

Introduction to Machine Learning

Rasika Bhalerao

Northeastern University

Climate Change AI Summer School 2023

Agenda

- **What is machine learning?**
- Supervised learning
- Unsupervised learning
- Reinforcement learning
- Generative models

Machine Learning:

“The science of getting computers to act without being explicitly programmed.”

- Andrew Ng

When should we use machine learning?

- Machine learning is not appropriate for every task!
 - If you can solve it analytically, that's better! (More explainable)
- When you have access to **data which “matches” the task**
- When **errors are allowable**
- When it's **cost effective**
 - Machine learning costs include (i) **dataset creation and processing** and (ii) **model development, deployment, and maintenance**

Machine learning strengths and limitations

Strengths

- Performing tasks at scale
- Modeling complex systems
- Generating derived data
- Integrating with other methods, e.g. domain and physical models

Limitations

- “Garbage in, garbage out”
- Inherits biases in data + human design/use
- Assumes patterns are persistent
- Finds correlation, not causation

Types of learning

- **Supervised learning**

Learning to predict or classify labels based on labeled input data
Performance feedback

- **Unsupervised learning**

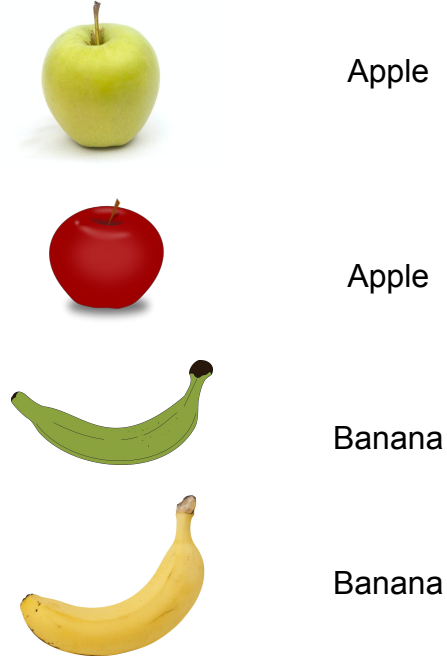
Finding patterns in unlabeled data
No performance feedback

- **Reinforcement learning**

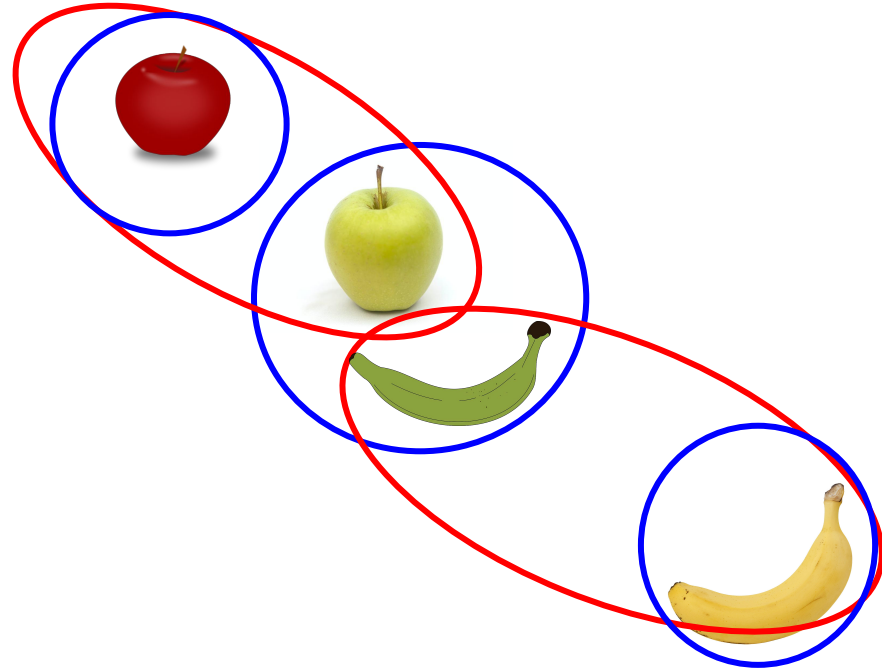
Learning well-performing behavior from state observations and rewards
Performance feedback

