

About descriptive analyses

Goal: Describe a set of data

- The first kind of data analysis performed
- Commonly applied to census data
- The description and interpretation are different steps
 - Descriptions can usually not be generalized without additional statistical modeling



Types of Data Science Questions

Jeffrey Leek
Johns Hopkins Bloomberg School of Public Health

3/17

Descriptive analysis



<http://www.census.gov/2010census/>

Types of Data Science Questions

In approximate order of difficulty

- Descriptive
- Exploratory
- Inferential
- Predictive
- Causal
- Mechanistic

4/17

2/17

About predictive analysis

Goal: To use the data on some objects to predict values for another object

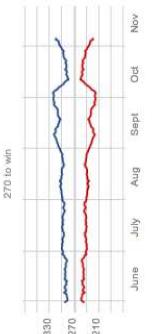
- If X predicts Y it does not mean that X causes Y
- Accurate prediction depends heavily on measuring the right variables
- Although there are better and worse prediction models, more data and a simple model works really well
- Prediction is very hard, especially about the future references

11/17

Predictive analysis

FiveThirtyEight Forecast

Updated 10:10 AM ET on Nov. 6	
President Nov. 6 Forecast	Senate Nov. 6 Forecast
Barack Obama	Mitt Romney
313.0 +14.0 since Oct. 30	Electoral vote 225.0 -14.0 since Oct. 30
270 to win	



<http://fivethirtyeight.blogs.nytimes.com/>

Inferential analysis

< Previous Article | Next Article >

Epidemiology:

January 2013 - Volume 24 - Issue 1 - p 23–31
doi: 10.1097/EDE.0b013e3182770237
Air Pollution

Effect of Air Pollution Control on Life Expectancy in the United States: An Analysis of 545 U.S. Counties for the Period from 2000 to 2007

Correia, Andrew W.^a; Pope, C. Arden III^b; Dockery, Douglas W.^c; Wang, Yun^a; Ezzati, Majid^d; Dominici, Francesca^a

FREE

SDC

Article Outline

[Correia et al. \(2013\) Epidemiology](#)

About inferential analysis

Goal: Use a relatively small sample of data to say something about a bigger population

- Inference is commonly the goal of statistical models
- Inference involves estimating both the quantity you care about and your uncertainty about your estimate
- Inference depends heavily on both the population and the sampling scheme

9/17

Inferential analysis

< Previous Article | Next Article >

Epidemiology:

January 2013 - Volume 24 - Issue 1 - p 23–31
doi: 10.1097/EDE.0b013e3182770237
Air Pollution

Effect of Air Pollution Control on Life Expectancy in the United States: An Analysis of 545 U.S. Counties for the Period from 2000 to 2007

Correia, Andrew W.^a; Pope, C. Arden III^b; Dockery, Douglas W.^c; Wang, Yun^a; Ezzati, Majid^d; Dominici, Francesca^a

FREE

SDC

Article Outline

[Correia et al. \(2013\) Epidemiology](#)

12/17

10/17

Causal analysis

Predictive analysis

The screenshot shows the homepage of The New England Journal of Medicine. At the top, there are links for 'SUBSCRIBE OR RENEW TODAY'. Below that, a search bar has 'SUBMIT YOUR MANUSCRIPT' and 'SEARCH' buttons. The main content area features an article titled 'Duodenal Infusion of Donor Feces for Recurrent Clostridium difficile' by Van Nood et al. The article summary states: 'Recurrence of C. difficile infection is difficult to treat, and failure rates for antibiotic therapy are high. We studied the effect of duodenal infusion of donor feces in patients with recurrent C. difficile infection.' The article is dated January 16, 2013, and has a DOI of 10.1053/j.nejmoa1206337. There are sections for 'BACKGROUND', 'METHODS', 'RESULTS', and 'CONCLUSIONS'. A figure titled 'FIGURE 1' is shown. At the bottom, there are links for 'Abstract', 'Article', 'References', and 'Comments'.

