

# Lecture 2: Introduction

Core 109S IDWT?, Spring 2017

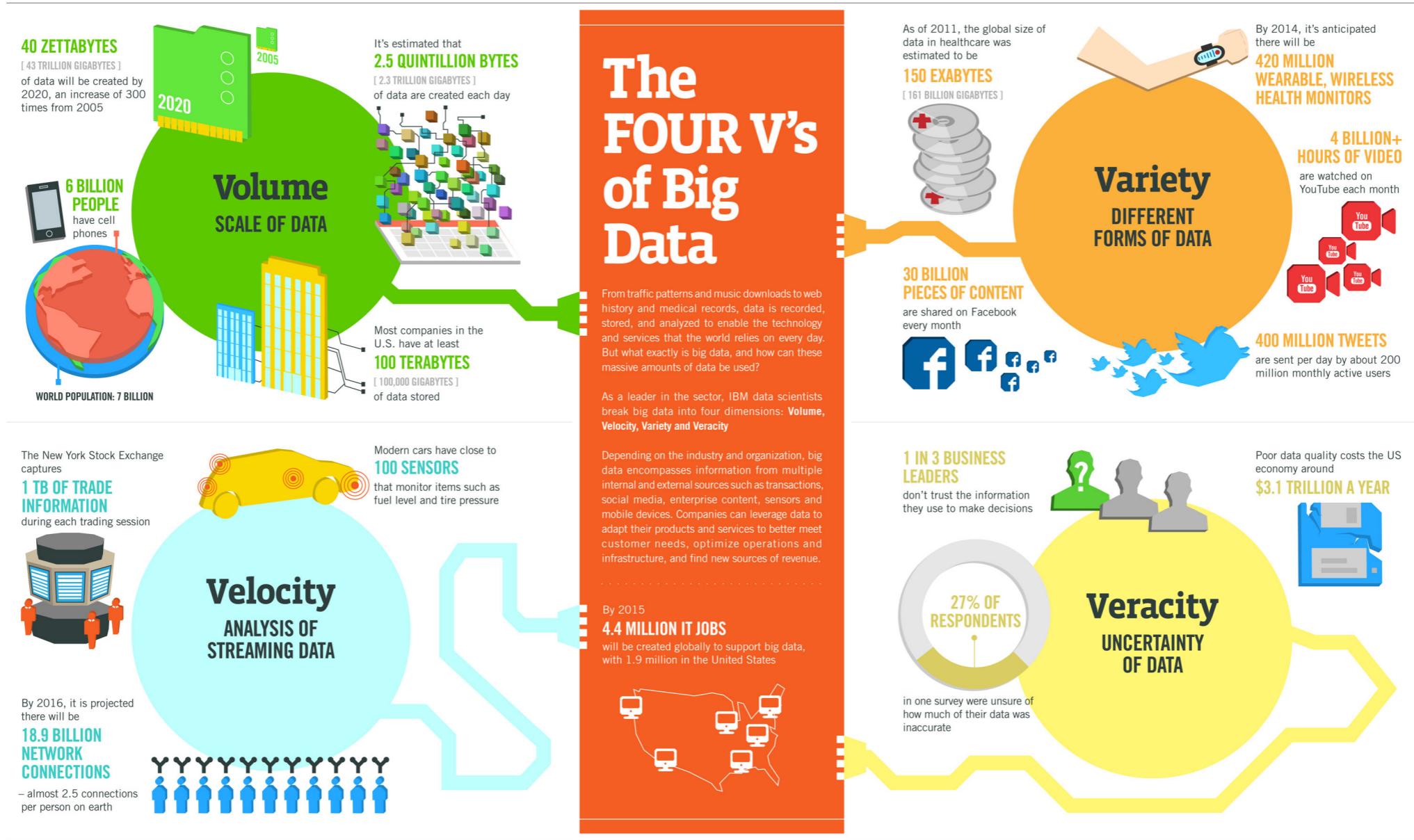
Michael Hay

# Reading Engagement

## Exercise #1



# Big Data

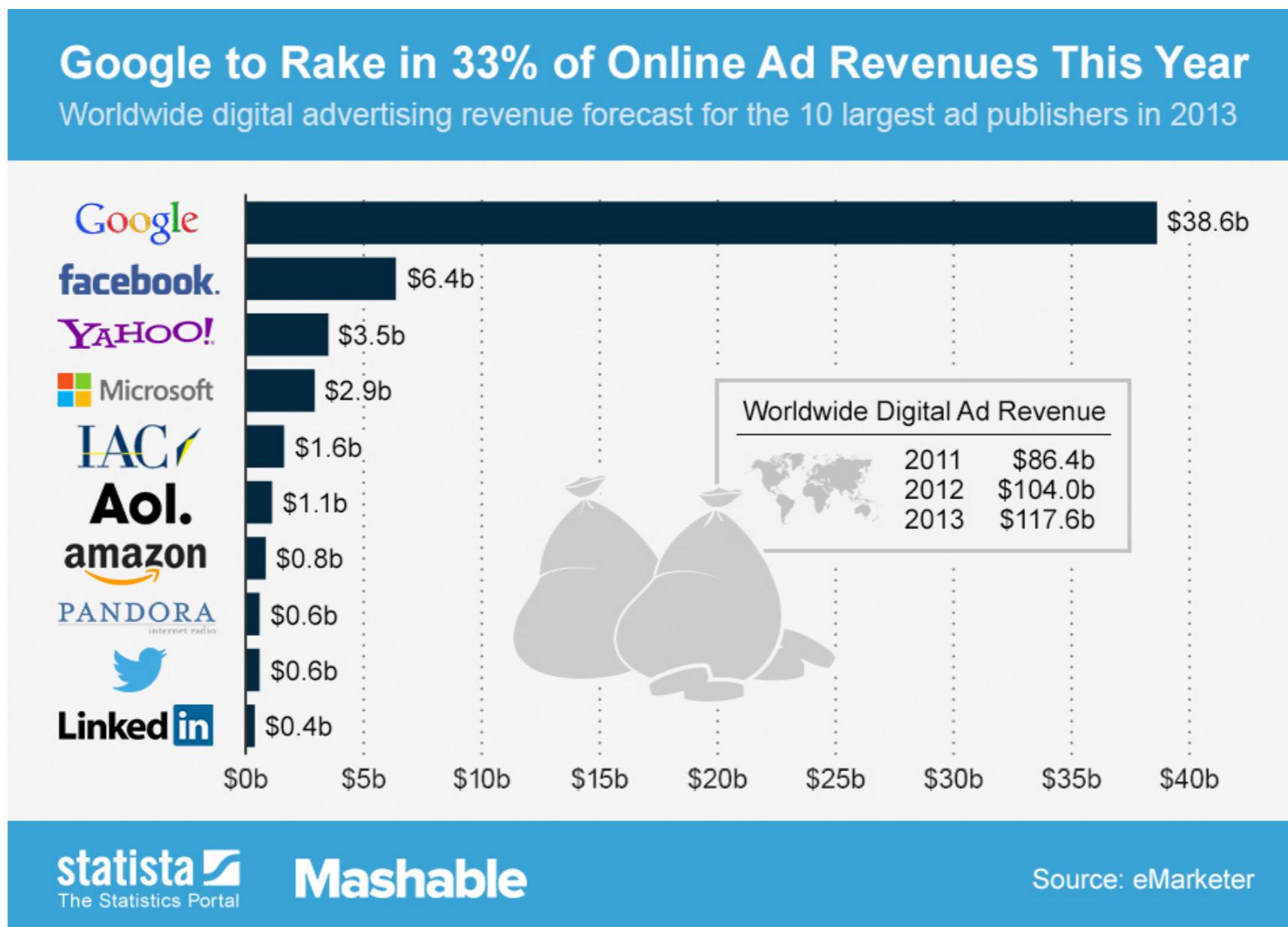


<http://www.ibmbigdatahub.com/infographic/four-vs-big-data>



# Applications of “Big Data”

# Let's talk about \$\$\$



Ad revenue is ≈ 89% of Google's revenue (as of 2016)

# Data and online content



Recommended links

**+79% clicks**  
vs. randomly selected

News Interests

**+250% clicks**  
vs. one size fits all

Top Searches

