

# Reproducible Research: Medication or Prevention?

Roger D. Peng  
*@rdpeng, @simplystats, simplystatistics.org*

UW SISBID  
July 2016

# Report Writing for Data Science in R



Roger D. Peng

The R Series

# Implementing Reproducible Research



Edited by  
**Victoria Stodden**  
**Friedrich Leisch**  
**Roger D. Peng**

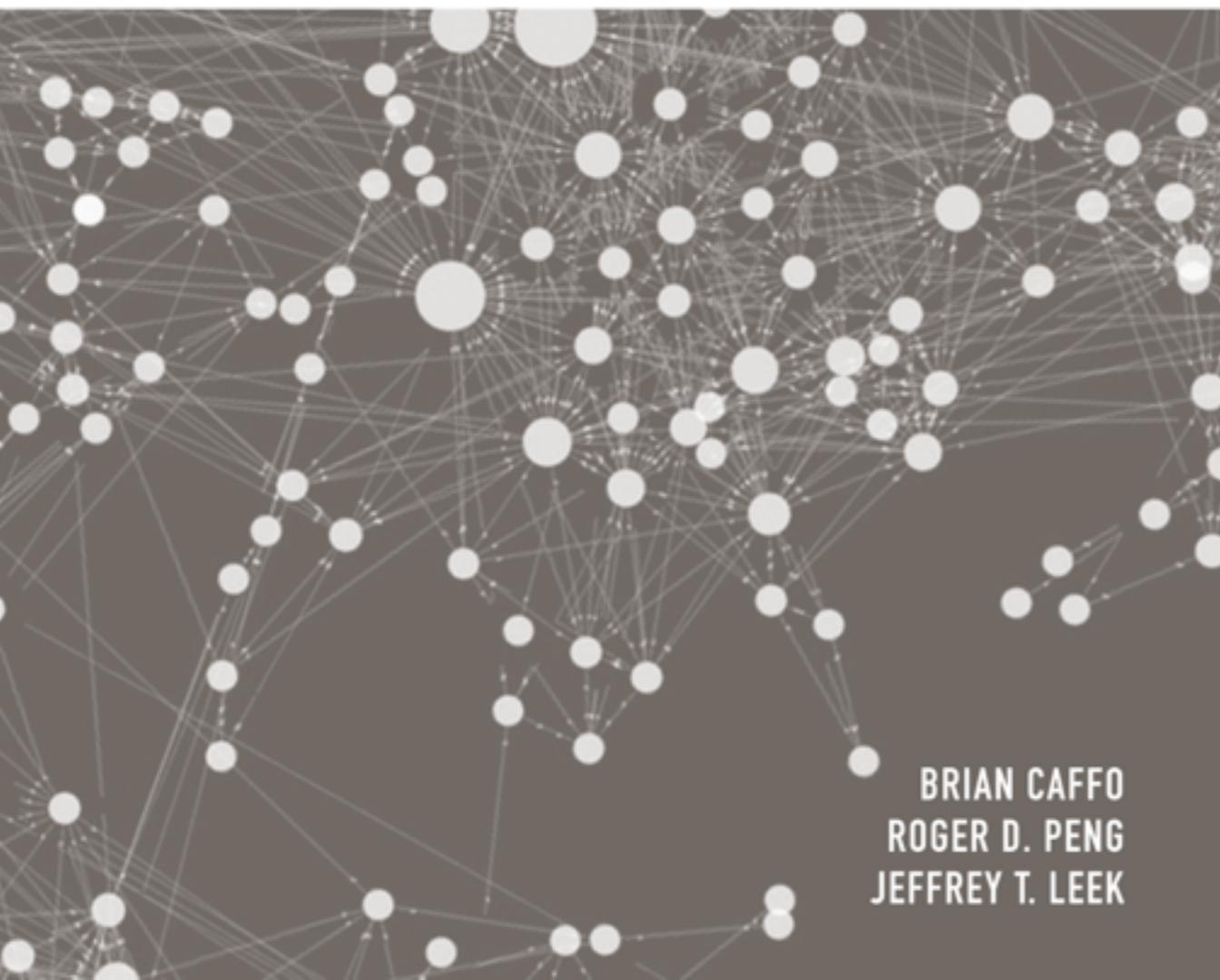
 CRC Press  
Taylor & Francis Group  
A CHAPMAN & HALL BOOK

[leanpub.com/reportwriting](https://leanpub.com/reportwriting)

Skip Ad ►

# EXECUTIVE DATA SCIENCE

A GUIDE TO TRAINING AND MANAGING THE BEST DATA SCIENTISTS



BRIAN CAFFO  
ROGER D. PENG  
JEFFREY T. LEEK

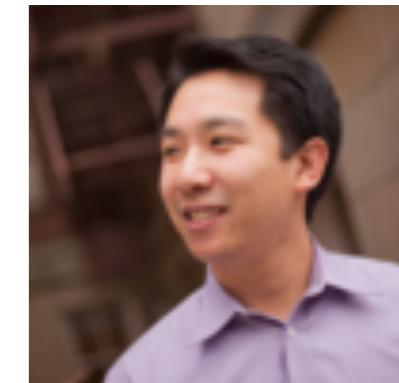
## Executive Data Science

A Guide to Training and Managing the Best Data Scientists

Brian Caffo, Roger D. Peng and Jeffrey Leek

---

This book teaches you how to assemble and lead a data science enterprise so that your organization can move towards extracting information from big data. This book is based on the acclaimed Johns Hopkins Executive Data Science Specialization. Printed copies of this book are available through Lulu (coming soon).



<https://leanpub.com/eds>

Skip Ad ►

# Parable

Scienceexpress

Report

## Genetic Signatures of Exceptional Longevity in Humans

Paola Sebastiani,<sup>1\*</sup> Nadia Solovieff,<sup>1</sup> Annibale Puca,<sup>2</sup> Stephen W. Hartley,<sup>1</sup> Efthymia Melista,<sup>3</sup> Stacy Andersen,<sup>4</sup> Daniel A. Dworkis,<sup>3</sup> Jemma B. Wilk,<sup>5</sup> Richard H. Myers,<sup>5</sup> Martin H. Steinberg,<sup>6</sup> Monty Montano,<sup>3</sup> Clinton T. Baldwin,<sup>6,7</sup> Thomas T. Perls<sup>4\*</sup>

<sup>1</sup>Department of Biostatistics, Boston University School of Public Health, Boston, MA 02118, USA. <sup>2</sup>IRCCS Multimedica, Milano, Italy; Istituto di Tecnologie Biomediche, Consiglio Nazionale delle Ricerche, Segrate, 20122, Italy. <sup>3</sup>Department of Medicine, Boston University School of Medicine, Boston, MA 02118, USA. <sup>4</sup>Section of Geriatrics, Department of Medicine, Boston University School of Medicine and Boston Medical Center, Boston, MA 02118, USA. <sup>5</sup>Department of Neurology, Boston University School of Medicine, Boston, MA 02118, USA. <sup>6</sup>Departments of Medicine and Pediatrics, Boston University School of Medicine and Boston Medical Center, Boston, MA 02118, USA. <sup>7</sup>Center for Human Genetics, Boston University School of Medicine, Boston, MA 02118, USA.

# Parable

## RETRACTION

*Post date 22 July 2011*

After online publication of our Report "Genetic signatures of exceptional longevity in humans" (1), we discovered that technical errors in the Illumina 610 array and an inadequate quality control protocol introduced false-positive single-nucleotide polymorphisms (SNPs) in our findings. An independent laboratory subsequently performed stringent quality control measures, ambiguous SNPs were then removed, and resultant genotype data were validated using an independent platform. We then reanalyzed the reduced data set using the same methodology as in the published paper. We feel the main scientific findings remain supported by the available data: (i) A model consisting of multiple specific SNPs accurately differentiates between centenarians and controls; (ii) genetic profiles cluster into specific signatures; and (iii) signatures are associated with ages of onset of specific age-related diseases and subjects with the oldest ages. However, the specific details of the new analysis change substantially from those originally published online to the point of becoming a new report. Therefore, we retract the original manuscript and will pursue alternative publication of the new findings.

