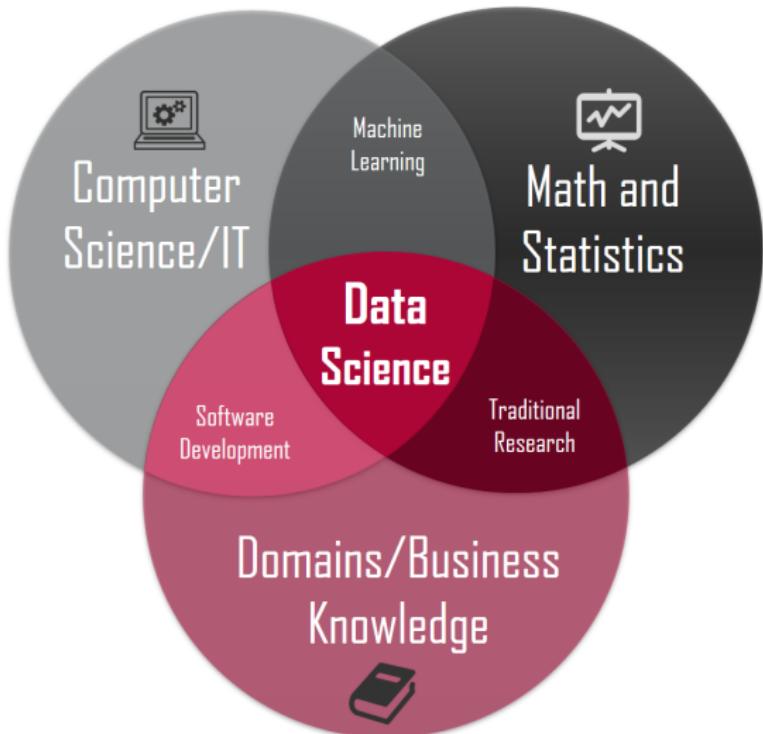


# Data Science 101: an introduction

Data Science Team

Stanford University, Department of Statistics

## One view



## A takehome message from this course:

Good Statistics and Data Science involves much more than running a Machine Learning algorithm on data.

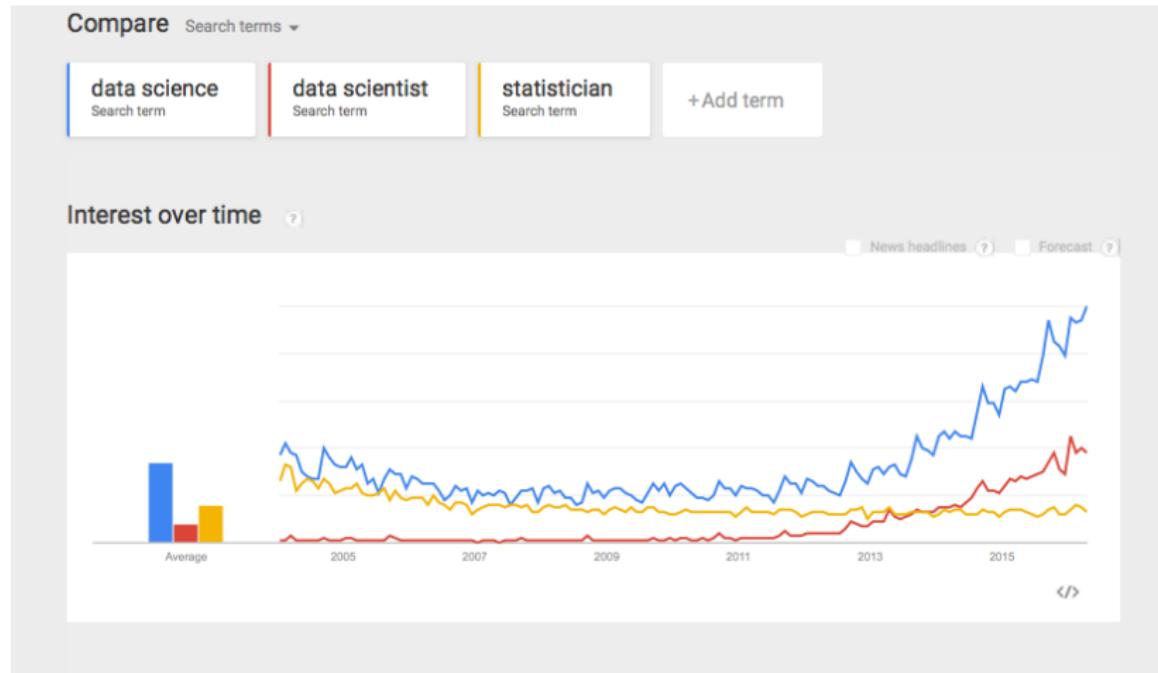
There are essential challenges like:

- ▶ How should I design my study to answer a question? Is it feasible and ethical to carry out such a study?
- ▶ What data should I collect? How much data?