**+43% clicks**  
vs. editor selected

Raghu Ramakrishnan,  
Yahoo! Research  
NSF Workshop on Social  
Networks and Mobility in  
the Cloud 2012

# Data and commerce



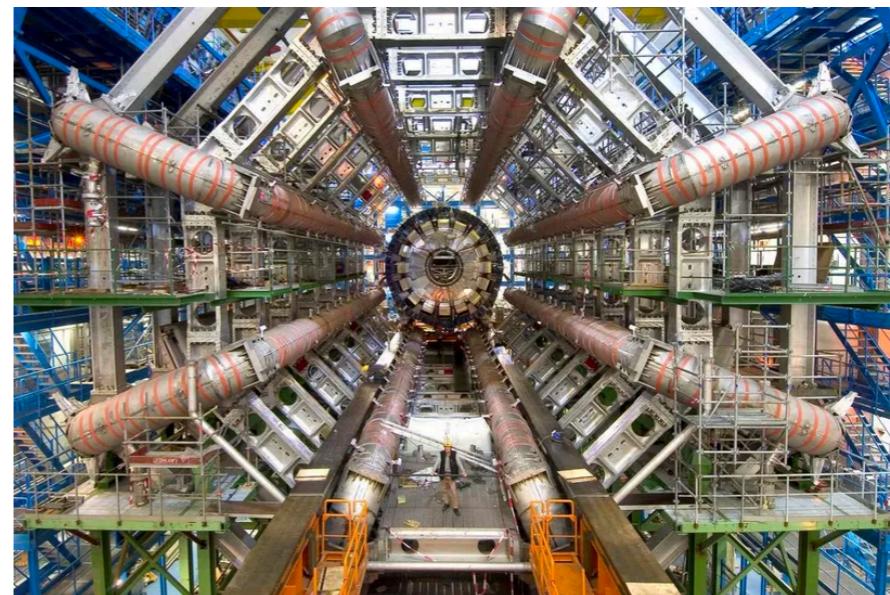
Target analyzed purchases to identify women in early-stage pregnancy. Exposed teen daughter's pregnancy.

# Data and science

- The world's largest particle collider at CERN—where the Higgs boson was confirmed—generates 30 petabytes of data per year



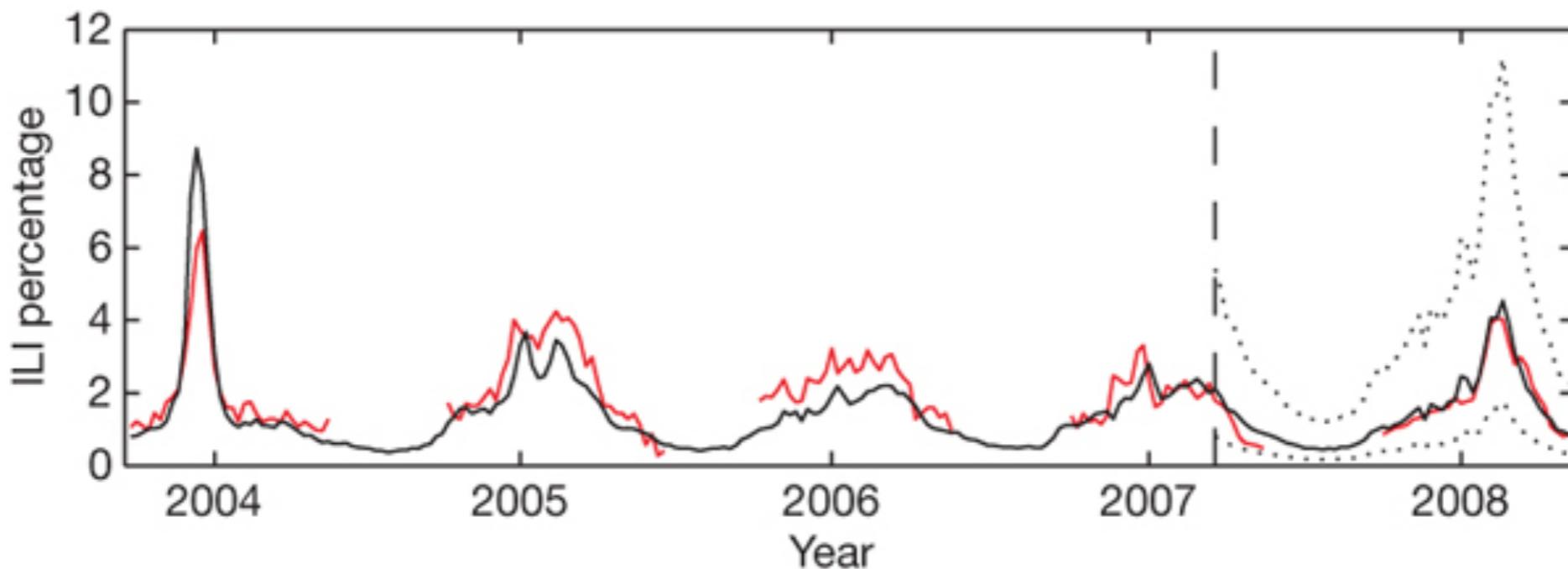
<http://home.web.cern.ch/about/computing>



