

CSCI 1070: Taming Big Data

Dr. Abby Stylianou

Logistics

- Are you on the waitlist? Talk to me after class!
- Class meetings: Tu/Th, 11AM–12:15PM, Ritter 117
 - Lecture, discussions, in class coding
 - Occasional quizzes about readings
 - Attendance not mandatory, but you're gonna have a bad time if you don't show up
- Office Hours: Tues, 2:30-3:30PM, Ritter 107

Logistics

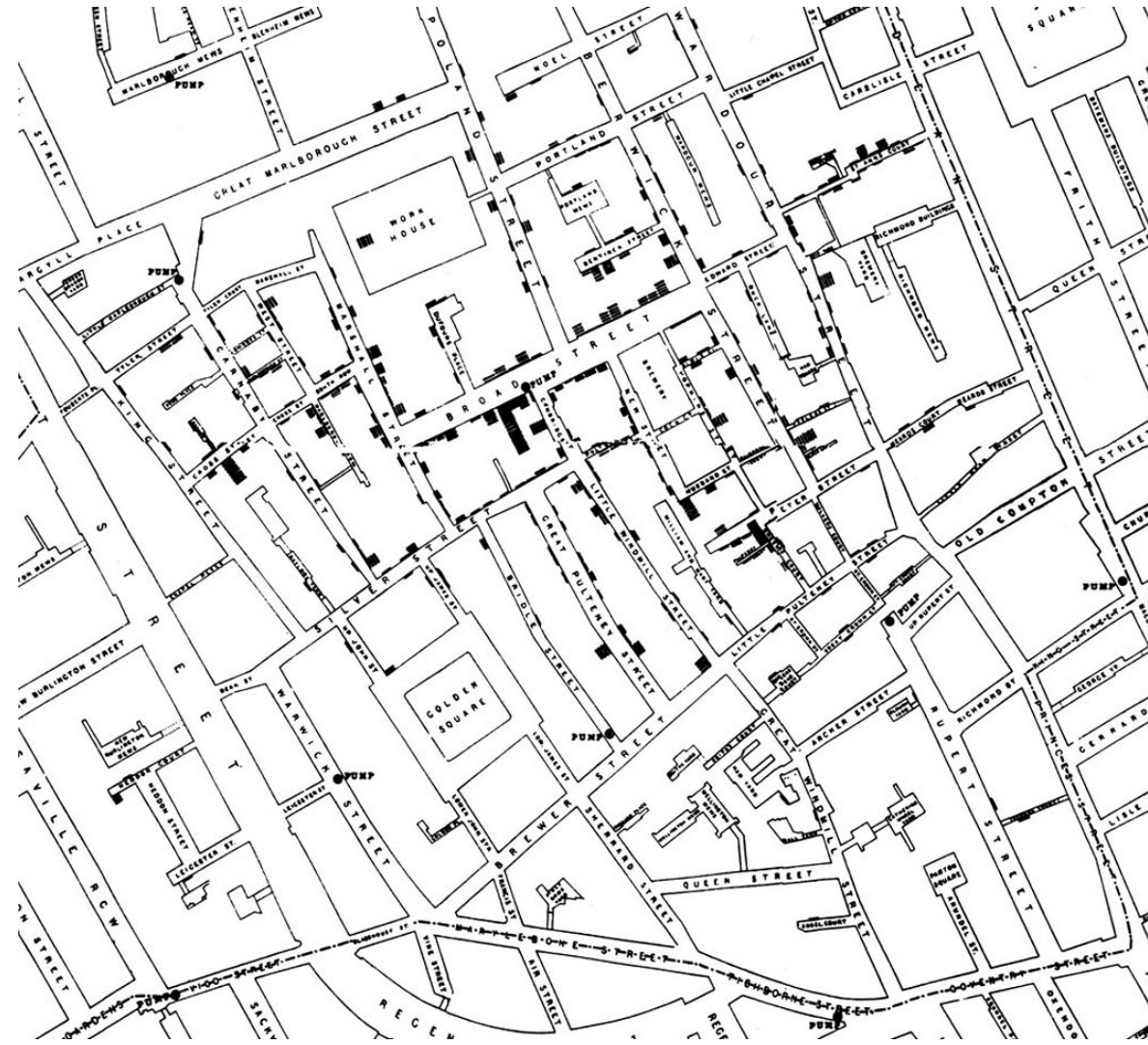
- Website/schedule/materials:
https://cs.slu.edu/~stylianou/1070_sp2020/
- Blackboard for turning in assignments and seeing grades
- Technologies: Python (w/ data science + scientific libraries), Jupyter Notebooks
- Class communications: Piazza! (*not email*)

