



INSTITUTE FOR DEFENSE ANALYSES

## Thinking About Data for Operational Test and Evaluation

Dean Thomas, *Project Leader*

Matthew Avery

November 2017

Approved for public release;  
distribution is unlimited.

IDA Non-Standard Document  
NS D-8729

Log: H 2017-000528



*The Institute for Defense Analyses is a non-profit corporation that operates three federally funded research and development centers to provide objective analyses of national security issues, particularly those requiring scientific and technical expertise, and conduct related research on other national challenges.*

#### About This Publication

While the human brain is a powerful tool for quickly recognizing patterns in data, it will frequently make errors in interpreting random data. Luckily, these mistakes occur in systematic and predictable ways. Statistical models provide an analytical framework that helps us avoid these error-prone heuristics and draw accurate conclusions from random data. This non-technical presentation highlights some tricks of the trade learned by studying data and the way the human brain processes. First, we introduce statistics as the science of data, and discuss how the popular conception of randomness differs from its technical definition. Later sections highlight the human brain as a pattern recognition machine. Examples from published literature and media highlight systematic and predictable errors in human cognition as well as how poor data analysis and graphical displays can cause critical errors in analysis. Finally, we'll talk about using statistical models for analysis, including how violations of model assumptions should affect our analyses.

#### For more information:

Dean Thomas, Project Leader  
[DThomas@ida.org](mailto:DThomas@ida.org) • (703) 845-6986

Robert R. Soule, Director, Operational Evaluation Division  
[rsoule@ida.org](mailto:rsoule@ida.org) • (703) 845-2482

#### Copyright Notice

© 2017 Institute for Defense Analyses  
4850 Mark Center Drive, Alexandria, Virginia 22311-1882 • (703) 845-2000

This material may be reproduced by or for the U.S. Government pursuant to the copyright license under the clause at DFARS 252.227-7013 (a)(16) [Jun 2013].

INSTITUTE FOR DEFENSE ANALYSES

IDA Non-Standard Document NS D-8729

## **Thinking About Data for Operational Test and Evaluation**

Dean Thomas, *Project Leader*

Matthew Avery



# Thinking About Data for Operational Test and Evaluation

---

## A. Introduction

The human brain is a powerful tool for quickly recognizing patterns, but it frequently makes errors in interpreting random data. Luckily, these mistakes occur in systematic and predictable ways. Statistical models provide an analytical framework that helps us avoid these error-prone heuristics and draw accurate conclusions from random data.

## B. Randomness

Statistics is the science of data in the same way that biology is the science of life and physics is the science of energy. This presentation highlights some of the tricks of the trade that studying and working with data for many years teaches you. One of the biggest differences in the way statisticians think about data is how they conceive of randomness. The technical definition varies in important and subtle ways from the colloquial way that randomness is thought of. These differences help explain why outcomes that are deterministic at a fundamental level may be profitably modeled as random. To that

end, understanding the strengths and weaknesses of statistical models is important, particularly how violations of their assumptions impact inference.

## C. Cognition & Data

Humans' built-in heuristics for processing data compel us to make errors in systematic and predictable ways. Understanding how these heuristics affect our information processing allows us to adjust and correct for known errors. The human brain is especially good at recognizing patterns. However, we frequently make Type I errors, mistaking randomness and noise for actual patterns. One example of this is the stock market. People regularly over-interpret market movements, attributing causal relationships to what later turns out to be nothing more than noise. "Charters" assure you that they can find specific patterns in stock charts when really they are just looking at noise.

Daniel Khaneman and Amos Tversky have demonstrated many ways in which humans err when considering randomness. For example, humans make systematic errors when attempting

to determine what a random sequence looks like, believing that chance is essentially fair and that things tend to balance out over time. We underestimate the likelihood of events that don't "look" random, such as a long sequence of Heads occurring over a large number of coin tosses. Prospect Theory explains the ways humans react to uncertainty, systematically overestimating the likelihood of unlikely events and underestimating the likelihood of likely events.

## D. Better Ways to Think About Data and Uncertainty

Even quantitatively literate members of the media and academia can fall victim to these pitfalls if they are not careful. Binned data (For example, Range to Target can be organized into bins of "Long Range" (20 to 30km), "Medium Range" (10 to 20 km) and "Short Range (0 to 10 km) for convenience in test planning) are frequently seen in operational tests, but when treating a continuous variable as an ordered categorical variable, analyzing binned data, it is important to understand how the data are distributed within the synthetic bins. This underlying distribution can have a substantial effect on the conclusions of an analysis. Anne Case and Angus Deaton provide a notable example in which failure to account for the within-age-bin distributions led to mis-estimation of mortality rates among middle-aged white Americans. Similarly, the way in which data are displayed can be misleading if care is not taken. Using

multiple Y-axes on a single plot can cause readers to misunderstand the magnitude of relative effects, and plotting a rate of change rather than the variable of interest can make mountains out of molehills.

## E. Statistical Models

Analysts can avoid many of these challenges by using an analytical framework that incorporates randomness, such as a statistical model. Statistical models rely on assumptions, but most popular models, such as simple linear regression, are robust to some violation of these assumptions. For example, regression estimates for the conditional mean are entirely robust to the assumption of equal variance. While there are no good hard-and-fast rules for how much a model's assumptions can be violated before the model should be discarded, knowing the right questions to ask about your data helps you know when you are okay with a small caveat and when you should start from scratch. Knowing your models thoroughly also helps you know how much confidence you should have in your results. Finally, it is important to understand the process by which data is collected. Often times, this process can bias subsequent analyses in unexpected ways if care is not taken.

## F. Conclusions

Thinking about data requires an understanding of the tricks the human brain can play as well as ways to avoid those pitfalls.

The origins of your data can impact the conclusions you draw. Although the human brain recognizes patterns well, know that it will often identify “patterns” that aren’t really there. These problems can be exacerbated if data is displayed carelessly in ways that misconstrue the actual underlying structure. Luckily, statistical models give us a framework to understand randomness and patterns, helping us learn about populations of interest. Understanding data and randomness is as challenging as any other science, so don’t be afraid to consult an expert!





# Thinking About Data for Operational Test and Evaluation

Perspectives from statistics, behavioral  
economics, and cognitive psychology

11 December 2017



# Thinking about data



The human brain misunderstands randomness in **systematic** and **predictable** ways. We often have to fight our instincts to not be misled by data.

Statistical models provide a **framework** that helps us **avoid** these **error-prone heuristics** and draw accurate conclusions from random data.

# Randomness



**“Randomness” is a technical term!**



# Per Wikipedia's entry on statistical randomness:

"A numeric sequence is said to be **statistically random** when it contains no recognizable **patterns** or **regularities**; sequences such as the results of an ideal dice roll, or the digits of  $\pi$  exhibit statistical randomness.

"Statistical randomness does not necessarily imply "true" randomness, i.e., **objective unpredictability**."

[https://en.wikipedia.org/wiki/Statistical\\_randomness](https://en.wikipedia.org/wiki/Statistical_randomness)

Randomness ≠ unpredictable

**Formally, randomness is thought of in terms of sequential of outcomes.**

**Note that the frequency of events may have a pattern  
(defined by the distribution function of a random  
variable).**

If it walks like a duck, and talks like a duck ...

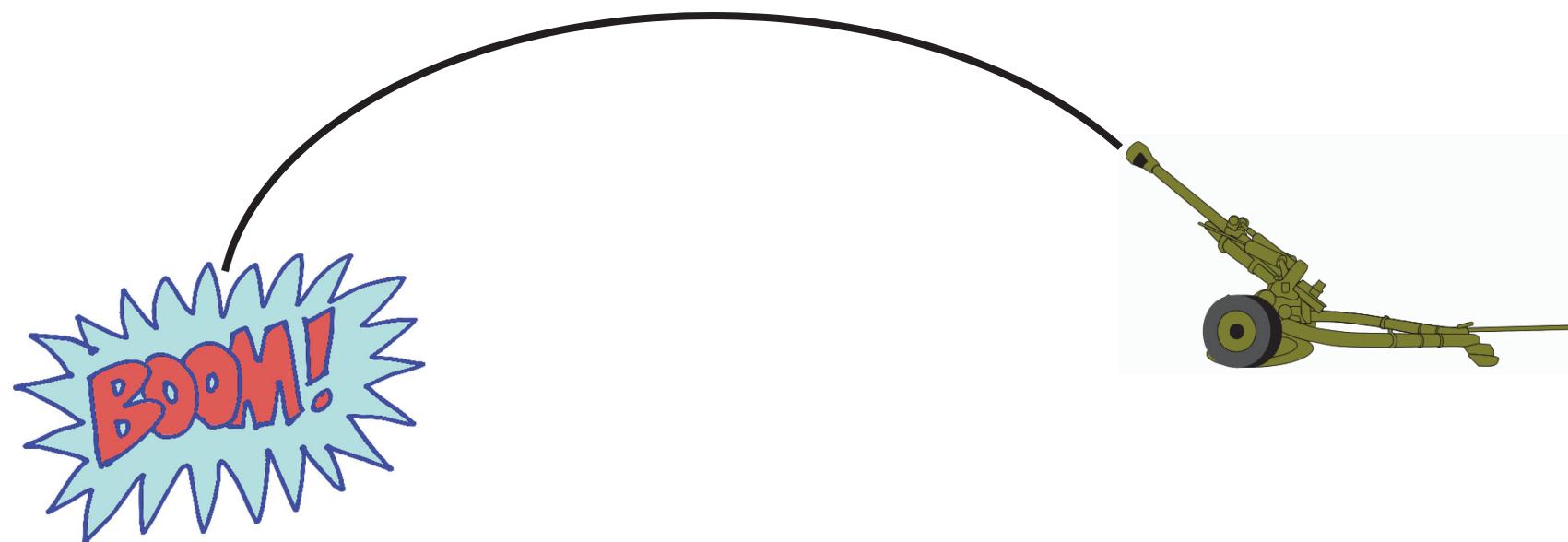


**... then a model built around the assumption that it is a duck may prove useful!**

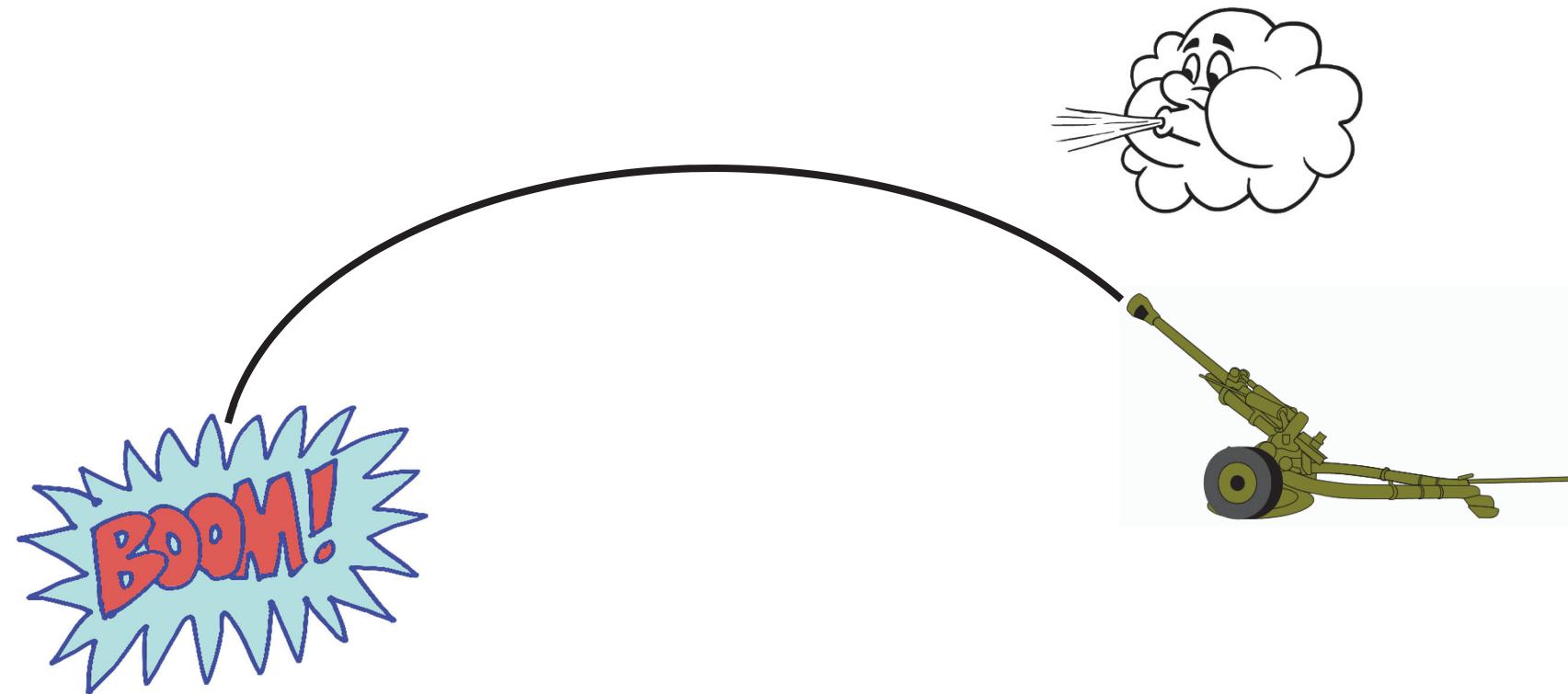
**Things that are deterministic at a fundamental level can often be modeled effectively as random.**

Example: How accurate is the newest Howitzer?