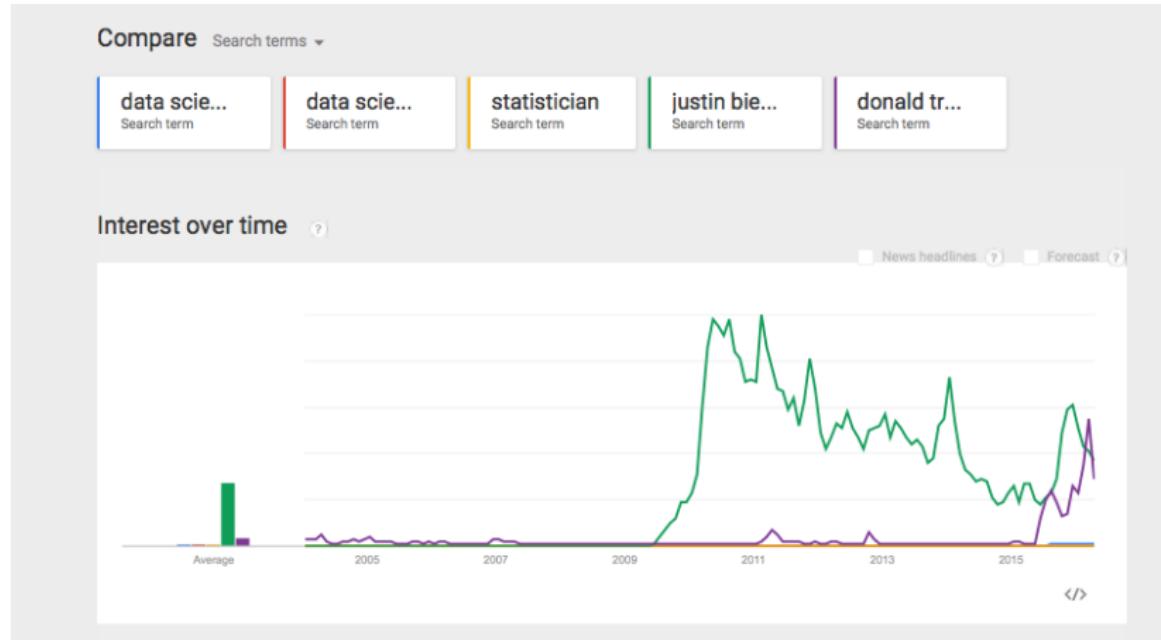
After collecting the data:

- ▶ What is the quality of my data?
- ▶ How should I model my data to try to answer the question at hand?
- ▶ What can I conclude from my study, and with how much certainty?