- **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”

© P

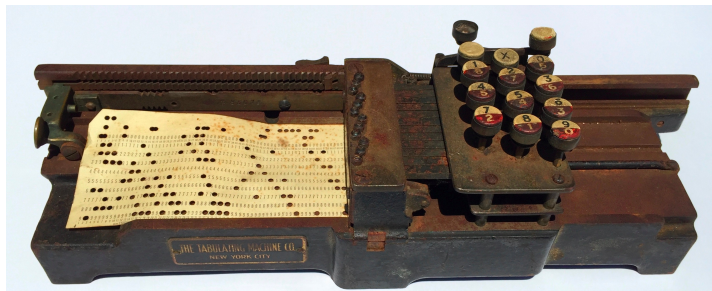
What is big data?

- **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



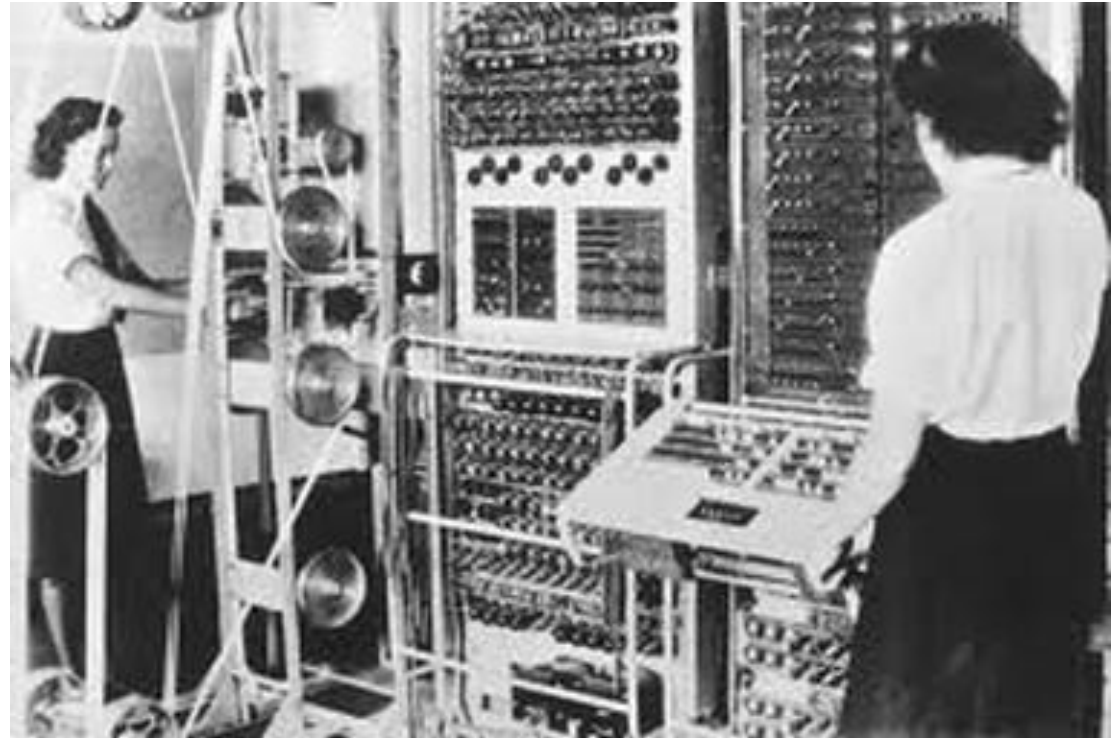
What is big data?

