



Experimental design

Jeffrey Leek
Johns Hopkins Bloomberg School of Public Health

Why you should care - an exciting result!

Genomic signatures to guide the use of chemotherapeutics






Anil Potti^{1,2}, Holly K Dressman^{1,3}, Andrea Bild^{1,3}, Richard F Riedel^{1,2}, Gina Chan⁴, Robyn Sayer⁴, Janiel Cragun⁴, Hope Cottrill⁴, Michael J Kelley², Rebecca Petersen⁵, David Harpole⁵, Jeffrey Marks⁵, Andrew Berchuck^{1,6}, Geoffrey S Ginsburg^{1,2}, Phillip Febbo^{1,2,3}, Johnathan Lancaster⁴ & Joseph R Nevins^{1,2,3}

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

ARTICLE LINKS

- Supplementary info

ARTICLE TOOLS

-  Send to a friend
-  Export citation
-  Export references
-  Rights and permissions
-  Order commercial reprints

SEARCH PUBMED FOR

- Anil Potti
- Holly K Dressman
- Andrea Bild
- Richard F Riedel
- Gina Chan
- Robyn Sayer

<http://www.nature.com/nm/journal/v12/n11/full/nm1491.html>

Why you should care - uh oh!

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

BY KEITH A. BAGGERLY* AND KEVIN R. COOMBES†

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 exercises 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 danger 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://arxiv.org/pdf/1010.1092.pdf>

Why you should care - serious trouble

Duke Lawsuit.pdf (page 1 of 90)

Previous Next Zoom Move Text Select Annotate Sidebar Search

NORTH CAROLINA DURHAM COUNTY	DURHAM COUNTY FILED SEP 7 2011 403 M	IN THE GENERAL COURT OF JUSTICE SUPERIOR COURT DIVISION 1 CVS 4121
Richard Aiken, Jean K. Carroll, Executrix of the Estate of Harold G. Carroll, Jean K. Carroll, Individually, Peggy Cox, as Administratrix of the Estate of Paul F. Cox, Peggy Cox, Individually, Helene L. Fligel, Jason Gannon, as Personal Representative of the Estate of Jennifer L. Gannon, John Haddock, as Executor of the Estate of Karen Heath, Walter Jacobs, as Executor of the Estate of Juliet J. Jacobs, Walter Jacobs, Individually, Polly Johnson, as Executor of the Estate of Malcom W. Johnson, and Polly Johnson, Individually, Plaintiffs		COMPLAINT (JURY TRIAL DEMANDED)
vs.		

Know and care about the analysis plan!

Abstract

Formula display: ☒ MathJax ?

Background

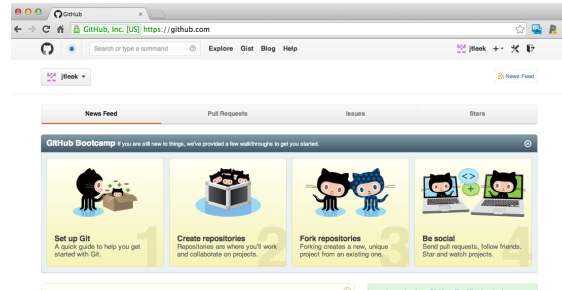
Many groups, including our own, have proposed the use of DNA methylation profiles as biomarkers for various disease states. While much research has been done identifying DNA methylation signatures in cancer vs. normal etc., we still lack sufficient knowledge of the role that differential methylation 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

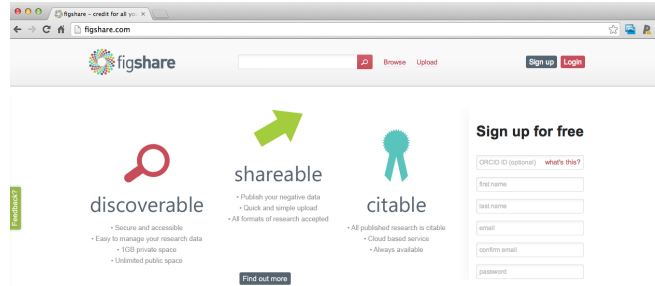
In this study, we have used (insert statistical method here) to compile unique DNA methylation signatures from normal human heart, lung, and kidney using the Illumina Infinium 27 K methylation arrays and compared those to gene expression by RNA sequencing. We have identified unique signatures of global 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.