Supervised vs. Unsupervised learning

Supervised



Unsupervised

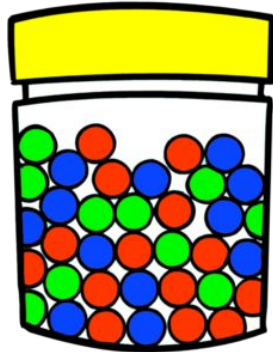


Data Types

- **Continuous**



- **Discrete**



- **Categorical**

- **Binary**

- Special case of categorical

- **Ordinal**

How do you feel today?

- ☒ 1 – Very Unhappy
- ☐ 2 – Unhappy
- ☐ 3 – OK
- ☐ 4 – Happy
- ☐ 5 – Very Happy

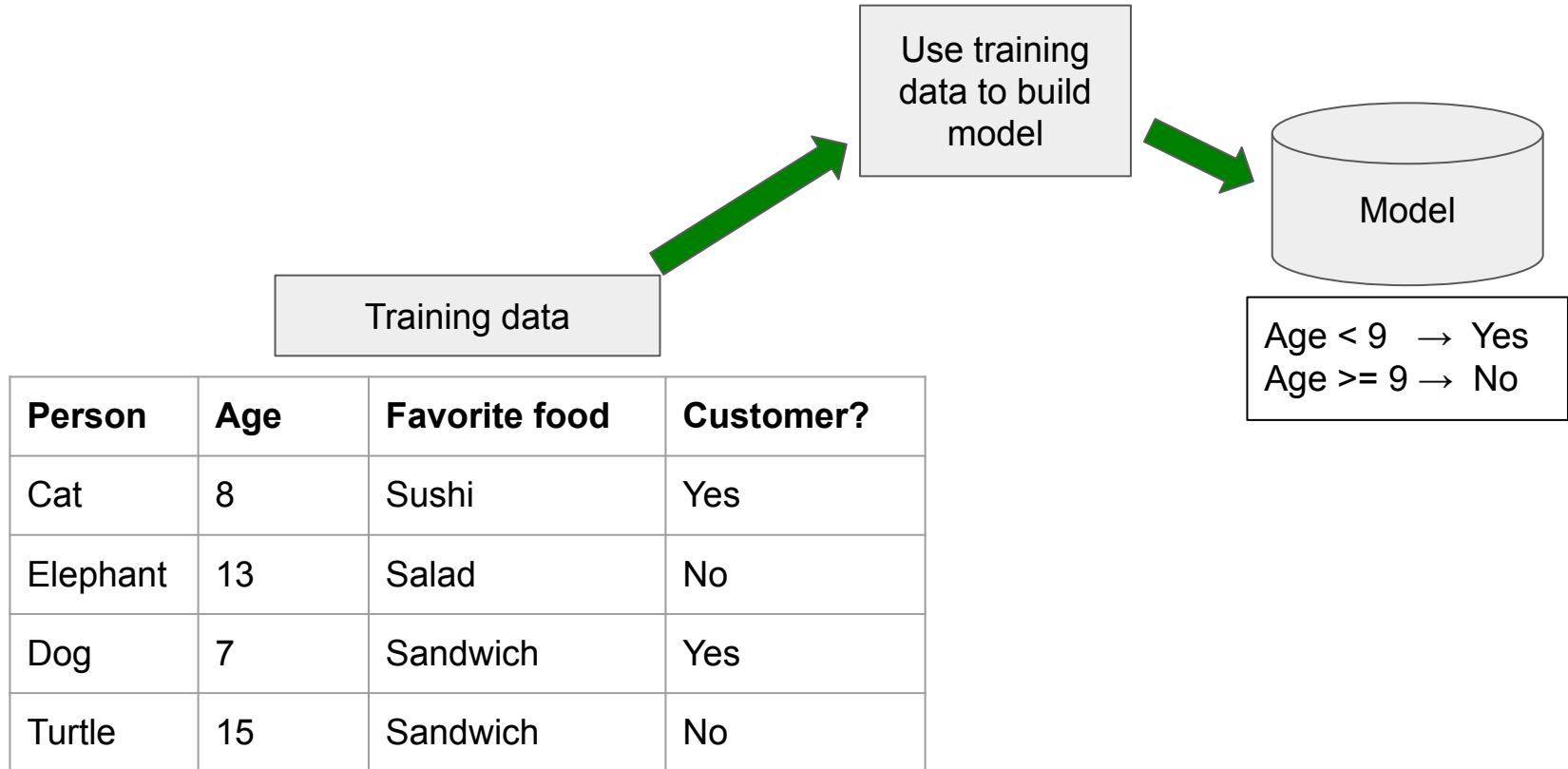
How satisfied are you with our service?