- **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?

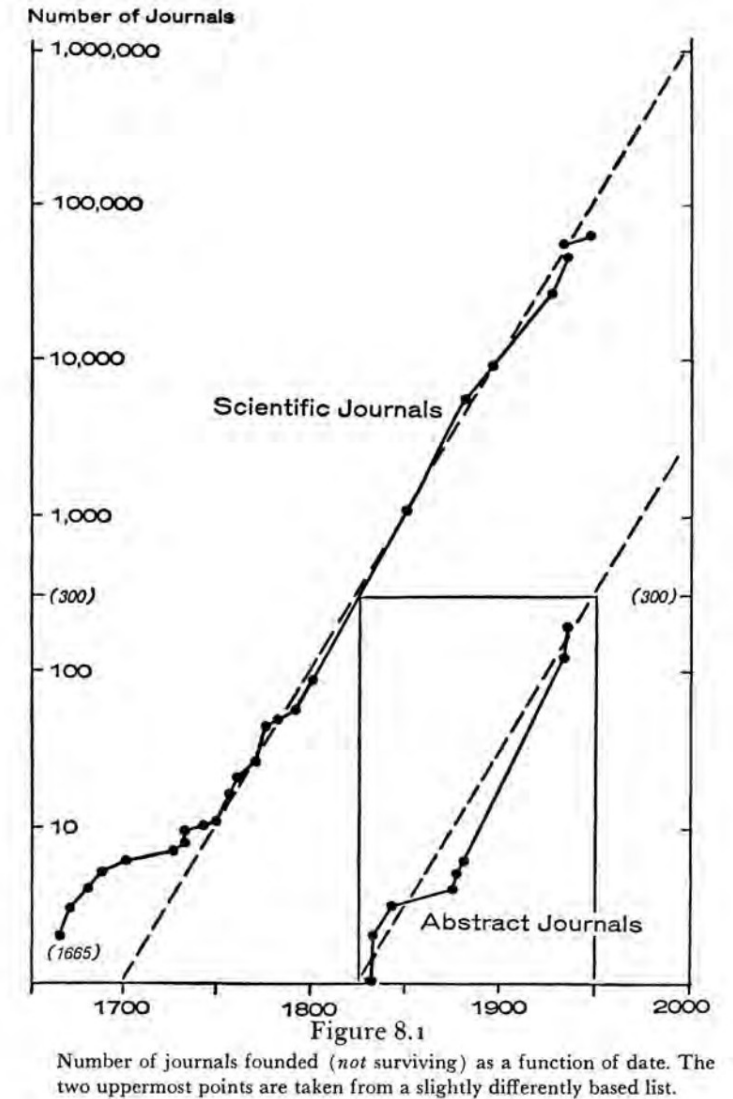
- 1944: Librarian Fremont Rider @ Wesleyan estimates American university libraries doubling in size every 16 years:

*“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

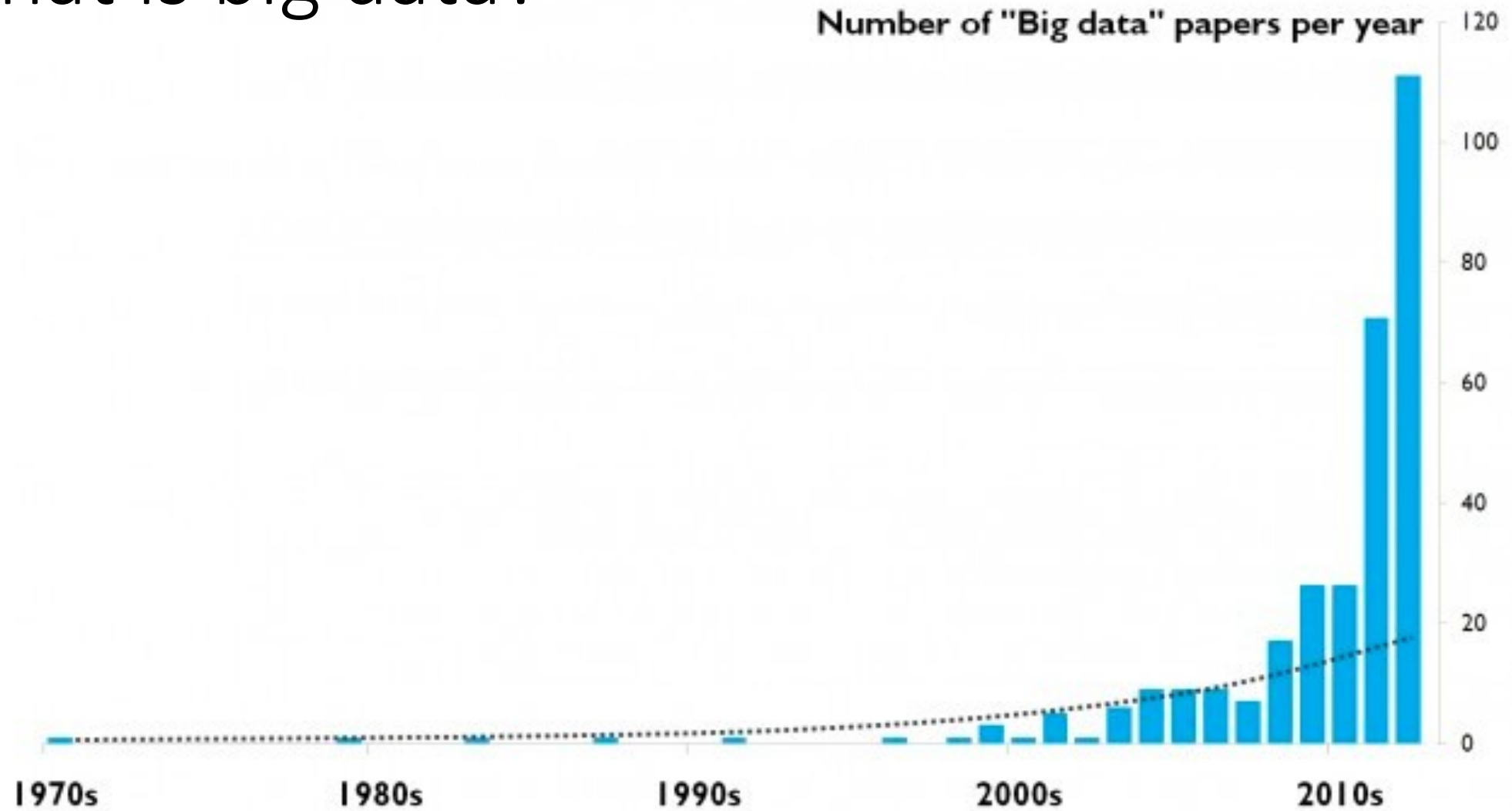
“each advance generates a new series of advances at a reasonably constant birth rate, so that the number of births is strictly proportional to the size of the population of discoveries at any given time”