van Nood et al. (2013) NEJM

15/17

The screenshot shows the Forbes website. The header includes a logo for 'Target' and links for 'New Posts', 'Most Popular', 'Video', and 'Lists'. The main article is titled 'How Target Figured Out A Teen Girl Was Pregnant Before Her Father Did' by Katherine Hill. The article discusses how Target used consumer data to identify pregnant customers. It includes a quote from a Target executive: 'Every time you go shopping, you share intimate details about your consumption patterns with retailers. And many of those retailers are mining that data to figure out what you are most likely to make you happy.' The article also mentions a study by Jennifer Bartz, which found that buying diapers is a good indicator of pregnancy. The article has 13.7k views and 5.6k likes. At the bottom, there are links for 'Share', 'Email', 'Print', 'Download Citation', 'Supplementary Material', 'Comments', and 'Topics'.

<http://www.forbes.com/sites/kashmirhill/2012/02/16/how-target-figured-out-a-teen-girl-was-pregnant-before-her-father-did/>

13/17

About mechanistic analysis

Goal: Understand the exact changes in variables that lead to changes in other variables for individual objects.

- Incredibly hard to infer, except in simple situations
- Usually modelled by a deterministic set of equations (physical/engineering science)
- Generally the random component of the data is measurement error
- If the equations are known but the parameters are not, they may be inferred with data analysis

About causal analysis

Goal: To find out what happens to one variable when you make another variable change.

- Usually randomized studies are required to identify causation
- There are approaches to inferring causation in non-randomized studies, but they are complicated and sensitive to assumptions
- Causal relationships are usually identified as average effects, but may not apply to every individual
- Causal models are usually the "gold standard" for data analysis

16/17

14/17

Definition of data

“ Data are values of qualitative or quantitative variables, belonging to a set of items.

<http://en.wikipedia.org/wiki/Data>

Definition of data

“ Data are values of qualitative or quantitative variables, belonging to a **set of items**.

<http://en.wikipedia.org/wiki/Data>

Set of items: Sometimes called the population; the set of objects you are interested in

Mechanistic analysis

The slide has a dark blue header with the FHWA logo and a yellow sidebar with the title 'Mechanistic - Empirical Pavement Design'. The main content area has a white background with a dark blue footer bar containing the URL 'http://www.fhwa.dot.gov/resourcecenter/teams/pavement/pave_3pdg.pdf'.

Mechanistic - Empirical Pavement Design

Deployment Process: After a successful pilot project, FHWA organized the Digital Guidebook Team (DGT) to form the FHWA, design offices, State highway agencies, industry partners, and other organizations and experts about the deployment process and to discuss implementation. To introduce the guide and to discuss implementation issues, the DGT has developed a one-day workshop.

Participants will be invited from across the Nation. More details will be provided at a later date. The meeting will be held at Novacor, West Valley City, Salt Lake City, UT, on May 29, 2007.

For more information, contact: Jeffrey Leek, INHS, Indianapolis, IN (July); Vicki Jancz, Kansas City, KS (August); Mystic, CT (August); and Phoenix, AZ (October).

http://www.fhwa.dot.gov/resourcecenter/teams/pavement/pave_3pdg.pdf

17/17



What is data?

Jeffrey Leek
Johns Hopkins Bloomberg School of Public Health

3/13

What do data look like?

http://brianknaus.com/software/srtoolbox/s_4_1_sequence80.txt

What do data look like?



<https://dev.twitter.com/docs/ani/1/get/blocks/blocking>

Definition of data

“ Data are values of qualitative or quantitative variables, belonging to a set of items. ”

<http://en.wikipedia.org/wiki/Data>

Definition of data

“ Data are values of qualitative or quantitative **variables**, belonging to a set of items.

<http://en.wikipedia.org/wiki/Data>

Variables: A measurement or characteristic of an item.