# Data Scientist vs Statistician in Google search



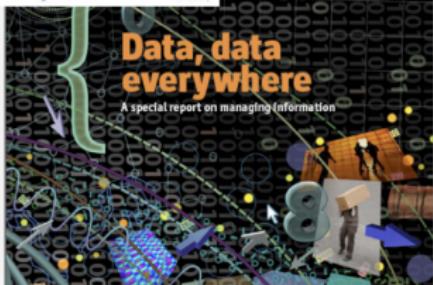
# Keeping things in perspective



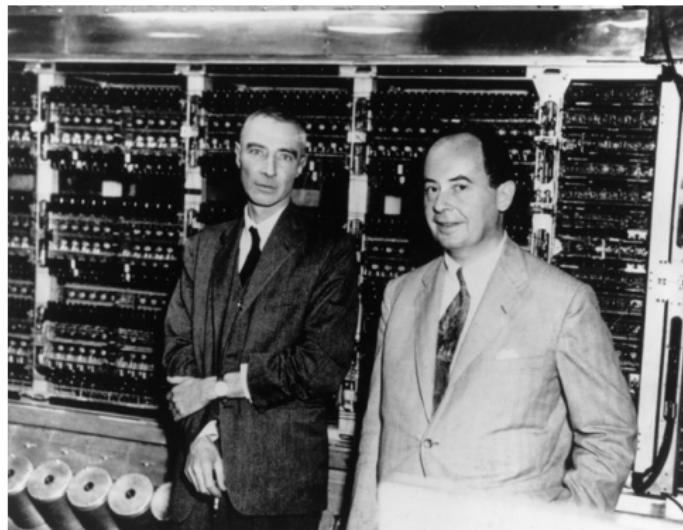
# Data Science 101

- ▶ Literacy for the digital age
- ▶ Data science 101 = Self-defense 101 in the modern workplace
- ▶ Data
  - ▶ What is it?
  - ▶ How to explore it?
  - ▶ Where does it come from?
- ▶ Science
  - ▶ What does it mean to learn from data?
  - ▶ Prediction vs Inference
  - ▶ How do we know when we are right or wrong?

# A lot of buzz about a lot of data



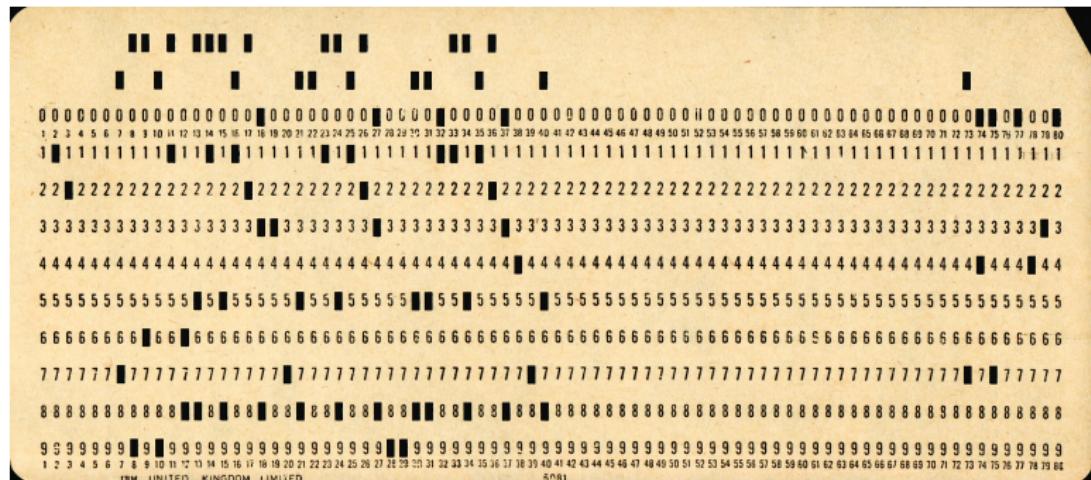
Things have changed a lot since the 1940s-1950s...



Robert Oppenheimer and John Von Neumann in front of the Electronic Numerical Integrator and Computer (ENIAC), [source](#).

The ENIAC was the first programmable, electronic, general-purpose digital computer (1945).

# And since the 1960s...



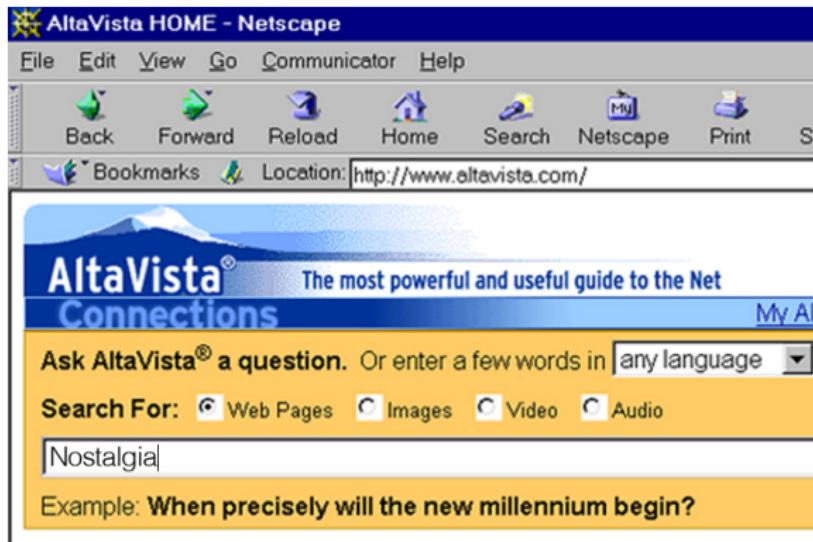
Punched card: Holds digital data represented by absence / presence of holes, [source](#).