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

Have a plan for data and code sharing



<https://github.com/>



<http://figshare.com/>

May I recommend?

The Leek group guide to data sharing — Edit

25 commits 1 branch 0 releases 8 contributors

branch: master datasharing / +

Merge pull request #9 from nikai3d/patch-1

jtleek authored 6 days ago latest commit e53857faa4

README.md fix typo 6 days ago

README.md

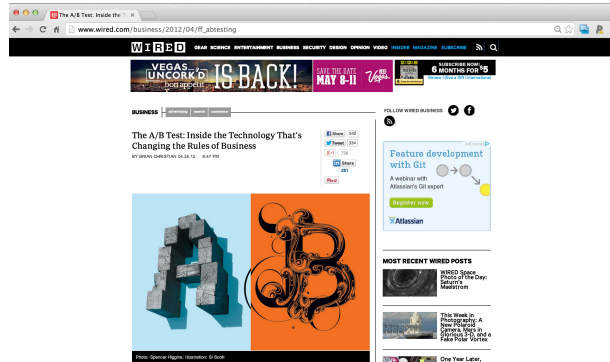
How to share data with a statistician

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

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

Formulate your question in advance



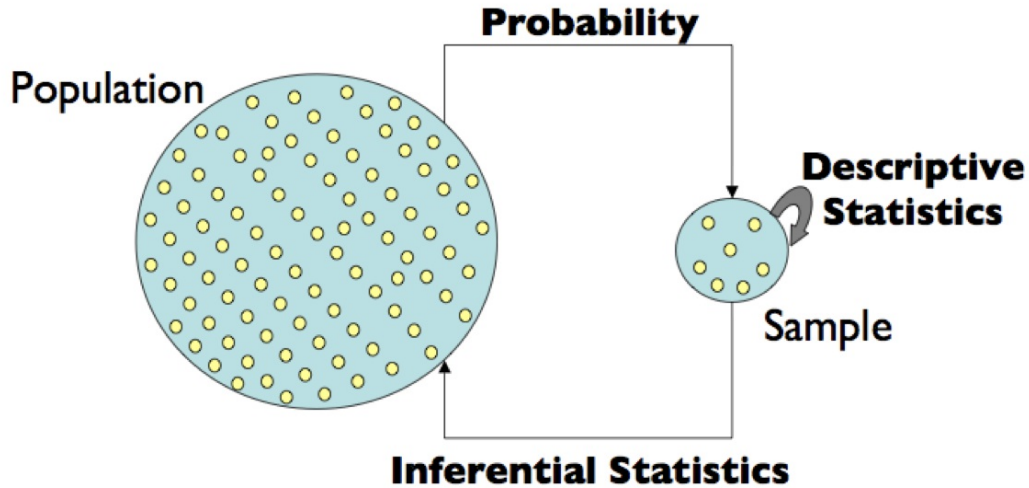
Question: Does changing the text on your website improve donations?

Experiment:

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

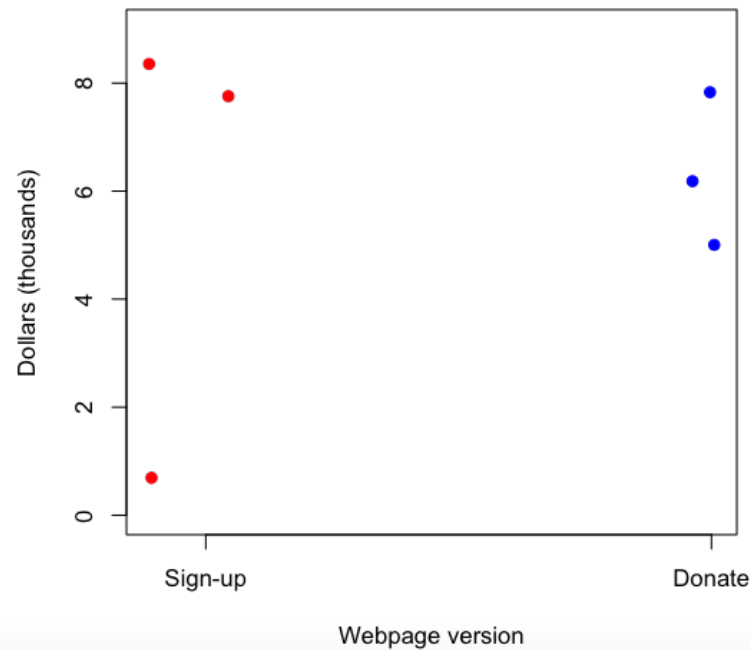
http://www.wired.com/business/2012/04/ff_abtesting