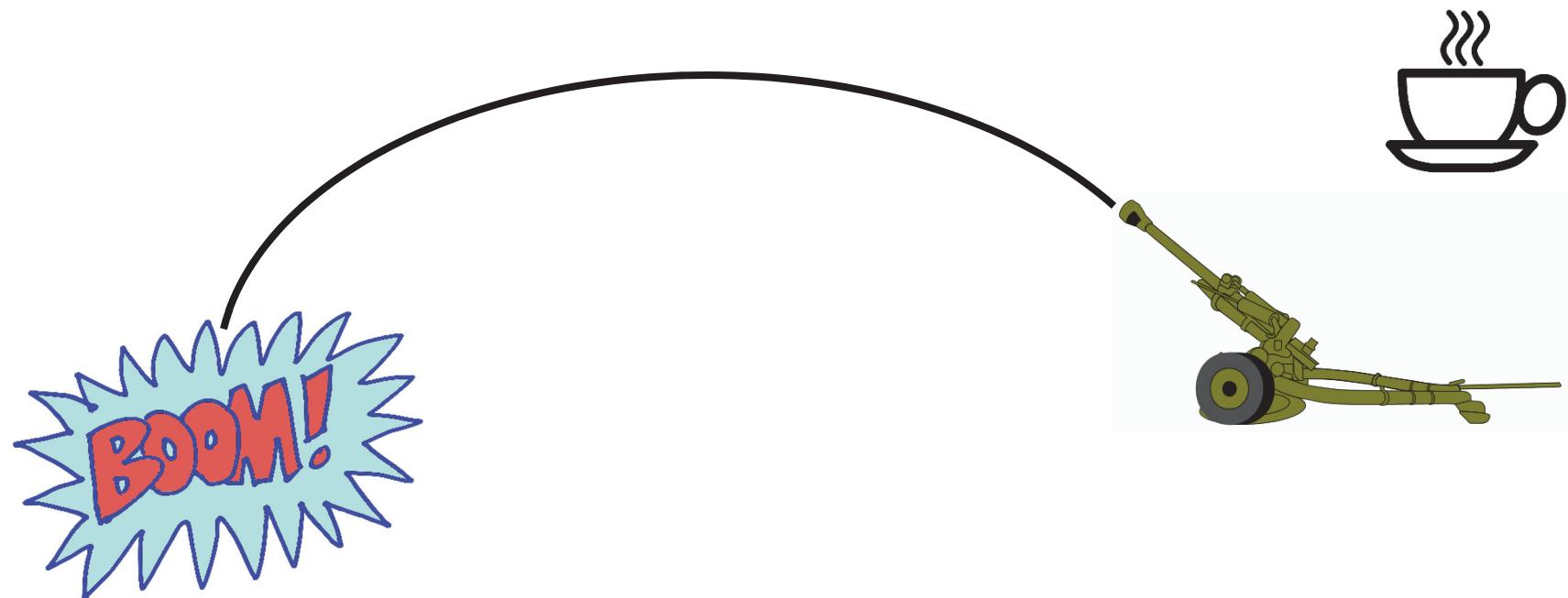
**Lots of things impact artillery accuracy beyond the factors we typically control:**



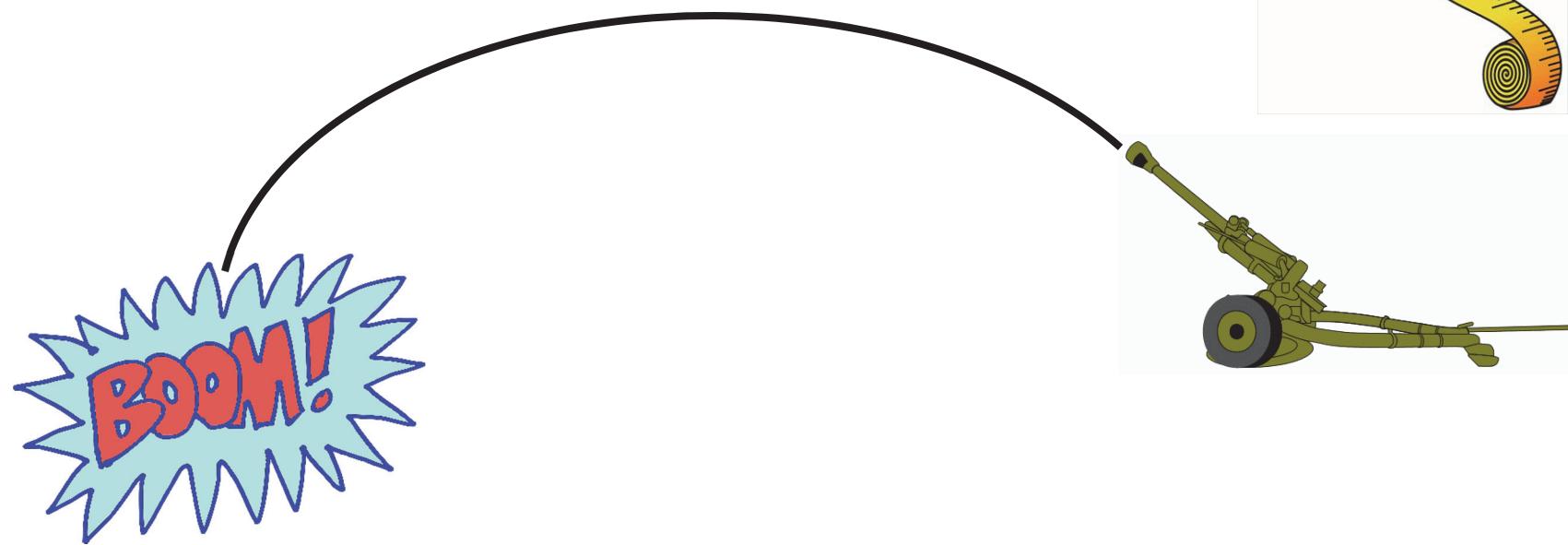
# Wind.



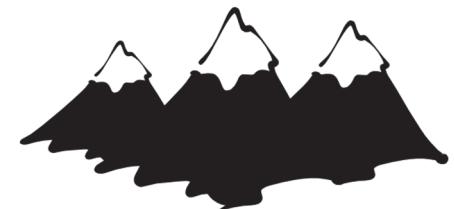
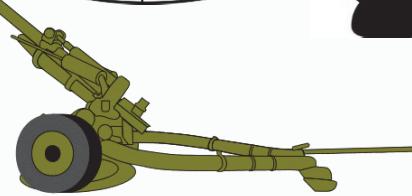
# How rested the operators are.



# The calibration of the gun.



# And other factors.



**If a model helps you understand the process that generated your data, then it's useful even if the process that generated it wasn't “random”.**

Think about where your data come from and how  
that impacts the conclusions you can draw.

# Cognition & Data



# Small Sample Size Theater





**Kevin Pelton**

@kpelton

Follow

The Warriors broadcast recently posted Steph Curry's best 3-point percentage by opponent, and I noticed four of the top five are East teams.

Opponent	G	3PA	3P%
Miami	12	101	51.5%
Charlotte	14	129	51.2%
New York	14	122	50.8%
Portland	26	225	48.9%
Washington	16	135	48.9%

1:59 PM - 6 Nov 2017

28 Retweets 99 Likes



**Kevin Pelton** @kpelton · 15h

Replying to @kpelton

If you just looked at that list, you might think Steph Curry shoots better against East opponents because they're weaker teams.



1



3



17



Kevin Pelton

@kpelton

Follow

But a funny thing happens when you look at Steph's worst 3-point percentage by opponent: 4 of the top 5 are East teams too:

Opponent	G	3PA	3P%
Chicago	15	95	31.6%
Philadelphia	15	122	32.8%
Atlanta	14	96	33.3%
LA Lakers	26	219	37.9%
Boston	14	98	38.8%

2:00 PM - 6 Nov 2017

5 Retweets 30 Likes



8

5

30



**Kevin Pelton**

@kpelton

Follow

With that information, the explanation becomes more obvious: Curry has more extreme 3P% against East opponents b/c he plays them less often.

2:01 PM - 6 Nov 2017

8 Retweets 54 Likes



2

8

54



**Kevin Pelton** @kpelton · 15h

Replies to @kpelton

This is a funny, trivial thing when we're talking about Steph Curry's 3-point percentage but it has important real-world implications.

1

3

20



**Kevin Pelton** @kpelton · 15h

If you look only at the best or worst performers in anything, those with smaller sample sizes will likely be overrepresented.

2

12

81



**Kevin Pelton** @kpelton · 15h

Howard Wainer of the Wharton School has written about the many mistakes created by ignorance of this fact (PDF): whr.tn/2Anw0Wp

1

12

121



**Kevin Pelton** @kpelton · 14h

BTW replies to this thread show the other issue: it's easy to construct post-hoc explanations for false "trends" that are simply randomness.

1

4

35



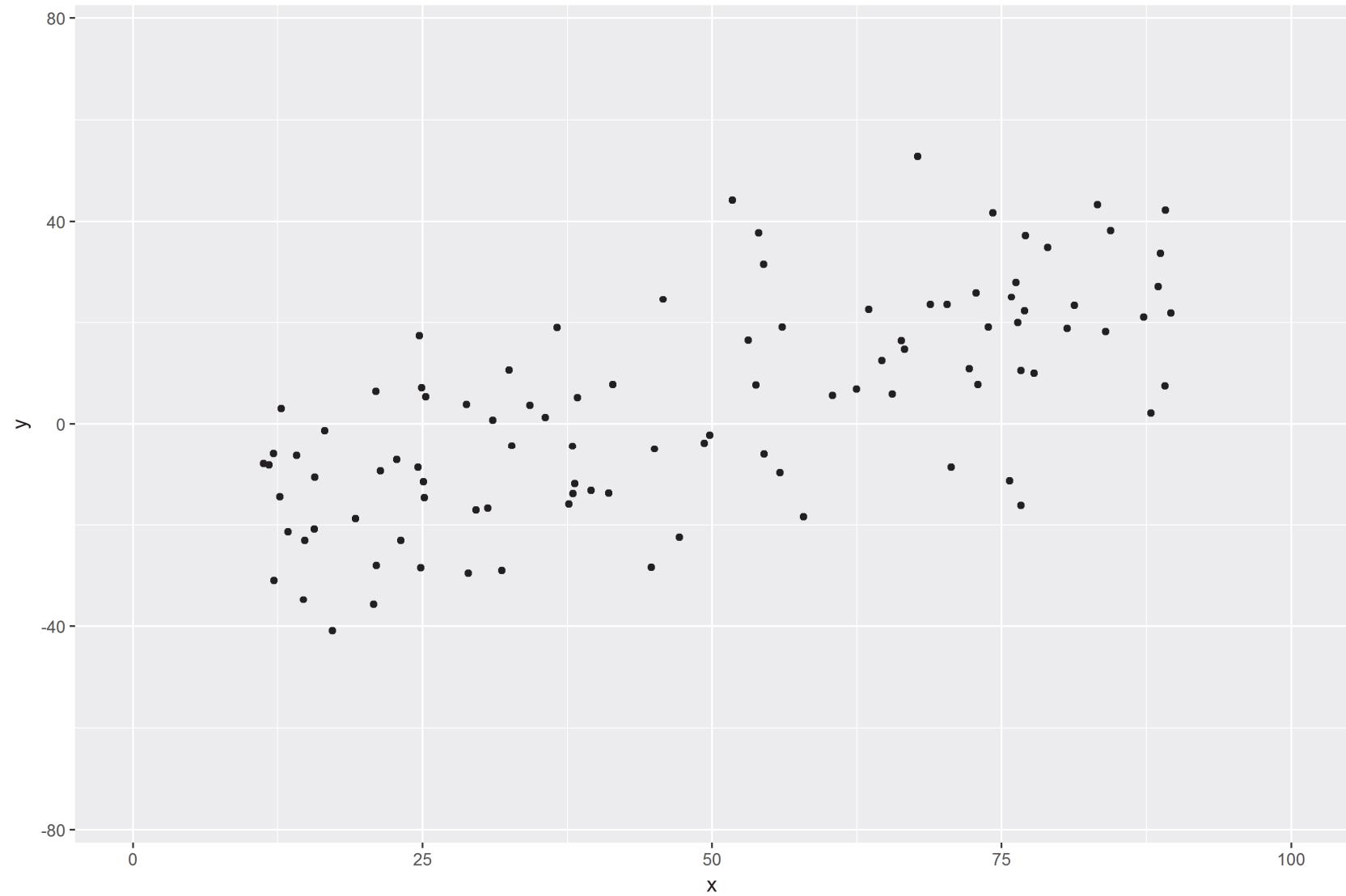
**Baby, we were born this way.**



**Humans have a lot of built-in heuristics for processing large amounts of information that make it possible for us to cope with the natural world.**

They're often **wrong** in systematic and **predictable** ways.