And since the 1980s...



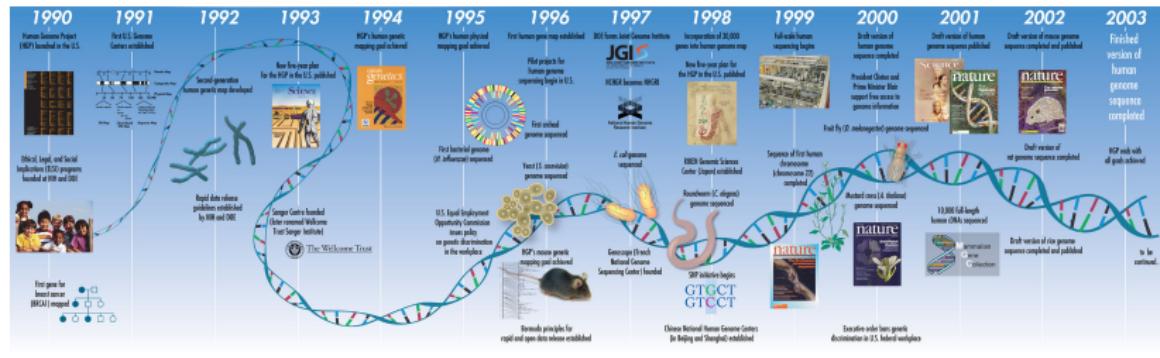
IBM PC 5150, [source](#). The IBM 5150 was the first [microcomputer](#) released in the IBM model line and had a strong influence on the personal computer market.

And since the 1990s...



AltaVista was established in 1995 and was one of the most-used early search engines, image source.

# And, yes, since the 2000s...



The Human Genome Project: 1990 - 2003, [image source](#).

## Some examples of Big Data

- ▶ **Genetics:** Now costs a few hundred dollars to sequence an entire genome. Typically there are many more measurements than subjects
- ▶ **Physics/CERN LHC:** LHC (Large Hadron Collider - particle accelerator) experiments produce about **90 petabytes** (1 petabyte =  $2^{50}$  bytes) of data per year (another 25 petabytes for non-LHC experiments)
- ▶ **Astronomy/Sloan Digital Sky Survey:** about 25.5TB telescope imaging data
- ▶ **Passively gathered data**
  - ▶ Products we buy
  - ▶ Topics that engage us
  - ▶ Our levels of physical activity
  - ▶ Who we talk to

## Some sources of Data

- ▶ [Stats for change](#) (public data sources related to social issues)
- ▶ [Data.gov](#) (U.S. Government open data)
- ▶ [dbGap](#) (database of Genotypes and Phenotypes)
- ▶ [CERN](#) (CERN particle physics data)
- ▶ [Google Dataset Search](#)
- ▶ Private data

### Google Dataset Search Beta

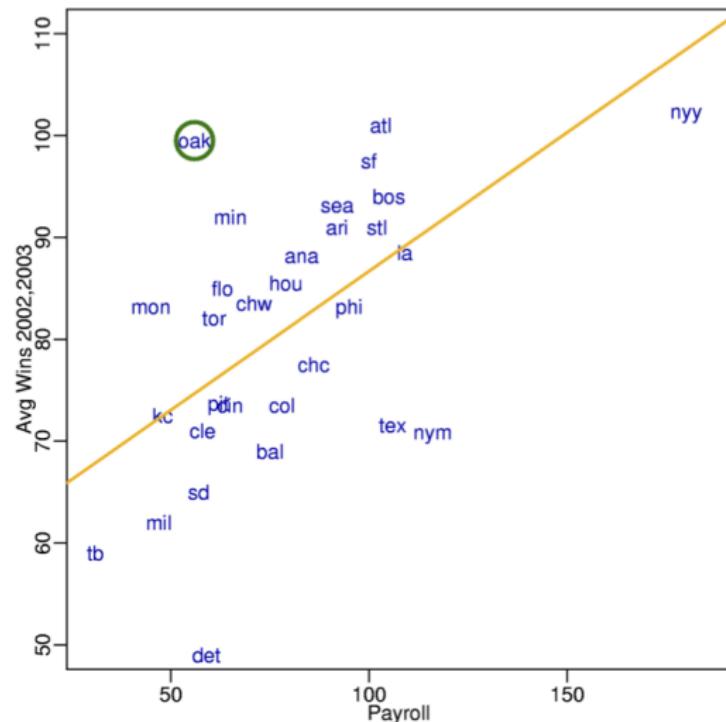
The screenshot shows the Google Dataset Search homepage. At the top, the title "Google Dataset Search" is displayed with "Beta" in red text next to it. Below the title is a search bar containing the placeholder text "Search for Datasets". To the right of the search bar is a blue magnifying glass icon. Below the search bar, there is sample search query text: "Try [boston education data](#) or [weather site:noaa.gov](#)".

