

# CSCI 1070: Taming Big Data

Prof. Abby Stylianou

# Logistics

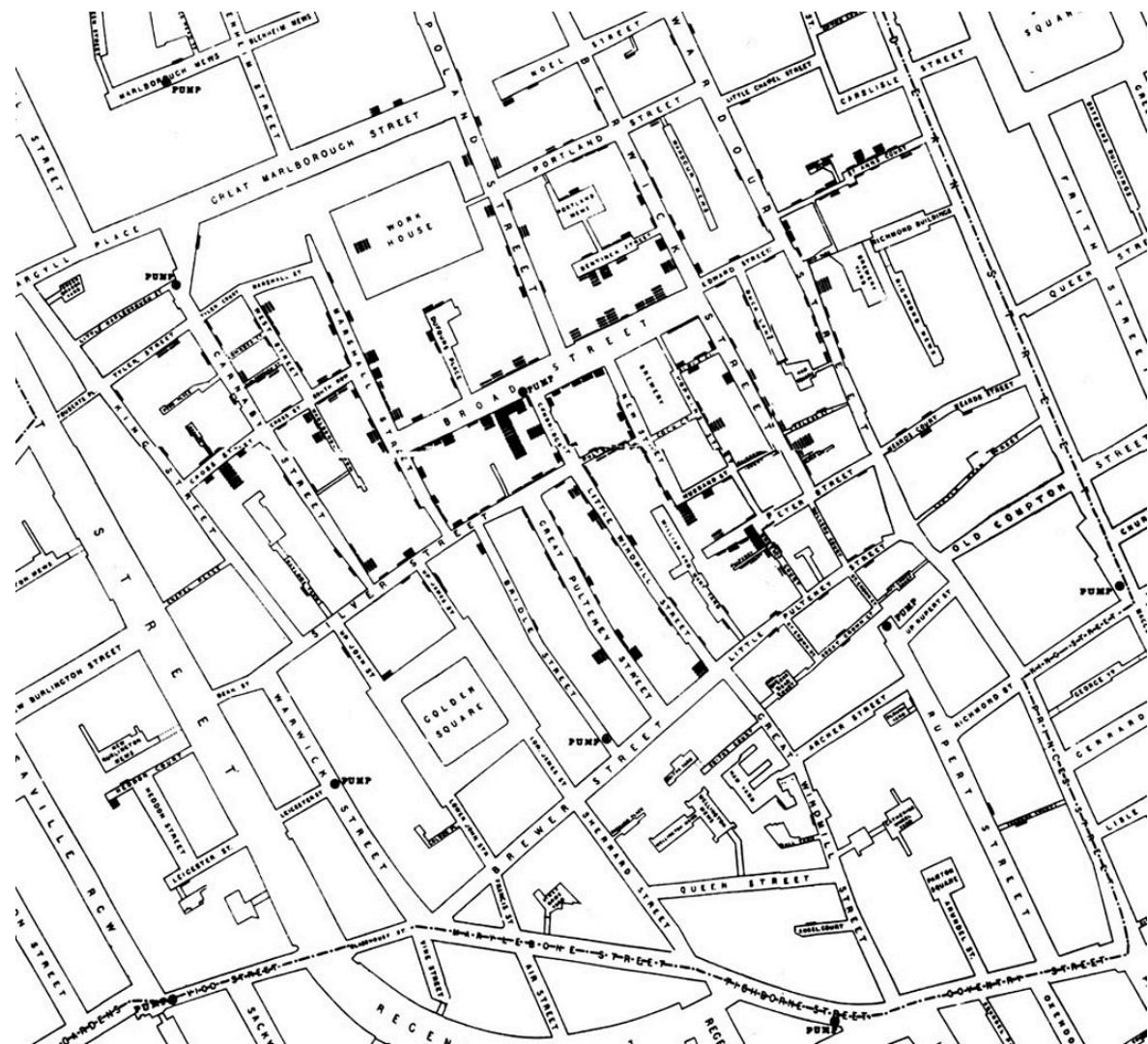
- Are you on the waitlist? Talk to me after class!
- Class meetings: Tu/Th, 11AM–12:15PM, Ritter 115
  - Lecture, discussions, in class coding
  - Occasional quizzes/journaling about readings
  - Attendance not mandatory, but you're gonna have a bad time if you don't show up.
- Office Hours: Wed, 1:30-2:30PM, Ritter 107
- Class communications: Piazza! (*not email*)
- Website/schedule/materials: <https://cs.slu.edu/~stylianou/1070/>
- Technologies: Python, Jupyter Notebooks, Git