# The human brain is a giant pattern recognition machine.



# This bread looks like a rabbit.



# People also see patterns when they aren't there, like the Virgin Mary in a piece of bread.

BBC NEWS

Watch One-Minute World News

News services Your news when you want it

News Front Page

Last Updated: Tuesday, 23 November, 2004, 11:54 GMT

E-mail this to a friend | Printable version

## 'Virgin Mary' toast fetches \$28,000

A decade-old toasted cheese sandwich said to bear an image of the Virgin Mary has sold on the eBay auction website for \$28,000.

An internet casino confirmed it had purchased the sandwich, saying it had become a "part of pop culture".

Goldenpalace.com says it will take the sandwich on world tour before selling it and donating the money to charity.

Diane Duyser, from Florida, says the sandwich has never gone mouldy since she made it 10 years ago.

By the time the sandwich auction closed on Monday the sale had received over 1.7 million hits on the auction site.

'Mystical power'



The toast is not intended for consumption

SEE ALSO:

- Woman 'blessed by the holy t 17 Nov 04 | Americas
- Ghanaians flock to see 'mirac 02 Nov 04 | Africa
- Virgin Mary 'seen in US hospit 17 Jun 03 | Americas
- Bangladeshis flock to 'weeping Virgin' 18 Feb 03 | South Asia

RELATED INTERNET LINKS:

- eBay
- GoldenPalace.com

The BBC is not responsible for content of external internet site

TOP AMERICAS STORIES

- US lifts lid on WikiLeaks probe
- Iran scientist heads home
- Argentina legalises gay marri

| News feeds

**Random doesn't always look random.**



**Humans are programmed to find patterns, and so they do, even when patterns aren't there.**

Since the 1930s (and probably before), people have been drawing pictures of noise and over-interpreting them.



From [www.stockcharts.com](http://www.stockcharts.com) (which is probably a good place to get advice on how to lose money).

**Sometimes, “random” doesn’t look random enough!**



"People view chance as unpredictable but essentially fair."

Kahneman & Tversky, 1972

When people are asked to generate their own “random” sequences of “coin flips”, the results are **locally representative** and don’t have enough **short “runs”**.

# People also don't like long runs.

For example, a run of six consecutive H or T will occur in over half of all sequences of 50 flips of a fair coin, but people will rarely include such a run when they generate a sequence.

“A major characteristic of apparent randomness  
is the absence of systematic patterns.”

Sequence	1	2	3	4	5	6	7	8
A	H	H	T	T	H	H	T	T
B	H	T	H	T	H	T	H	T
C	H	H	T	H	T	T	T	H
D	H	T	H	H	H	T	H	H

## The “Law” of Small Numbers:

Even numerate scientists act as if any sample, no matter how small, is representative of the population.

Note: This is a **fallacy**.

**On each round of a game, 20 marbles are distributed at random among five children. Consider the following distributions:**

Child	1	2
Alan	4	4
Ben	4	4
Carl	5	4
Dan	4	4
Ed	3	4

In many rounds of the game, will there be more results of Type 1 or Type 2?

**Both Type 1 and Type 2 are very unlikely, but Type 2 is more likely (25% more likely, in fact!).**

$$P(\text{Type 1}) = 0.00256$$

$$P(\text{Type 2}) = 0.00320$$

*In practice, 36 of 52 subjects said  
Type 1 was more likely.*

**My hypothesis: People are answering the question they want to answer rather than the one that was asked.**

$$P(\text{Type 1}) = 0.00256$$

$$P(\text{Type 2}) = 0.00320$$

$$P(\text{Type 1 or something that looks like Type 1}) = 0.05120$$

**People are bad enough at self-generating random numbers that forensic accountants can often catch people “cooking the books” because the fake numbers don’t have the same distribution as real numbers.**

### Benford's Law

Probabilities					
Digit	1st Digit	2nd Digit	3rd Digit	4th Digit	5th Digit
0	NA	0.11968	0.10178	0.10018	0.10002
1	0.30103	0.11389	0.10138	0.10014	0.10001
2	0.17609	0.10882	0.10097	0.10010	0.10001
3	0.12494	0.10433	0.10057	0.10006	0.10001
4	0.09691	0.10031	0.10018	0.10002	0.10000
5	0.07918	0.09668	0.09979	0.09998	0.10000
6	0.06695	0.09337	0.09940	0.09994	0.09999
7	0.05799	0.09035	0.09902	0.09990	0.09999
8	0.05115	0.08757	0.09864	0.09986	0.09999
9	0.04576	0.08500	0.09827	0.09982	0.09998

# To be fair to people, though, differentiating true randomness is pretty tough.

Though to be fair to humans, randomness is hard to see:  
Williams, J. J., & Griffiths, T. L. (2008). "Why are people bad at detecting randomness? Because it is hard."

**So you're telling me there's a chance!**



**Long-tail events are rare, not impossible.**

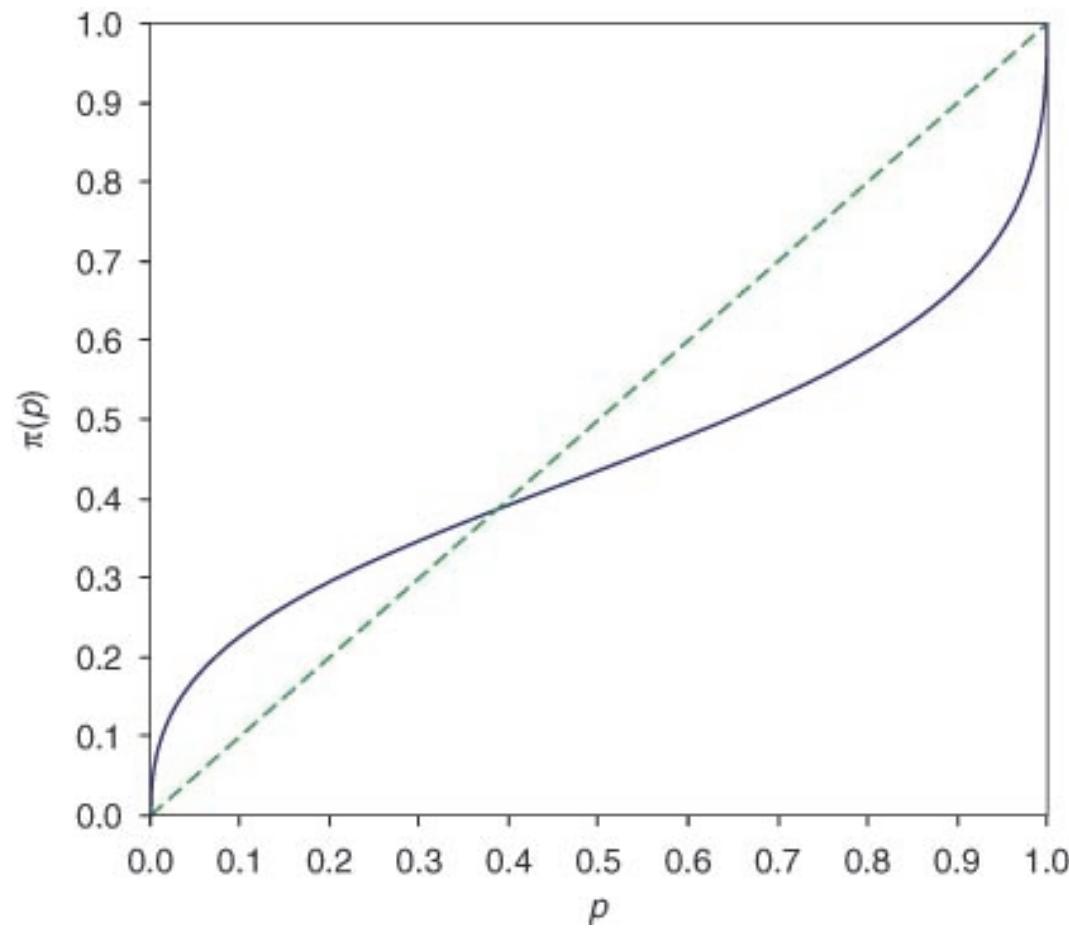
Think 1-in-10 or 1-in-100-type events, not 1-in-a-million-type events.

# The difference between “kinda rare” and “really, really rare”.

Systematically differentiating between these two is the focus of a large part of statistical research.

# People tend to make decisions irrationally in the face of high- and low-probability events.

How people perceive risk, according to Prospect Theory:



See Kahneman, Daniel; Tversky, Amos (1979). ["Prospect Theory: An Analysis of Decision under Risk"](#) (PDF). *Econometrica*. **47** (2):263

Don't be fooled by the pattern recognition box in  
your head.

# Better Ways to Think About Data and Uncertainty



# Misadventures with binning



# Doing analysis on binned continuous data can lead to incorrect conclusions.

Continuous Variable	Values	Binned Values						
Range	4,000 – 21,000 ft	4,000-9,000	9,000-16,000	16,000-21,000				
Age	30-64 yrs	30-34	35-39	40-44	45-49	50-54	55-59	60-64

# Mortality among middle-aged white Americans is increasing!



## Rising morbidity and mortality in midlife among white non-Hispanic Americans in the 21st century

Anne Case<sup>1</sup> and Angus Deaton<sup>1</sup>

Woodrow Wilson School of Public and International Affairs and Department of Economics, Princeton University, Princeton, NJ 08544

Contributed by Angus Deaton, September 17, 2015 (sent for review August 22, 2015; reviewed by David Cutler, Jon Skinner, and David Weir)