Data science can do a lot: sports



Moneyball was a 2011 movie. Using “data science / statistics techniques”, the Oakland Athletics baseball team build a team of “undervalued” talent, [source](#).

## Data science can do a lot: sports



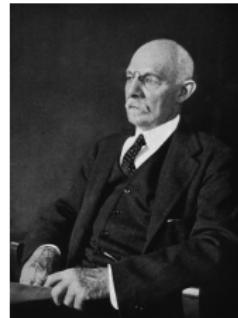
Starting around 2001, the Oakland A's picked players that scouts thought were no good but data said otherwise.

## The problem with bad data



- ▶ Picture of Truman holding the Chicago Daily Tribune with an [incorrect](#) banner headline on front.
- ▶ Truman won over his opponent Dewey in the 1948 presidential election.
- ▶ The paper relied on their analyst who predicted the winner in previous presidential contests. The conventional wisdom (supported by various polls) was that Dewey would win the election by a “landslide”, see [here](#).

## Data science can do a lot: medicine



- ▶ William Stewart Halsted of Johns Hopkins University introduced radical mastectomy for the treatment of Breast Cancer in the 1920s
- ▶ It became the standard treatment in the field, and became more radical (and disfiguring) as time went on
- ▶ Few doctors dared to challenge the conventional wisdom, partly because of Halsted's stature

## Fast forward to 1958



- ▶ Finally in 1958 Dr. Bernard Fisher of the University of Pittsburgh managed to run the first clinical trial comparing radical mastectomy to less invasive treatments.
- ▶ His work showed that early-stage breast cancer could be more effectively treated by lumpectomy, in combination with radiation therapy, chemotherapy, and/or hormonal therapy

# Data science can do a lot: medicine

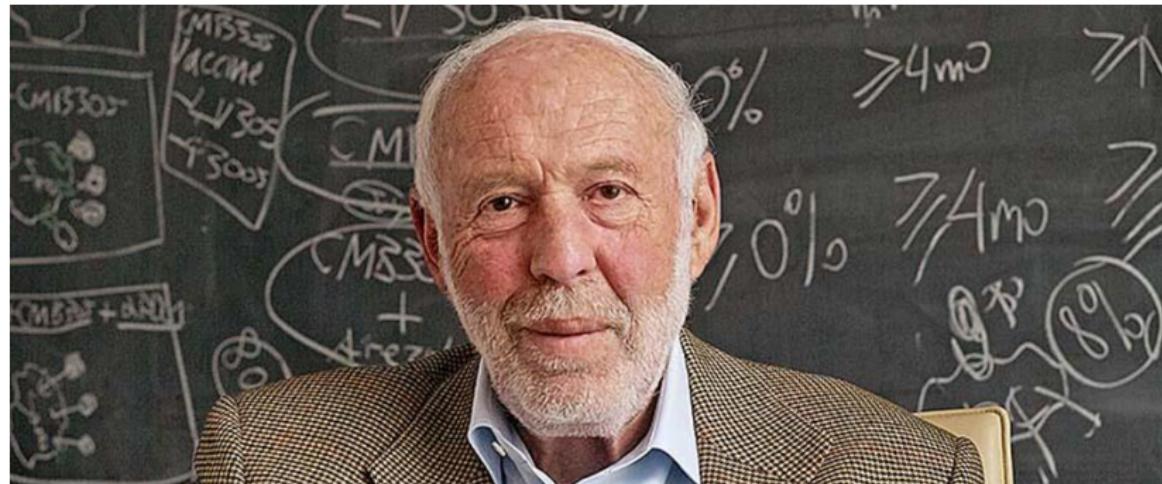


The National Institutes of Health has announced a new opportunity for organizations interested in helping engage volunteers in the [All of Us Research Program](#), part of the Precision Medicine Initiative. This funding opportunity, open to national and regional organizations, as well as local community groups, will support activities to promote enrollment and retention in the *All of Us* Research Program across diverse communities.