# What is big data?

- **1663:** John Graunt is the first person credited w/ statistical data analysis in his studies of the bubonic plague in Europe, dealing with what he referred to as “overwhelming amounts of information”

# What is big data?

- **1854:** John Snow maps London Cholera outbreaks and finds that they are clustered around a single pump; it was found that a cesspit was leaking into that pump



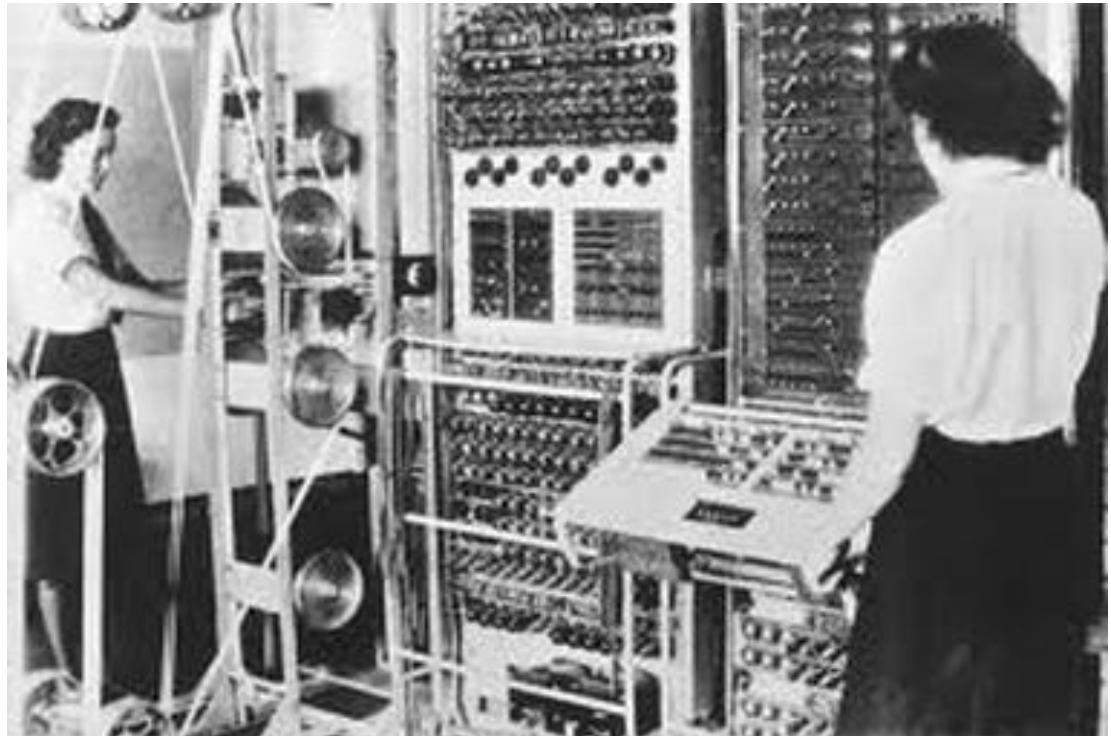
# What is big data?

- **1880:** US Census Bureau estimates it will take eight years to process the data collected in the 1880 census, and over 10 years to process 1890 census data
- **Hollerith Tabulating Machine**  
(punch card tabulation) reduces to ~3 months



# What is big data?

- **WW2:** British invent the Colossus machine to scan for patterns in intercepted Nazi codes. Scans 5,000 characters a second, reducing workload from weeks to hours



# What is big data?

- 1944: Librarian Fremont Rider @ Wesleyan estimates American university libraries doubling in size every 16 years:

*“the Yale Library in 2040 will have approximately 200,000,000 volumes,  
which will occupy over 6,000 miles of shelves...  
[requiring] a cataloging staff of over six thousand persons”*

# What is big data?

- 1961: Derek Price shows # of new scientific journals growing exponentially rather than linearly

“each advance generates a new series of advances at a reasonably constant birth rate, so that the number of births is strictly proportional to the size of the population of discoveries at any given time”