# The money plot

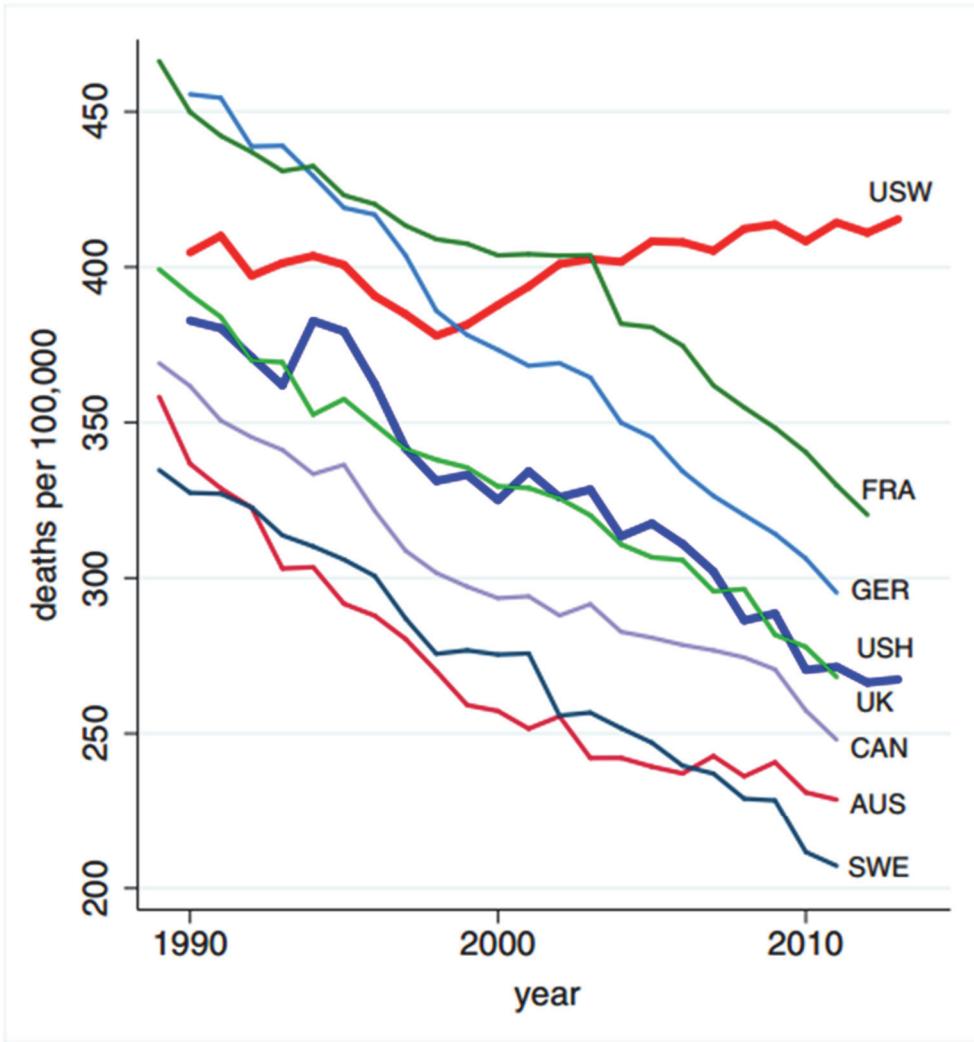


Fig. 1. All-cause mortality, ages 45–54 for US White non-Hispanics (USW), US Hispanics (USH), and six comparison countries: France (FRA), Germany (GER), the United Kingdom (UK), Canada (CAN), Australia (AUS), and Sweden (SWE).

Middle-aged  
despair?

Loss of meaning  
from work/faith/  
family?

Effects of the  
Great Recession?

Except ...

SCIENCE

THE STATE OF THE UNIVERSE.

NOV. 11 2015 12:38 PM

Slate

# Is the Death Rate Really Increasing for Middle-Aged White Americans?

I ran the numbers, and the story isn't as simple as it seems.

By Andrew Gelman



490



108

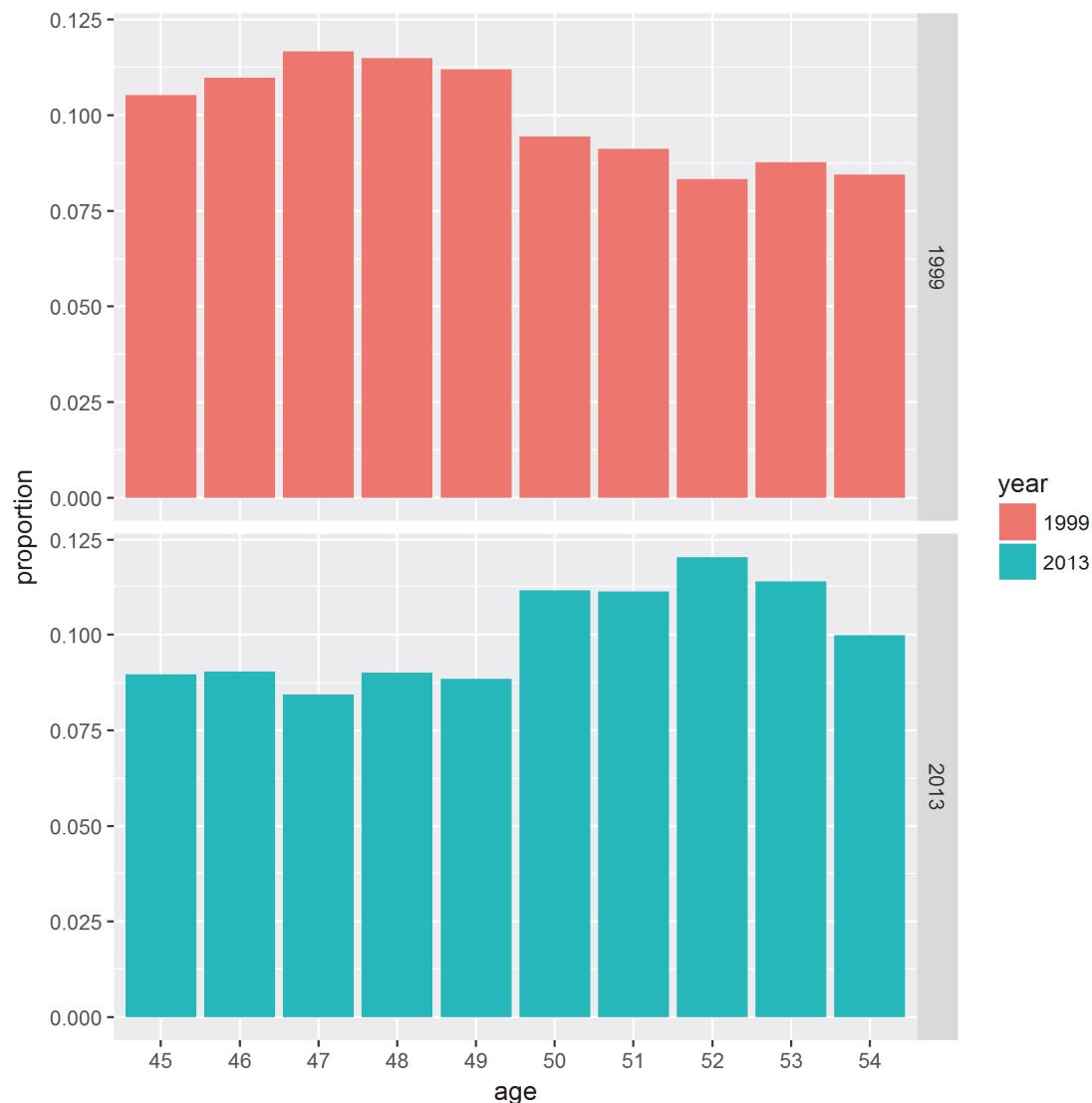


**Case and Deaton used binned data, but the composition of the bins had changed.**

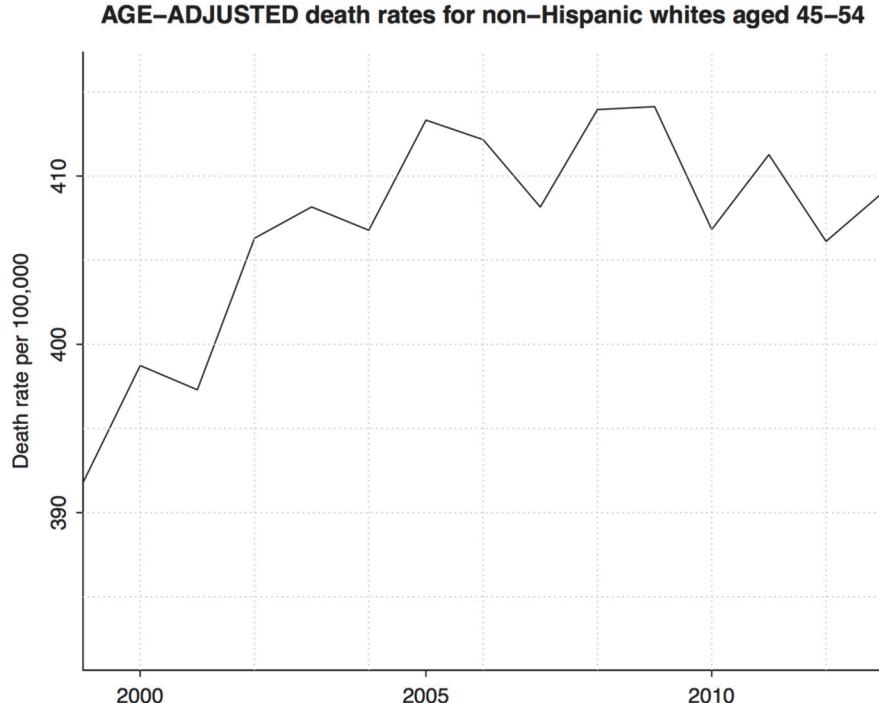
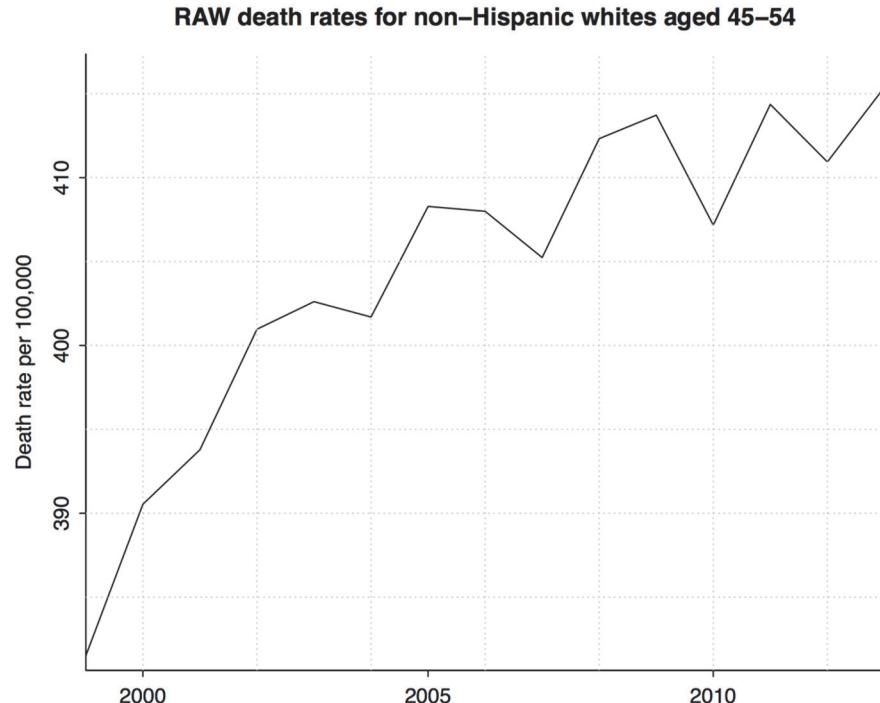
White Americans age 45-54 were on average 0.4 years older in 2013 than they were in 1999!

(Thanks, baby boomers!)