- ☒ 1 – Very Unsatisfied
- ☐ 2 – Somewhat Unsatisfied
- ☐ 3 – Neutral
- ☐ 4 – Somewhat Satisfied
- ☐ 5 – Very Satisfied

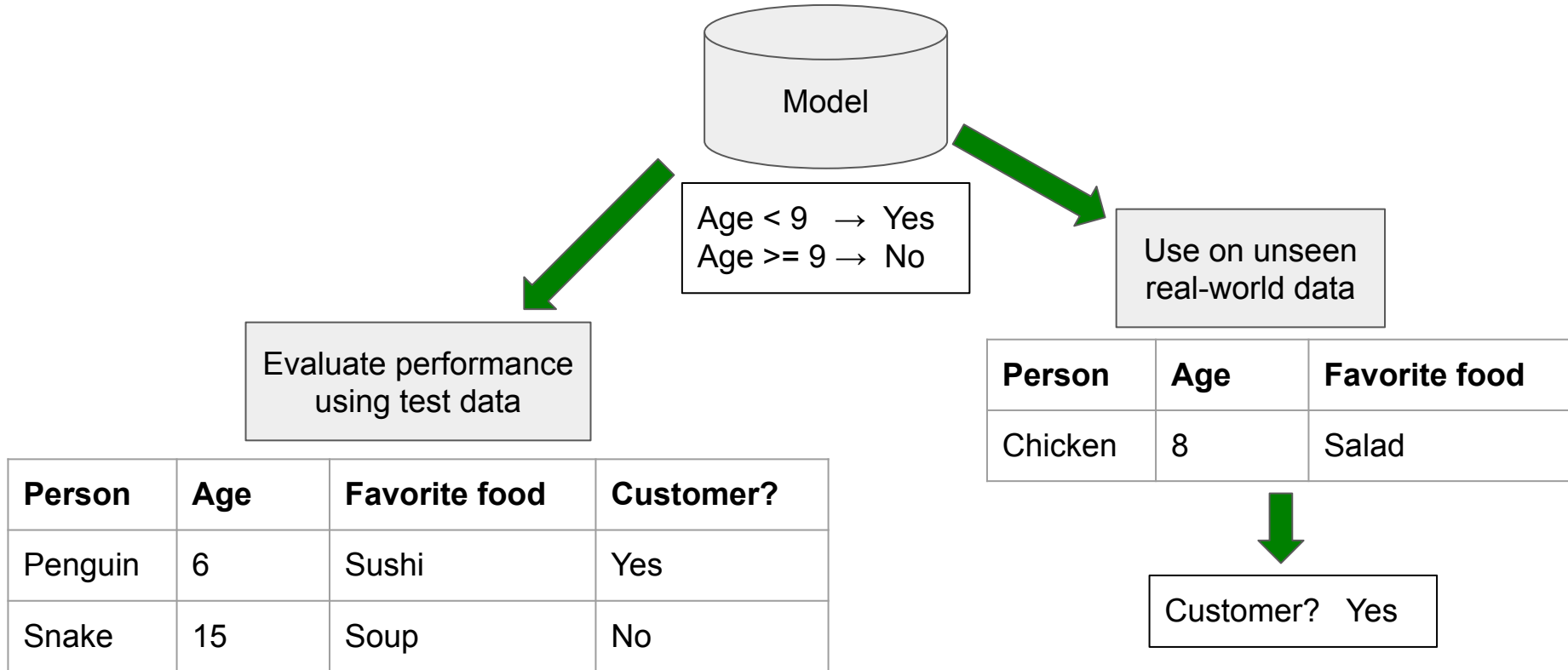
Agenda

- What is machine learning?
- **Supervised learning**
- Unsupervised learning
- Reinforcement learning
- Generative models

Supervised learning: training



Supervised learning: test / prediction



Supervised learning: categorical versus continuous labels

- Classification: **categorical labels**
 - Examples: pregnant or not, from which country, which type of road sign
- Regression: **continuous labels**
 - Examples: future stock price, life expectancy, distance to obstacle

Example: predicting bicycle counts

<https://www.climatechange.ai/papers/iclr2023/15>

Given: historical data of the number of bicycles in certain locations per hour

Want to predict: number of bicycles in future times at those locations

Which type are the labels? Categorical or continuous?

Supervised models

| Model | When to use it? |
|---|---|
| KNN | <ul style="list-style-type: none">➤ Little / no training time, large prediction time➤ Given small / medium dataset size |