*All of Us* is an ambitious effort to gather data over time from 1 million or more people living in the United States, with the ultimate goal of accelerating research and improving health. Unlike research studies that are focused on a specific disease or population, *All of Us* will serve as a national research resource to inform thousands of studies, covering a wide variety of health conditions. Researchers will use data from the program to learn more about how individual differences in lifestyle, environment and biological make-up can influence health and disease. By taking part, people will be able to learn more about their own health and contribute to an effort that will advance the health of generations to come. NIH plans to launch the program later this year.

- ▶ The *All of Us* research program wants to build the most diverse health database in history.
- ▶ The goal is to investigate how biology, lifestyle and environment affect health in order to help find ways to treat and prevent disease.

## Data science can do a lot: finance



Jim Simons, Renaissance capital, see [here](#). He is a mathematician and founder of a quantitative hedge fund. He is known for his studies on pattern recognition.

# Data science can do a lot: recommendations

Netflix: Movies You'll Love

http://www.netflix.com/RecommendationsHome?linkctr=mh2

Chris Volinsky | Your Account | Buy / Redeem

Browse DVDs | Browse Instant | Your Queue | Movies You'll Love | Friends & Community | DVD Sale \$5.99 | Movies, actors, directors, genres

Suggestions (663) | Suggestions by Genre | Rate Movies | Rate Genres | Movies You've Rated (103)

Movies You'll Love

Suggestions based on your ratings

You have 6 suggestions from 103 rated.

### INDEPENDENT SUGGESTIONS (19) [See all 19 >](#)

**Wristcutters: A Love Story**  
Because you enjoyed:  
*Lost in Translation*  
*Garden State*  
*Children of Men*

**DEAD MAN**  
Because you enjoyed:  
*Taxi Driver*  
*Being John Malkovich*  
*Harold and Maude*

**Trainspotting: Collector's Edition**  
Because you enjoyed:  
*Pulp Fiction*  
*Reservoir Dogs*  
*Taxi Driver*

**STRANGER THAN PARADISE**  
Because you enjoyed:  
*Annie Hall*  
*This Is Spinal Tap*  
*Taxi Driver*

### DOCUMENTARY SUGGESTIONS (107) [See all 107 >](#)

**The King of Kong**  
Because you enjoyed:  
*This Is Spinal Tap*  
*Spellbound*  
*Children of Men*

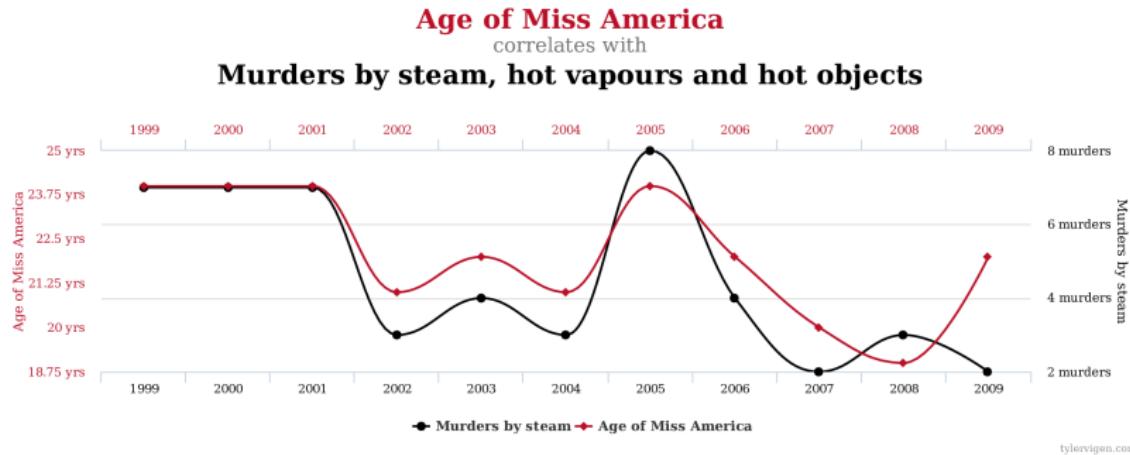
**The Business of Being Born**  
Because you enjoyed:  
*Life Is Beautiful*  
*Spellbound*  
*Super Size Me*

