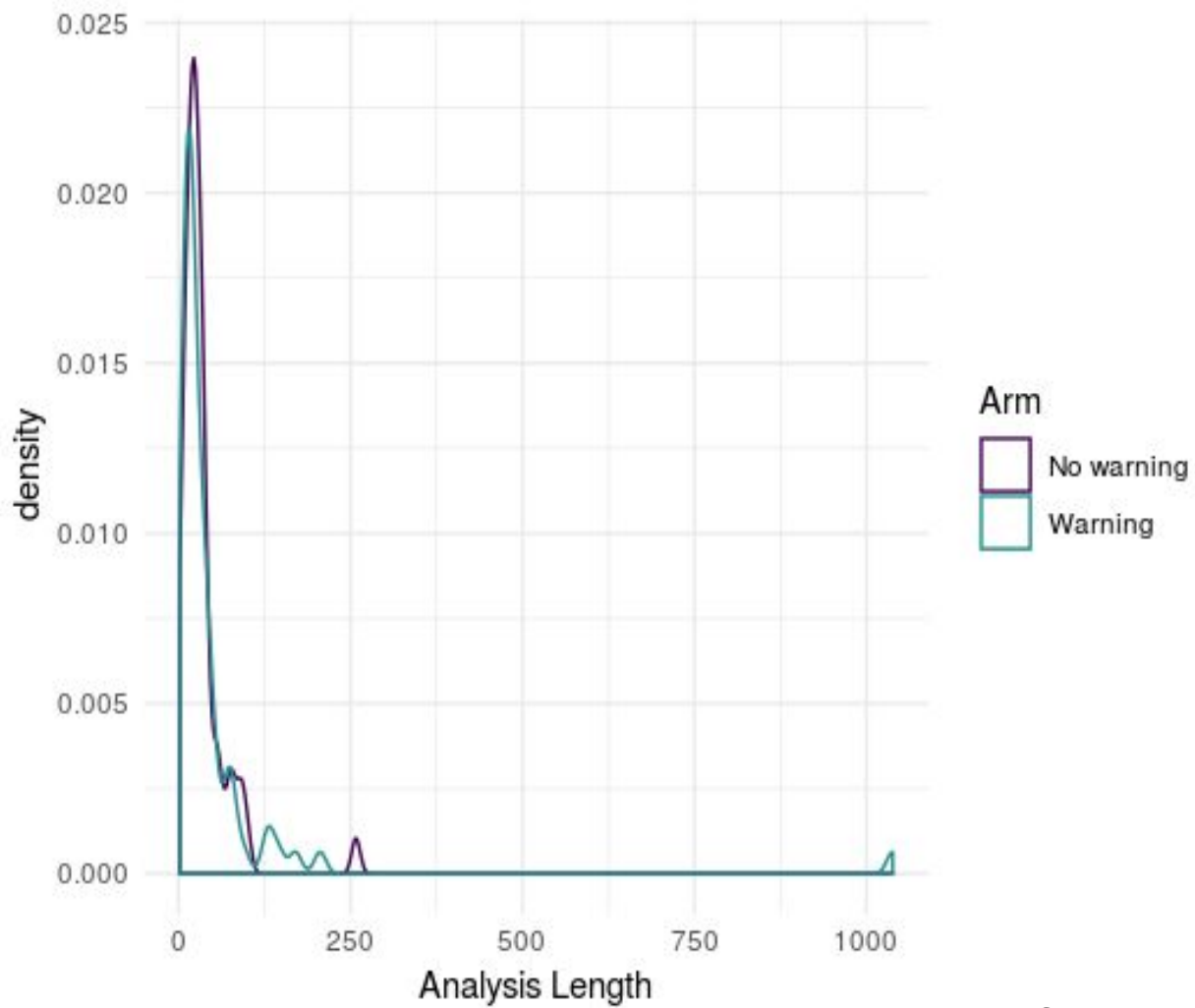
summary(fit5)
tail(hatvalues(fit5)[order(hatvalues(fit5, decreasing = T))], 5)
tail(cooks.distance(fit5)[order(cooks.distance(fit5, decreasing = T))], 5)