# Binning data can obscure changes within bins, resulting in unfounded conclusions.

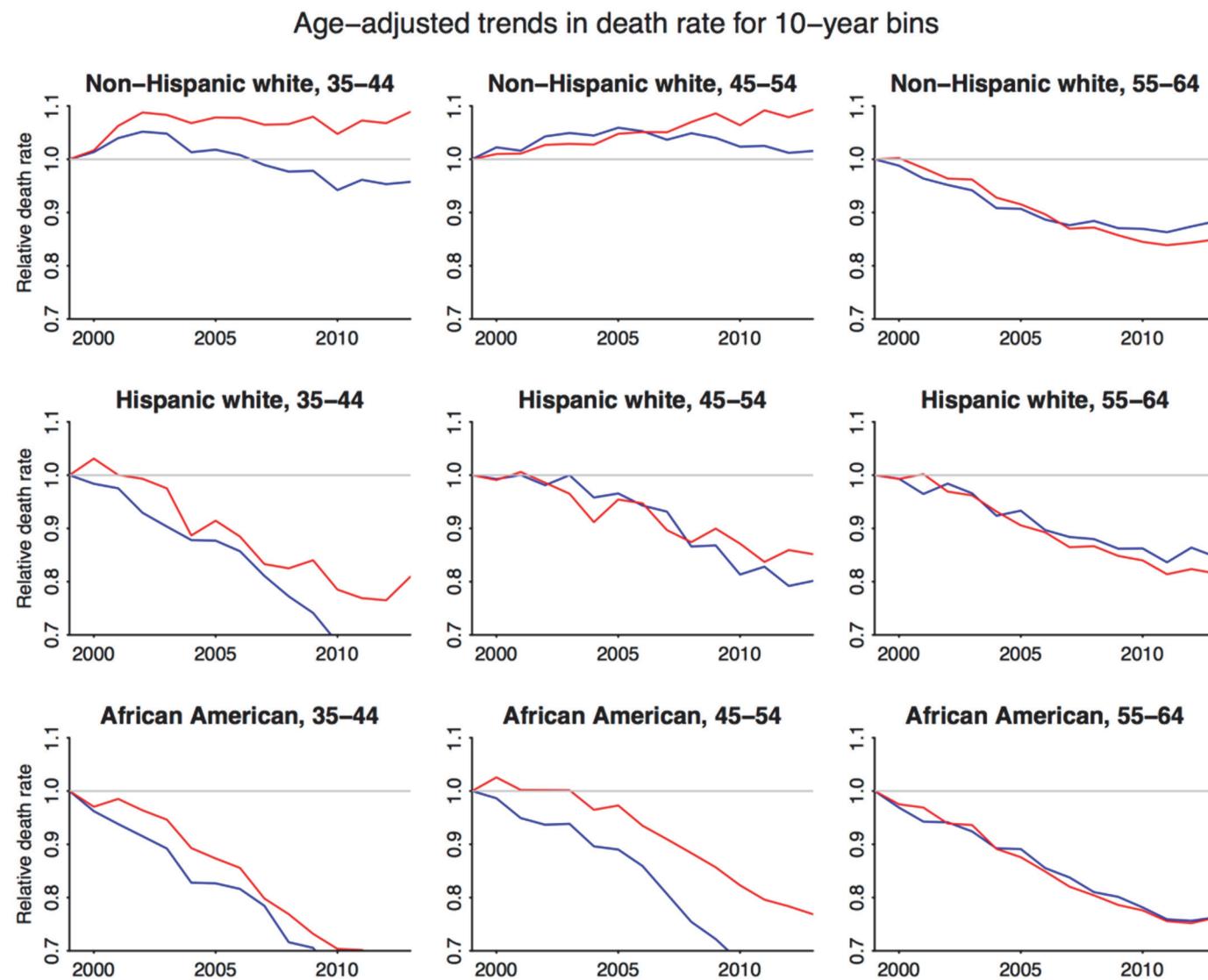


# Taking that difference in age into account ...



“Op-ed writers: Go back to your notebooks! Time to come up with new explanations.”

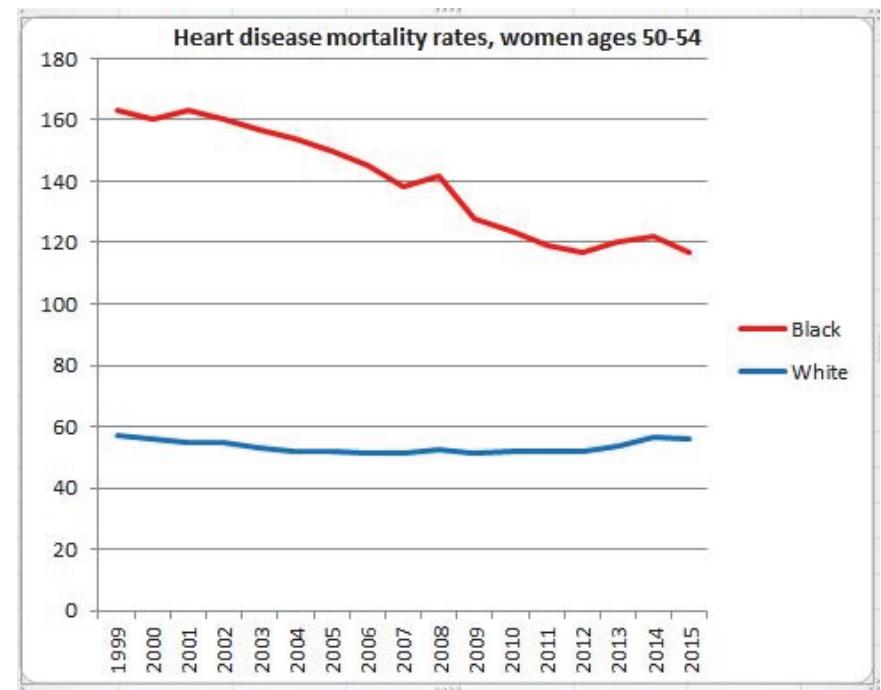
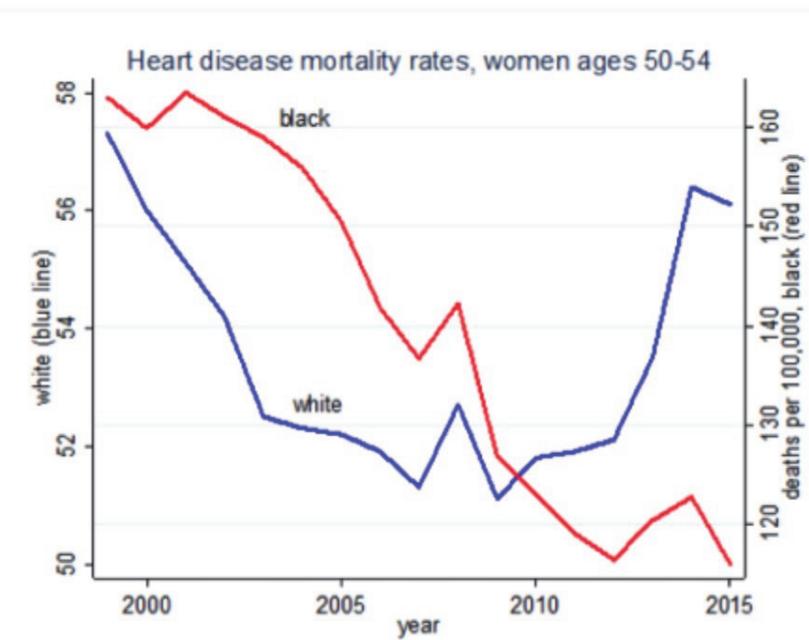
# So what's the real story with this data?



**Don't believe every picture you read.**



# Not to keep beating up Case and Deaton, but ...

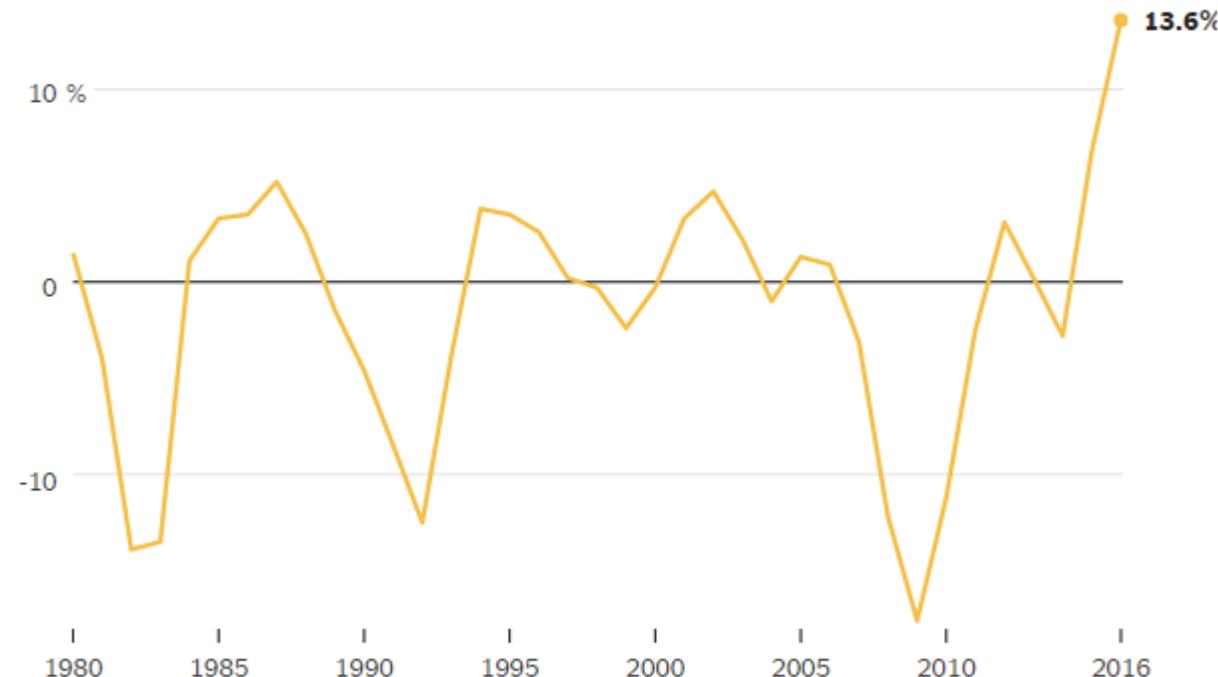


On the left, Figure 2.1 from Case and Deaton. On the right, the same trends shown with a single y-axis.

# Rates of traffic fatalities are spiking like never before!

## A Surge in Fatalities

Change in U.S. vehicle deaths over successive two-year periods.

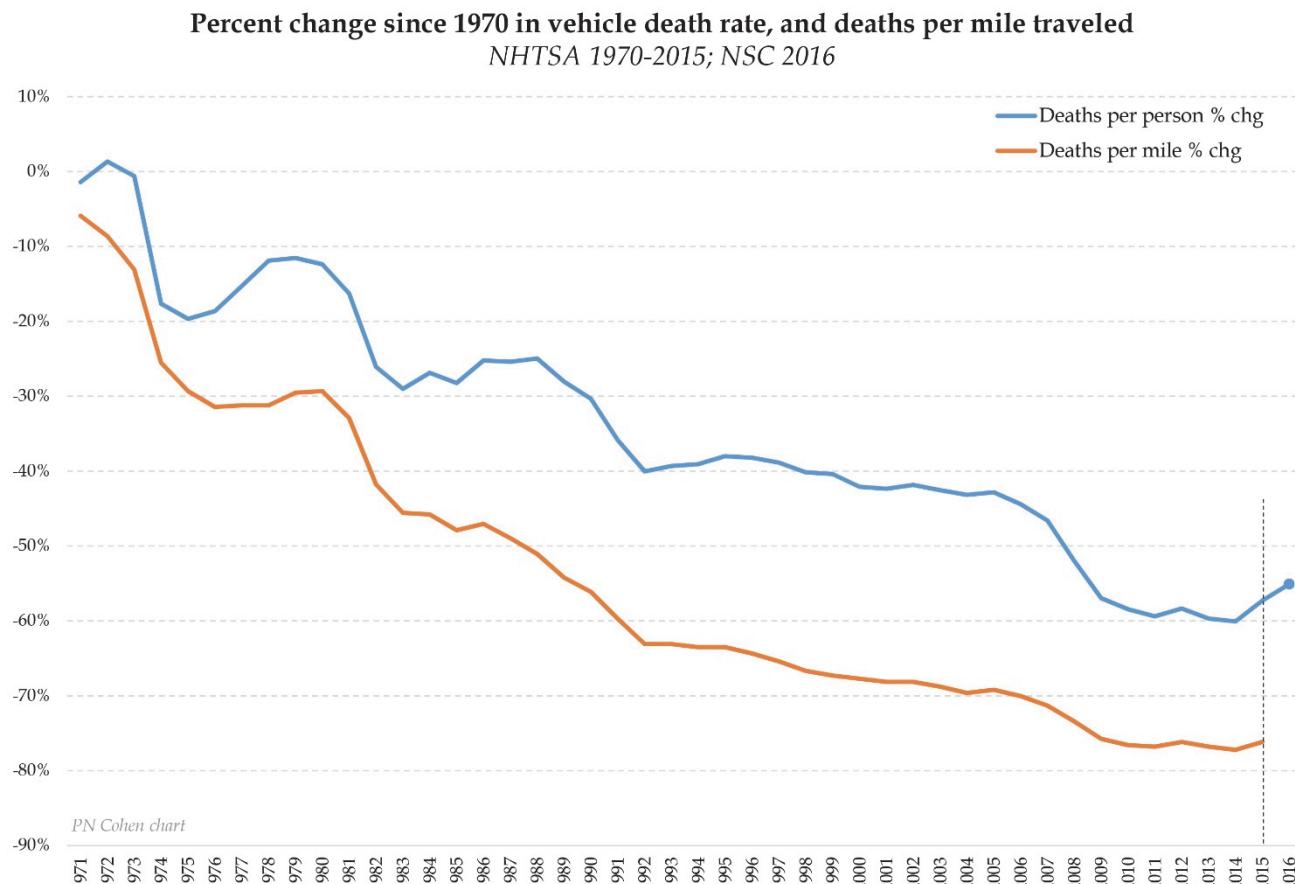


Deaths decline with recessions (such as 2007–9) because driving declines.