PAOLA SEBASTIANI,<sup>1\*</sup> NADIA SOLOVIEFF,<sup>1</sup> ANNIBALE PUCA,<sup>2</sup> STEPHEN W. HARTLEY,<sup>1</sup> EFTHYMIA MELISTA,<sup>3</sup> STACY ANDERSEN,<sup>4</sup> DANIEL A. DWORKIS,<sup>3</sup> JEMMA B. WILK,<sup>5</sup> RICHARD H. MYERS,<sup>5</sup> MARTIN H. STEINBERG,<sup>6</sup> MONTY MONTANO,<sup>3</sup> CLINTON T. BALDWIN,<sup>6,7</sup> THOMAS T. PERLS<sup>4\*</sup>

# Parable

## Reproducibility!

What needs to happen next? For a start, **the authors should release the raw intensity data for their genotyping experiments**, which would allow independent investigators to spot obvious problems. Doing so immediately on a public database would go a long way towards showing they're not trying to cover up any methodological flaws. Ideally, they should also validate their putative associated SNPs using an independent platform and release those raw data as well

More broadly, **this is an important lesson for the increasing number of investigators wandering into the GWAS arena**: they need to be aware that the genotype data they're working with aren't just clean, digital data points, but best-guess estimates (typically very reliable, but sometimes badly flawed) based on a noisy fluorescent intensity signal. There's a reason why researchers working on GWAS spend so much of their time on a regimented series of upstream "data cleaning" steps and careful downstream validation of new associations – it's all too easy for noisy data to introduce bias that produces a false association signal. So, kids, don't end up in Newsweek for all the wrong reasons: talk to someone who really knows what they're doing when it comes to GWAS data.

# What Happens Everyday



# What's Next for Reproducibility?

- Reproducibility is critical for communicating a data analysis
- One cannot sufficiently describe an analysis in journal pages or supplementary materials
- General consensus about its importance
- No credible plan (yet) for how to implement such a requirement (hint: this is what's next)

# What Problem Does Reproducibility Solve?

- Transparency / Improved information transfer
- Data availability
- Software / Methods

# Reproducible Research at *Biostatistics*

*Biostatistics* (2009), **10**, 3, pp. 409–423  
doi:10.1093/biostatistics/kxp010  
Advance Access publication on April 17, 2009

R

*Biostatistics* (2012), **13**, 1, pp. 166–178  
doi:10.1093/biostatistics/kxr013  
Advance Access publication on June 17, 2011

R

## Significance analysis and statistical dissection of variably methylated regions

ANDREW E. JAFFE

Departments of Epidemiology and Biostatistics,  
Johns Hopkins Bloomberg School of Public Health,  
Baltimore, MD 21205, USA

ANDREW P. FEINBERG

Center for Epigenetics, Johns Hopkins University, Baltimore,  
MD 21205, USA

RAFAEL A. IRIZARRY, JEFFREY T. LEEK\*

Department of Biostatistics,  
Johns Hopkins Bloomberg School of Public Health,  
Baltimore, MD 21205, USA  
jleek@jhsph.edu

## Air pollution and health in Scotland: a multicity study

DUNCAN LEE\*, CLAIRE FERGUSON

Department of Statistics, University of Glasgow, Glasgow, G12 8QQ UK  
duncan@stats.gla.ac.uk

RICHARD MITCHELL

Public Health and Health Policy, University of Glasgow, Glasgow, G12 8QQ UK

*Biostatistics* (2009), **10**, 4, pp. 756–772  
doi:10.1093/biostatistics/kxp029  
Advance Access publication on July 27, 2009

C

## Second-order estimating equations for the analysis of clustered current status data

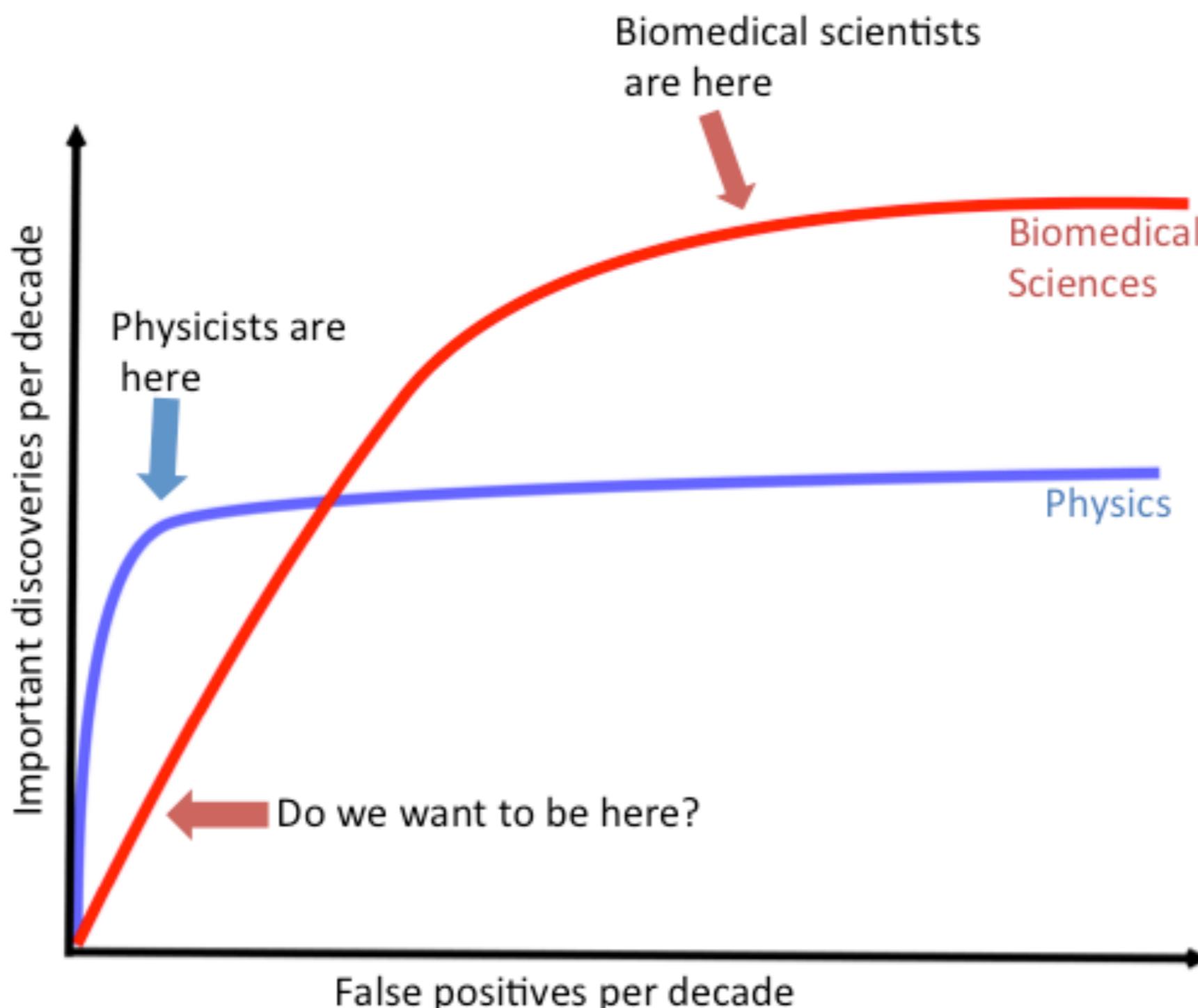
RICHARD J. COOK\*, DAVID TOLUSSO

Department of Statistics and Actuarial Science, University of Waterloo,  
Waterloo, ON, Canada N2L 3G1  
rjcook@uwaterloo.ca

# What Problem Does Reproducibility Solve?

- An analysis can be reproducible and still be wrong
- We want to know “can we trust this analysis?”
- Does requiring reproducibility deter bad analysis?
- While reproducibility is a key part of science, what people really want is *correct* or *well-done* research