4/13

Definition of data

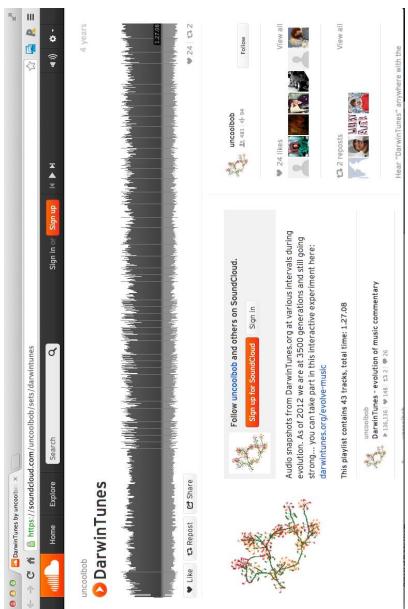
“ Data are values of qualitative or quantitative variables, belonging to a set of items. ”

<http://en.wikipedia.org/wiki/Data>

Quantitative: Height, weight, blood pressure

What do data look like?

What do data look like?



<https://www.pnas.org/content/109/30/12081.full> | <https://soundcloud.com/uncoolbob/sets/darwintunes>

10/13

ALLERGIES	
Last Updated:	01 Dec 2011 @ 0855
Medication Name:	TRIMETHOPRIN
location:	DAVITZ9
date Entered:	09 Mar 2011
action:	initial
Allergy Type:	DRUGS
A Drug Class:	ANTI-INFECTIVES, OTHER
Observed/Historical:	HISTORICAL
Comments:	The reaction to this allergy was MILD (NO SQUELAE)
Medication Name:	TRANZOOL
location:	DAVITZ9
are Entered:	initial
action:	URINARY RETENTION
Allergy Type:	DRUGS
A Drug Class:	MUS-CPTOID ANALGESICS
Observed/Historical:	HISTORICAL
Comments:	gradually worsening difficulty emptying bladder

<http://blue-button.github.com/challenge/>

8/13

What do data look like?

What do data look like?



<http://www.data.gov/>

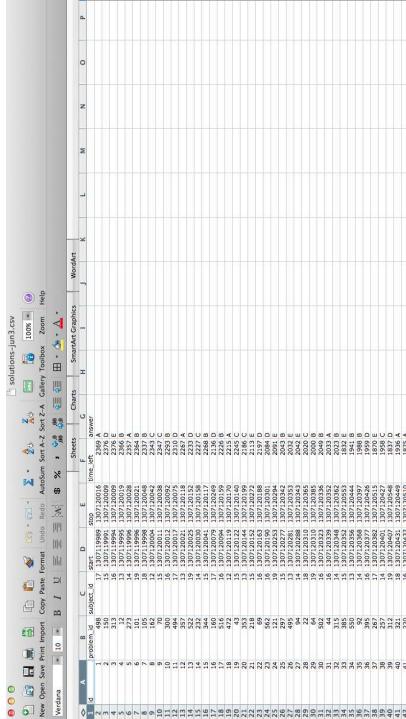
<http://www.nytimes.com/2012/06/26/technology/in-a-big-network-of-computers-evidence-of-machine-learning.html?pagewanted=all&r=0>



11/13

9/13

What do data look like? Rarely



What about big data?

Jeffrey Leek
Johns Hopkins Bloomberg School of Public Health



How much is there?

THE WORLD'S INFORMATION IS DOUBLING
EVERY TWO YEARS, WITH A COLOSSAL
1.8 zettabytes
to be created & replicated in
2011



<http://mashable.com/2011/06/28/data-infographic/>

The data is the second most important thing

- The most important thing in data science is the question
- The second most important is the data
- Often the data will limit or enable the questions
- But having data can't save you if you don't have a question

Why big data now?

So what about big data?