Source: National Safety Council

Per David Leonhardt at NYT

# ... except overall traffic fatalities are way down over the same time frame?



h/t Philip Cohen

Think carefully about how you process and display  
your data. It matters a lot!

# Statistical Models



# But Why Statistical Models?



Statistical models provide a framework for evaluating data  
that accounts for randomness.

**Statistical models can help us avoid the “naked eye” data interpretation errors described in the previous sections.**

... but if you use a statistical model, you need to be aware of its limitations as well!



You know what they say about assumptions ...

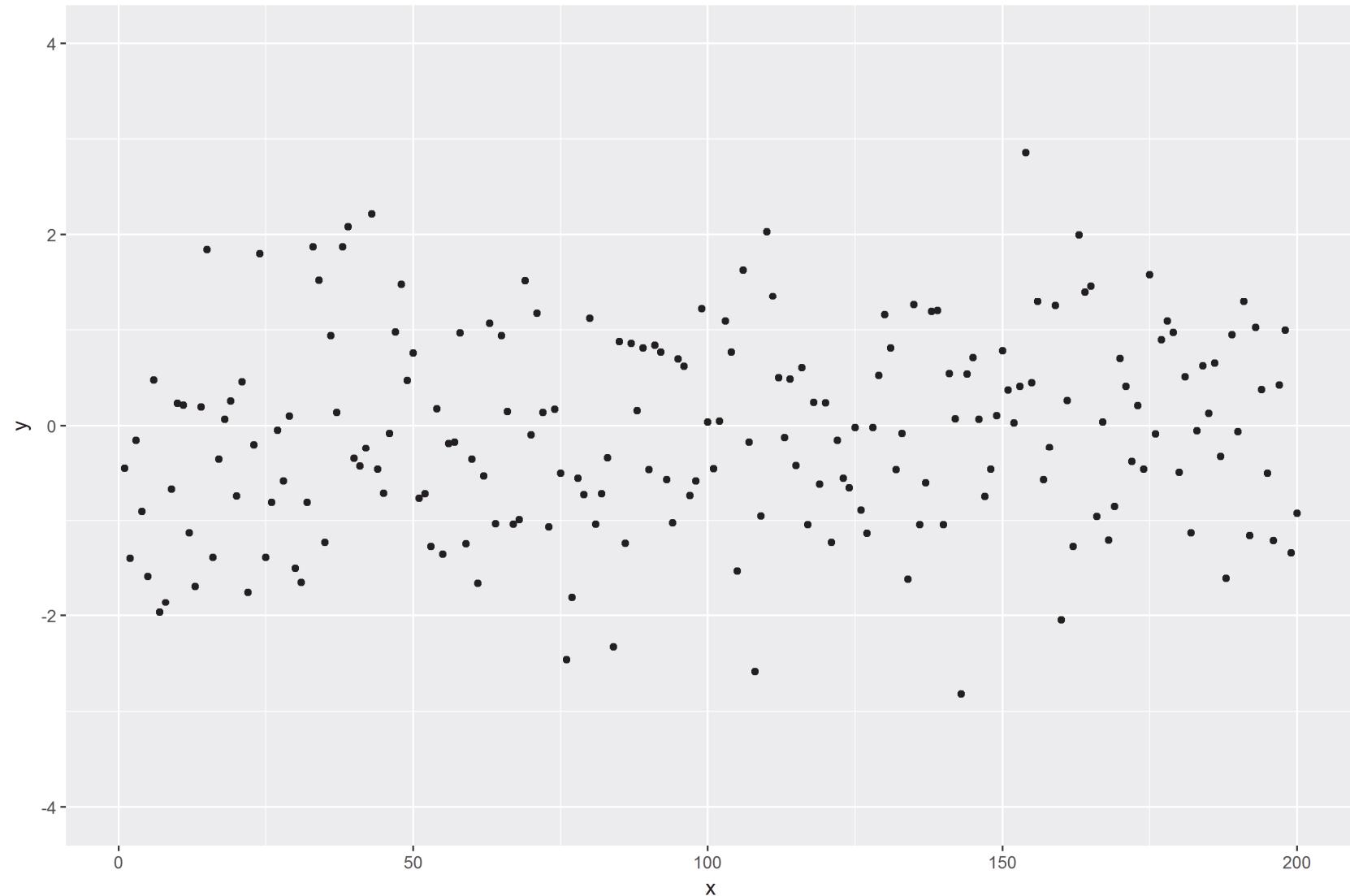


# **Model robustness is an underrated quality that most common statistical methods have.**

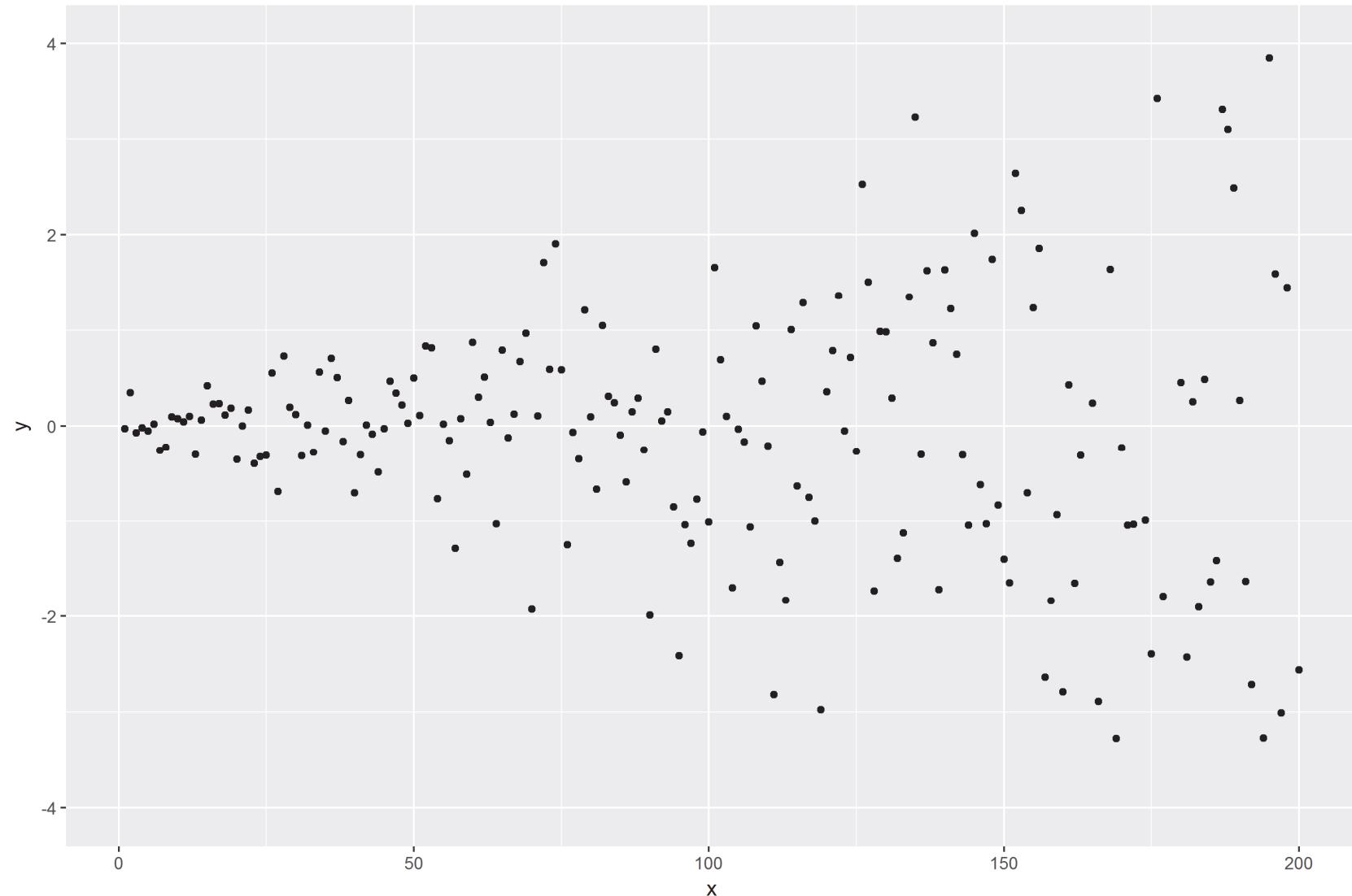
All models rely on assumptions, but many common statistical models are pretty robust to failures in these assumptions.

**Take the heteroskedasticity assumption ...**

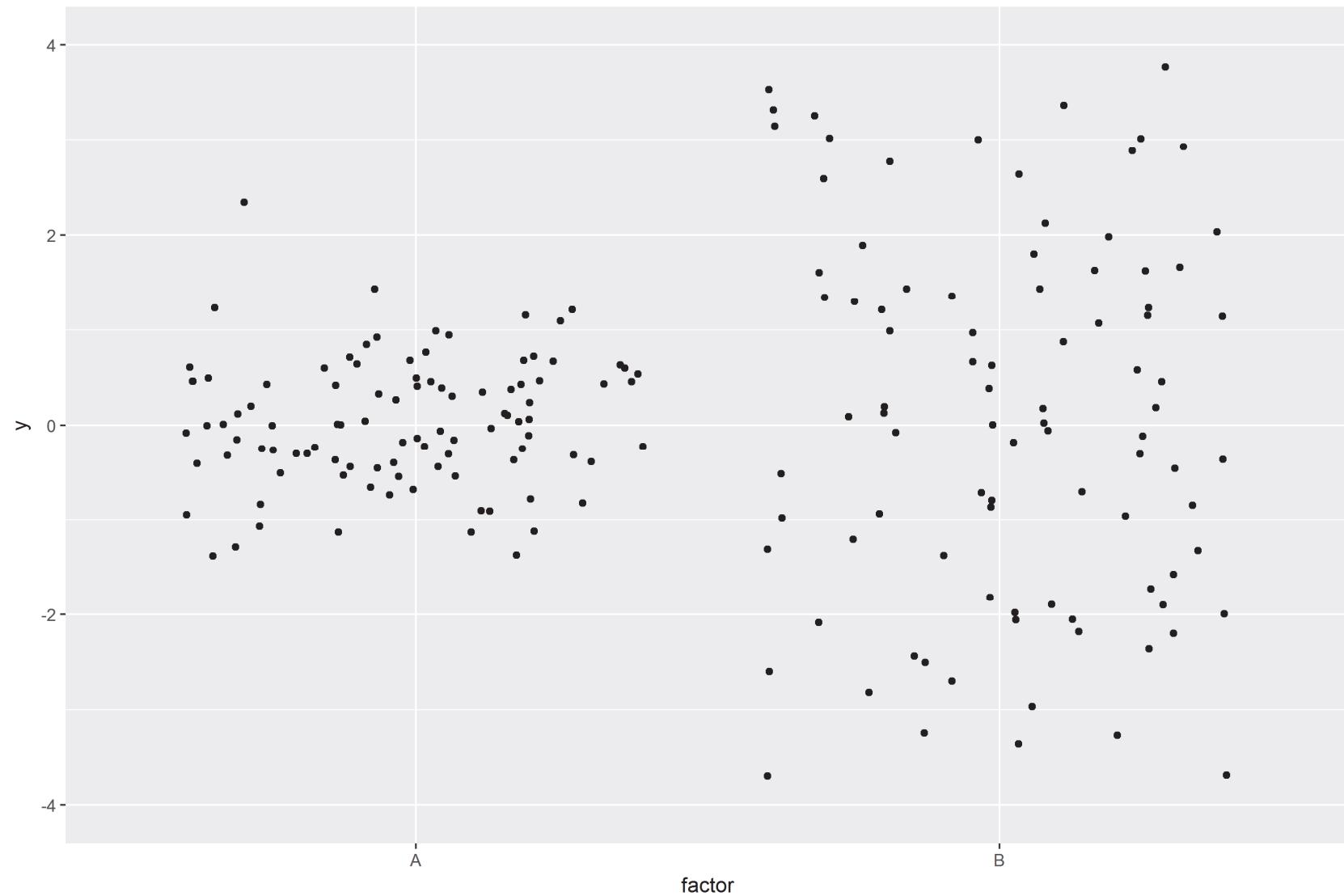
# Homoskedastic variance



# Heteroskedastic variance



# Heteroskedastic variance (categorical factor)



**If the variance of our response variable is heteroskedastic, much of our analysis is still valid!**

**Estimates of the mean are still unbiased and have good statistical properties (BLUE) ...**

**... but inference (CIs, p-values about parameter estimates) is no longer valid.**

# How bad is too bad?



**In the real world (and especially the world of OT!), data rarely follows all of the assumptions we make when we build models for it:**

Linearity

Normality

Constant variance (as in the example above)

Independence of runs/observations

Independence of factors

**Determining hard-and-fast rules for when a violation of these assumptions is acceptable is difficult and not productive.**

What are you using your model for?