[http://www.theverge.com/2016/4/25/11501078/  
cern-300-tb-lhc-data-open-access](http://www.theverge.com/2016/4/25/11501078/cern-300-tb-lhc-data-open-access)

- CERN's data center has 11,000 servers with 100,000 cores... yet it still can't crunch all data!

# Data and health

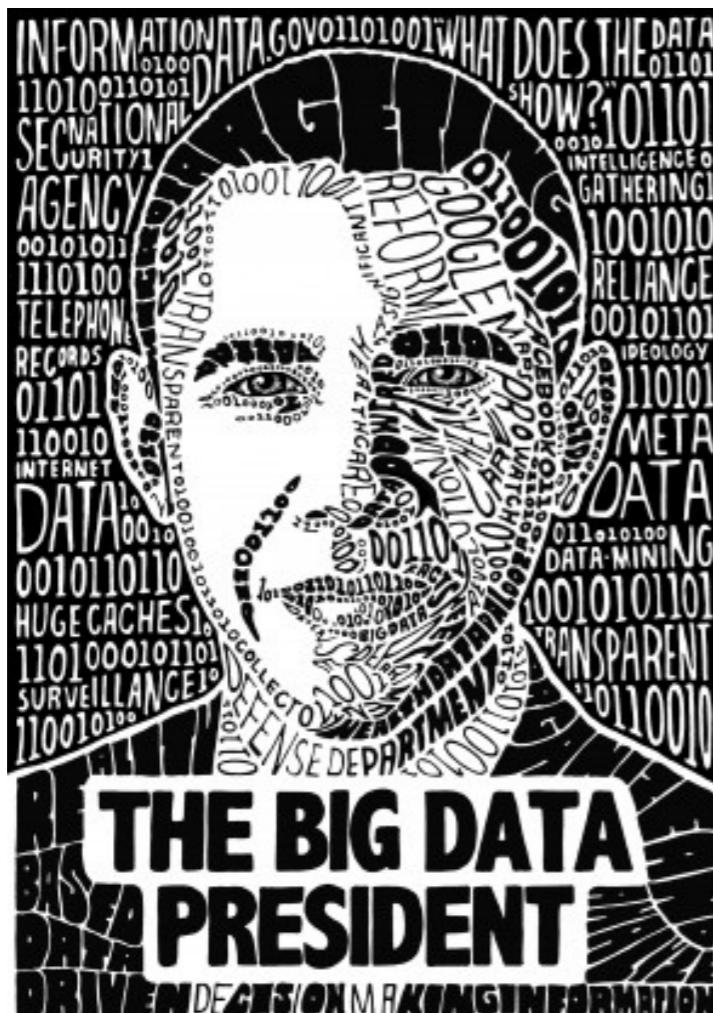


**Red:** official numbers from Center for Disease Control and Prevention; weekly  
**Black:** based on Google search logs; daily (potentially instantaneously)

## Detecting influenza epidemics using search engine query data

<http://www.nature.com/nature/journal/v457/n7232/full/nature07634.html>

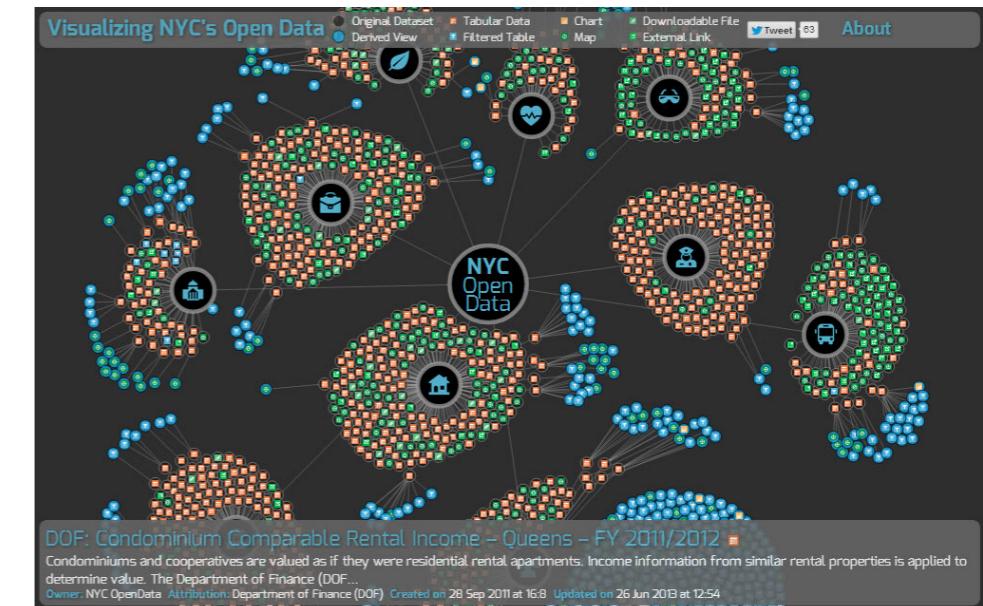
# Data and government



[http://www.washingtonpost.com/opinions/obama-the-big-data-president/2013/06/14/1d71fe2e-d391-11e2-b05f-3ea3f0e7bb5a\\_story.html](http://www.washingtonpost.com/opinions/obama-the-big-data-president/2013/06/14/1d71fe2e-d391-11e2-b05f-3ea3f0e7bb5a_story.html)



