CSCI 3360 | Spring 2024
# Data Science I

Jin Sun, PhD
School of Computing

Week 1: Introduction

# Outline

- Background
- Class logistics and policy
- Class topics
- Q&A



"students learning data science in Athens"

# Data Science Motivating Examples

# Data Science Motivating Examples

# Data Science Motivating Examples



**Data Science in Finance**

- Risk Analytics
- Managing Customer Data
- Fraud Detection
- Real-time Analytics
- Consumer Analytics
- Algorithmic Trading

# Data Science Motivating Examples

# What is "Data Science"?

Data science (DS) is a **multidisciplinary** field that combines techniques from statistics, mathematics, computer science, and domain knowledge to extract insights and knowledge from data. It involves collecting, cleaning, analyzing, and interpreting large volumes of data to uncover patterns, trends, and relationships. Data scientists use various tools and techniques, such as data visualization, machine learning, and statistical modeling, to make data-driven decisions and solve complex problems. In this course, we will explore the fundamentals of data science and learn how to apply these techniques to real-world datasets.

# What is "Data Science"?

Data science (DS) is a **multidisciplinary** field that combines techniques from statistics, mathematics, computer science, and domain knowledge to extract insights and knowledge from data. It involves collecting, cleaning, analyzing, and interpreting large volumes of data to uncover patterns, trends, and relationships. Data scientists use various tools and techniques, such as data visualization, machine learning, and statistical modeling, to make data-driven decisions and solve complex problems. In this course, we will explore the fundamentals of data science and learn how to apply these techniques to real-world datasets.

-- Generated by Copilot

# DS vs others

- Statistics

Statistics is a branch of mathematics dealing with data collection, analysis, interpretation, presentation, and organization. It provides methodologies to design experiments and surveys, and techniques to analyze the results to draw conclusions.

# DS vs others

- Machine learning

Machine Learning is a subset of artificial intelligence that uses statistical techniques to enable machines to improve with experience. It involves the creation of algorithms that can modify themselves without human intervention to produce desired outputs by feeding itself through structured data.



10

# DS vs others

- Deep learning

Deep Learning is a subfield of machine learning that uses algorithms inspired by the structure and function of the brain's neural ne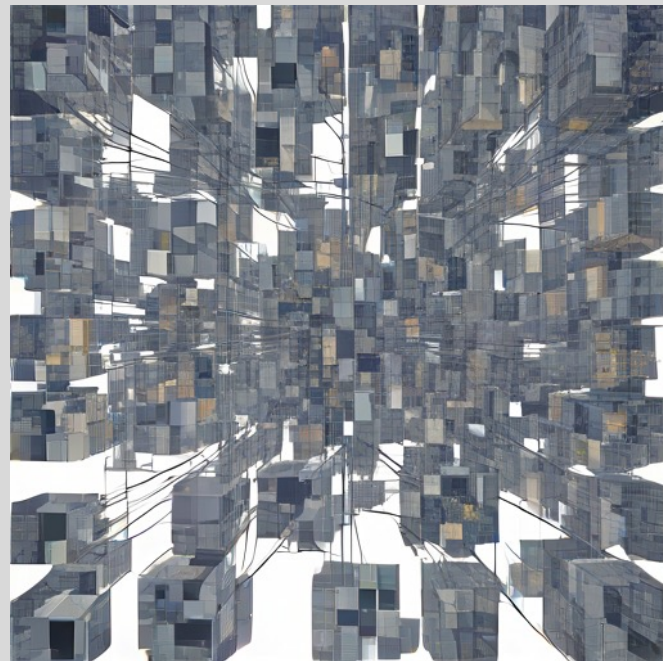tworks. These artificial neural networks are designed to simulate human decision-making and are becoming essential in areas where we are inundated with data. They provide a means to extract and learn complex patterns from massive amounts of data, and are often used in image recognition, speech recognition, and natural language processing tasks.

# DS vs others

- "A.I."

At this point, I don't even know what AI means anymore.

It could be: knowledge base system, chat bot, simple statistical model, or an image filter.

# What this class is about?

- Introduction to data science
- Fundamental understanding of data science pipeline
- Python + DS programming
- Hands-on examples

There is Data Science II for more advanced algorithms and techniques.

# Learning objectives

- Familiar with Python and relevant libraries
- Familiar with data science pipeline
- Understand the fundamentals of data learning
- Can formulate a learning problem from raw data
- Can train simple models to learn from data
- Can validate the performance of a model

# Outline

- Background
- <span style="color:red">Class logistics and policy</span>
- Class topics
- Q&A

# Class format

Lectures                                                    Lab

# Class format

~~Lectures~~ Theory

- Concepts
- Maths
- Slides and Whiteboard

~~Lab~~ Practice

- Coding
- Development
- Toy examples and playground

For each week's topic, we will start with the fundamental concepts of the topic and then learn to realize those concepts by programming.

The learning objectives of each week will be evaluated by in-class quiz.

# Textbooks

- An introduction to statistical learning (https://www.statlearning.com/)

- Python data science handbook
  (https://jakevdp.github.io/PythonDataScienceHandbook/)

Both books are free online.

# Teaching team

Instructor:

Prof. Jin Sun

Office Hours:  Thursdays 4-5pm or by appointment

Office: 804 Boyd

Email: jinsun@uga.edu

Teaching assistant:

TBD

# Evaluation and grading

The final course grade will be weighted as follows:

| | |
|---|---|
| Quiz: | 10% |
| Homework: | 40% |
| Midterm exam: | 10% |
| Final exam: | 15% |
| Project: | 25% |

Late policy: 10% of total score deduction for each late day (including partial day).

# Project

You will work in a team on a course project. Each team should have 2-3 members.
You are encouraged to design the project to solve <span style="color:red">a real-world problem</span>.

<span style="color:red">Project Proposal (5%):</span> What do you plan to do? What's the learning problem? Data?

<span style="color:red">Project Milestone (5%):</span> Preliminary results and progress report.

<span style="color:red">Project Report and Presentation (15%):</span> All results and findings.

# Assessment and feedback

For each learning objective, we will have methods for you to do self-assessment.

Advanced teaching techniques might be explored in this class.

Feedbacks are welcome!

# Outline

- Background
- Class logistics and policy
- <span style="color:red">Class topics, aka a whole semester in a day</span>
- Q&A

# Programming tools for DS

- Python

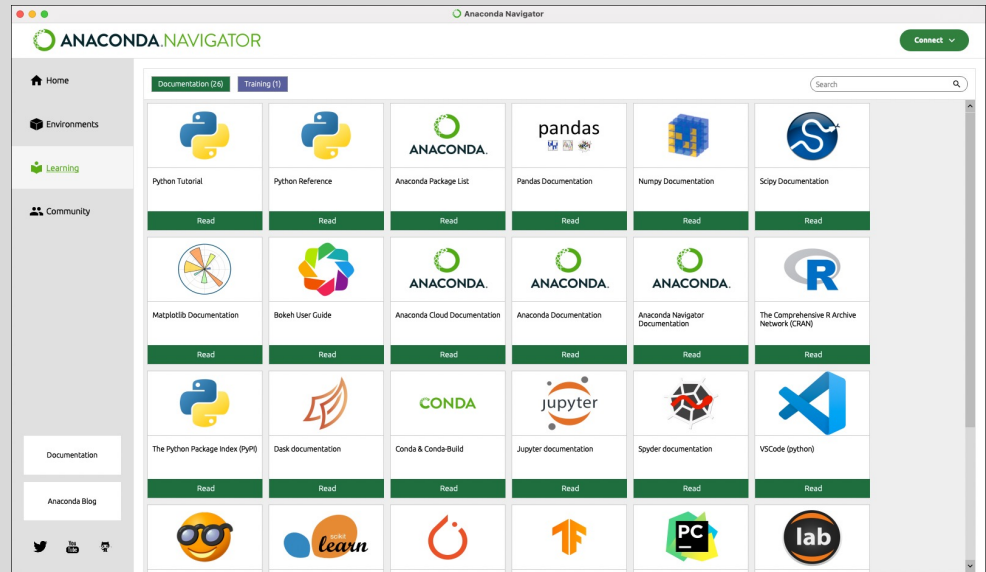"THE" programming language of data science. It is simple to config, develop, and run. Most importantly, there is a huge community.

```python
class Car:

    def __init__(self, speed=0):
        self.speed = speed
        self.odometer = 0
        self.time = 0

    def say_state(self):
        print(f"I'm going {my_car}kph!".format(self.speed))

    def accelerate(self):
        self.speed += 5
```

# Programming tools for DS

● Conda

Package manager for Python. Make it easy for users to manage environments and packages. It is cross-platform.