How does the violated assumption impact your results?

How badly wrong is the assumption?

How robust is your model to violations of this assumption?

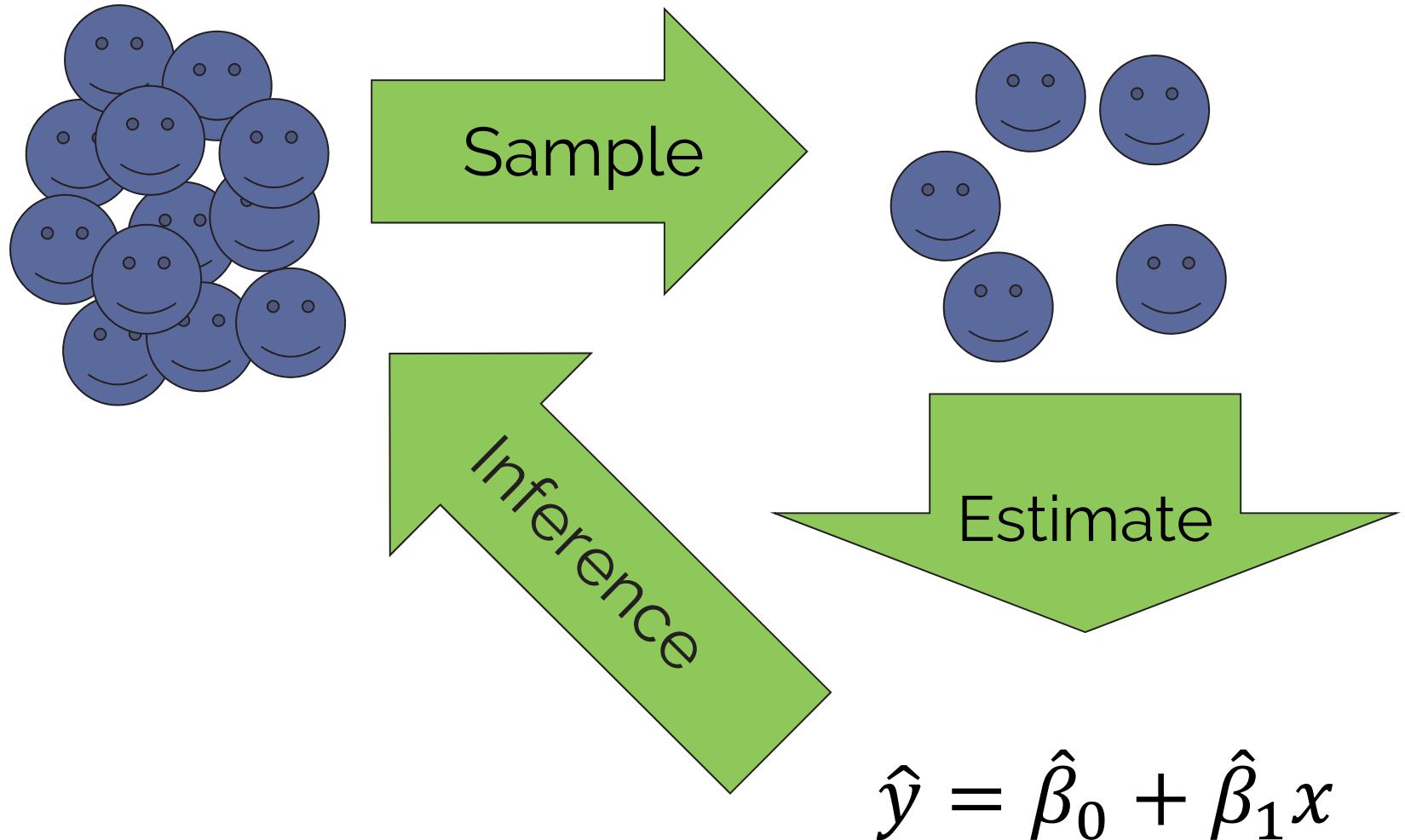
What alternative models or analysis approaches are available?



**“The model is the result!”**



When we analyze operational test data, we're interested in  
the **population**, not the **sample**.



Statistical models provide a framework to avoid errors commonly encountered by humans processing random data.



# Conclusions



Think about where your data come from and how that impacts the conclusions you can draw.

Don't be fooled by the pattern recognition box in your head.

Think carefully about how you process and display your data. It matters a lot!

Statistical models provide a framework to avoid errors commonly encountered by humans processing random data.

Thinking about data can be hard. Don't be afraid to consult an expert!



# Backup



# When summary statistics go wrong



# Notional example of how aggregating over factors (“roll ups”) can give you incorrect results.

P(Detect) for 360-degree Mode  
66.2%

P(Detect) for 360-degree Mode by Munition

Artillery	Mortar	Rocket
45.4%	65.8%	94.8%

P(Detect) for 360-degree Mode by Munition and Fire Rate

	Standard	Artillery	Mortar	Rocket
Mode	72.8%	78.0%	94.8%	
Volley	Artillery	Mortar	Rocket	
	7.8%	27.6%	N/A	

P(Detect) for 360-degree Mode by Munition and QE (Standard Fire only)

High QE	Artillery	Mortar	Rocket
88.4%	72.3%	N/A	
Low QE	Artillery	Mortar	Rocket
60.1%	95.6%	94.8%	



**Statisticians take p-values less seriously than you  
do.**



**When you know a tool well, you understand its strengths and weaknesses and feel more comfortable discounting it when appropriate.**

# Analogy: Advanced metrics in sports

The folks who create metrics have an understanding of what those metrics do well and what they don't. This provides important context when trying to make sense of the ratings they spit out.

**PER is useful but not the only thing worth considering when evaluating NBA players.**

Hollinger Stats - Player Efficiency Rating - Qualified Players													
RK	PLAYER	GP	MPG	TS%	AST	TO	USG	ORR	DRR	REBR	PER	VA	EWA
1	Russell Westbrook, OKC	81	34.6	.554	23.4	12.2	42.5	5.4	28.8	17.1	<b>30.70</b>	823.7	27.5
2	Kevin Durant, GS	62	33.4	.651	18.4	8.5	26.7	2.2	23.6	13.6	<b>27.68</b>	530.7	17.7
3	Kawhi Leonard, SA	74	33.4	.610	13.3	7.9	29.6	3.7	15.7	9.8	<b>27.62</b>	632.4	21.1
4	Anthony Davis, NO	75	36.1	.580	7.3	8.4	29.8	6.7	27.9	17.2	<b>27.59</b>	650.4	21.7
5	James Harden, HOU	81	36.4	.613	27.6	14.1	35.1	3.5	20.9	12.2	<b>27.43</b>	744.7	24.8
6	LeBron James, CLE	74	37.8	.619	25.6	12.0	30.1	4.0	20.7	12.6	<b>27.11</b>	692.7	23.1
7	Isaiah Thomas, BOS	76	33.8	.625	18.5	8.7	32.8	1.9	7.0	4.4	<b>26.59</b>	597.9	19.9
8	Nikola Jokic, DEN	73	27.9	.640	24.2	11.5	24.0	11.6	27.2	19.5	<b>26.40</b>	453.3	15.1
9	Chris Paul, LAC	61	31.5	.614	35.0	9.1	25.8	2.4	14.9	8.8	<b>26.25</b>	437.3	14.6
10	Giannis Antetokounmpo, MIL	80	35.6	.599	19.8	10.7	27.4	5.9	22.6	14.3	<b>26.13</b>	663.7	22.1

# Things like p-values should also be understood in the proper context.

P-values are contingent on a specific hypothesis (including a statistical model) and a particular set of test data. Consider how much data you have, the size of your observed effect, and what you know about the factor you're looking at in addition to the p-value.

**Conclusions first, analysis later!**



**When there are lots of ways to measure something,  
“motivated” researchers can prove anything they want.**

# Democrats are good for the economy!

## Hack Your Way To Scientific Glory



You're a social scientist with a hunch: **The U.S. economy is affected by whether Republicans or Democrats are in office.** Try to show that a connection exists, using real data going back to 1948. For your results to be publishable in an academic journal, you'll need to prove that they are "statistically significant" by achieving a low enough p-value.

### 1 CHOOSE A POLITICAL PARTY

Republicans

Democrats

### 2 DEFINE TERMS

Which politicians do you want to include?

- Presidents
- Governors
- Senators
- Representatives

How do you want to measure economic performance?

- Employment
- Inflation
- GDP
- Stock prices

Other options

- Factor in power  
Weight more powerful positions more heavily
- Exclude recessions  
Don't include economic recessions

### 3 IS THERE A RELATIONSHIP?

Given how you've defined your terms, does the economy do better, worse or about the same when more Democrats are in power? Each dot below represents one month of data.



### 4 IS YOUR RESULT SIGNIFICANT?

If there were no connection between the economy and politics, what is the probability that you'd get results at least as strong as yours? That probability is your p-value, and by convention, you need a p-value of 0.05 or less to get published.



### Result: Publishable

You achieved a p-value of 0.03 and showed that **Democrats** have a **positive effect** on the economy. Get ready to be published!

If you're interested in reading real (and more rigorous) studies on the connection between politics and the economy, see the work of Larry Bartels and Alan Blinder and Mark Watson.

Data from The @unitedstates Project, National Governors Association, Bureau of Labor Statistics, Federal Reserve Bank of St. Louis and Yahoo Finance.

Check out <https://projects.fivethirtyeight.com/p-hacking/> for more!

# Err, I mean Republicans!

## Hack Your Way To Scientific Glory



You're a social scientist with a hunch: **The U.S. economy is affected by whether Republicans or Democrats are in office.** Try to show that a connection exists, using real data going back to 1948. For your results to be publishable in an academic journal, you'll need to prove that they are "statistically significant" by achieving a low enough p-value.

### 1 CHOOSE A POLITICAL PARTY

Republicans

Democrats

### 2 DEFINE TERMS

Which politicians do you want to include?

- Presidents
- Governors
- Senators
- Representatives

How do you want to measure economic performance?

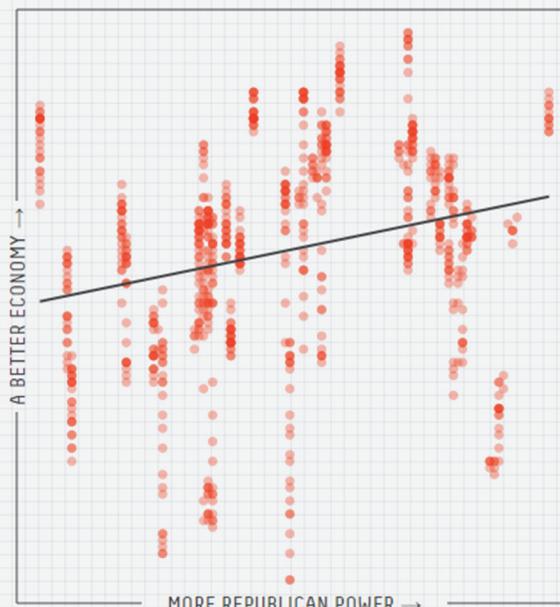
- Employment
- Inflation
- GDP
- Stock prices

Other options

- Factor in power  
Weight more powerful positions more heavily
- Exclude recessions  
Don't include economic recessions

### 3 IS THERE A RELATIONSHIP?

Given how you've defined your terms, does the economy do better, worse or about the same when more Republicans are in power? Each dot below represents one month of data.



### 4 IS YOUR RESULT SIGNIFICANT?

If there were no connection between the economy and politics, what is the probability that you'd get results at least as strong as yours? That probability is your p-value, and by convention, you need a p-value of **0.05 or less** to get published.