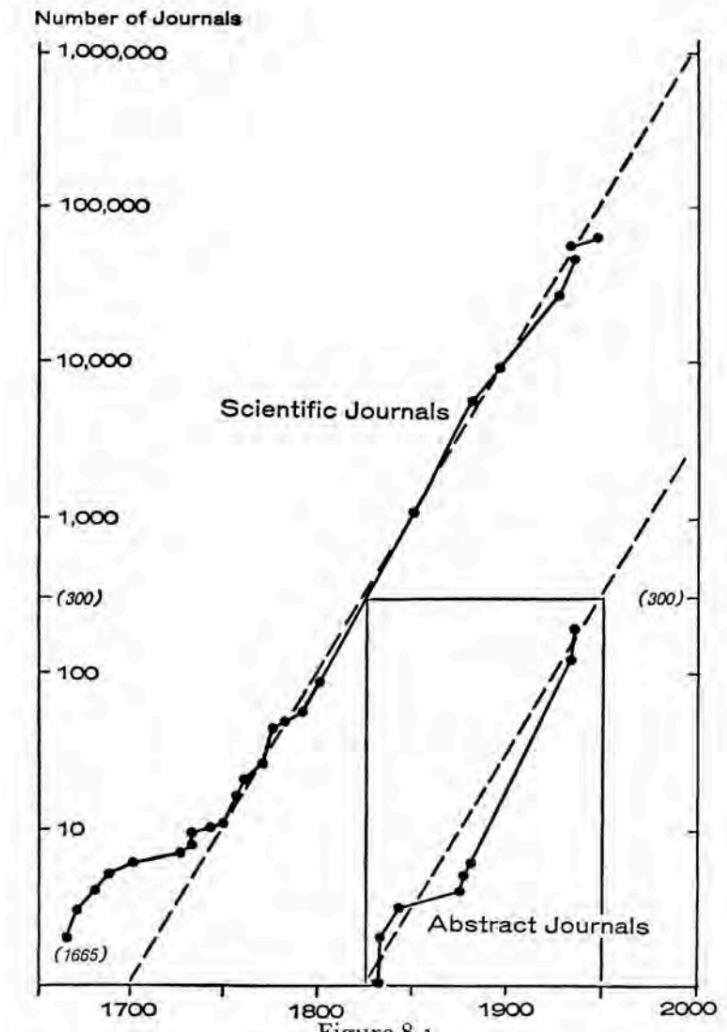


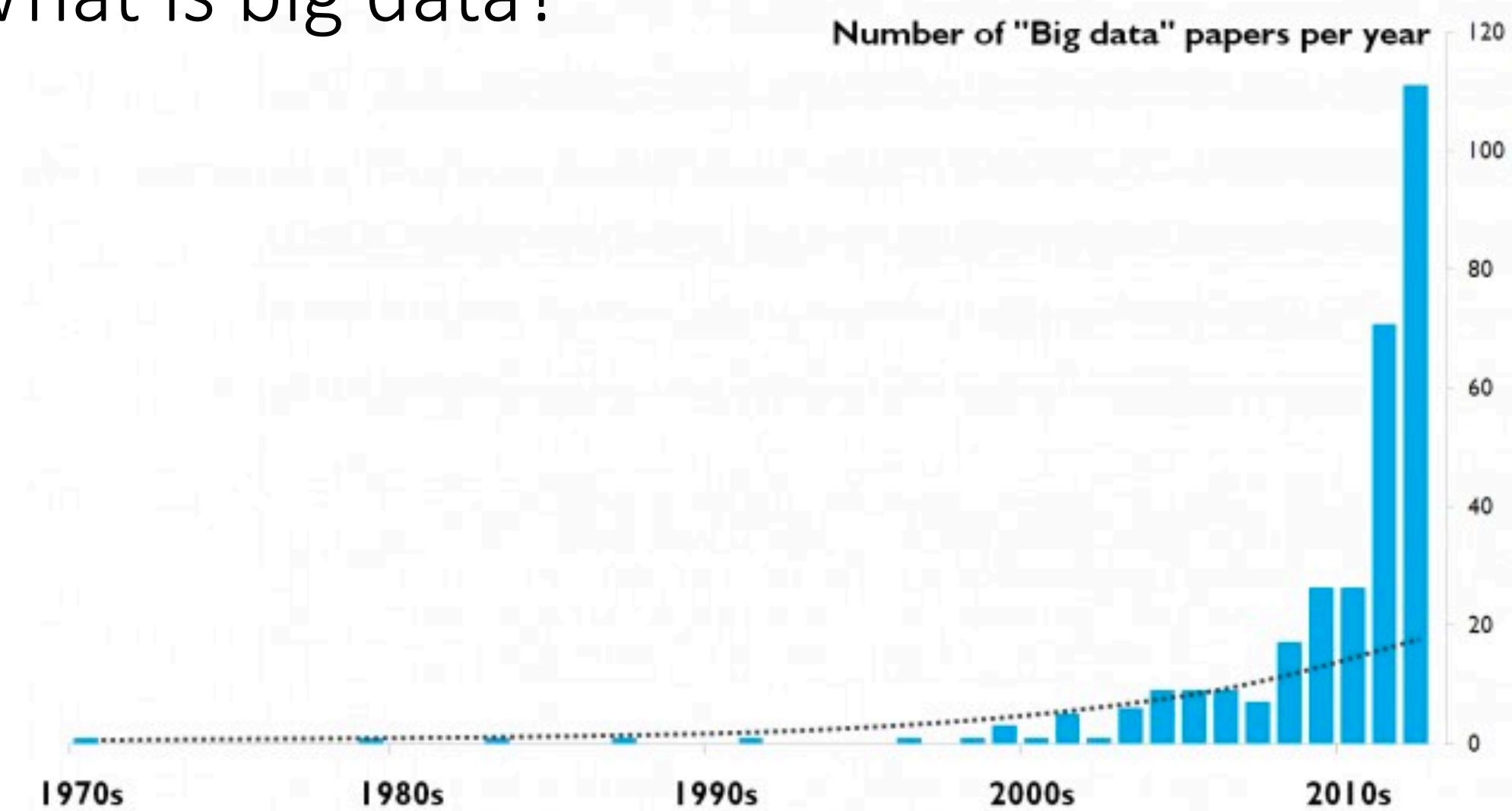
Figure 8.1

Number of journals founded (*not surviving*) as a function