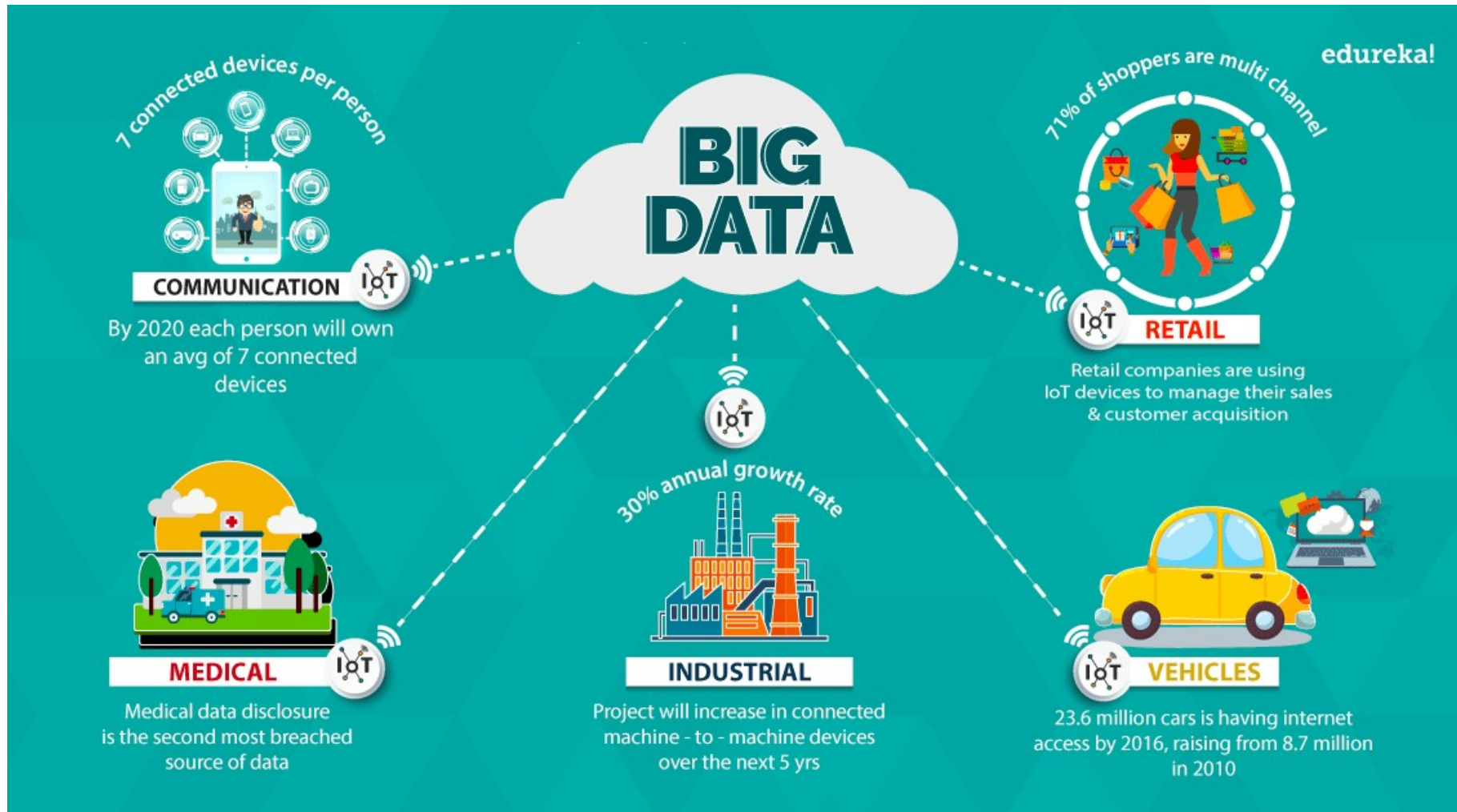
What is big data?

- **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 Mougals coins term 'Big Data', referring to the scale of data that is nearly impossible to manage and process w/ available tools