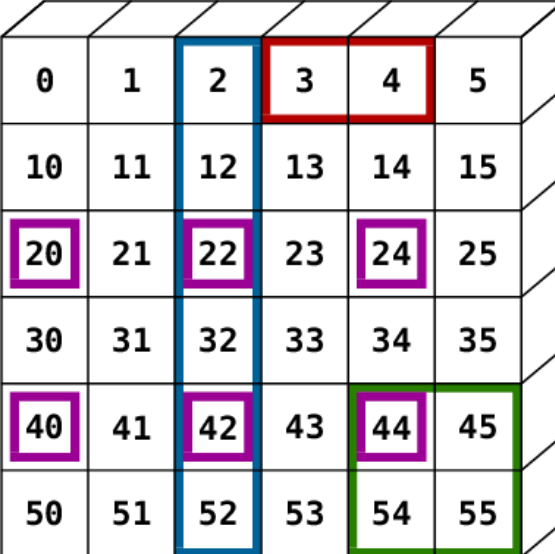# Programming tools for DS

- Numpy

Python library with Matlab-like syntax for matrices and vectors operations. And much more.

```
>>> a[0, 3:5]
array([3, 4])

>>> a[4:, 4:]
array([[44, 55],
       [54, 55]])

>>> a[:, 2]
a([2, 12, 22, 32, 42, 52])

>>> a[2::2, ::2]
array([[20, 22, 24],
       [40, 42, 44]])
```

| 0 | 1 | 2 | 3 | 4 | 5 |
| 10 | 11 | 12 | 13 | 14 | 15 |
| 20 | 21 | 22 | 23 | 24 | 25 |
| 30 | 31 | 32 | 33 | 34 | 35 |
| 40 | 41 | 42 | 43 | 44 | 45 |
| 50 | 51 | 52 | 53 | 54 | 55 |

# Programming tools for DS

- Pandas

Python library to handle data.

```python
df = pd.DataFrame({'City': ['Singapore','London','Hong Kong','Paris','Moscow'],
                   'City Population': [563, 898, 745, 215, 1192],
                   'City Area': [721.5, 1572, 1106, 105.4, 2511],
                   'Currency':['SGD','GBP','HKD','EUR','RUB'],
                   'Continent':['Asia','Europe','Asia','Europe','Europe'],
                   'Main Language': ['English','English','Chinese','French','Russian']})

df
```

|   | City | City Population | City Area | Currency | Continent | Main Language |
|---|------|-----------------|-----------|----------|-----------|---------------|
| 0 | Singapore | 563 | 721.5 | SGD | Asia | English |
| 1 | London | 898 | 1572.0 | GBP | Europe | English |
| 2 | Hong Kong | 745 | 1106.0 | HKD | Asia | Chinese |
| 3 | Paris | 215 | 105.4 | EUR | Europe | French |
| 4 | Moscow | 1192 | 2511.0 | RUB | Europe | Russian |

# Programming tools for DS

- Jupyter notebook

Intuitive IDE. Nice interface and flexibility. You can run all coding part of this class on Jupyter notebook, or on [Google Colab](#).

# Data

Of course, the first thing we should really be talking about is **data**.
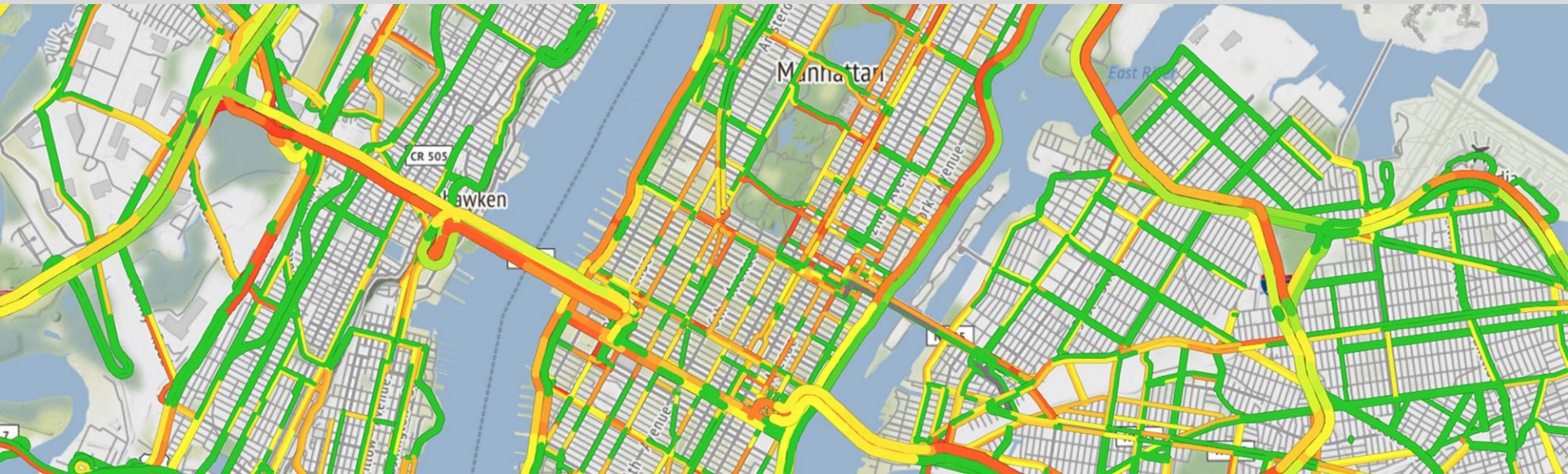
Data is the most important thing in your data science project. Period.