of date. The two uppermost points are taken from a slightly differently based list.

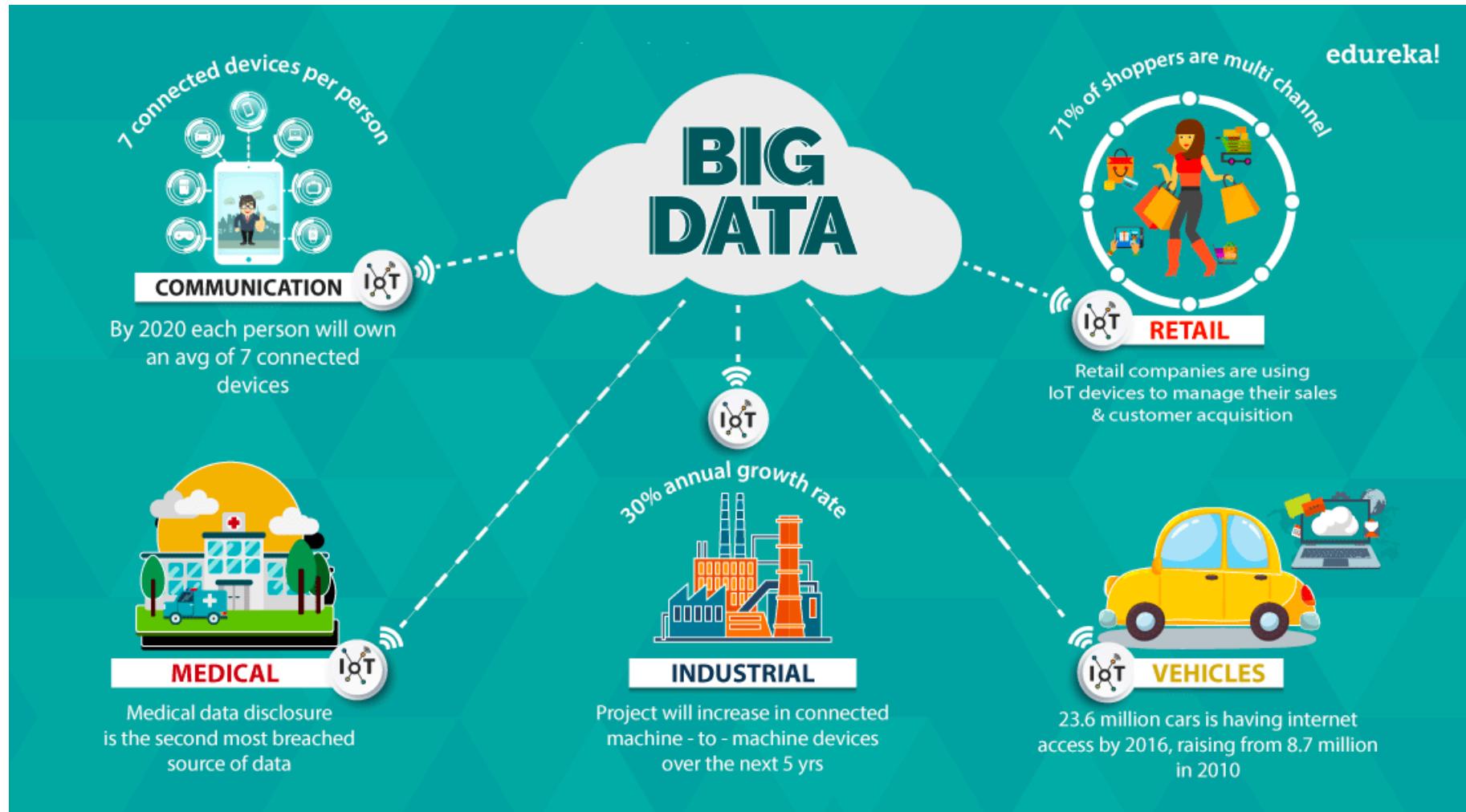
# What is big data?