# ROC Curves of Science



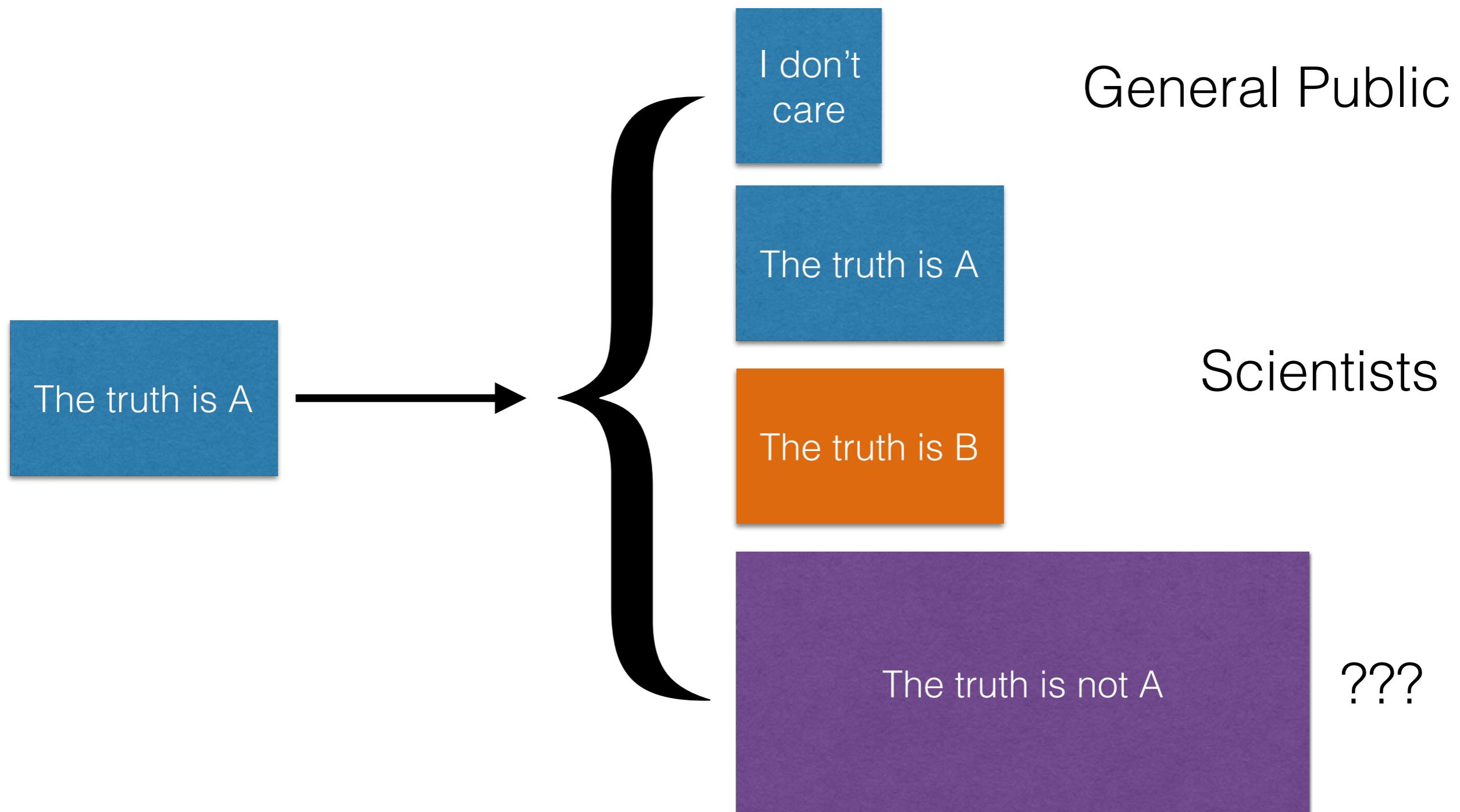
# Limitations of Reproducibility

- The premise of reproducible research is that with data/code available, people can check each other and the whole system is self-correcting
- Addresses the most “downstream” aspect of the research process – post-publication
- Assumes everyone is capable of doing analysis and wants to achieve the same goals (i.e. scientific discovery, search for truth)

# Who Reproduces Research?

- Someone needs to *do* something
  - Re-run the analysis; check results match
  - Check the code for bugs/errors
  - Try alternate approaches; check sensitivity
- The need for someone to do something is inherited from traditional notion of replication
- Who is “someone” and what are their goals?

# Who Reproduces Research?



# Primary Prevention

- Once bad research is published, it can be a long road to rectify (Duke episode is a prime example)
- Even reproducible research can be difficult to untangle unless examined by knowledgeable people (e.g. Baggerly and Coombes)
- How do we prevent shoddy / fraudulent research from appearing **in the first place?**

# Scary Bedtime Stories

ARTICLES

• Retracted •

nature  
medicine

## Genomic signatures to guide the use of chemotherapeutics

Anil Potti<sup>1,2</sup>, Holly K Dressman<sup>1,3</sup>, Andrea Bild<sup>1,3</sup>, Richard F Riedel<sup>1,2</sup>, Gina Chan<sup>4</sup>, Robyn Sayer<sup>4</sup>, Janiel Cragun<sup>4</sup>, Hope Cottrill<sup>4</sup>, Michael J Kelley<sup>2</sup>, Rebecca Petersen<sup>5</sup>, David Harpole<sup>5</sup>, Jeffrey Marks<sup>5</sup>, Andrew Berchuck<sup>1,6</sup>, Geoffrey S Ginsburg<sup>1,2</sup>, Phillip Febbo<sup>1–3</sup>, Johnathan Lancaster<sup>4</sup> & Joseph R Nevins<sup>1–3</sup>

RETRACTED 22 JULY 2011; SEE LAST PAGE

Scienceexpress Report

### Genetic Signatures of Exceptional Longevity in Humans

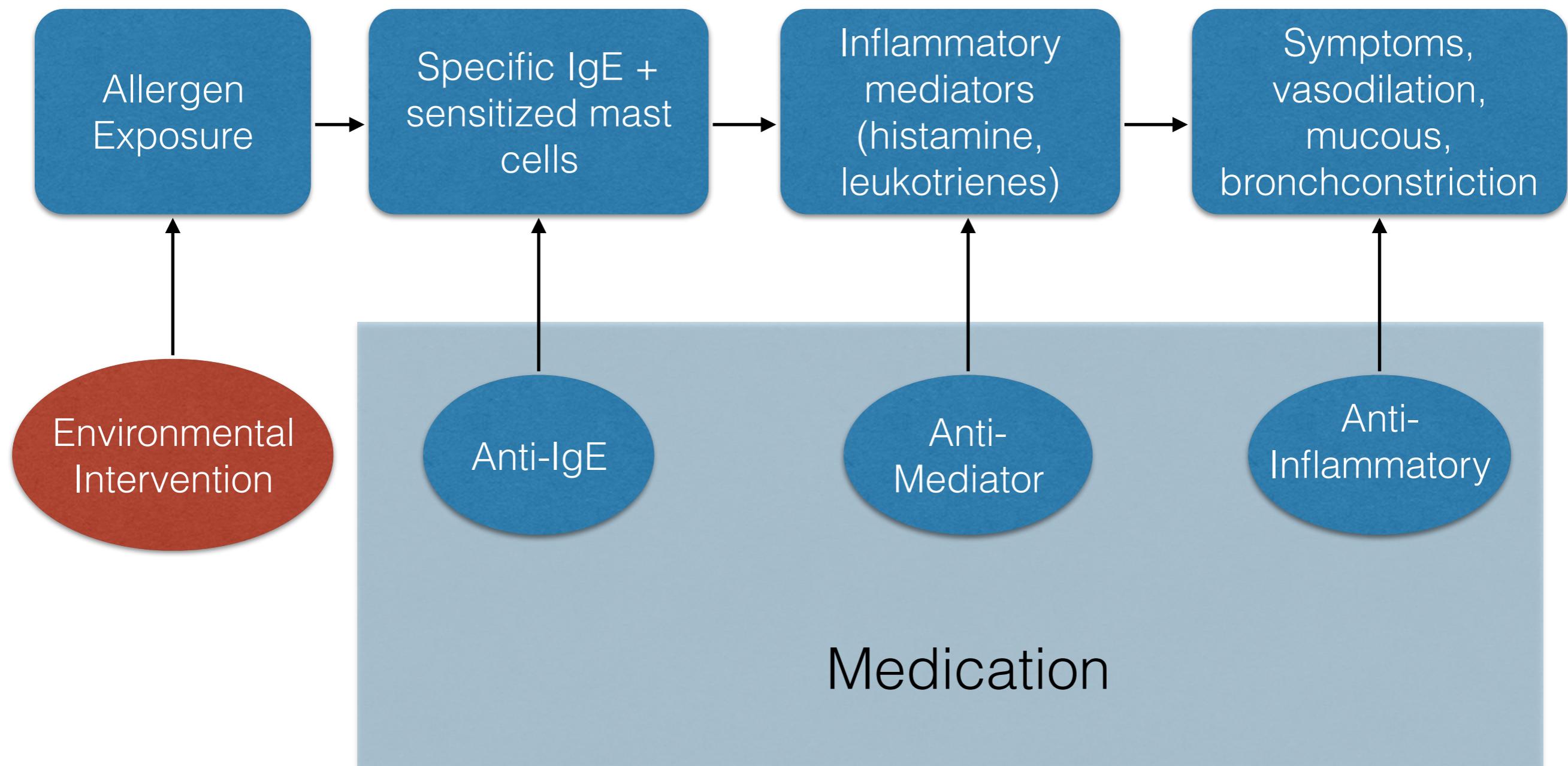
Paola Sebastiani,<sup>1\*</sup> Nadia Solovieff,<sup>1</sup> Annibale Puca,<sup>2</sup> Stephen W. Hartley,<sup>1</sup> Efthymia Melista,<sup>3</sup> Stacy Andersen,<sup>4</sup> Daniel A. Dworkis,<sup>3</sup> Jemma B. Wilk,<sup>5</sup> Richard H. Myers,<sup>5</sup> Martin H. Steinberg,<sup>6</sup> Monty Montano,<sup>3</sup> Clinton T. Baldwin,<sup>6,7</sup> Thomas T. Perls<sup>4\*</sup>