Without high quality data, you cannot run meaningful analysis or learn useful models from it.
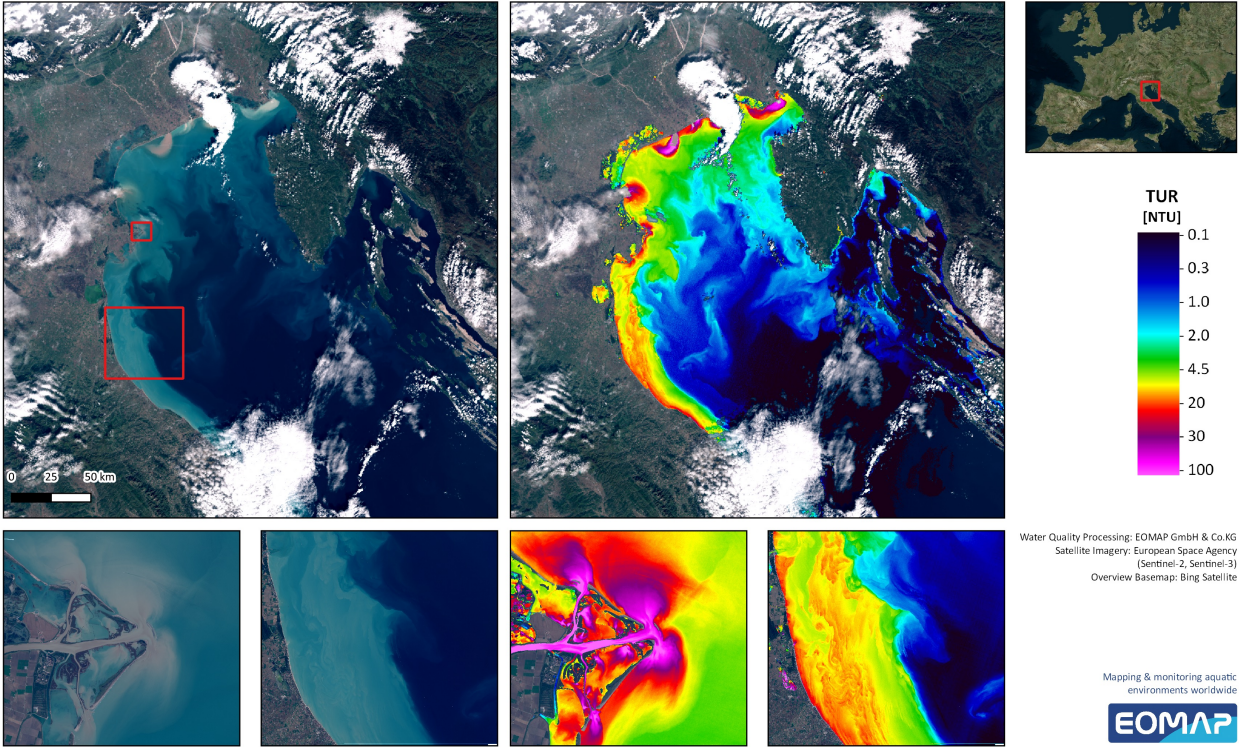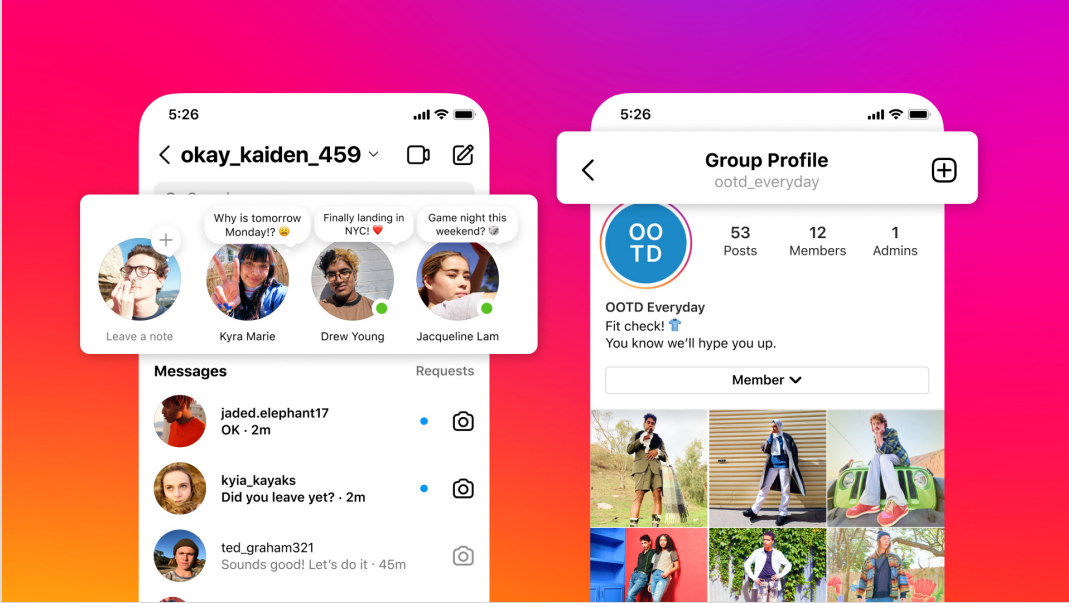
# Example data streams

Traffic

# Example data streams

Satellite



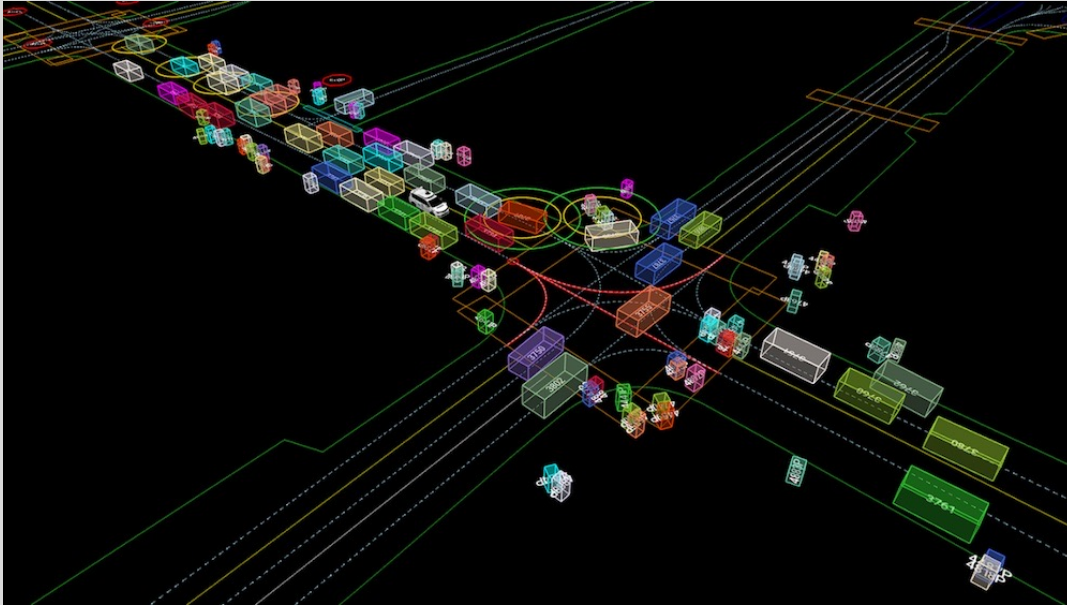Turbidity in the Adriatic Sea (2018/10/31)

# Example data streams

## Images and videos

# Example data streams

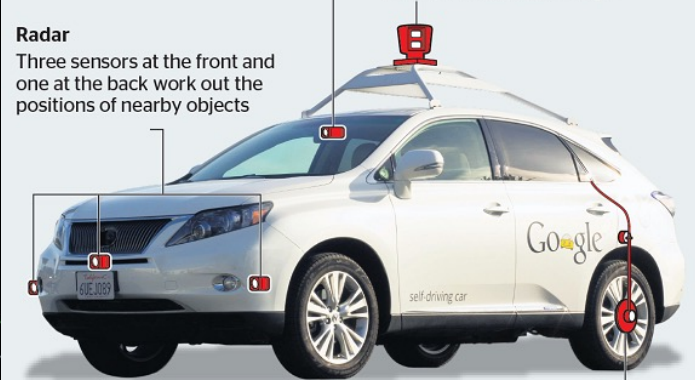Images and videos



**Look — no driver**

**Video camera**
Detects traffic lights, oncoming vehicles and other obstacles

**Lidar**
A rotating sensor on the roof scans 200ft in all directions to create a 3D map of its surroundings

**Radar**
Three sensors at the front and one at the back work out the positions of nearby objects

**Position estimator**
A sensor on the left rear wheel measures the car's movements so that its position can be mapped with accuracy