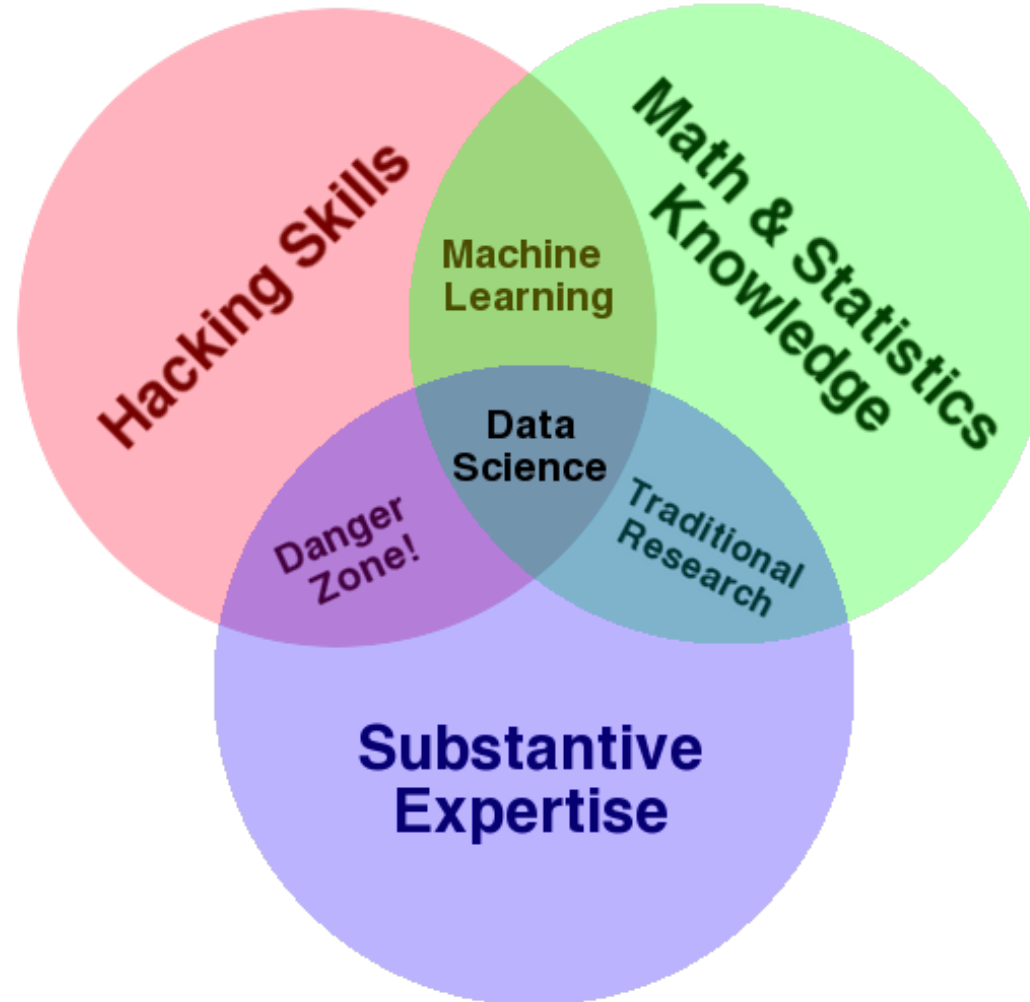
What is big data?



Where is big data?



What is data science?



The Big Data that I Work With



Google search results for "traffickcam".

Search results include:

- 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. TraffickCam has built a database of more than 4.5 million images of hotel ...

What's the plan for this class?

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

What's the plan for this class?

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

Intro to Python, Data Types, Sequences



Melissa Kline

@melissaekline



My strongest memory of learning to code is sitting thru my very first lecture silently panicking because while the FOR loop concept made perfect sense, I had no idea what I was supposed to do with it. I knew 'type it into Word' was wrong, but didn't know how to ask the Q...

Programming / Coding / Hacking

- Instruct a computer to carry out tasks
 - “Add 2 + 4”
 - “Round 3.14159265”

Programming / Coding / Hacking

- Instruct a computer to carry out tasks
 - “Add 2 + 4”
 - “Round 3.14159265”
- Combine multiple tasks into programs
 - “Add 2 + 4, then round 3.14159265, add the outputs together”
 - “Read in a file with the stock prices for a stock over the last 5 years and compute the average return”

Programming / Coding / Hacking

- Instruct a computer to carry out tasks
 - “Add 2 + 4”
 - “Round 3.14159265”
- Combine multiple tasks into programs
 - “Add 2 + 4, then round 3.14159265, add the outputs together”
 - “Read in a file with the stock prices for a stock over the last 5 years and compute the average return”
- Write these instructions in a language the computer can understand