<http://www.whitehouse.gov/blog/Democratizing-Data>



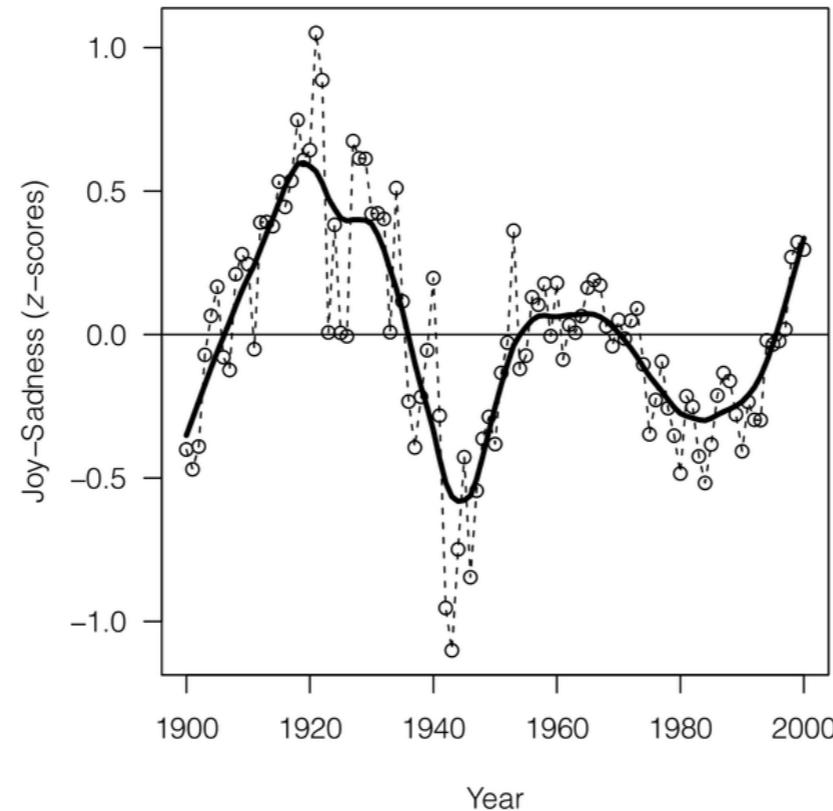
<https://nycopendata.socrata.com/>

# Data and culture (“culturomics”)



## Quantitative Analysis of Culture Using Millions of Digitized Books

<http://science.sciencemag.org/content/331/6014/176>

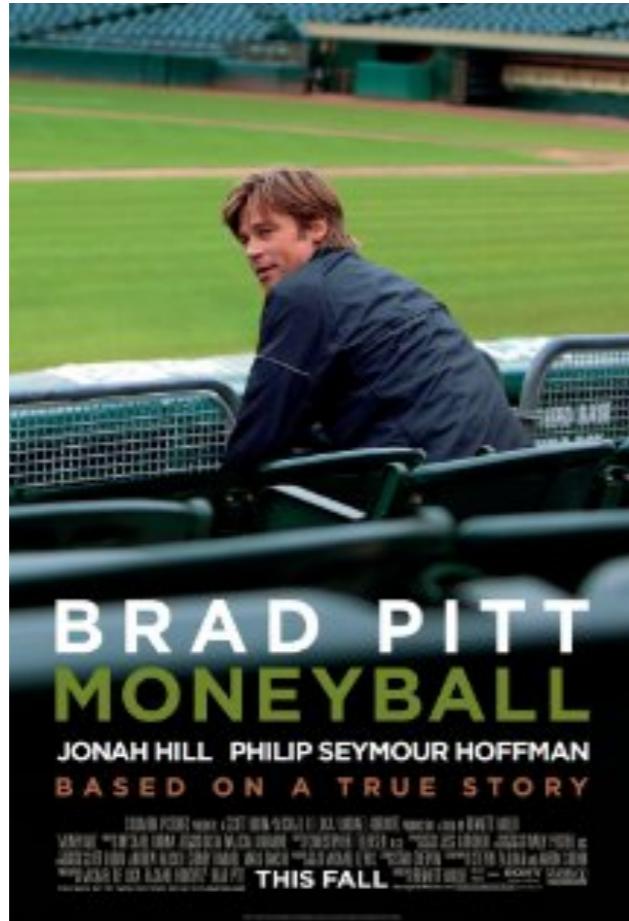


- Frequencies of emotion words in English-language books in Google's database

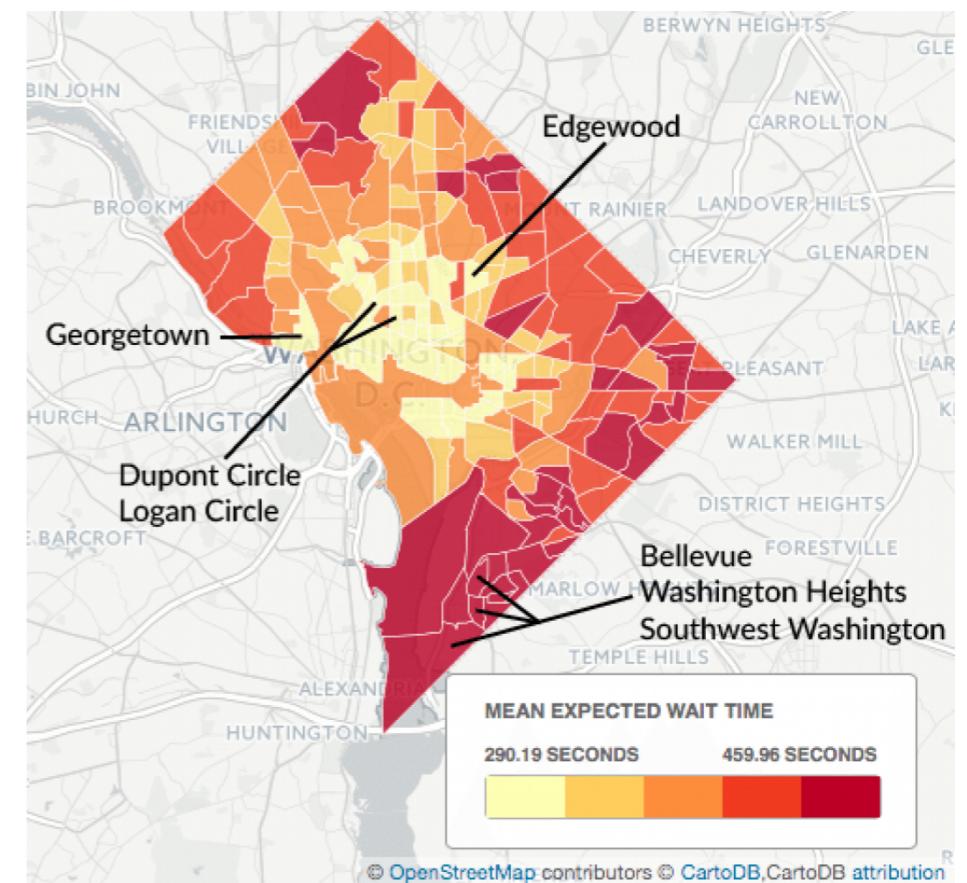
<http://blogs.plos.org/everyone/2013/03/20/what-are-you-in-the-mood-for-emotional-trends-in-20th-century-books/>