# Example data streams

Images and videos

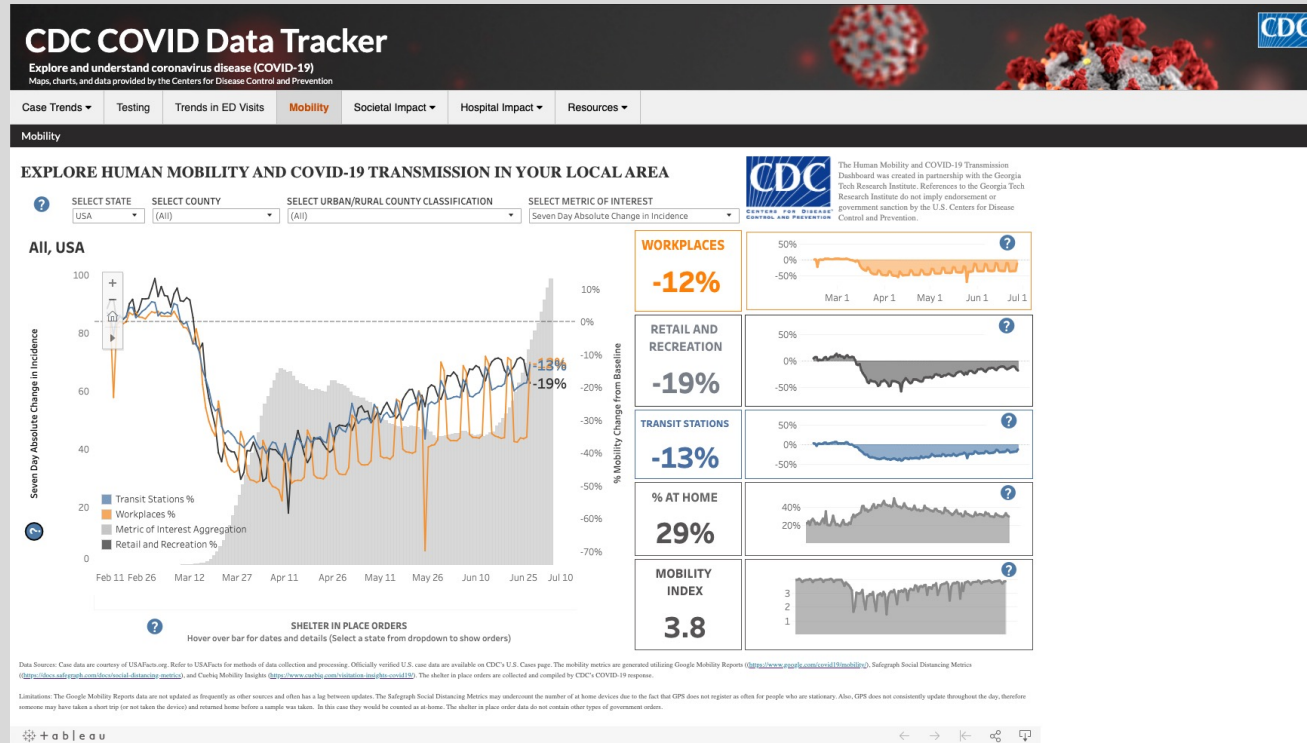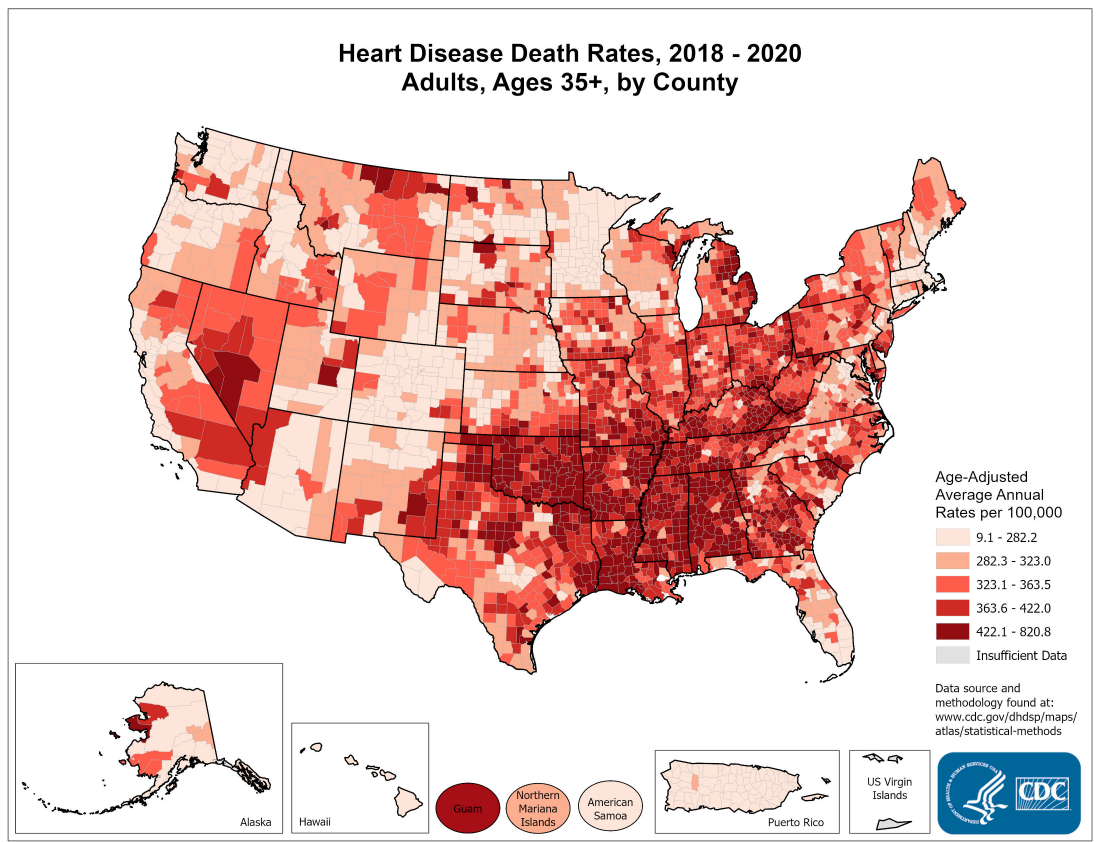- Surveillance

# Example data streams

Health

- Disease tracker

# Example data streams

Health

- Disease tracker



Heart Disease Death Rates, 2018 - 2020
Adults, Ages 35+, by County

Age-Adjusted
Average Annual
Rates per 100,000

- 9.1 - 282.2
- 282.3 - 323.0
- 323.1 - 363.5
- 363.6 - 422.0
- 422.1 - 820.8
- Insufficient Data

Data source and
methodology found at:
www.cdc.gov/dhdsp/maps/
atlas/statistical-methods