Parts of a Program

- Expressions

| Expression Type | Operator | Example | Value |
|-----------------|----------|----------|---------|
| Addition | + | 2 + 3 | 5 |
| Subtraction | - | 2 - 3 | -1 |
| Multiplication | * | 2 * 3 | 6 |
| Division | / | 7 / 3 | 2.66667 |
| Remainder | % | 7 % 3 | 1 |
| Exponentiation | ** | 2 ** 0.5 | 1.41421 |

Parts of a Program

- Variables

| NAME | VALUE | TYPE |
|---------|--------|--------|
| number | 123 | int |
| sum | -456 | int |
| pi | 3.1416 | double |
| average | -55.66 | double |

A variable has a name, stores a value of the declared type

Parts of a Program

- Operations

a = 1

b = 2

c = a + b

Parts of a Program

- Operations

```
a = 1  
b = 2  
c = a + b
```

- Functions (named operations)

```
max(2, 2 + 3, 4)
```

```
give_everyone_in_class_an_a(class_list)
```

Python

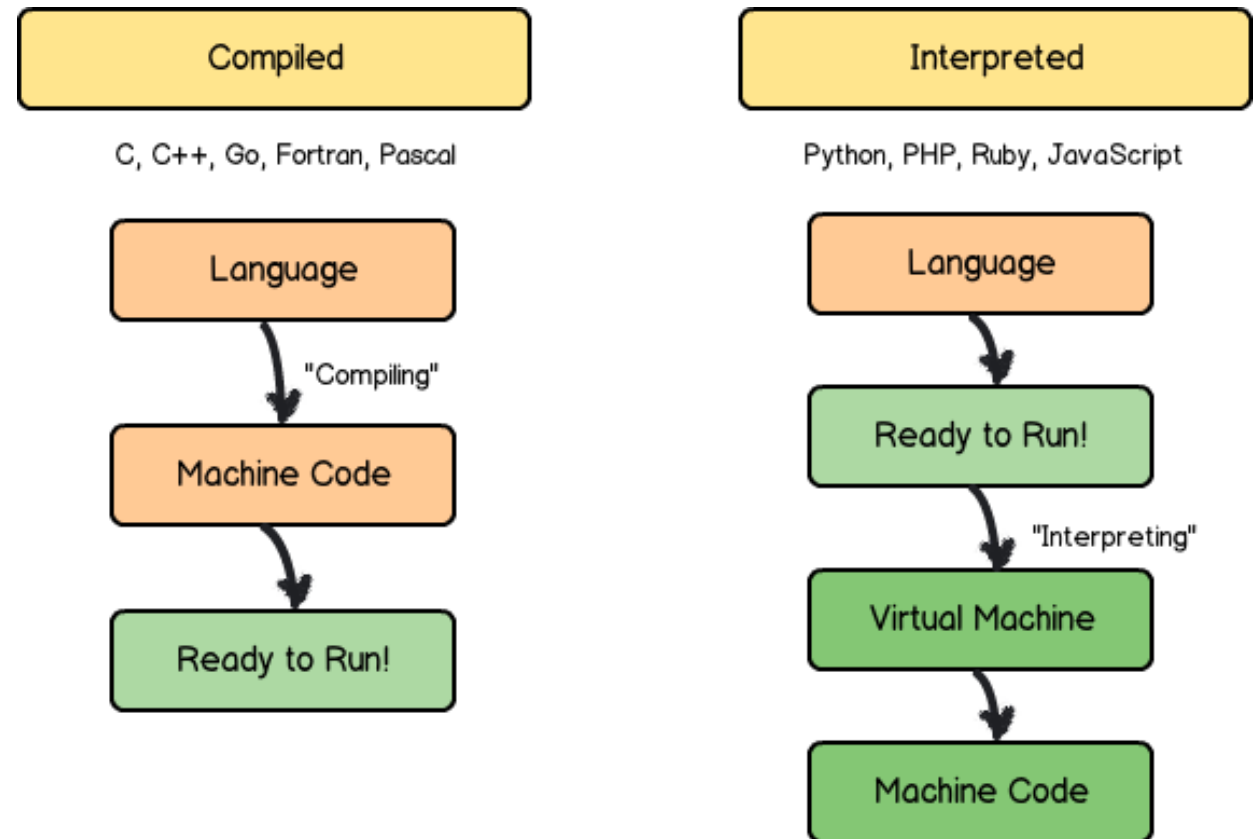
- Free & widely used (including in CSCI 1300)