### Result: Publishable

You achieved a p-value of **less than 0.01** and showed that **Republicans have a positive effect on the economy**. Get ready to be published!

If you're interested in reading real (and more rigorous) studies on the connection between politics and the economy, see the work of Larry Bartels and Alan Blinder and Mark Watson.

Data from The @unitedstates Project, National Governors Association, Bureau of Labor Statistics, Federal Reserve Bank of St. Louis and Yahoo Finance.

Check out <https://projects.fivethirtyeight.com/p-hacking/> for more!

# No, no, Democrats!!

## Hack Your Way To Scientific Glory



You're a social scientist with a hunch: **The U.S. economy is affected by whether Republicans or Democrats are in office.** Try to show that a connection exists, using real data going back to 1948. For your results to be publishable in an academic journal, you'll need to prove that they are "statistically significant" by achieving a low enough p-value.

### 1 CHOOSE A POLITICAL PARTY

Republicans

Democrats

### 2 DEFINE TERMS

Which politicians do you want to include?

- Presidents
- Governors
- Senators
- Representatives

How do you want to measure economic performance?

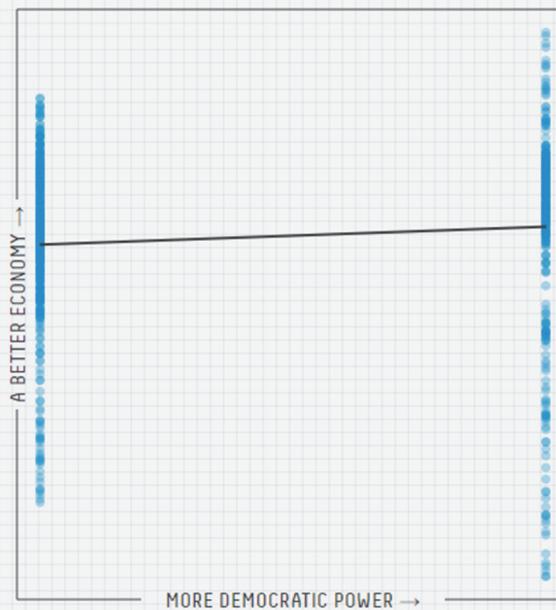
- Employment
- Inflation
- GDP
- Stock prices

Other options

- Factor in power
  - Weight more powerful positions more heavily
- Exclude recessions
  - Don't include economic recessions

### 3 IS THERE A RELATIONSHIP?

Given how you've defined your terms, does the economy do better, worse or about the same when more Democrats are in power? Each dot below represents one month of data.



### 4 IS YOUR RESULT SIGNIFICANT?

If there were no connection between the economy and politics, what is the probability that you'd get results at least as strong as yours? That probability is your p-value, and by convention, you need a p-value of 0.05 or less to get published.



### Result: Publishable

You achieved a p-value of less than 0.01 and showed that **Democrats** have a **positive effect** on the economy. Get ready to be published!

If you're interested in reading real (and more rigorous) studies on the connection between politics and the economy, see the work of Larry Bartels and Alan Blinder and Mark Watson.

Data from The @unitedstates Project, National Governors Association, Bureau of Labor Statistics, Federal Reserve Bank of St. Louis and Yahoo Finance.

Check out <https://projects.fivethirtyeight.com/p-hacking/> for more!

# Wait, no, Republicans!

## Hack Your Way To Scientific Glory



You're a social scientist with a hunch: **The U.S. economy is affected by whether Republicans or Democrats are in office.** Try to show that a connection exists, using real data going back to 1948. For your results to be publishable in an academic journal, you'll need to prove that they are "statistically significant" by achieving a low enough p-value.

### 1 CHOOSE A POLITICAL PARTY

Republicans

Democrats

### 2 DEFINE TERMS

Which politicians do you want to include?

- Presidents
- Governors
- Senators
- Representatives

How do you want to measure economic performance?

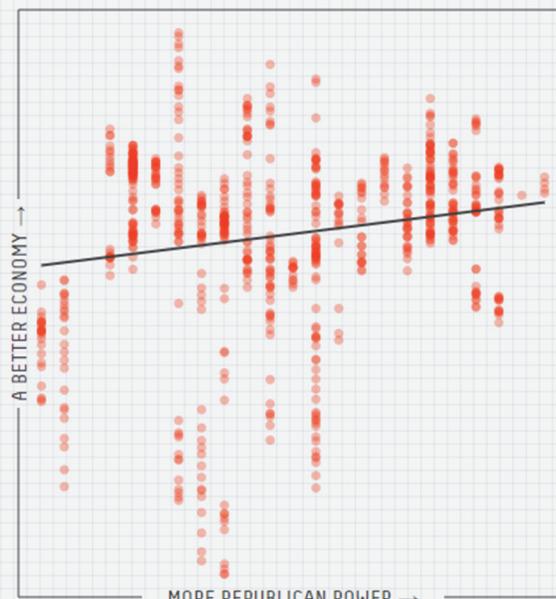
- Employment
- Inflation
- GDP
- Stock prices

Other options

- Factor in power
  - Weight more powerful positions more heavily
- Exclude recessions
  - Don't include economic recessions

### 3 IS THERE A RELATIONSHIP?

Given how you've defined your terms, does the economy do better, worse or about the same when more Republicans are in power? Each dot below represents one month of data.



### 4 IS YOUR RESULT SIGNIFICANT?

If there were no connection between the economy and politics, what is the probability that you'd get results at least as strong as yours? That probability is your p-value, and by convention, you need a p-value of 0.05 or less to get published.



### Result: Publishable

You achieved a p-value of less than 0.01 and showed that **Republicans** have a **positive effect** on the economy. Get ready to be published!

If you're interested in reading real (and more rigorous) studies on the connection between politics and the economy, see the work of Larry Bartels and Alan Binder and Mark Watson.

Data from The @unitedstates Project, National Governors Association, Bureau of Labor Statistics, Federal Reserve Bank of St. Louis and Yahoo Finance.

Check out <https://projects.fivethirtyeight.com/p-hacking/> for more!

**And this can happen even when a researcher isn't actively trying to produce a biased result!**

# Researchers make decisions about ...

What data to collect

What response variables to consider

What factors to compare

What analysis methods they should use

What results to publish

# The garden of forking paths:

Why multiple comparisons can be a problem, even when there is no “fishing expedition” or “p-hacking” and the research hypothesis was posited ahead of time.

Andrew Gelman and Eric Loken

“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.”

-John Tukey

“A well-ordered pile of numbers offers nearly unlimited scope for storytelling.”

-Andrew Gelman

We're really bad at assessing success for long-tail events.



**“Kinda rare” to a statistician is another way of saying  
“random noise”.**

# Evaluating predictions or models that produce probabilistic outcomes is challenging.

Model	Prediction	Outcome
A	20%	Yes
B	49%	Yes
C	51%	Yes

Which model is right? Which is wrong?

# The problem is even worse with rare events.

Model	Prediction	Outcome
A	5%	Yes
B	1%	Yes
C	0.1%	Yes

Without substantial repetition, it's hard to differentiate, and in operational testing we may not have enough repetition to see such rare events multiple times.



# The Winner's Curse



**Determining whether past success was due to efficacy or randomness is hard.**

# Or Allah in a fish.

**BBC NEWS**

Watch One-Minute World News

News services Your news when you want it

News Front Page

Last Updated: Tuesday, 31 January 2006, 19:32 GMT

E-mail this to a friend | Printable version

## Tropical fish 'has Allah marking'

A pet shop owner has found the markings on one of his tropical fish appear to spell the word "Allah" in Arabic.

An Asian customer spotted the markings on the astronotus ocellatus at Walker Aquatics in Waterfoot, Rawtenstall, Lancs, and offered to buy it for £10.

The customer bought the most expensive tank in the shop to house the fish, at a cost of £700.

Shop owner Tony Walker said he would honour the buyer's £10 offer but said he expected more interest in the fish.

The fish is also believed to have the word "Mohammed" in its markings on its other side.



A customer has offered £10 for the fish

**BBC Lancashire**  
Sport, travel, weather, things to do, features and much more

**SEE ALSO:**

- Scientists find 'smallest fish'  
25 Jan 06 | Science/Nature
- World's biggest fish 'shrinking'  
17 Jan 06 | Science/Nature
- Tropical fish 'may be UK first'  
19 Sep 05 | Norfolk
- Boy struck by giant tropical fish  
28 Aug 05 | Wales

**RELATED BBC LINKS:**

- Nature - sea life

**TOP LANCASHIRE STORIES**

- Police 'solve' 1975 murder case

**This phenomena is not limited to fans of fish and enthusiasts of grilled cheese.**

# People tried to read a lot into the way the stock market moved at various points in the 2016 election.

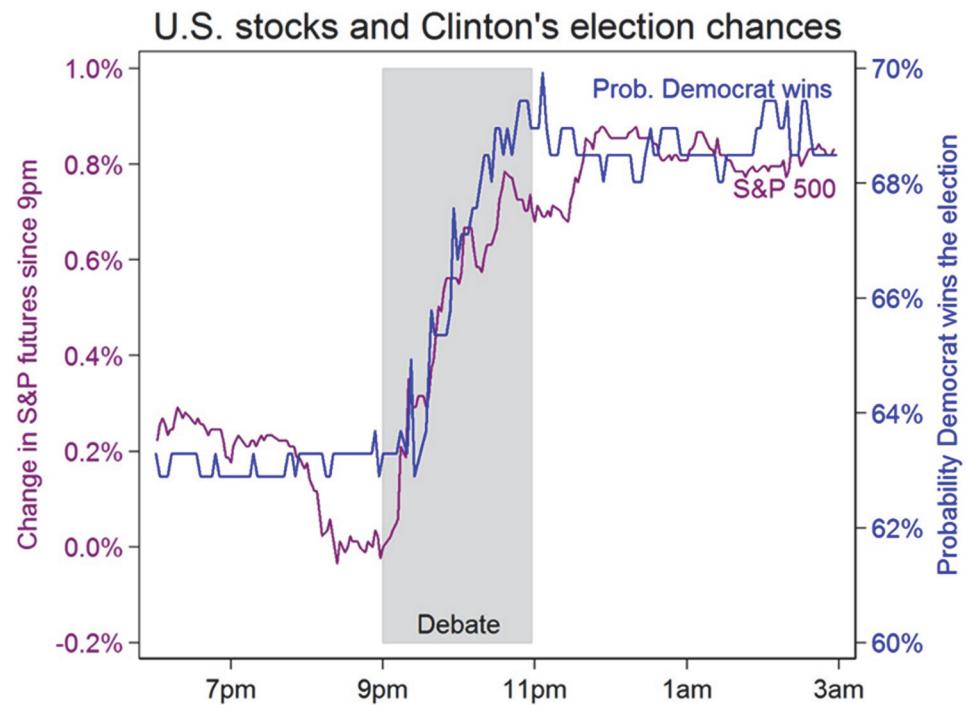
## Debate Night Message: The Markets Are Afraid of Donald Trump

Economic View  
By JUSTIN WOLFERS SEPT. 30, 2016



"Market movements over the October 7-9 weekend, during which a tape was released ... tell a largely consistent story."

– 20 October 2016



# They were generally wrong.

Markets Sent a Strong Signal on Trump ... Then Changed Their Minds



Justin Wolfers @JustinWolfers NOV. 18, 2016



## US MARKETS

[MARKET MOVERS](#) | [DOW 30](#) | [NASDAQ 100](#) | [IQ 100](#) | [SECTORS](#) | [WORLD HEAT MAP](#)

# Dow closes above 20,000 for first time as Trump orders send stocks flying

Fred Imbert | @foimbert  
Wednesday, 25 Jan 2017 | 4:23 PM ET



**This is because trying to predict the market is a fool's errand, though people still try to do it.**

Easy as 2, 4, 6



I'm thinking of a rule that governs this sequence of three numbers:

2, 4, 6

What is my rule? You may submit triplicates, and I'll tell you if they agree or don't agree with the rule.

# Only 21% of subjects guessed the rule correctly.

People tend to generate guesses with expected returns of “yes” rather than “no”, seeking to learn via confirmation rather than contradiction.

Peter Watson, "On the failure to eliminate hypotheses in a conceptual task", Quarterly Journal of Experimental Psychology, 12: 129-140, 1960