# Data and           ← your favorite subject

## Sports



## Journalism



<https://www.washingtonpost.com/news/wonk/wp/2016/03/10/uber-seems-to-offer-better-service-in-areas-with-more-white-people-that-raises-some-tough-questions/>

# Hal Varian, Chief Economist at Google

- “I keep saying ***the sexy job in the next ten years will be statisticians.***

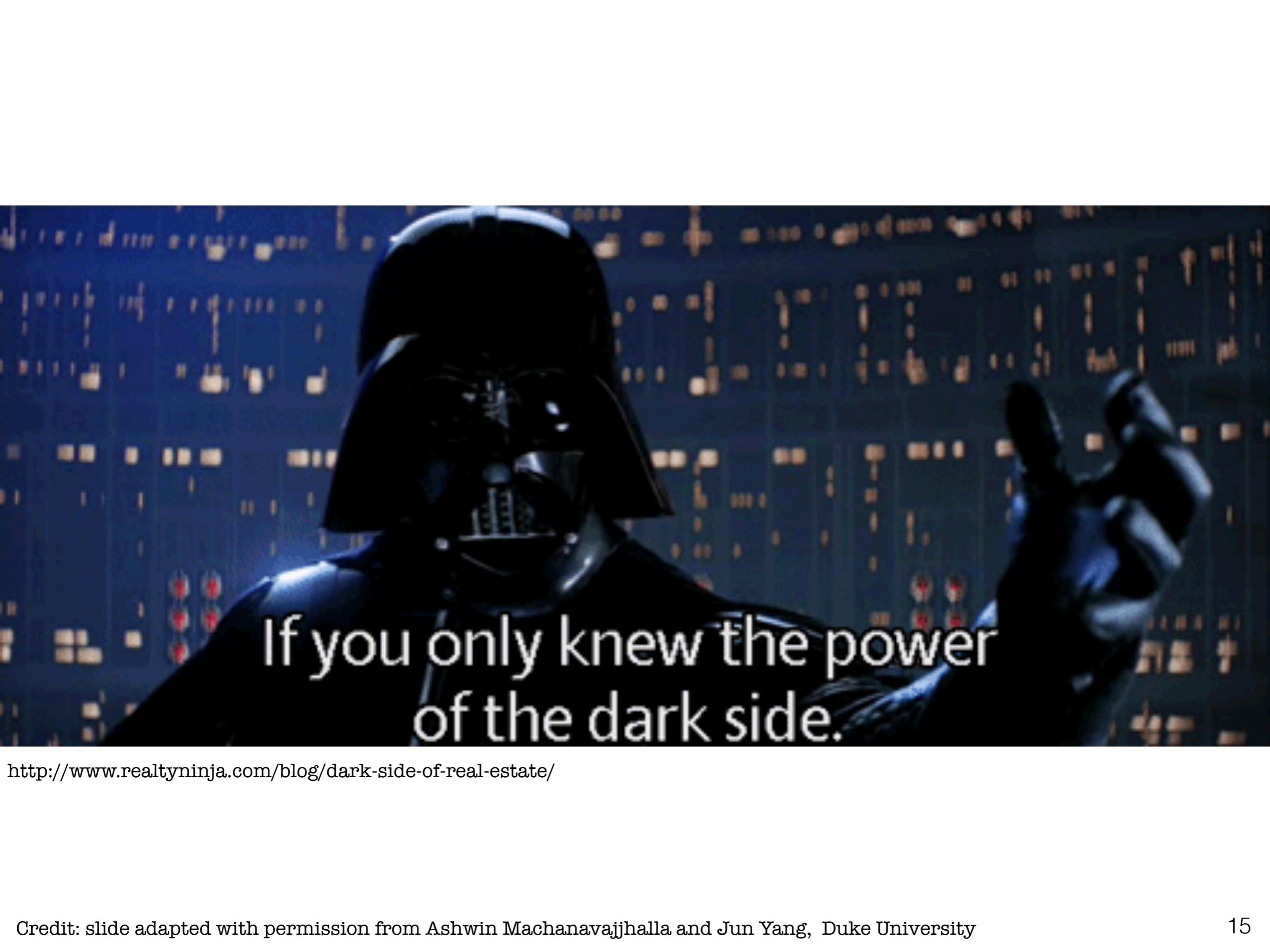
People think I’m joking, but who would’ve guessed that computer engineers would’ve been the sexy job of the 1990s? The ability ***to take data***—to be able to understand it, to process it, ***to extract value from it***, to visualize it, to communicate it—that’s going to be a hugely important skill in the next decades...” (2009)



<http://www.mckinsey.com/industries/high-tech/our-insights/hal-varian-on-how-the-web-challenges-managers>