<sup>1</sup>Department of Biostatistics, Boston University School of Public Health, Boston, MA 02118, USA. <sup>2</sup>IRCCS Multimedica, Milano, Italy; Istituto di Tecnologie Biomediche, Consiglio Nazionale delle Ricerche, Segrate, 20122, Italy. <sup>3</sup>Department of Medicine, Boston University School of Medicine, Boston, MA 02118, USA. <sup>4</sup>Section of Geriatrics, Department of Medicine, Boston University School of Medicine and Boston Medical Center, Boston, MA 02118, USA. <sup>5</sup>Department of Neurology, Boston University School of Medicine, Boston, MA 02118, USA. <sup>6</sup>Departments of Medicine and Pediatrics, Boston University School of Medicine and Boston Medical Center, Boston, MA 02118, USA. <sup>7</sup>Center for Human Genetics, Boston University School of Medicine, Boston, MA 02118, USA.

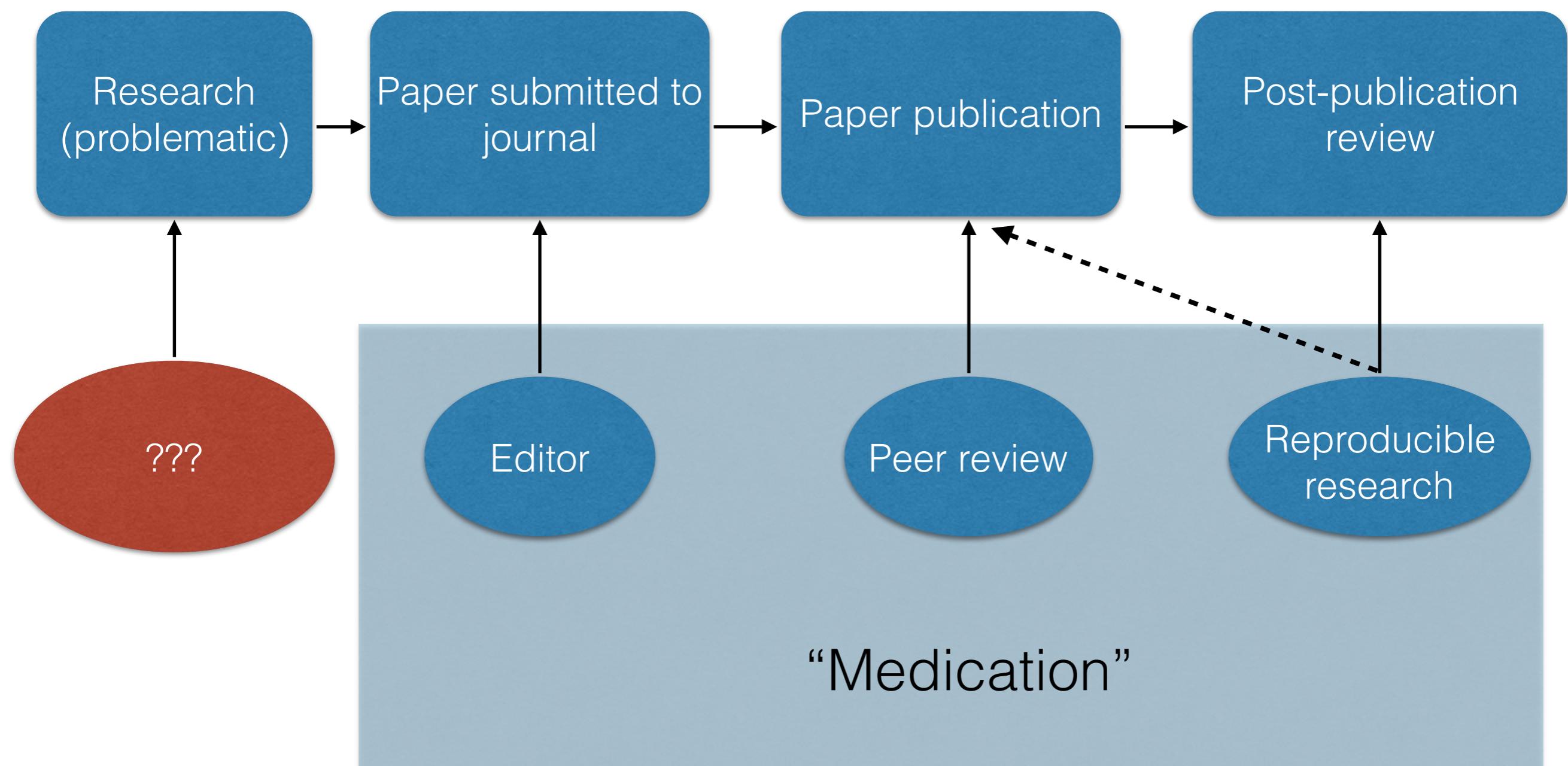
# Fundamental Problem in Statistics

**How is it that we have all this  
knowledge and no one seems  
to know about it?**

# An Analogy from Asthma



# Scientific Dissemination



# **THE FUTURE OF STATISTICAL SOFTWARE**

**Proceedings of a Forum**

**Panel on Guidelines for Statistical Software  
Committee on Applied and Theoretical Statistics**

**Board on Mathematical Sciences**

**Commission on Physical Sciences, Mathematics, and Applications  
National Research Council**

# Incorporating Statistical Expertise Into Software (1991)

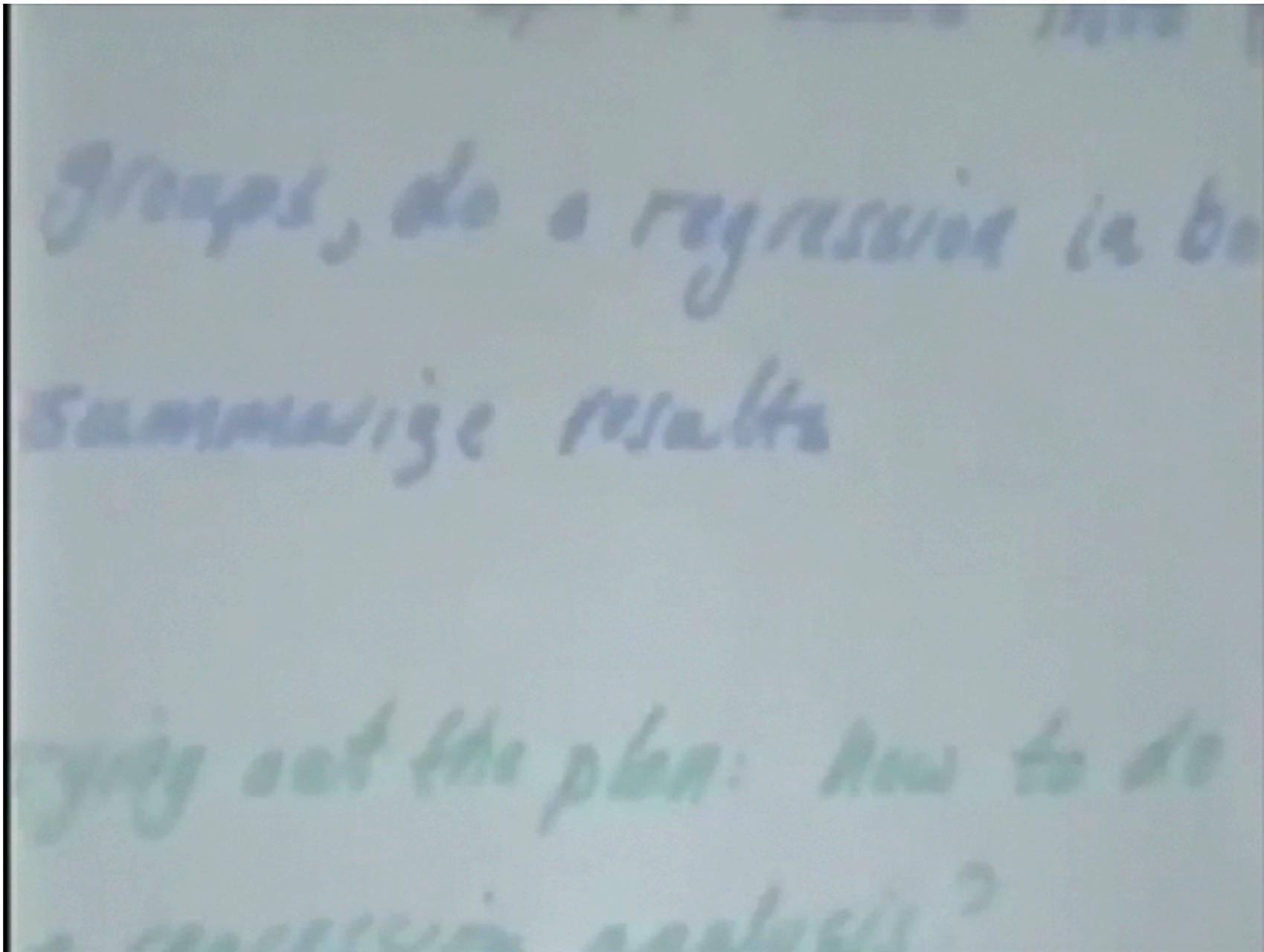
“Throughout American or even global industry, there is much advocacy of statistical process control and of understanding processes. **Statisticians have a process they espouse but do not know anything about.** It is the process of putting together many tiny pieces, the process called data analysis, and is not really understood.”