Alaska

Hawaii
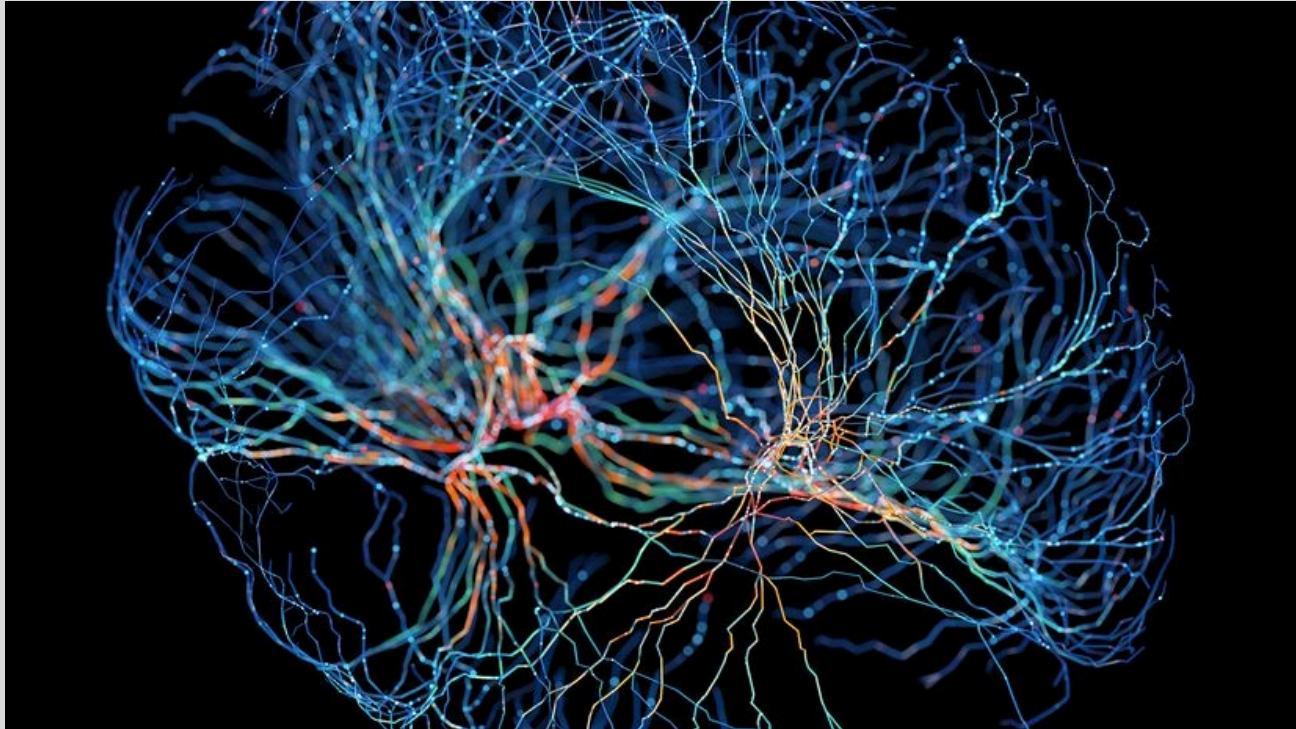
Guam

Northern
Mariana
Islands
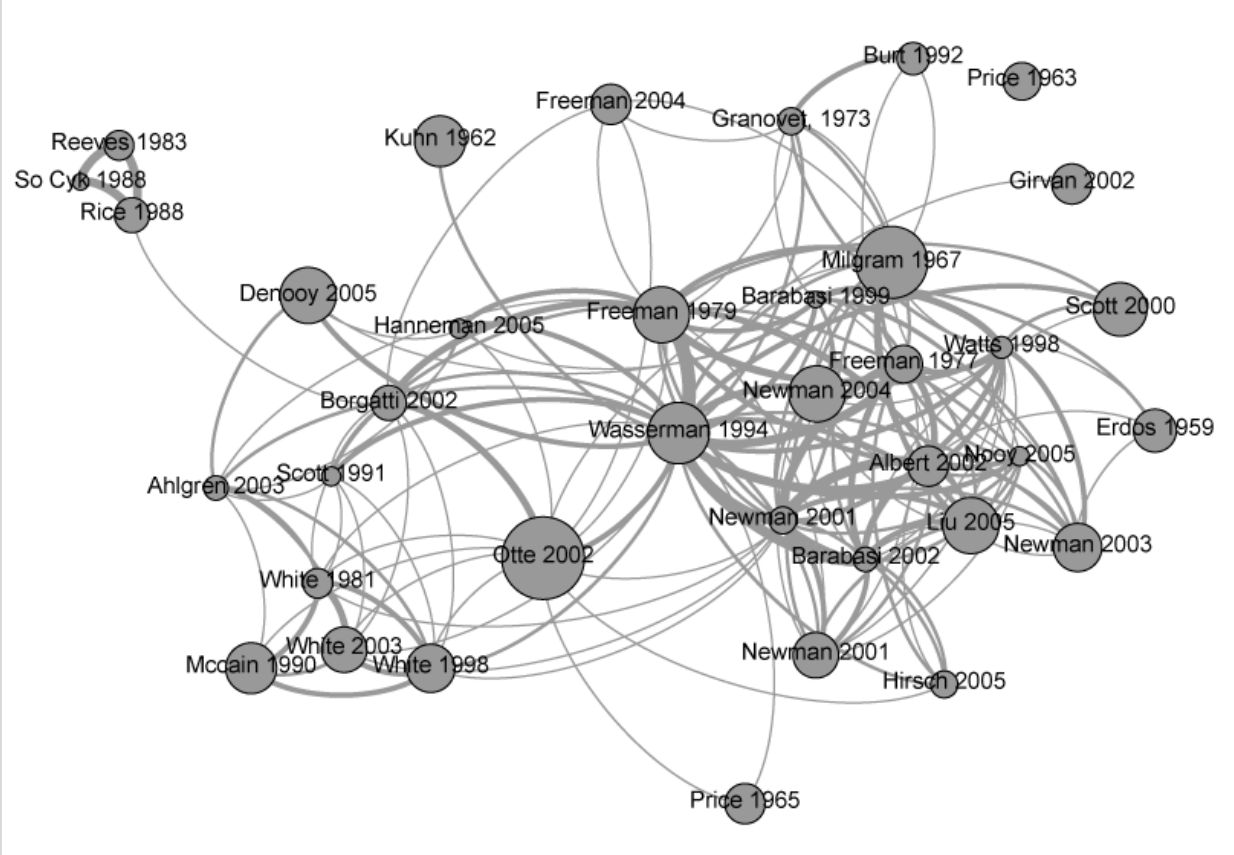
American
Samoa

Puerto Rico

US Virgin
Islands

CDC

# Example data streams

Brain

# Example data streams

Scientific research

# Data formats

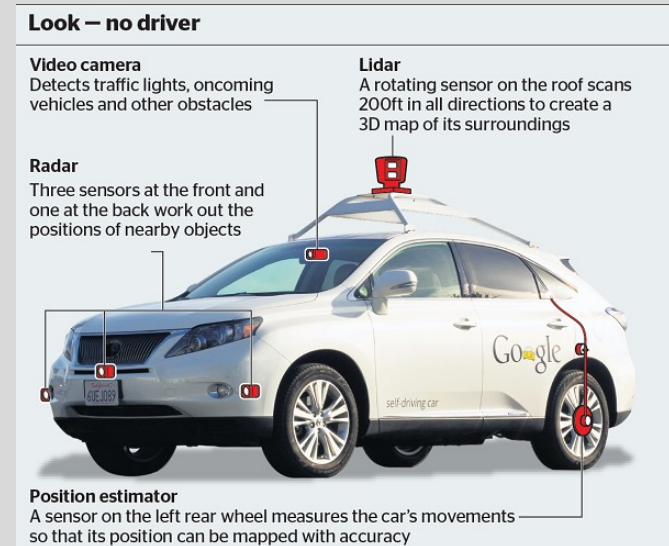Exact format of data depends on its domain of applications.

- Text strings
- Images
- Geo-tagged sensor readings
- Structured forms

# Data formats

In most cases, raw data are in mixed formats.

For example, in self-driving cars, data contains:

- RGB images
- LiDAR
- GPS
- Vehicle speed and orientation
- etc



**Look — no driver**

**Video camera**
Detects traffic lights, oncoming vehicles and other obstacles

**Lidar**
A rotating sensor on the roof scans 200ft in all directions to create a 3D map of its surroundings

**Radar**
Three sensors at the front and one at the back work out the positions of nearby objects

**Position estimator**
A sensor on the left rear wheel measures the car's movements so that its position can be mapped with accuracy

40

# (Almost) Infinite amount of data

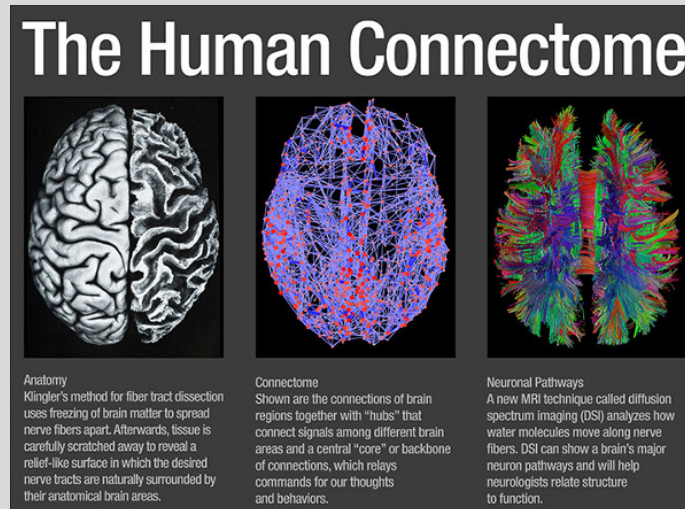In many scenarios, we have more data than we can handle.

How much data YouTube generates?