An Experimental Study of the Small World Problem*

JEFFREY TRAVERS

Harvard University

AND

STANLEY MILGRAM

The City University of New York

Arbitrarily selected individuals ($N=396$) in Nebraska and Boston are asked to generate acquaintance chains "through a person in one's social network, employing the small world method" (Milgram, 1967). Sixty-four chains reach the target person. Within this sample, the mean sum of intermediaries between starters and targets is 5.2. Boston starting chains reach the target

Travers and Milgram (1969) Sociometry

5/9



3/9

Why big data now?

arXiv.org > physics > arXiv:0803.0939

Physics > Physics and Society

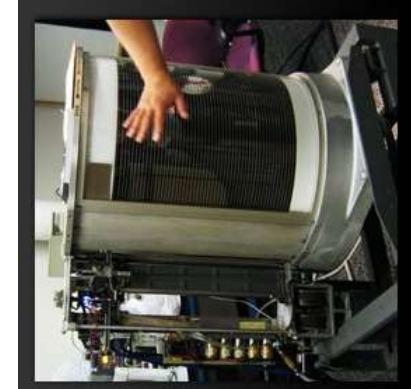
Planetary-Scale Views on an Instant-Messaging Network

Jure Leskovec, Eric Horvitz

(Submitted on 6 Mar 2008)

We present a study of anonymized data capturing a month of high-level communication activities within the whole of the Microsoft Messenger instant-messaging system. We examine characteristics and patterns that emerge from the collective dynamics of over 1 billion numbers of people, rather than the individual characteristics of individuals. The dataset contains summary properties of 30 billion conversations among 1240 million people. From the data, we construct a communication graph with 180 million nodes and 1.7 billion undirected edges, averaging over 3 edges per node. The social network constructed and analyzed to date. We report on multiple aspects of the dataset and synthesized graph. We find that the graph is well-connected and robust to node removal. We investigate on a planetary-scale the oft-cited "six degrees of separation" fact that people are separated by ``six degrees of separation'' and find that the average path length among Messenger users is 6.6. We also find that people tend to communicate more with each other when they have similar age, language, and location. Interestingly, cross-gender conversations are both more frequent and of longer duration than conversations with the same gender.

Leskovec and Horvitz WWW '08



Depends on your perspective

6/9

4/9

Big or small - you need the right data

Big or small - you need the right data

“ ...no matter how big the data are. ”

[Leek](#)



Chris Stucchio

Home Blog Code Work

Don't use Hadoop - your data isn't that big

Posted Mon, 16 Sept 2013
big data buzzwords, hadoop
Follow @stucchio | Tweet | 2/169
Like Share 1055 prepaid I bet this Sign Up to see what your friends like.
81 · 457 Recommended this on Google
http://www.chrissstucchio.com/blog/2013/hadoop_hatred.html

So, how much experience do you have with Big Data and "Hadoop"? They asked me, I told them that I use Hadoop all the time, but rarely for jobs larger than a few TB. I'm basically a big data neophyte - I know the concepts, I've written code, but never at scale. The next question they asked me: "Could you use Hadoop to do a simple group by and sum?" Of course I could, and just told them I needed to see an example of the file format.

http://www.chrissstucchio.com/blog/2013/hadoop_hatred.html

9/9

Big or small - you need the right data



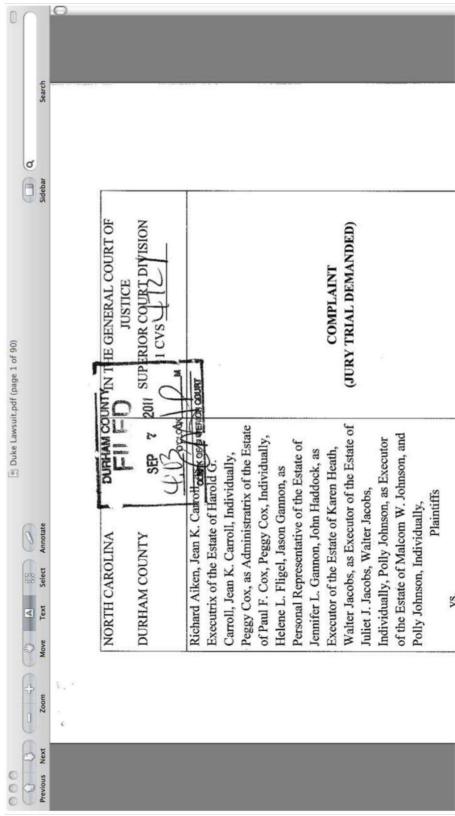
“ The data may not contain the answer. The combination of some data and an aching desire for an answer does not ensure that a reasonable answer can be extracted from a given body of data... ”

[Tukey](#)

Experimental design

Why you should care - serious trouble

Why you should care - an exciting result!