Daryl Pregibon, NRC Report 1991

<http://goo.gl/vsBUKn>

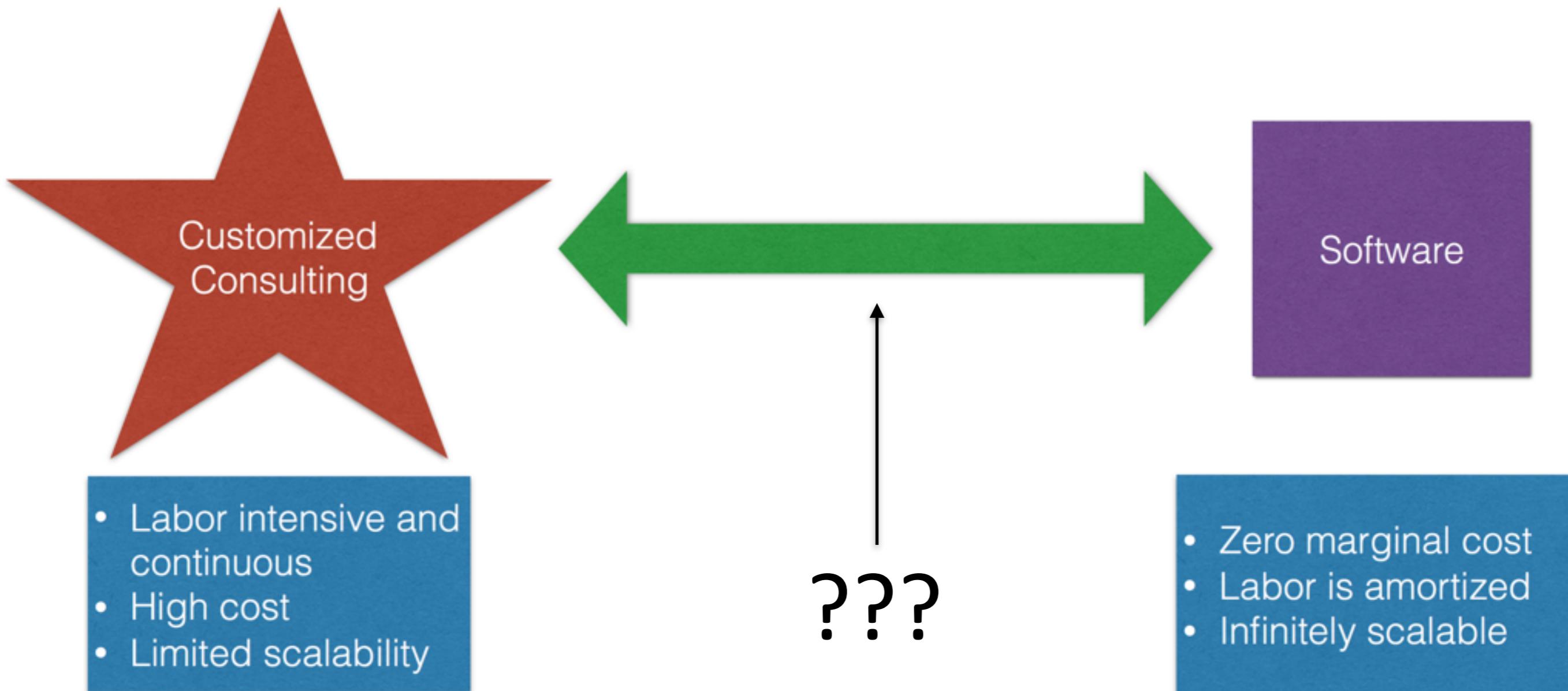
<http://youtu.be/8h96LgVpUrl>

# Data Analysis as an Algorithm?

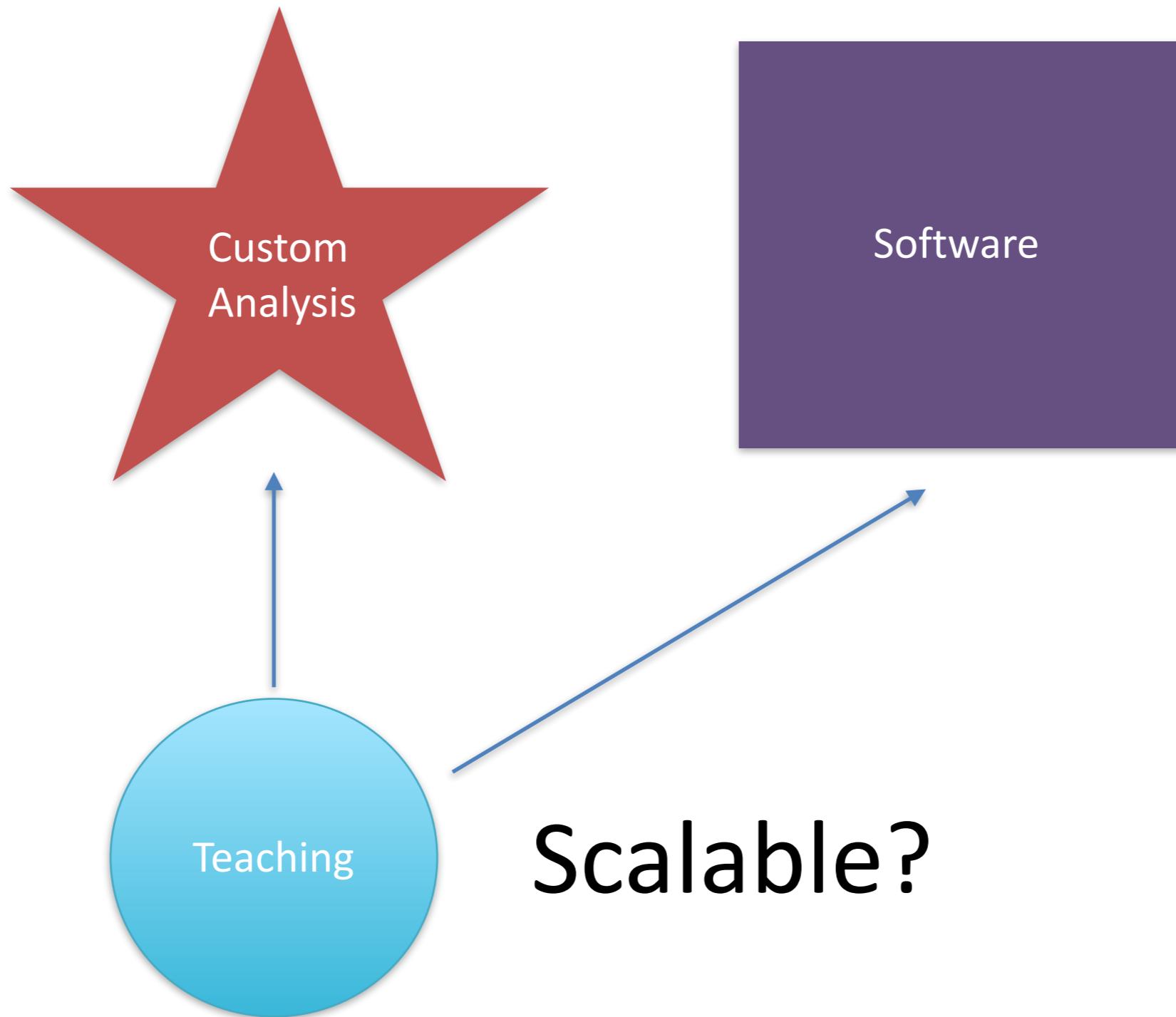


Daryl Pregibon (1987)

# Data Analysis Spectrum



# Scaling Data Analysis: The Third Option



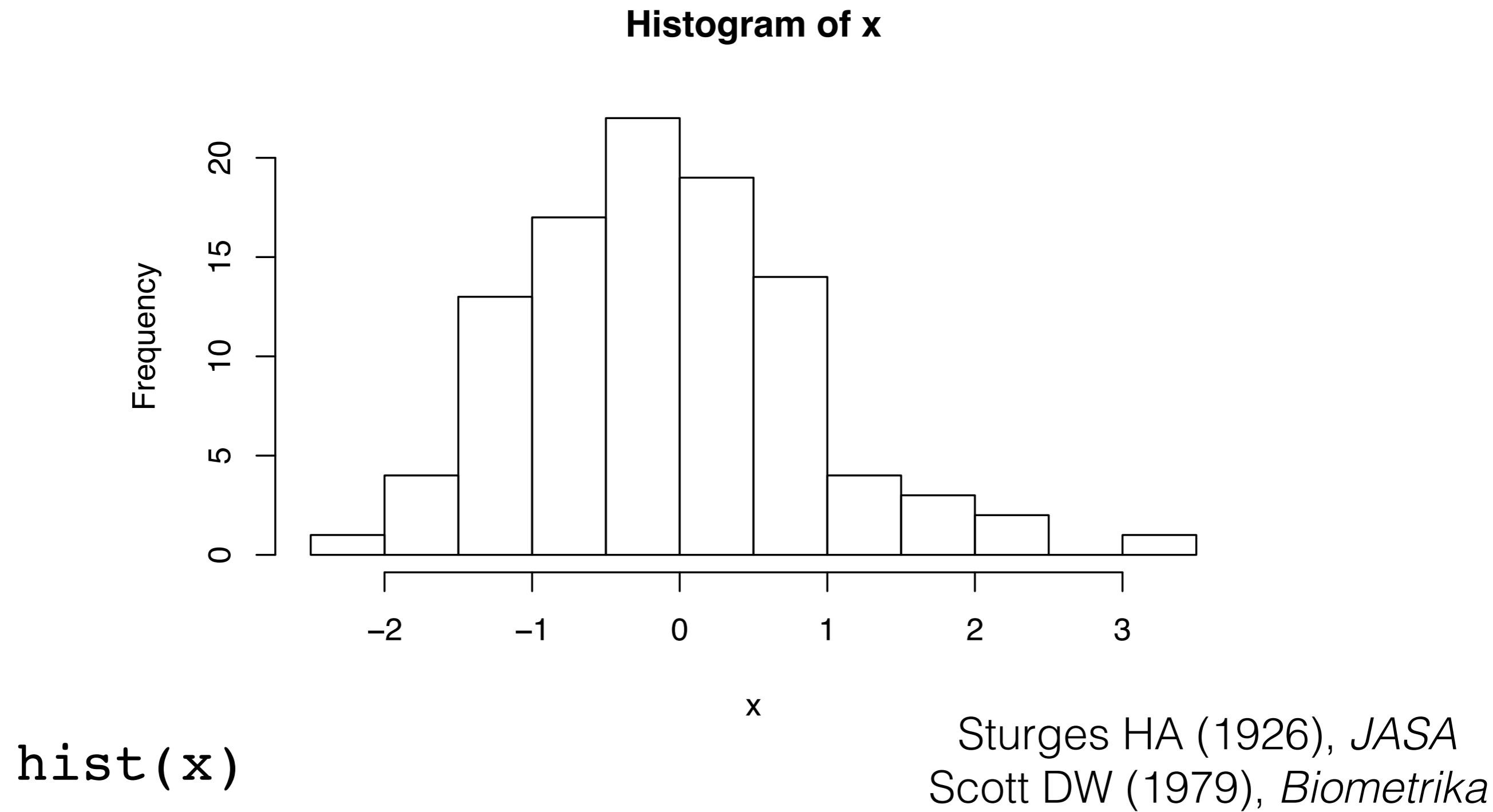
# “Evidence-based Data Analysis”

- Most data analyses involve stringing together many different tools and methods (“data science”)
- Some methods are standard for a given field, others are often applied ad hoc
- We should apply thoroughly studied, mutually agreed upon methods to analyze data whenever possible
- There should be evidence to justify the use of a given analysis method
- Many methods are used “off-label”; often okay, but....