"Every minute, over 500 hours (about 3 weeks) of video are uploaded to YouTube. According to Global Media Insight, that's equivalent to over 300,000 hours (about 34 years) of video every day. This represents around 30,000 hours of new video content being uploaded per hour." -- src

# Not enough data

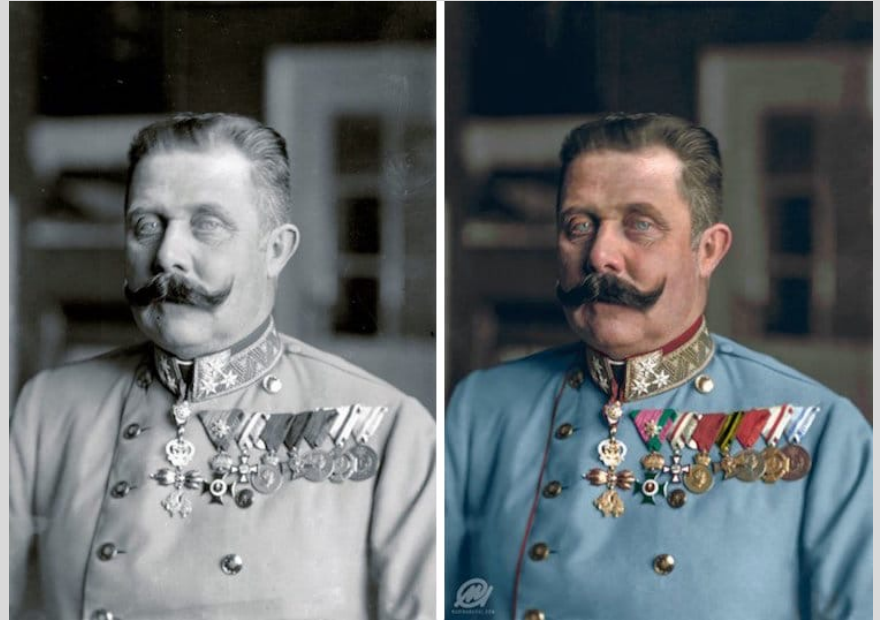On the other hand, in many cases we simply don't have enough data or no data at all.

- Limited data: e.g., We cannot record the activity of every single neuron in human brain.



**The Human Connectome**

**Anatomy**
Klingler's method for fiber tract dissection uses freezing of brain matter to spread nerve fibers apart. Afterwards, tissue is carefully scratched away to reveal a relief-like surface in which the desired nerve tracts are naturally surrounded by their anatomical brain areas.

**Connectome**
Shown are the connections of brain regions together with "hubs" that connect signals among different brain areas and a central "core" or backbone of connections, which relays commands for our thoughts and behaviors.

**Neuronal Pathways**
A new MRI technique called diffusion spectrum imaging (DSI) analyzes how water molecules move along nerve fibers. DSI can show a brain's major neuron pathways and will help neurologists relate structure to function.

# Not enough data

On the other hand, in many cases we simply don't have enough data or no data at all.

- No data: e.g., We do not have colored photo of the world before it was invented.

# Data vs label

'labels' are usually the 'meaning' attached to data.

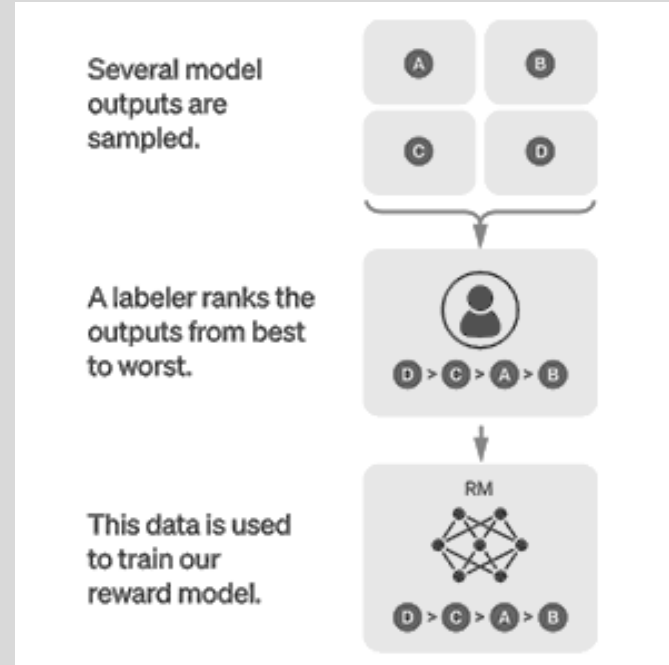In many applications, raw data are already labeled. Can you think of such cases?

In many applications, raw data recorded are 'unlabeled'.
Can you think of such cases?

# Data vs label

The process of obtaining the labels via human effort is called annotation.

For example, ChatGPT was trained using human annotated conversations.



Several model outputs are sampled.

A labeler ranks the outputs from best to worst.

This data is used to train our reward model.

# Data bias and ethics



Doctor walks in a bar

Nurse walks in a bar

# Data bias and ethics

# Understanding data

Visualization

# Understanding data

Analysis

# Learning from the data

What is 'learning'?

"... learning is about <span style="color:red">predicting</span> the <span style="color:blue">future</span> based on the <span style="color:green">past</span>."

-- CIML Book

<span style="color:green">Past</span>:  Training data

Fully labeled, unsupervised, semi-supervised

Classification, regression

<span style="color:blue">Future</span>:  Testing data

Out-of-distribution, domain shift, life-long learning

<span style="color:red">Predicting</span>:  Model

Linear model, kernel method, nearest neighbor, decision trees, neural networks

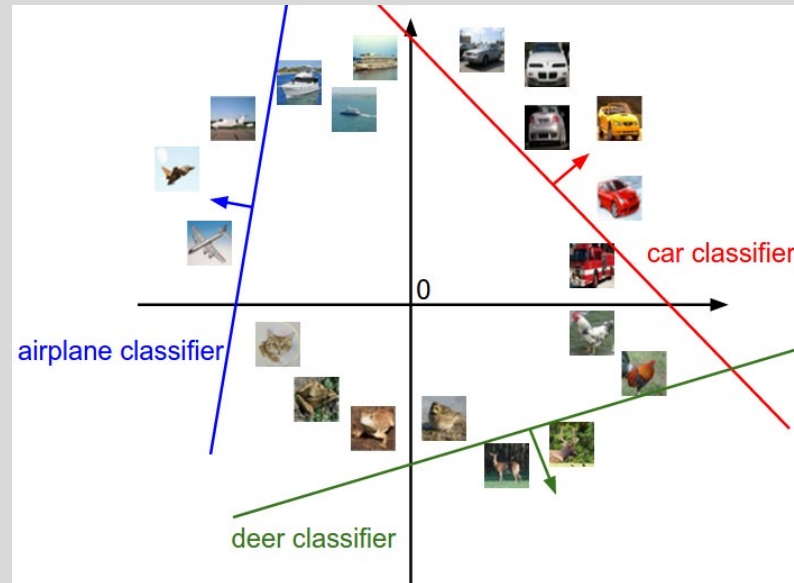<span style="color:orange">Evaluation</span>  Task-specific metrics, user studies

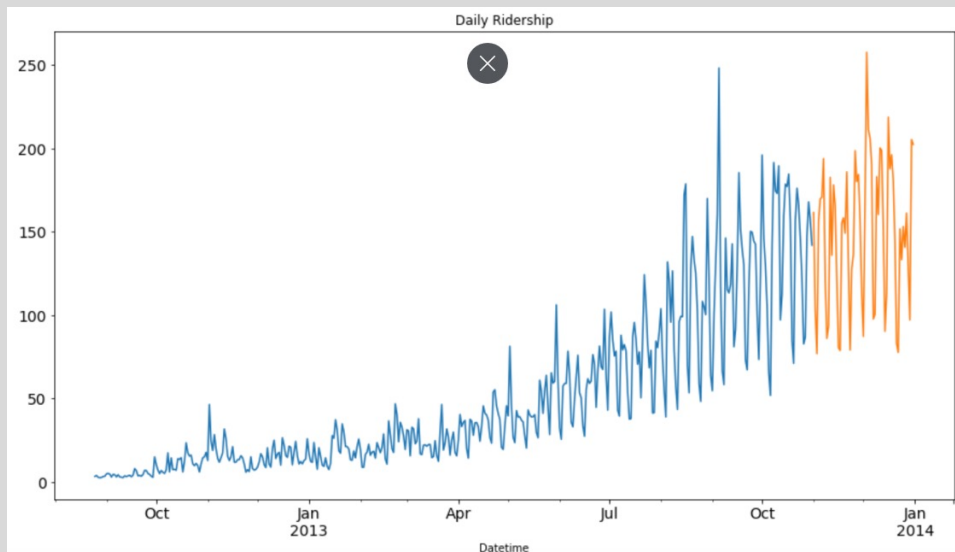# Learning from the data

- Classification

    e.g., rain or no rain? Cat or dog?