4/23

Know and care about the analysis plan!

Abstract

Background

Many groups, including our own, have proposed the use of DNA methylation profiles as biomarkers for various disease states. While much research has been done identifying DNA methylation signatures in cancer vs. normal etc., we still lack sufficient knowledge of the role that differential methylation plays during normal cellular differentiation and tissue specification. We also need thorough, genome level studies to determine the meaning of methylation of individual CpG dinucleotides in terms of gene expression.

Results

In this study, we have used ([insert statistical method here](#)) to compile unique DNA methylation signatures from normal human heart, lung, and kidney using the Illumina Infinium 27 K methylation array and compared those to gene expression by RNA sequencing. We have identified unique signatures of global DNA methylation for human heart, kidney and liver, and showed that DNA methylation data can be used to correctly classify various tissues. It indicates that DNA methylation reflects tissue specificity and may play an important role in tissue differentiation. The integrative analysis of methylation and RNA-Seq data showed that gene methylation and its transcriptional levels were comprehensively correlated. The location of methylation markers in terms of distance to transcription start site and CpG island showed no effects on the regulation of gene expression by DNA methylation in normal tissues.

Why you should care - uh oh!

Formula display: [MathJax](#)

DERIVING CHEMOSENSITIVITY FROM CELL LINES:
FORENSIC BIOINFORMATICS AND REPRODUCIBLE
RESEARCH IN HIGH-THROUGHPUT BIOLOGY

By KEITH A. BAGGERLY* AND KEVIN R. COOMBEST†

U.T. M.D. Anderson Cancer Center
High-throughput biological assays such as microarrays let us ask very detailed questions about how diseases operate, and promise to let us personalize therapy. Data processing, however, is often not described well enough to allow for exact reproduction of the results, leading to increases in "forensic bioinformatics," where aspects of raw data and reported results are used to infer what methods must have been employed. Unfortunately, poor documentation can shift from an inconvenience to an active driller when it obscures not just methods but errors. In this report, we examine several related papers purporting to use microarray-based signatures of drug sensitivity derived from cell lines to predict patient response. Patients in clinical trials are currently being allocated to treatment arms on the basis of these results. However, we show in five case studies that the results incorporate several simple errors that may be putting patients at risk. One theme that emerges is that the most common errors are simple (e.g., row or column offsets); conversely, it is our experience that the most simple errors are common. We then discuss steps we are taking to avoid such errors in our own investigations.

Annals of Applied Statistics

<http://nsaunder.wordpress.com/2012/07/23/we-really-dont-care-what-statistical-method-you-used/>

<http://arxiv.org/pdf/1010.1092.pdf>

2/23

5/23

3/23

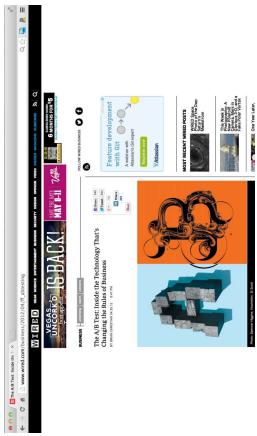
The screenshot shows a web page from the Nature journal website. At the top, there's a header with the journal name and some navigation links. Below the header, there's a sidebar with sections for 'ARTICLE LINKS' and 'ARTICLE TOOLS'. The main content area contains the abstract of a paper. The abstract discusses the development of gene expression signatures to predict drug sensitivity using Affymetrix microarray data. It highlights that these signatures can predict individual therapeutic drugs and potentially predict response to multidrug regimens. The paper is cited as 'Potti A, Dressman HK, Sayer R, et al. (2011) Using in vitro drug sensitivity data coupled with Affymetrix microarray data, we developed gene expression signatures that predict sensitivity to individual chemotherapeutic drugs. Each signature was validated with response data from an independent set of cell line studies. We further show that many of these signatures can accurately predict clinical response in individuals treated with these drugs. Notably, signatures developed to predict response to individual agents, when combined, could also predict response to multidrug regimens. Finally, we integrated the chemotherapy response signatures with signatures of oncogenic pathway deregulation to identify new therapeutic strategies that make use of all available drugs. The development of gene expression profiles that can predict response to'.