# “Evidence-based Data Analysis”

- Create analytic pipelines from evidence-based components
- Create a standard by which we can judge deviations
- *Deterministic Statistical Machine* (<http://goo.gl/Qvlhuv>)
- Don’t mess with evidence-based pipeline (“transparent box” analysis)
- Reduce the “researcher degrees of freedom”
- Analogous to a pre-specified clinical trial protocol

# Evidence-based Histogram



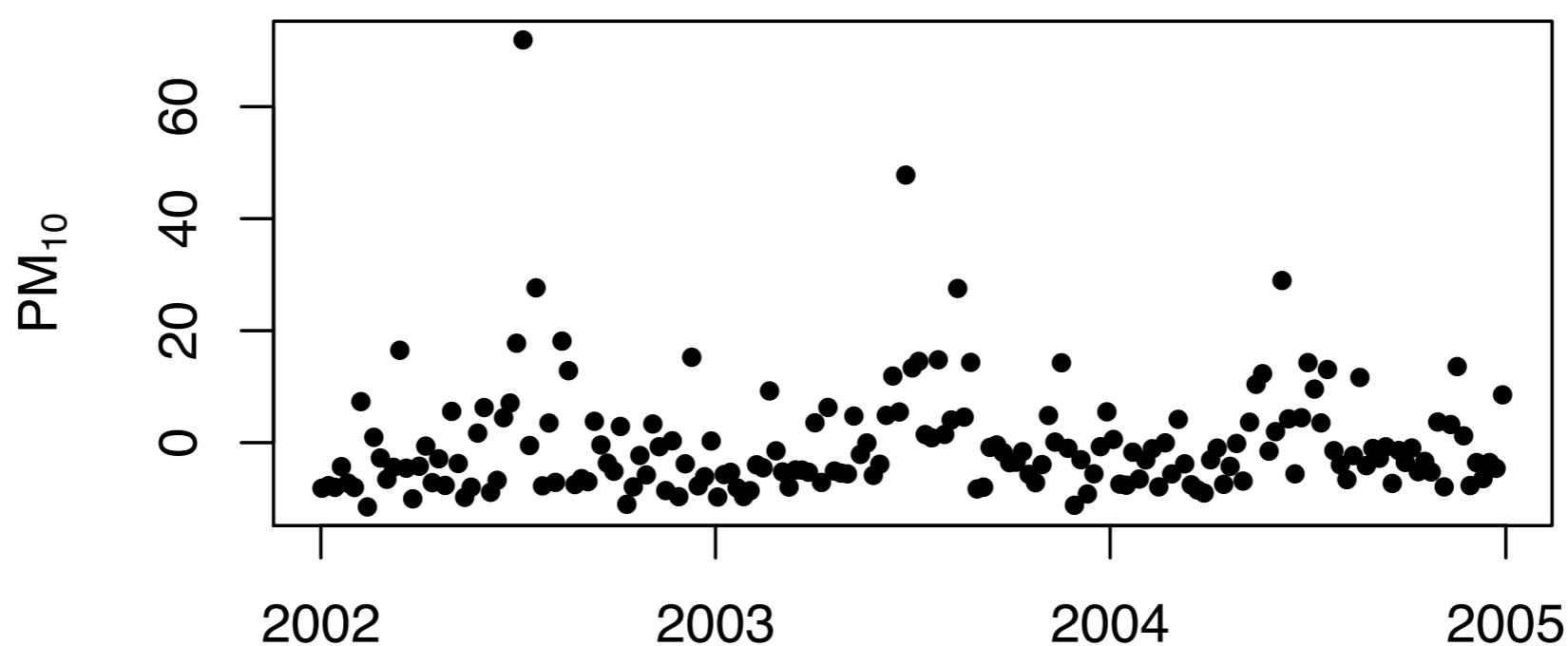
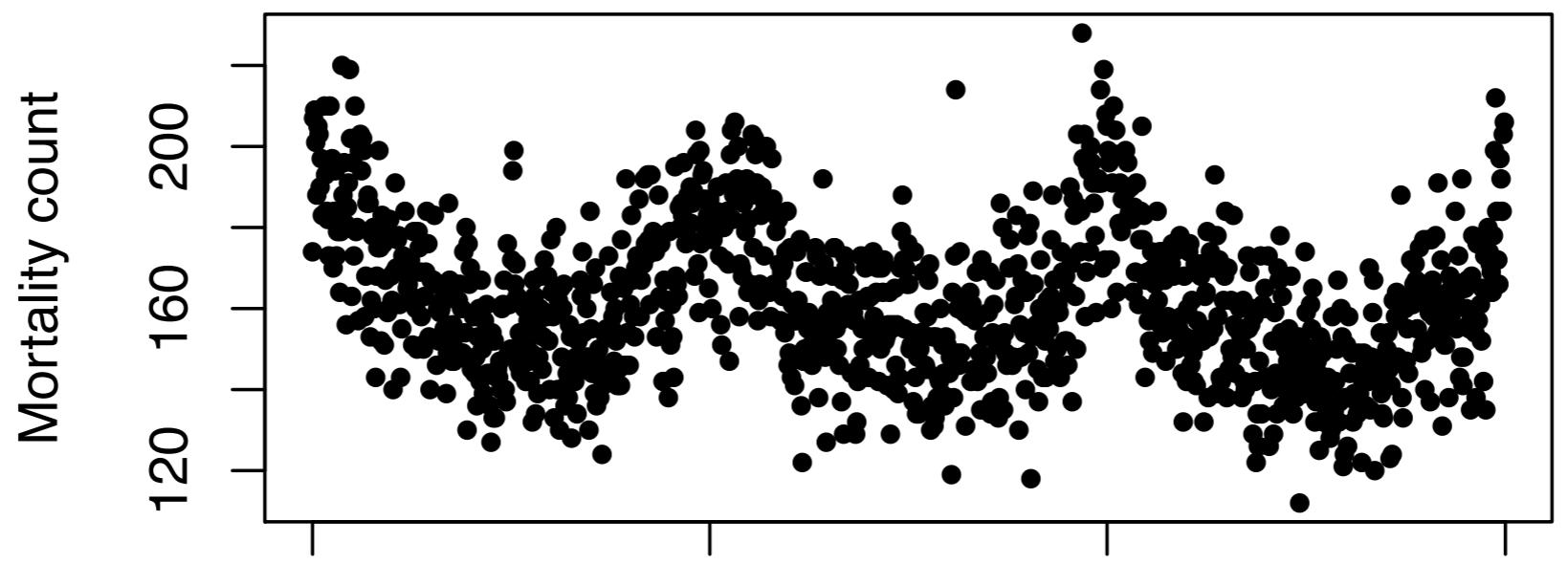
# Evidence-based Regression

- How do you solve the matrix equation  $X^T X \beta = X^T y$ ?
- Textbooks tell us  $\beta = (X^T X)^{-1} X^T y$  but we *never actually compute the solution this way*
- Matrix inversion is unstable; we know that because numerical analysts always yell at us
- R uses the QR decomposition; other decompositions like LU can be used too

# Case Study: Estimating Acute Effects of Ambient Air Pollution Exposure

- Acute/short-term effects typically estimated via panel studies or time series studies
- Work originated in late 1970s early 1980s
- Key question: “Are short-term changes in pollution associated with short-term changes in a population health outcome?”
- Studies usually conducted at community level
- Long history of statistical research investigating proper methods of analysis

# New York Data



# Case Study: Estimating Acute Effects of Ambient Air Pollution Exposure

- Can we encode everything that we have found in statistical/epidemiological research into a single package?
- Time series studies do not have a huge range of variation; typically involves similar types of data and similar questions
- We can create a deterministic statistical machine for this area?

# DSM Modules for Time Series Studies of Air Pollution and Health

1. Check for outliers, high leverage, overdispersion
2. Fill in missing data? NO!
3. Model selection: Estimate degrees of freedom to adjust for unmeasured confounders
  - Other aspects of model not as critical
4. Multiple lag analysis
5. Sensitivity analysis wrt
  - Unmeasured confounder adjustment
  - Influential points

Dominici, McDermott, Hastie (2004) JASA; Peng, Dominici, Louis (2006) JRSS-A

# Not Enough Geeks

06/28  
2011

## Critical Shortage Of “Data Geek” Talent Predicted By 2018



shortage of 1.5 million managers and analysts who have the ability to understand and make decisions using big data.