# Opportunities & challenges of big data

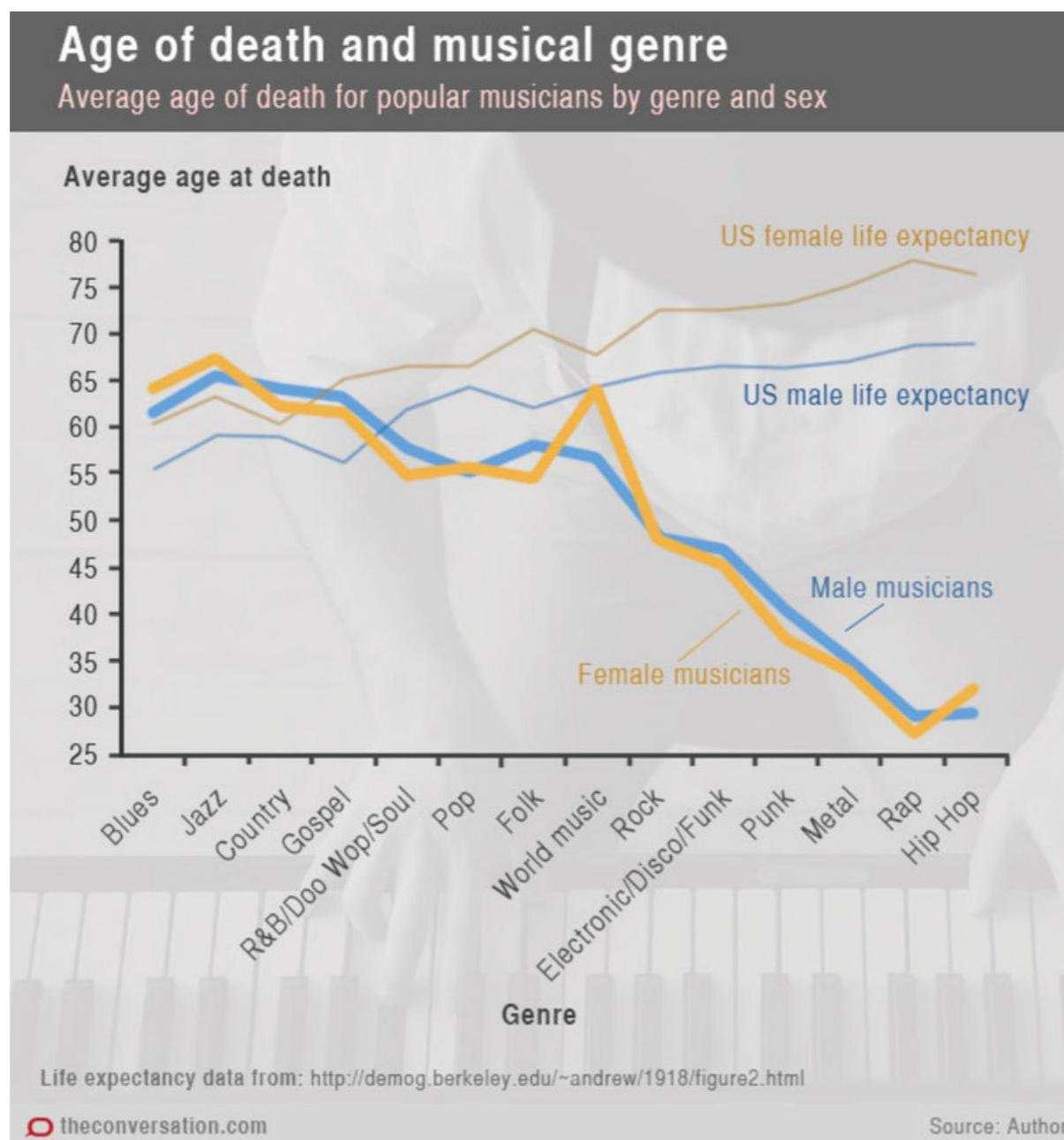
- **Big**
  - Quantitative shift → qualitative difference
  - Technical: *How to process large volumes of data efficiently?*
- **Messy**
  - Data exhaust. Measurement error, ambiguity, bias. Combining disparate data sources.
  - Technical: *How to link data from different sources*
- **Correlation**
  - A shift away from causal understanding
  - Technical: *How to find correlations automatically; how ensure pattern is real and not spurious?*
- **Societal implications**
  - Explore the dark side of big data...
  - Technical: *Explore some technological responses to social issues (e.g., privacy-preserving data analysis)*



If you only knew the power  
of the dark side.