Below the abstract, there's a link to the full article: <http://www.nature.com/nm/journal/v12/n1/full/nm1491.html>.

2/23

Formulate your question in advance

Have a plan for data and code sharing



Question: Does changing the text on your website improve donations?

Experiment:

1. Randomly show visitors one version or the other
2. Measure how much they donate
3. Determine which is better

http://www.wired.com/business/2012/04/ff_abtesting

8/23



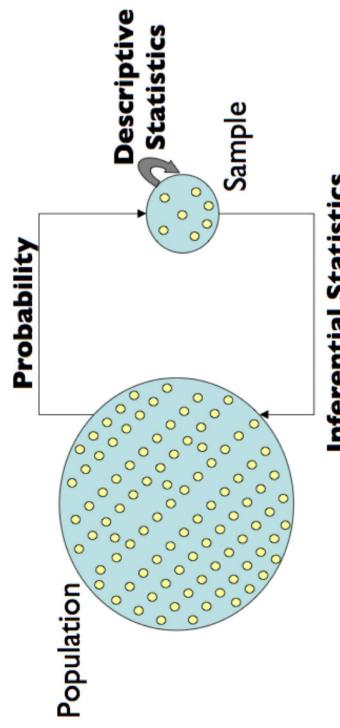
<https://github.com/>



<http://figshare.com/>

6/23

Statistical inference



May I recommend?

The Leek group guide to data sharing — Edit

A screenshot of a GitHub repository page for 'datasharing'. The repository has 25 commits, 1 branch, and 0 releases. It has 8 contributors. The repository page shows several files: 'Merge pull request #9 from nikai3d/patch-1', 'jleek authored 6 days ago', 'README.md', and 'README.md'. The commit history shows a merge pull request and a fix for a typo.

How to share data with a statistician

This is a guide for anyone who needs to share data with a statistician. The target audiences I have in mind are:

- Scientific collaborators who need statisticians to analyze data for them
- Students or postdocs in scientific disciplines looking for consulting advice
- Junior statistics students whose job it is to collate/clean data sets