- **1969:** ARPANET kicks off the Internet
- **1977:** First personal computers come on the market
- **1989:** Tim Berners-Lee introduces the concept of the World Wide Web and the underlying protocols that support it (HTML, URL, HTTP)
- **1993:** CERN announces WWW will be free for everyone to develop and use
- **1990s– 2000s:** The explosion of the internet
- **2005:** Roger Moughalas coins term ‘Big Data’, referring to the scale of data that is nearly impossible to manage and process w/ available tools

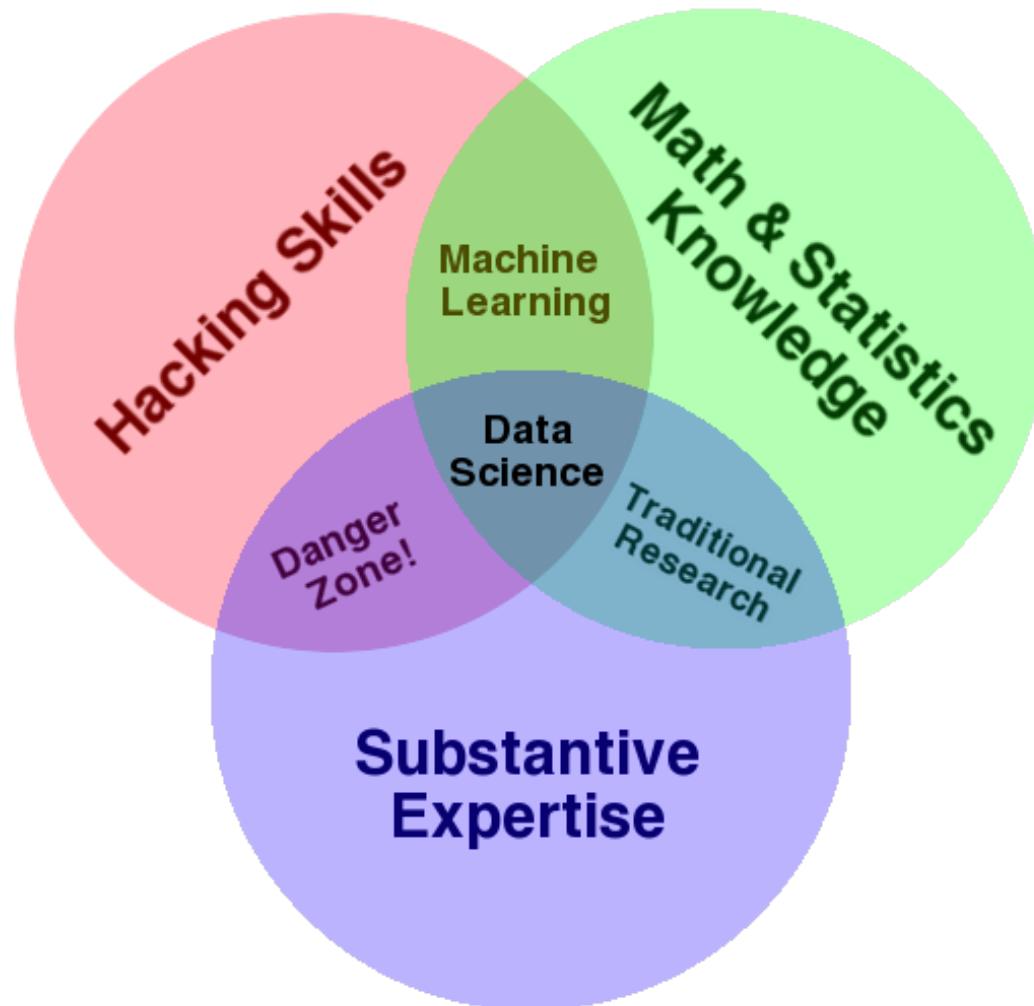
# What is big data?