<http://www.realtyninja.com/blog/dark-side-of-real-estate/>

# Easy to get it wrong...



**Data censoring  
(and other statistical traps)**

# Easy to get it wrong...

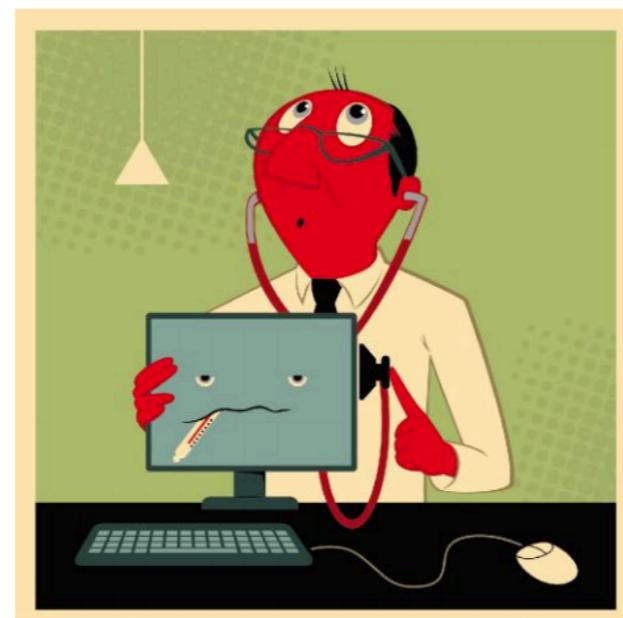
BIG DATA

## The Parable of Google Flu: Traps in Big Data Analysis

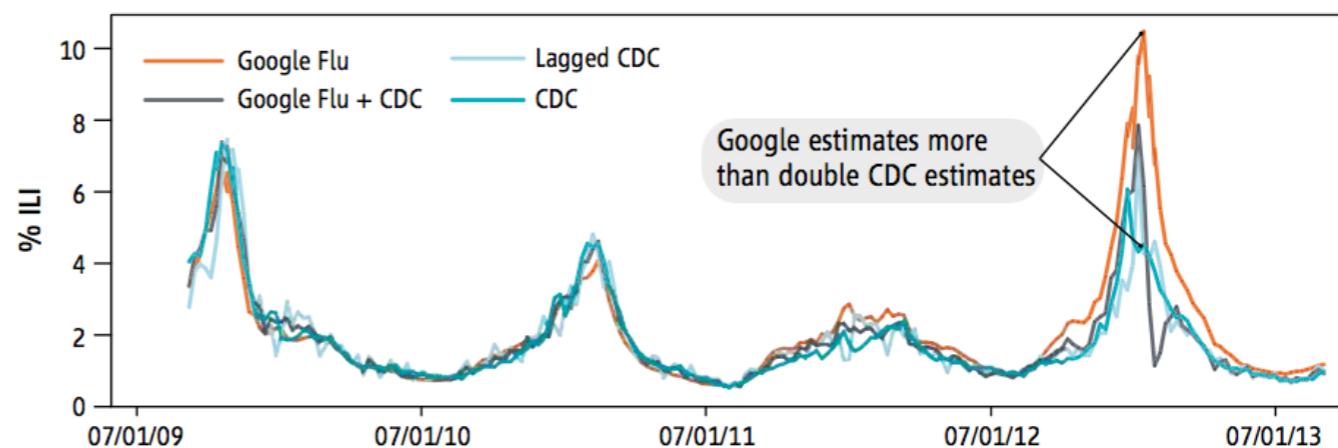
David Lazer,<sup>1,2\*</sup> Ryan Kennedy,<sup>1,3,4</sup> Gary King,<sup>3</sup> Alessandro Vespiagnani<sup>3,5,6</sup>

In February 2013, Google Flu Trends (GFT) made headlines but not for a reason that Google executives or the creators of the flu tracking system would have hoped. *Nature* reported that GFT was predicting more than double the proportion of doctor visits for influenza-like illness (ILI) than the Centers for Disease Control and Prevention (CDC), which bases its estimates on surveillance reports from laboratories across the United States (1, 2). This happened despite the fact that GFT was built to predict CDC reports. Given that GFT is often held up as an exemplary use of big data (3, 4), what lessons can we draw from this error?

The problems we identify are not limited to GFT. Research on whether search or social media can

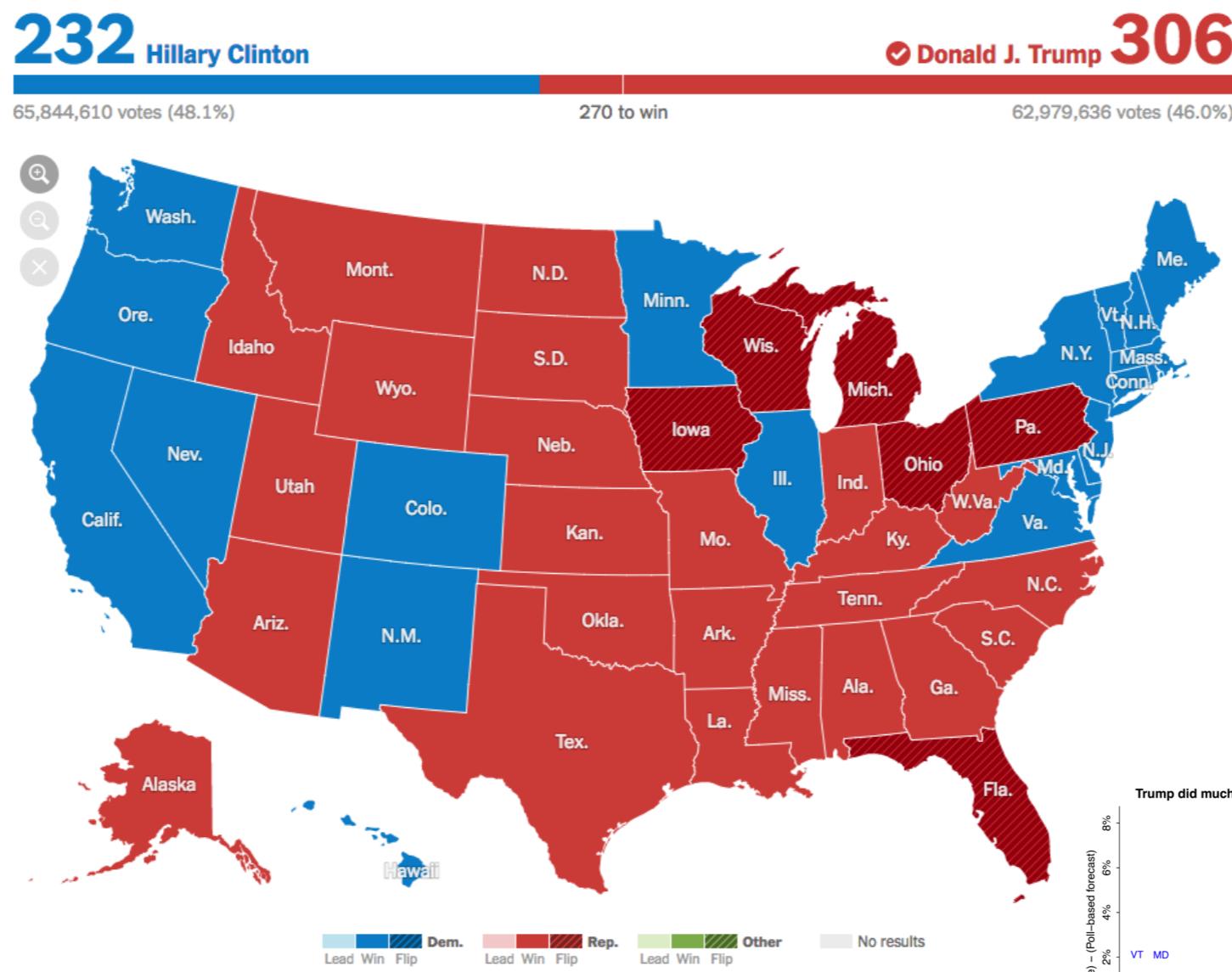


Large errors in flu prediction were largely avoidable, which offers lessons for the use of big data.