<http://www.gs.washington.edu/academics/courses/akey/56008/lecture/lecture2.pdf>

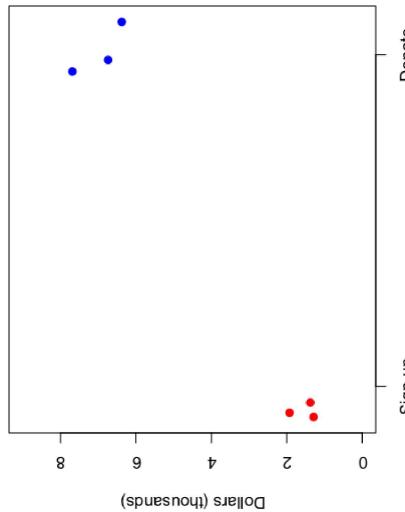
<https://github.com/jtleek/datasharing>

9/23

7/23

Variability - Scenario 3

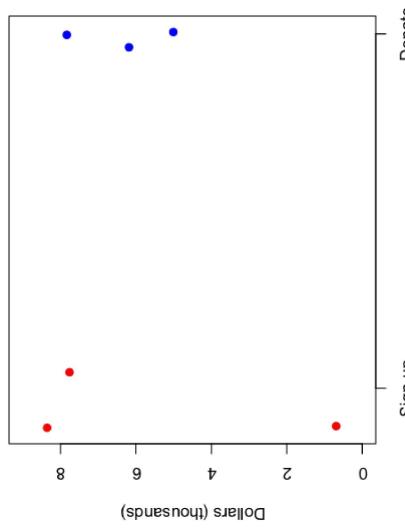
Variability - Scenario 1



12/23

10/23

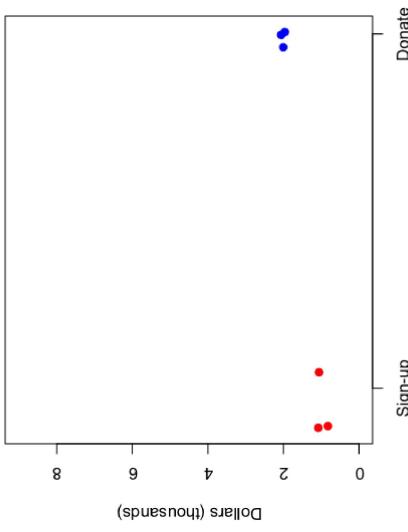
Confounding



Webpage version

Sign-up
Donate

Variability - Scenario 2



Webpage version

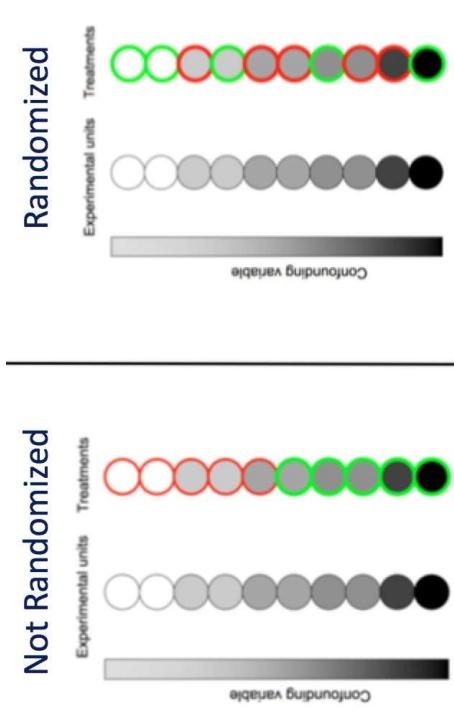
Sign-up
Donate

13/23

11/23

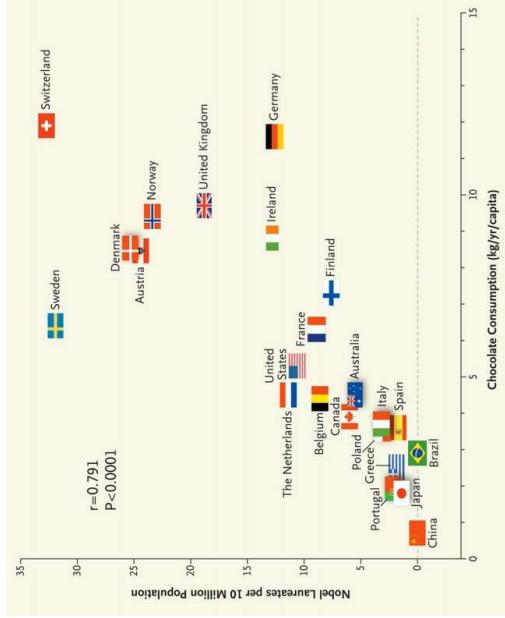
Why does randomization help?