```

Exploratory

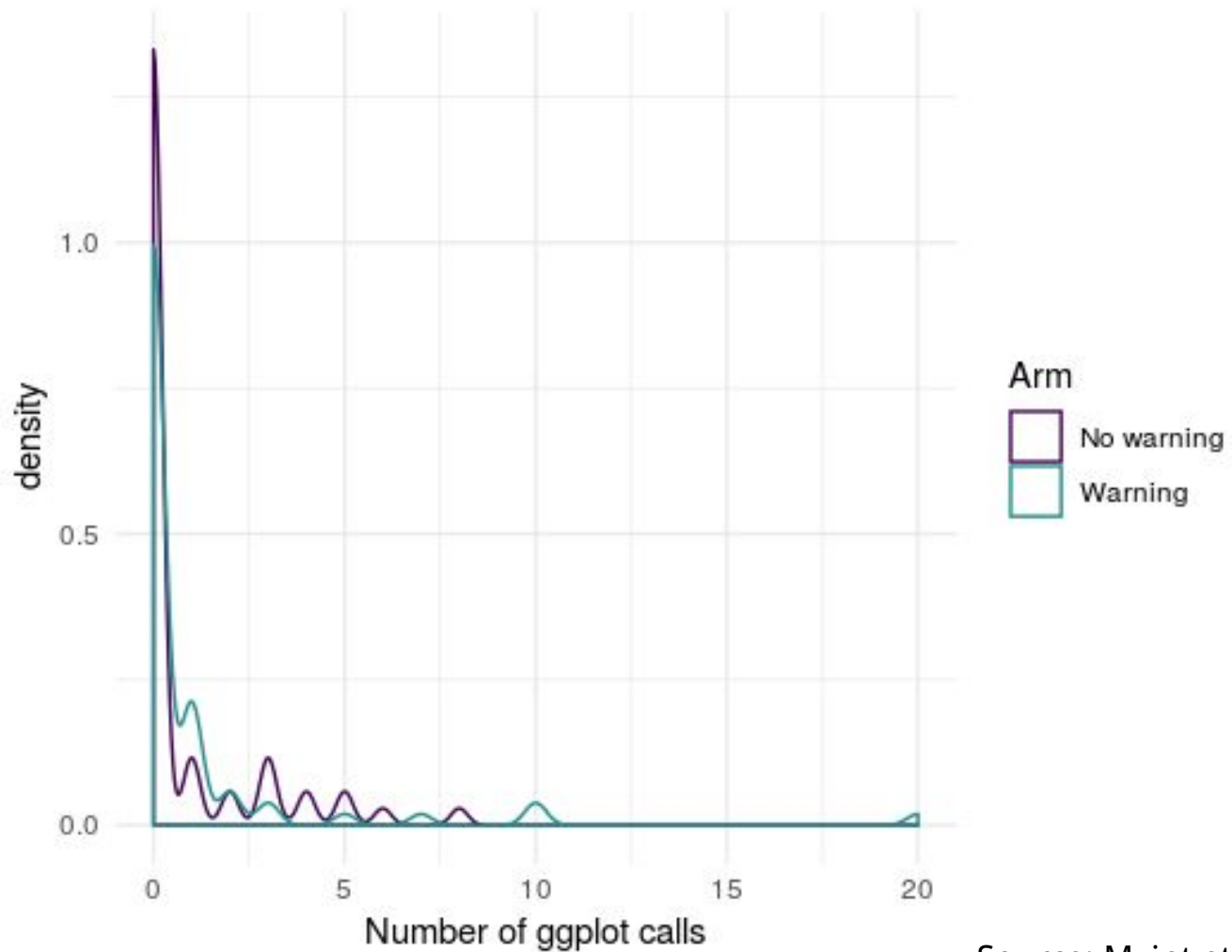
Munging

Modeling

Evaluation



Source: Myint et al (in prep)



Source: Myint et al (in prep)