Statistical inference

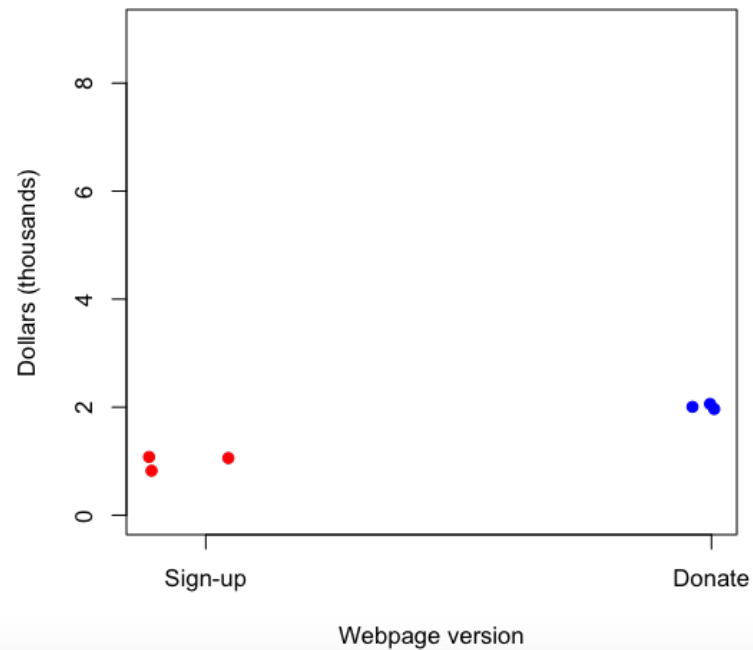


<http://www.gs.washington.edu/academics/courses/akey/56008/lecture/lecture2.pdf>

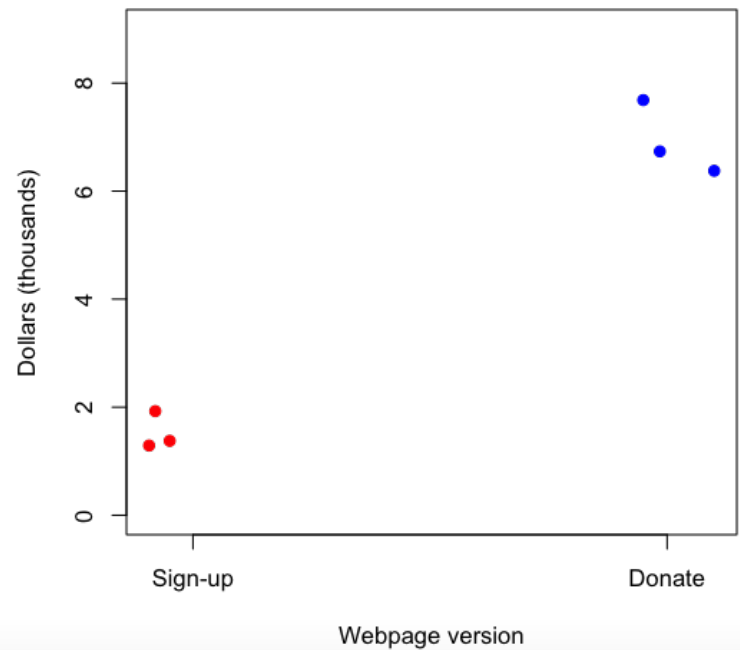
Variability - Scenario 1



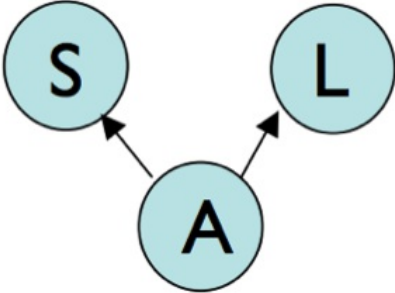
Variability - Scenario 2



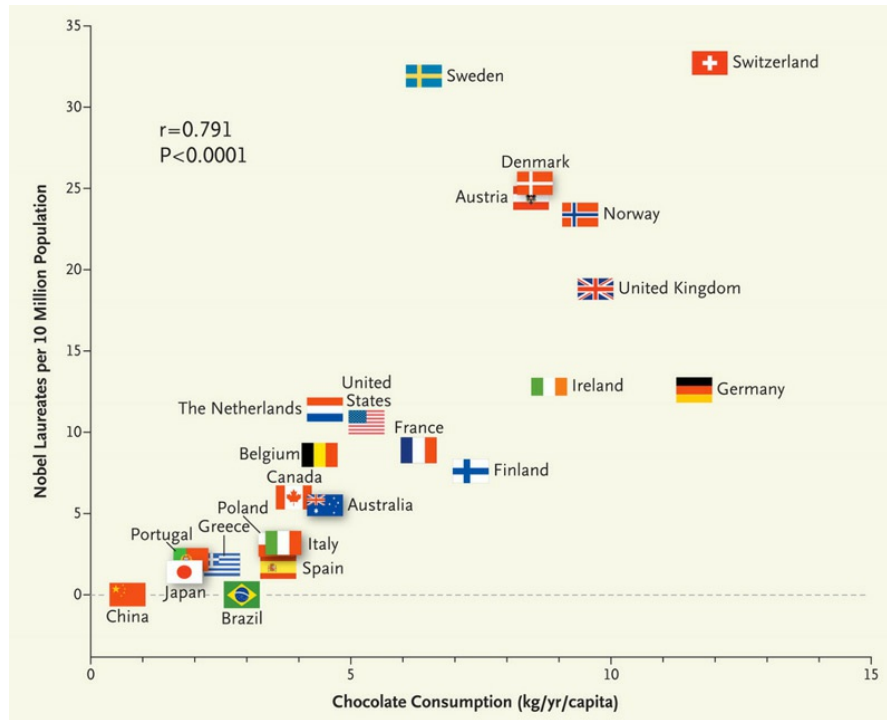
Variability - Scenario 3



Confounding



Correlation is not causation*



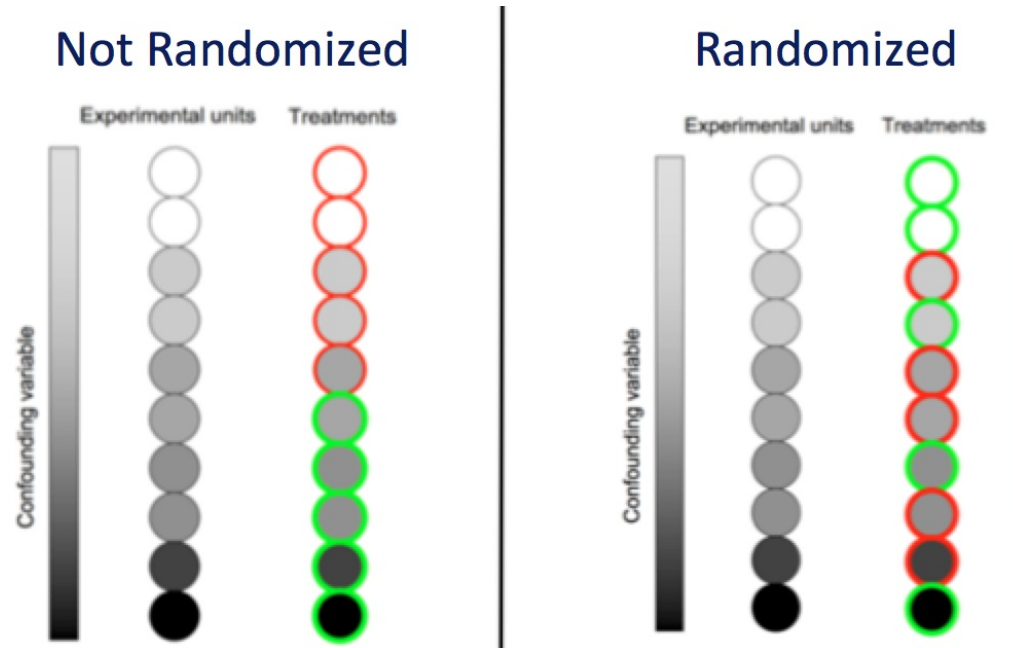