Correlation is not causation*



<http://www.gs.washington.edu/academics/courses/akey/56008/lecture/lecture1.pdf>

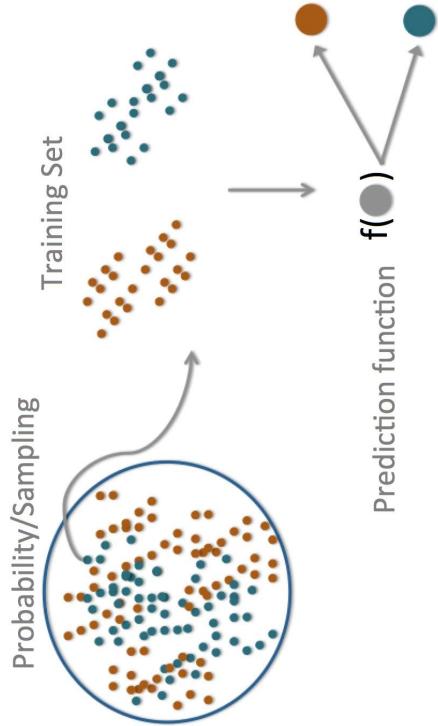
16/23



<http://www.nejm.org/doi/full/10.1056/NEJMMon1211064>

14/23

Prediction



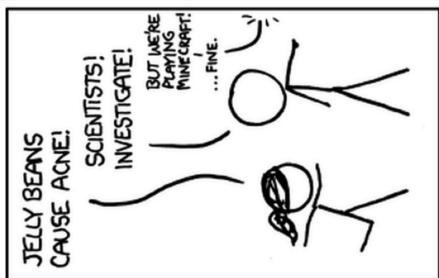
Randomization and blocking

- If you can (and want to) fix a variable
 - Website always says Obama 2014 on it
- If you don't fix a variable, stratify it
 - If you are testing sign up phrases and have two website colors, use both phrases equally on both.
- If you can't fix a variable, randomize it

17/23

15/23

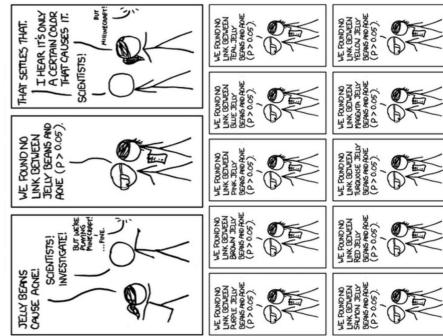
Beware data dredging



<http://xkcd.com/882/>

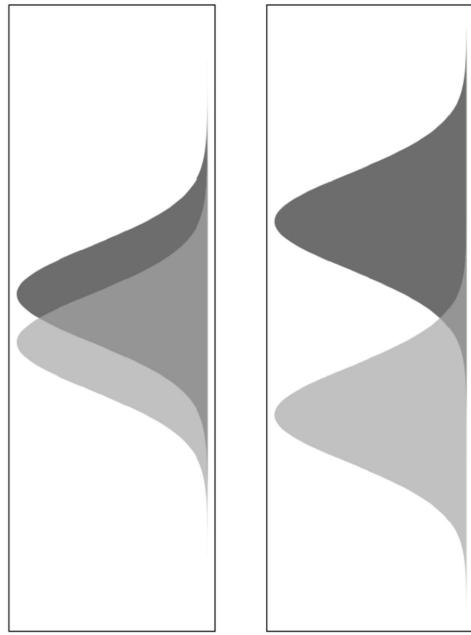
20/23

Beware data dredging



<http://xkcd.com/882/>

Prediction versus inference



18/23

Prediction key quantities

		DISEASE	
		+	-
TEST	+	TP	FP
	-	FN	TN

- Sensitivity → $\Pr(\text{positive test} \mid \text{disease})$
- Specificity → $\Pr(\text{negative test} \mid \text{no disease})$
- Positive Predictive Value → $\Pr(\text{disease} \mid \text{positive test})$
- Negative Predictive Value → $\Pr(\text{no disease} \mid \text{negative test})$
- Accuracy → $\Pr(\text{correct outcome})$

<http://www.biostat.jhsph.edu/~iruczins/teaching/140.615/>

21/23

19/23

Beware data dredging



<http://xkcd.com/882/>

22/23

Summary

- Good experiments
 - Have replication
 - Measure variability
 - Generalize to the problem you care about
 - Are transparent
- Prediction is not inference
 - Both can be important
- Beware data dredging

23/23