Python

- Free & widely used (including in CSCI 1300)
- Compatible w/ different systems

Python

- Free & widely used (including in CSCI 1300)
- Compatible w/ different systems
- Interpreted rather than compiled



Python

- Free & widely used (including in CSCI 1300)
- Compatible w/ different systems
- Interpreted rather than compiled
- Dynamically typed

Static typing:

```
String name;
```

```
name = "John";
```

```
name = 34;
```

Variables have types

Values have types

Variables cannot change type

Dynamic typing:

```
var name;
```

```
name = "John";
```

```
name = 34;
```

Variables have no types

Values have types

Variables change type dynamically

@jordanwilder

Python

- Free & widely used (including in CSCI 1300)
- Compatible w/ different systems
- Interpreted rather than compiled
- Dynamically typed
- Built in memory management

Python

- Free & widely used (including in CSCI 1300)
- Compatible w/ different systems
- Interpreted rather than compiled
- Dynamically typed
- Built in memory management
- Lots of built in modules + tons of very useful, well developed external libraries (NumPy, Pandas)

Python

- Free & widely used (including in CSCI 1300)
- Compatible w/ different systems
- Interpreted rather than compiled
- Dynamically typed
- Built in memory management
- Lots of built in modules + tons of very useful, well developed external libraries (NumPy, Pandas)
- Readable!

Hello World: Python vs. Java

A **"Hello, World!" program** generally is a [computer program](#) that outputs or displays the message "Hello, World!". Such a program is very simple in most [programming languages](#), and is often used to illustrate the basic [syntax](#) of a programming language. It is often the first program written by people learning to code.^{[1][2]}


https://en.wikipedia.org/wiki/%22Hello,_World!%22_program

```
def hello_world():  
    print("Hello, world!")  
  
hello_world()
```

```
public class HelloWorld {  
    public static void main (String[] args) {  
        System.out.println("Hello, world!");  
    }  
}
```

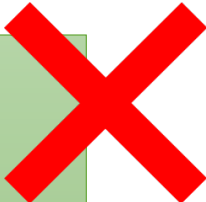
Whitespace in Python

```
def hello_world():  
    print("Hello, world!")
```




```
hello_world()
```

```
def hello_world()  
print("Hello, world!")
```




```
hello_world()
```

```
public class HelloWorld {  
    public static void main (String[] args) {  
        System.out.println("Hello, world!");  
    }  
}
```




```
public class HelloWorld {  
    public static void main (String[] args) {  
        System.out.println("Hello, world!");  
    }  
}
```



```
public class HelloWorld { public static void main (String[]  
args) {System.out.println("Hello, world!"); }}
```

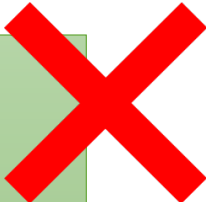


Whitespace in Python



```
def hello_world():  
    print("Hello, world!")
```

```
hello_world()
```



```
def hello_world()  
print("Hello, world!")
```

```
hello_world()
```

```
>>> def hello_world():  
...     print("Hello, world!")  
...
```

```
>>> hello_world()  
Hello, world!
```

```
[>>>
```

```
[>>>
```

```
>>> def hello_world():  
...     print("Hello, world!")  
      File "<stdin>", line 2  
        print("Hello, world!")
```

^

```
IndentationError: expected an indented block
```

```
>>> █
```

Python