<http://www.nejm.org/doi/full/10.1056/NEJMon1211064>

*Sometimes called spurious correlation**

Randomization and blocking

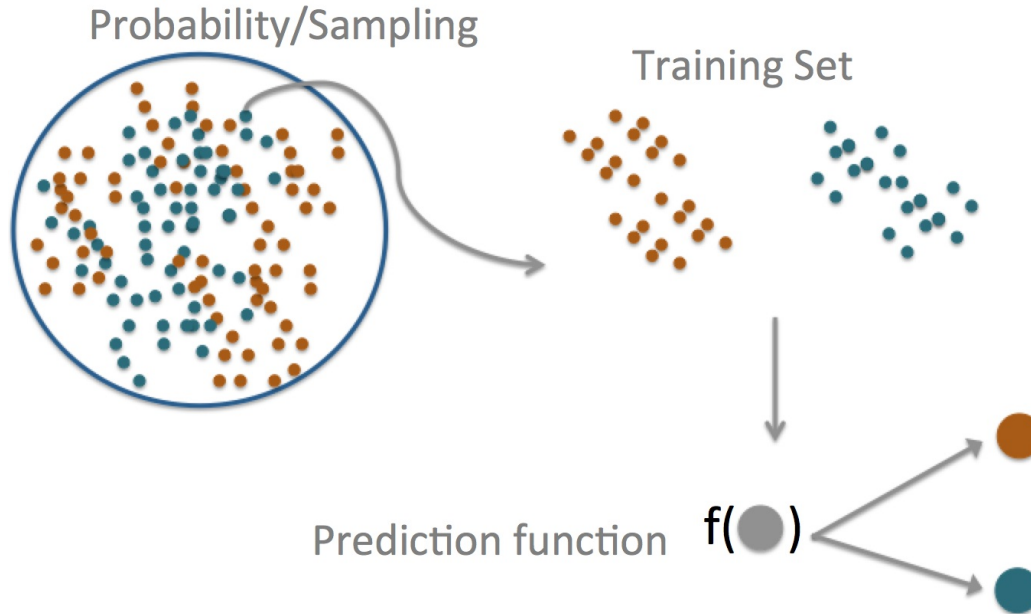
- If you can (and want to) fix a variable
 - Website always says Obama 2014 on it
- If you don't fix a variable, stratify it
 - If you are testing sign up phrases and have two website colors, use both phrases equally on both.
- If you can't fix a variable, randomize it

Why does randomization help?