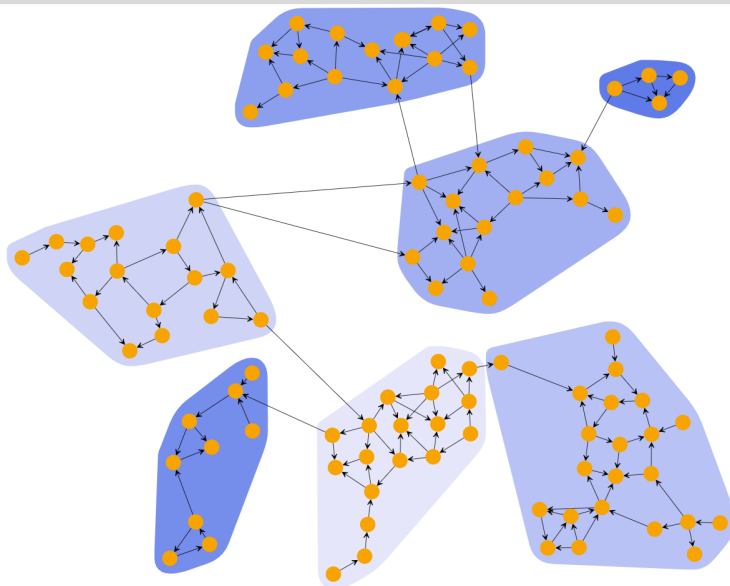# Learning from the data

- Regression

   e.g., what is the stock price for TSLA tomorrow?

# Learning from the data
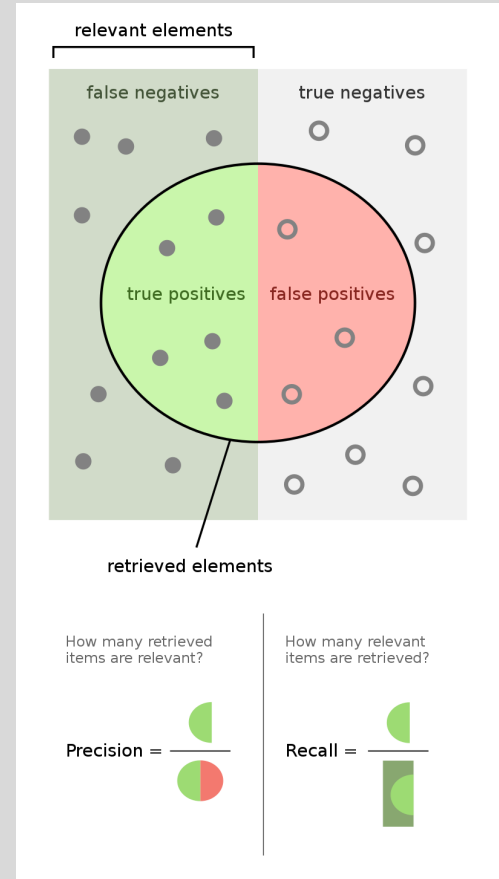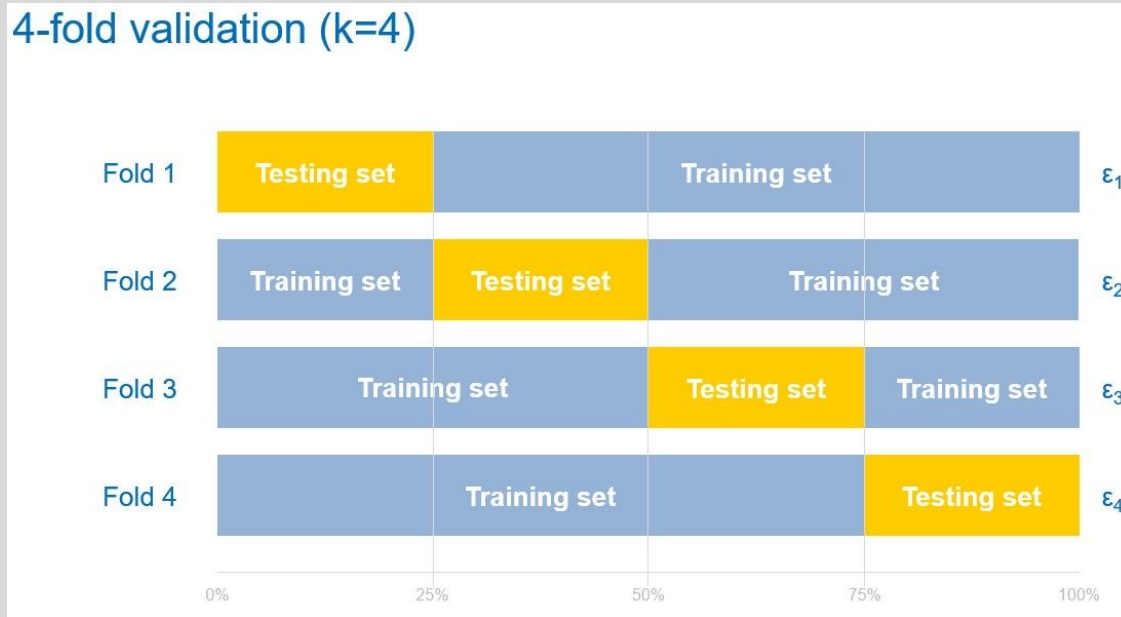
- Unsupervised learning

  e.g., how many sub-communities exist in this social network?
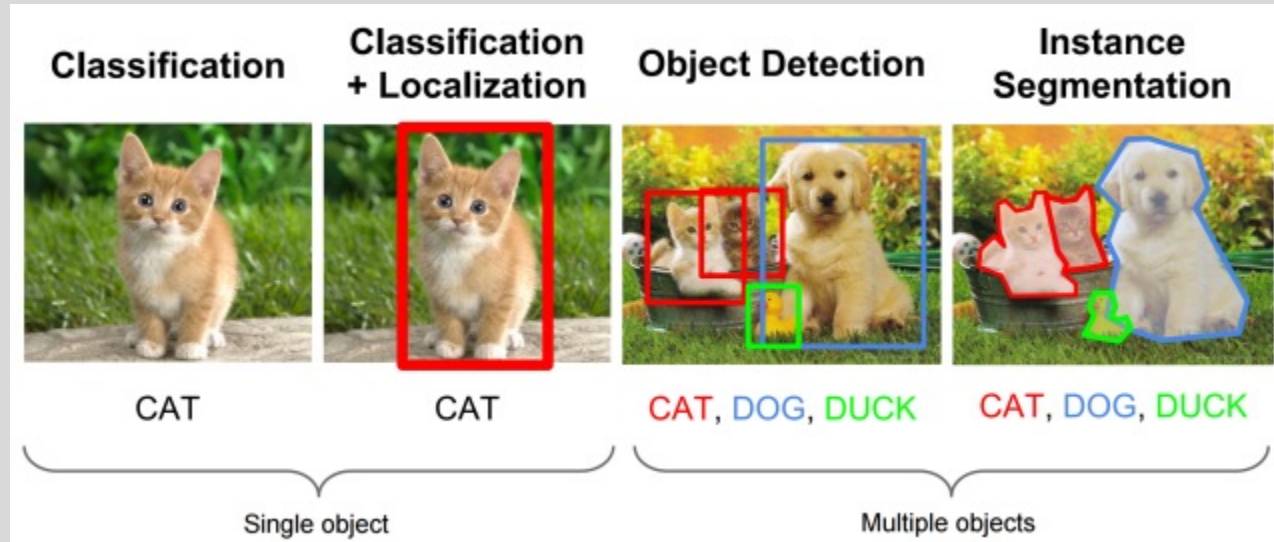
# Learning from the data