**Jimmy Carter: Man from Plains**  
Because you enjoyed:  
*Annie Hall*  
*Being John Malkovich*  
*Lost in Translation*

**Lake of Fire**  
Because you enjoyed:  
*Annie Hall*  
*Fargo*  
*The Graduate*

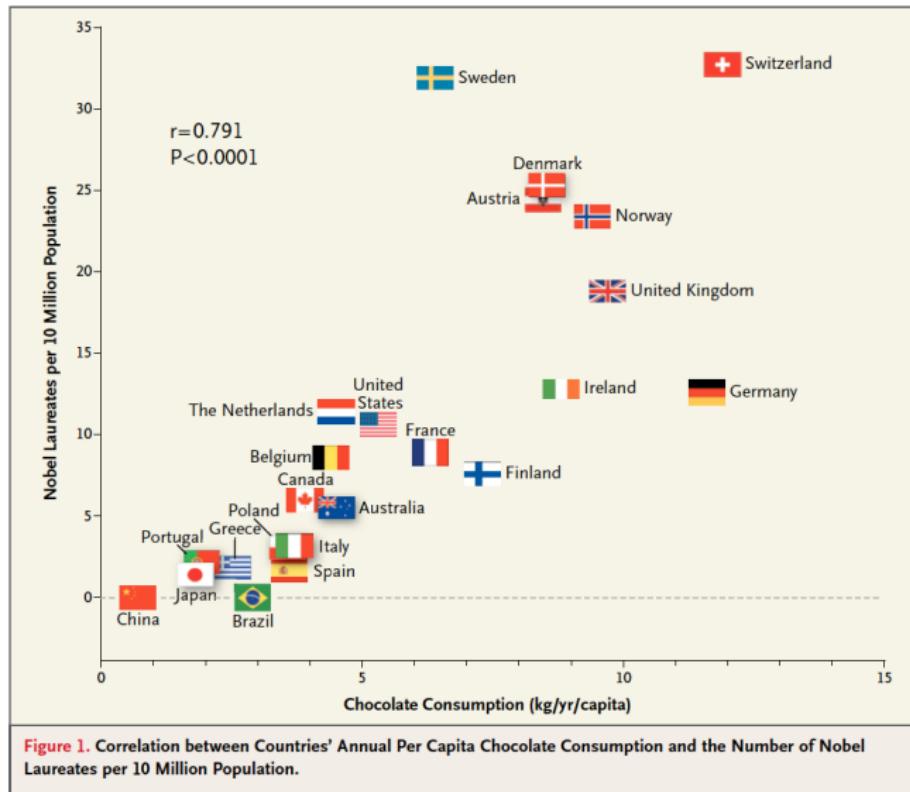
Find:  Next  Highlight all Match case Phrase not found

# Is data all we need?



- ▶ It is easy to find “interesting” patterns where none exist! (image)
- ▶ How should we judge whether a “pattern” is interesting?
- ▶ When should we worry about falsely labeling patterns “interesting”? (E.g. Google mistranslates a sentence vs. incorrect cancer diagnosis... )

# Is data all we need?



Chocolate Consumption, Cognitive Function, and Nobel Laureates

# Is data all we need?



Kari Käasper  
@karikas

Follow

Google translate gender bias also present when translating from English to Estonian and back. How could this be averted?

The screenshot shows two Google Translate windows side-by-side. The top window has English as the source language and Estonian as the target language. It translates "He is a babysitter." to "Ta on lapsehoidja." and "She is a lawyer." to "Ta on advokaat." The bottom window has Estonian as the source language and English as the target language. It translates "Ta on lapsehoidja." back to "She's a babysitter." and "Ta on advokaat." back to "He is a lawyer." This demonstrates that the machine learning algorithm used by Google Translate preserves gender bias when translating between English and Estonian.

5:51 AM - 5 Oct 2017

9 Retweets 32 Likes



25

Blindly used, machine learning algorithms can reinforce “hidden” biases.

## Questions From the Pandemic

- To mask or not to mask?
- Should young children get vaccinated?
- Hydroxychloroquine? Ivermectin? Molnupiravir?
- Long covid?

## Learn how to use data

- ▶ **Explore:** identify patterns in data
- ▶ **Predict:** make informed guesses about the future, or about subjects we haven't gathered data on yet
- ▶ **Infer:** quantify what you know/understand the limits of what you can know based on the data you have