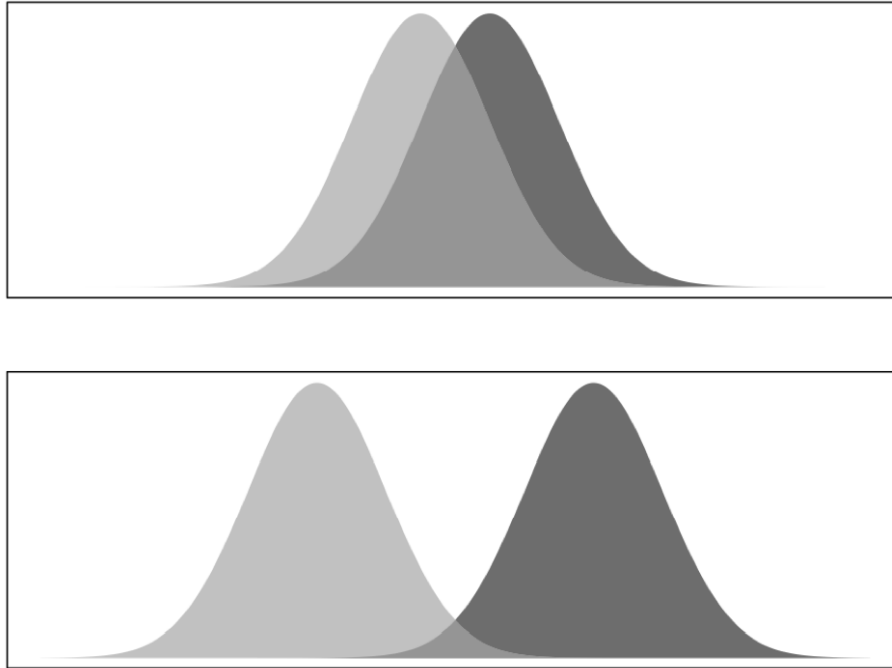


<http://www.gs.washington.edu/academics/courses/akey/56008/lecture/lecture1.pdf>

Prediction



Prediction versus inference



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

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

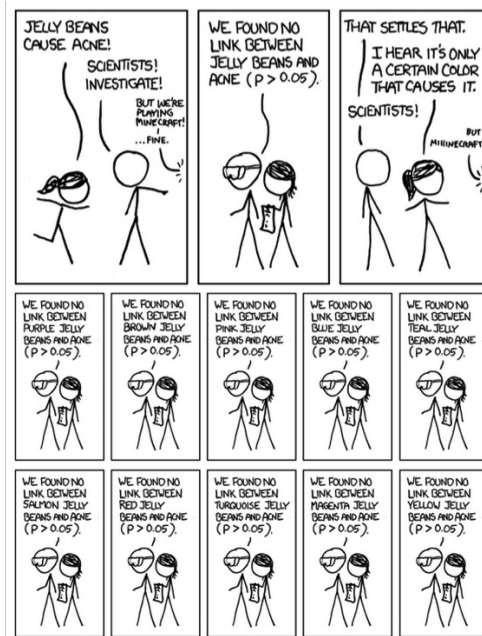
→ $\Pr(\text{correct outcome})$

Beware data dredging



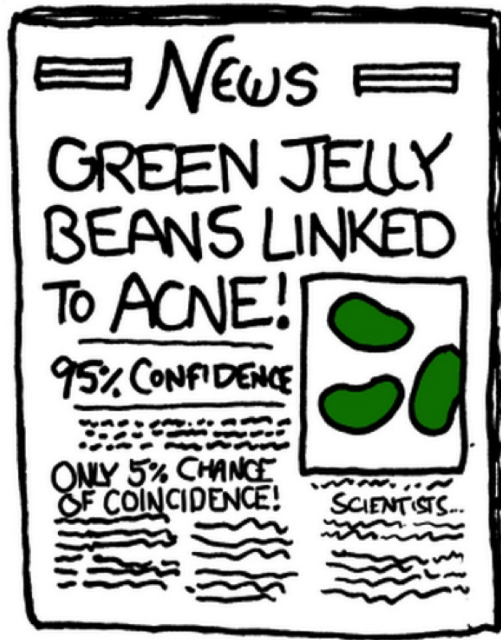
<http://xkcd.com/882/>

Beware data dredging



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

Beware data dredging



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

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