Evaluation and metrics
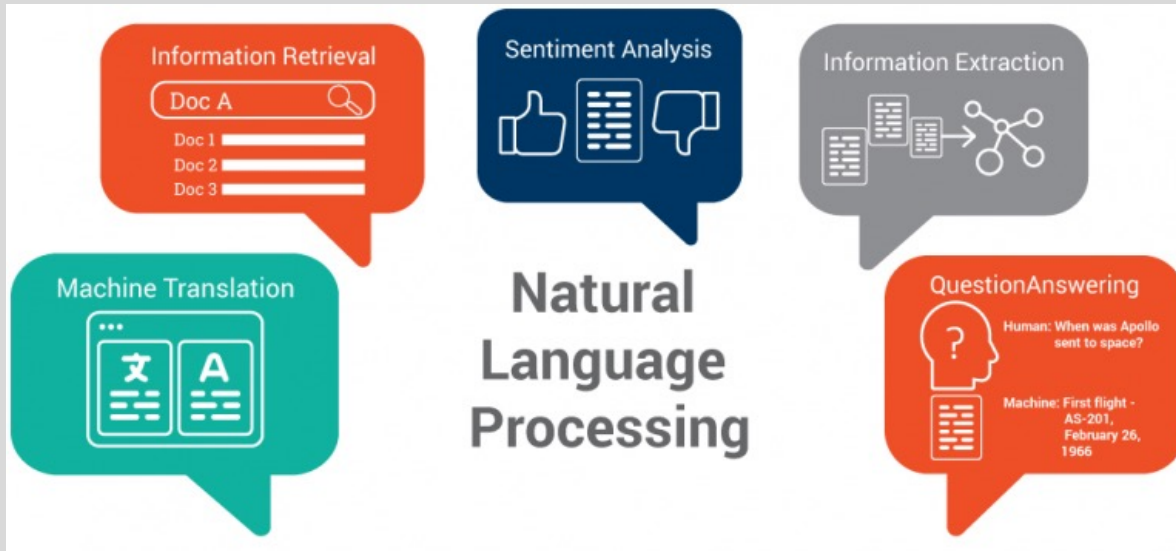


4-fold validation (k=4)

# Applications and practices

Vision



| Classification | Classification + Localization | Object Detection | Instance Segmentation |
|---|---|---|---|
| CAT | CAT | CAT, DOG, DUCK | CAT, DOG, DUCK |
| Single object | | Multiple objects | |

Cool things in computer vision now: text-to-image, 2d-to-3d, text-to-video, 3d reconstruction
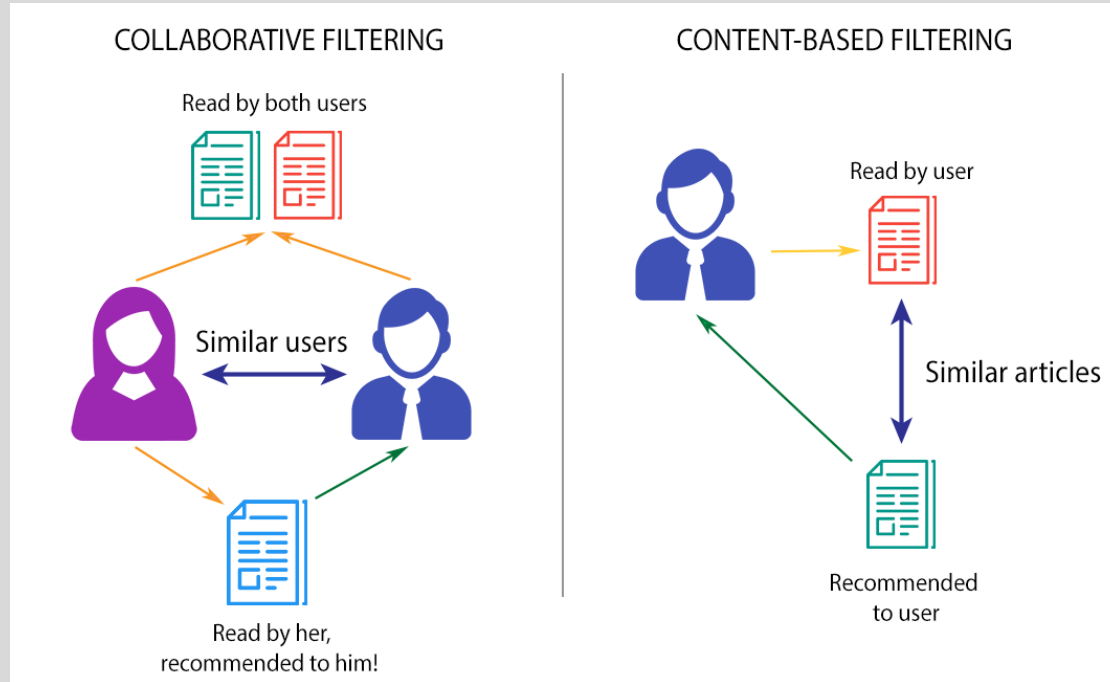
# Applications and practices

Languages

# Applications and practices
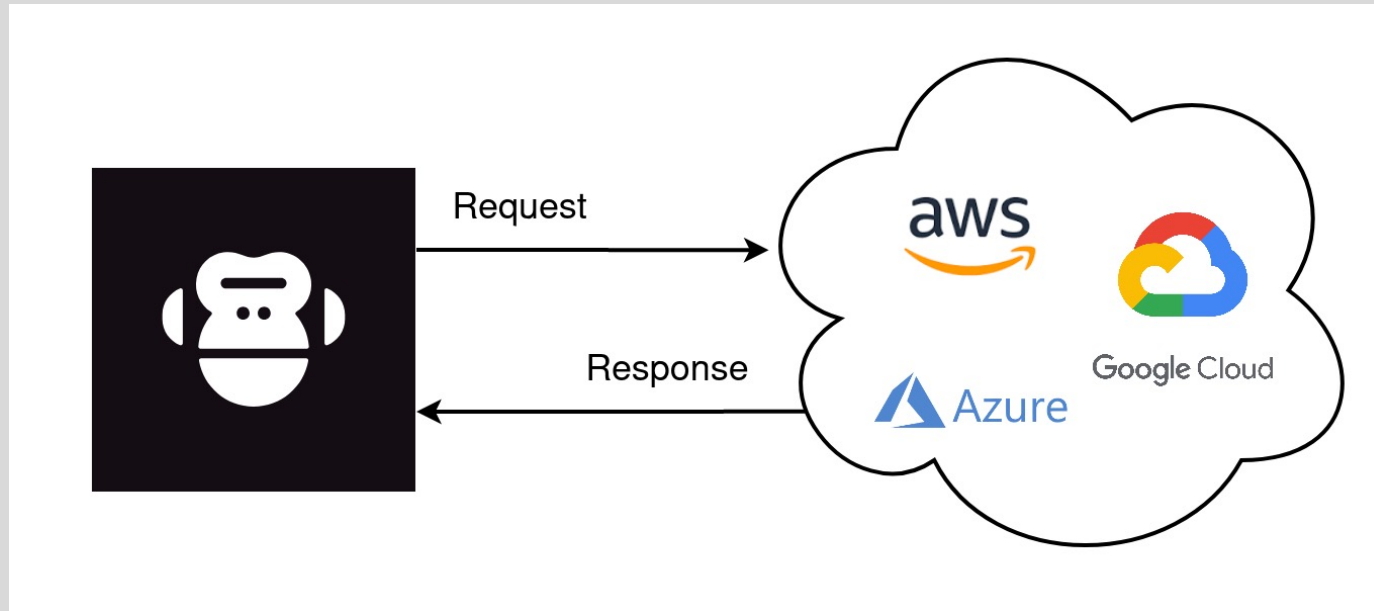
Recommendation systems



COLLABORATIVE FILTERING — Read by both users — Similar users — Read by her, recommended to him!

CONTENT-BASED FILTERING — Read by user — Similar articles — Recommended to user

# Applications and practices

Real-world ready development

# Outline

- Background
- Class logistics and policy
- Class topics
- Q&A

# What's next?

We will walk through example data science project pipelines and work on your first Kaggle challenge.

Start to think about your project and team!