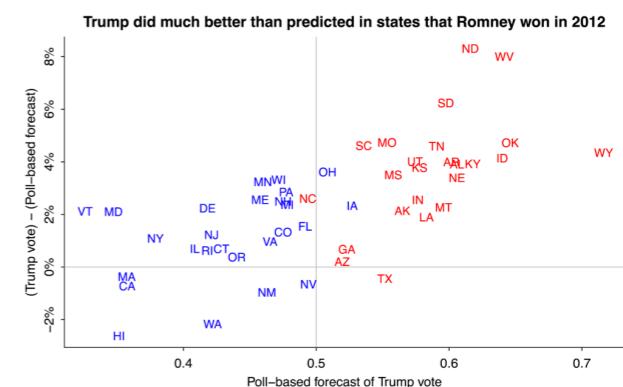


<http://science.sciencemag.org/content/343/6176/1203>

Easy to get it wrong...



<http://www.nytimes.com/elections/results/president>



<http://andrewgelman.com/2016/11/09/explanations-shocking-2-shift/>

# Easy to abuse... ethics



## Experimental evidence of massive-scale emotional contagion through social networks

Adam D. I. Kramer<sup>a,1</sup>, Jamie E. Guillory<sup>b</sup>, and Jeffrey T. Hancock<sup>c,d</sup>

<sup>a</sup>Core Data Science Team, Facebook, Inc., Menlo Park, CA 94025; <sup>b</sup>Center for Tobacco Control Research and Education, University of California, San Francisco, CA 94143; and Departments of <sup>c</sup>Communication and <sup>d</sup>Information Science, Cornell University, Ithaca, NY 14853

Edited by Susan T. Fiske, Princeton University, Princeton, NJ, and approved March 25, 2014 (received for review October 23, 2013)

Emotional states can be transferred to others via emotional contagion, leading people to experience the same emotions without their awareness. Emotional contagion is well established in laboratory experiments, with people transferring positive and negative emotions to others. Data from a large real-world social network, collected over a 20-y period suggests that longer-lasting moods (e.g., depression, happiness) can be transferred through networks [Fowler JH, Christakis NA (2008) *BMJ* 337:a2338], although the results are controversial. In an experiment with people who use Facebook, we test whether emotional contagion occurs outside of in-person interaction between individuals by reducing the amount of emotional content in the News Feed. When positive expressions were reduced, people produced fewer positive posts and more negative posts; when negative expressions were re-

demonstrated that (i) emotional contagion occurs via text-based computer-mediated communication (7); (ii) contagion of psychological and physiological qualities has been suggested based on correlational data for social networks generally (7, 8); and (iii) people's emotional expressions on Facebook predict friends' emotional expressions, even days later (7) (although some shared experiences may in fact last several days). To date, however, there is no experimental evidence that emotions or moods are contagious in the absence of direct interaction between experiencer and target.

On Facebook, people frequently express emotions, which are later seen by their friends via Facebook's "News Feed" product (8). Because people's friends frequently produce much more content than one person can view, the News Feed filters posts, stories, and activities undertaken by friends. News Feed is the

<http://www.pnas.org/content/111/24/8788.full>

## Facebook emotion study breached ethical guidelines, researchers say

Lack of 'informed consent' means that Facebook experiment on nearly 700,000 news feeds broke rules on tests on human subjects, say scientists

<https://www.theguardian.com/technology/2014/jun/30/facebook-emotion-study-breached-ethical-guidelines-researchers-say>

# Easy to abuse... privacy

## A Face Is Exposed for AOL Searcher No. 4417749

By MICHAEL BARBARO and TOM ZELLER Jr.  
Published: August 9, 2006

 SIGN IN TO E-  
THIS



## Why 'Anonymous' Data Sometimes Isn't

By Bruce Schneier  12.13.07

Last year, Netflix published 10 million movie rankings by 500,000 customers, as part of a challenge for people to come up with better recommendation systems than the one the company was using.

The Scientist » The Nutshell

## “Anonymous” Genomes Identified

The names and addresses of people participating in the Personal Genome Project can be easily tracked down despite such data being left off their online profiles.

By Dan Cossins | May 3, 2013



# Easy to abuse... liberty



[http://www.theguardian.com/world/2013/jun/23/  
edward-snowden-nsa-files-timeline](http://www.theguardian.com/world/2013/jun/23/edward-snowden-nsa-files-timeline)

# Easy to misuse... bias



Bernard Parker, left, was rated high risk; Dylan Fugett was rated low risk. (Josh Ritchie for ProPublica)

## Machine Bias

There's software used across the country to predict future criminals.  
And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica  
May 23, 2016

<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

# 39% of experts agree...

*Thanks to many changes, including the building of “the Internet of Things,” human and machine analysis of **Big Data will cause more problems than it solves** by 2020. The existence of huge data sets for analysis will **engender false confidence in our predictive powers** and will lead many to make **significant and hurtful mistakes**. Moreover, analysis of Big Data will be **misused by powerful people and institutions with selfish agendas** who manipulate findings to make the case for what they want. And the advent of Big Data has a harmful impact because it **serves the majority (at times inaccurately) while diminishing the minority** and ignoring important outliers. Overall, the rise of **Big Data is a big negative for society in nearly all respects.***

2012 Pew Research Center Report

<http://pewinternet.org/Reports/2012/Future-of-Big-Data/Overview.aspx>

# But it's here, now!

- Learn to...
  - Think critically about technology's promise and limitations
  - Help yourself and other avoid being taken advantage of



<http://rosemarynonnyknight.com/use-the-force-your-name-here/>

# Introductions

- Please tell us...
  - your name
  - one hobby/interest/extra-curricular activity
  - possible major
  - experience with computer science / programming  
*(no experience needed for this course!)*
  - what are you hoping to get out of this course?

# Core SP

- Core SP classes unified by some common goals
- Colgate wants to assess progress towards those goals
- Two surveys:
  - One now, one at end of semester
  - Not graded, but please take seriously