# Where is big data?



# What is data science?



# The Big Data that I Work With



Google

traffickcam

Q All News Maps Videos Images More Settings Tools

About 459 results (0.17 seconds)

TNW How to help AI find human trafficking victims  
Three years ago an app hit the market called “Traffickcam.” It’s a simple, free app published by Exchange Initiative that lets you upload a picture ...  
Feb 12, 2019

CNN Your hotel room photos could help catch sex traffickers  
That’s where TraffickCam comes in. It’s a simple phone app that uses crowdsourced snapshots of hotel rooms to help law enforcement locate ...  
Mar 20, 2017

TechCrunch You can help stop human trafficking with the TraffickCam app  
In a world where the phrase “oh god, not another app” often springs to mind, along with “Yeah, yeah, I’m sure you want to make a world a better ...  
Jun 25, 2016

Washington Post An incredibly simple way your phone may help save sex trafficked ...  
She said the TraffickCam app, which is available for iOS and Android, isn’t going to solve the problem, but it’s one more tool that could help.  
Jul 1, 2016

Gizmodo Researchers Create Hotel-Recognition System to Aid Human Trafficking Investigations  
... as crowdsourced images from the mobile app TraffickCam, which asks users who are traveling to take photos of their hotel rooms and submit ...  
Feb 5, 2019

snopes.com Taking Pictures of Your Hotel Room Could Help Stop Human Trafficking

# What's the plan for this class?

- Learn the data science tools to collect, clean, understand, manipulate big data
- Exploration, Inference and Prediction

# What's the plan for this class?

- Learn the data science tools to collect, clean, understand, manipulate big data
- Exploration, Inference and Prediction
  - Exploration: Identify patterns in information
    - Tools: Visualization + Descriptive Statistics

# What's the plan for this class?

- Learn the data science tools to collect, clean, understand, manipulate big data
- Exploration, Inference and Prediction
  - Exploration: Identify patterns in information
    - Tools: Visualization + Descriptive Statistics
  - Prediction: make informed guesses about what we want to know, based on the patterns that we identified
    - Tools: Machine Learning + Optimization

# What's the plan for this class?

- Learn the data science tools to collect, clean, understand, manipulate big data
- Exploration, Inference and Prediction
  - Exploration: Identify patterns in information
    - Tools: Visualization + Descriptive Statistics
  - Prediction: make informed guesses about what we want to know, based on the patterns that we identified
    - Tools: Machine Learning + Optimization
  - Inference: quantify our certainty about our predictions
    - Tools: Statistical Tests + Models

# Login to Systems

- Class website: <https://cs.slu.edu/~stylianou/1070>
- JupyterHub: <http://cs1070.com>
  - Login: Your SLUNetID
  - Password: Set it the first time you login
- Piazza: <https://piazza.com/class/jzbep3yp13k3pr>
  - Should have gotten an email to active – let me know if you didn't!
- GitLab: [http://git.cs.slu.edu/courses/fall19/csci\\_1070/<your-SLUNetID>](http://git.cs.slu.edu/courses/fall19/csci_1070/<your-SLUNetID>)
  - Login: SLUNetID
  - Password:
    - either you've done this already and should know it
    - or do a password reset the first time

# Demo

# Next class

- No class on Thursday 8/29 (Mass of the Holy Spirit)
- For next Tuesday, 9/3:
  - Assignment 0 ('Getting to Know You') due via git at beginning of class
  - Do readings that are listed for 8/27 and 8/29
  - We will discuss at the beginning of class