New research by the McKinsey Global Institute (MGI) forecasts a 50 to 60 percent gap between the supply and demand of people with deep analytical talent. These “data geeks” have advanced training in statistics machine learning as well as the ability to analyze data sets. The study projects there will be approximately 140,000 to 190,000 unfilled positions for data analytics experts in the U.S. by 2018 and a



# Scale and Reach?

	<b>Statistics Master's Degrees</b>	<b>2011</b>	<b>2012</b>	<b>2013</b>	<b>2003-2013</b>
		2011	2012	2013	
1	Columbia University in the City of New York	242	288	294	1943
2	Rutgers University-New Brunswick	47	62	79	576
3	Ohio State University-Main Campus	45	43	25	486
4	Stanford University	39	30	54	414
5	University of Michigan-Ann Arbor	44	47	55	407
6	University of Illinois at Urbana-Champaign	46	36	61	373
7	California State University-East Bay	49	43	55	354
8	Cornell University	35	51	54	346
9	Michigan State University	44	36	25	341
10	North Carolina State University at Raleigh	29	38	28	329

# Johns Hopkins Data Science Specialization

The Johns  
Hopkins Data  
Science  
Specialization



Powered by Tumblr  
Colorfall Theme by Paul Mackenzie



1. [The Data Scientist's Toolbox](#) - Get yourself set up.
2. [R programming](#) - Learn to code.
3. [Getting and Cleaning Data](#) - You need data. Get some.
4. [Exploratory Data Analysis](#) - What's that in my data?
5. [Reproducible Research](#) - Did you do what you think you did?
6. [Statistical Inference](#) - You don't have infinite money. Try sampling.
7. [Regression Models](#) - The duct tape of data science.
8. [Practical Machine Learning](#) - Predict the future with data. Easy.
9. [Developing Data Products](#) - There better be an app for that data.

# Johns Hopkins Data Science Specialization

- Data science toolbox
- Probability / math stat, statistical inference
- Getting + cleaning data
- R programming
- Regression modeling / machine learning
- **Reproducible research tools**
- Exploratory data analysis
- Data products
- Capstone project with industry partners

# Genomic Data Science

About This Specialization

Courses

Pricing

Creators

FAQs

Genomic Data  
Science  
Specialization

From \$49 USD

Enroll

Starts Aug 1



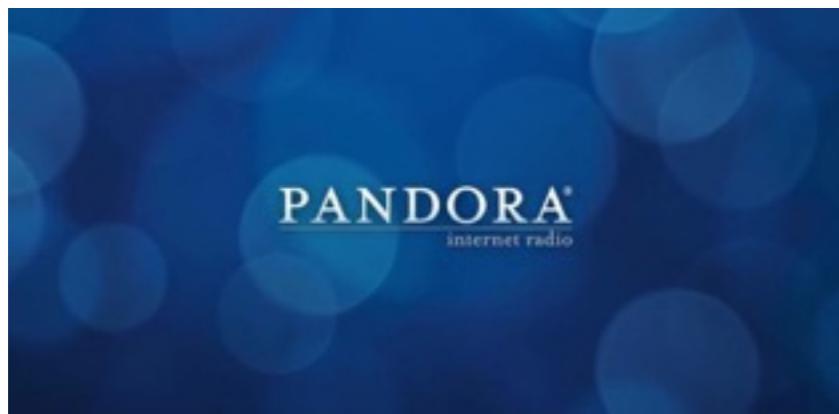
## Become a next generation sequencing data scientist

Master the tools and techniques at the forefront of the sequencing data revolution.

About This Specialization

This specialization covers the concepts and tools to understand, analyze, and interpret data from next generation sequencing experiments. It teaches the most common tools used in genomic data science including how to use the command line, Python, R, Bioconductor, and Galaxy. The sequence is a stand alone introduction to genomic data science or a perfect compliment to a primary degree or postdoc in biology, molecular biology, or genetics.

# Where's the Curation?



# Curation for Data Analysis

- Provide packages that encode data analysis pipelines for given problems, technologies, questions
- Curated by experts knowledgeable in the various areas (i.e. statisticians)
- Documentation given supporting each module in the pipeline
- Changes introduced after passing relevant benchmarks and unit tests
- We already have the tools and much of the knowledge

# Cochrane Summaries



**COCHRANE SUMMARIES**

Independent high-quality evidence for health care decision making

English | 简体中文 | Hrvatski  
| Français | Português  
| Español

[to Cochrane.org](#)  
[to The Cochrane Library](#)

Enter words like "aspirin for headaches" or "vaccines for influenza"

[Browse health topics](#) | [New and updated](#)

A product of The Cochrane Collaboration

[How to use this site](#)

*If you are looking for information you can trust,  
to help you make choices about health care, this site is for you.*

---

<http://summaries.cochrane.org>

# Cochrane Summaries



**COCHRANE SUMMARIES**

Trusted evidence. Informed decisions. Better health.

English | 简体中文 | 繁體中文

Hrvatski | Français | Português | Español

[to Cochrane.org](#)

[to The Cochrane Library](#)

ibuprofen for headaches



Search

[Browse health topics](#) |  [Retain current filters](#) | [Reset](#)

A product of The Cochrane Collaboration

[How to use this site](#)

## Filter by

### Search results: 6

Current search:  ibuprofen for headaches  Dentistry & oral health

Sort by: Best match

Results per page: 10

[Subscribe to this search](#) [More...](#)

[ShareThis](#)



[Like](#) { 0 }

#### Health topics

- ▶ Child health (1)
- ▼ Dentistry & oral health
  - ▶ Antibiotic therapy (1)
  - ▶ Craniofacial anomalies (2)
  - ▶ Oral pain (3)
- ▶ Rheumatology (1)

#### ▼ Cochrane group topics ?

- ▶ Child Health (1)
- ▶ Complementary Medicine (1)
- ▶ Nursing Care (1)
- ▶ Oral Health (6)
- ▶ Pre-hospital and Emergency Care (1)

#### Ibuprofen versus paracetamol (acetaminophen) for pain relief after surgical removal of lower wisdom teeth

12 Dec 2013

Review question This review, carried out by the Cochrane Oral Health Group, seeks to compare the effectiveness of two commonly used painkillers, paracetamol and ibuprofen and the combination of both in a single tablet in the relief of pain following surgical removal of lower wisdom teeth.

Background Worldwide the number of surgical operations ...

#### Does giving children painkillers such as paracetamol and ibuprofen before dental treatment help reduce pain after the treatment?

12 Sep 2012

Dental pain is common after dental procedures and can be linked with increase in fear of dental treatment, avoidance of dental treatment and other associated problems. Reducing pain is important, particularly in children and adolescents. One way of managing this might be to give painkillers before treatment so that they can start to work right away. ...

#### Antibiotic use for severe toothache (irreversible pulpitis)

18 Dec 2013

Review question Are oral antibiotics effective and safe for treating pain in irreversible pulpitis (inflammation of the nerve inside the tooth/nerve damage)?

# Cochrane Summaries

## **Ibuprofen versus paracetamol (acetaminophen) for pain relief after surgical removal of lower wisdom teeth**

Bailey E, Worthington HV, van Wijk A, Yates JM, Coulthard P, Afzal Z

Published Online: 12 December 2013

### **Key results**

Ibuprofen is more effective than paracetamol at all doses studied in this [review](#). On limited evidence, the combination of ibuprofen and paracetamol appeared to be no more effective than the single drugs when measured two hours after surgery. However, again on limited evidence, it was found to be more effective than the drugs taken singly when measured at six hours after surgery. Participants taking the combined drug also had a smaller chance of requiring rescue medication.

### **Quality of the evidence**

All of the results (outcomes) comparing ibuprofen to paracetamol are of high quality. This means that further research is very unlikely to change our confidence in the estimates of the effect.

<http://goo.gl/2Lz3ax>

# We Need to Avoid...

Lots of data

+

Not sure what to  
do with data

=

Whatever works or  
is easiest!

# Summary

- Reproducibility is important, but likely would not have prevented recent notorious examples
- Reproducibility focuses on the most “downstream” aspect of the research dissemination process
- Evidence-based data analysis would provide standardized, best practices for given scientific areas and questions
- With development of personalized medicine, poor data analysis has the potential to seriously harm people
- More effort should be put into improving the quality of “upstream” aspects of scientific research