| Linear / polynomial regression | <ul style="list-style-type: none">➤ Linear / polynomial relationship between input and output➤ Small training time➤ Given small / medium dataset size |
| Logistic regression, SVM, decision tree | <ul style="list-style-type: none">➤ Categorical output➤ Given small / medium training time |
| Neural network | <ul style="list-style-type: none">➤ Large training time, large computer➤ Given large dataset size |

k-Nearest Neighbors Algorithm

Training set: n instances, each with a feature vector and an output category

Now, given another (unseen) instance, we want to determine its category

Check the k instances in the training data that are closest to your new instance

- Categorical: choose the majority of those values
- Continuous: choose the mean/median of those values

Training set:

$(1,2) \rightarrow \text{red}$

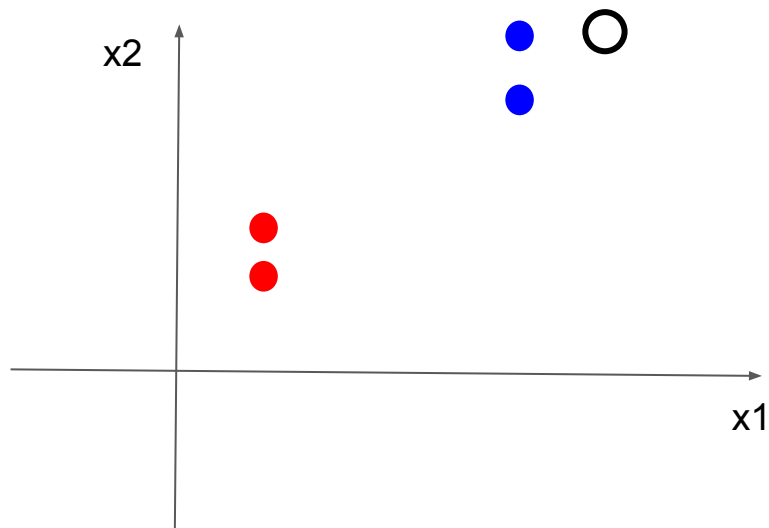
$(1,3) \rightarrow \text{red}$

$(5,5) \rightarrow \text{blue}$

$(5,6) \rightarrow \text{blue}$

New instance:

$(6,6) \rightarrow \text{blue}$



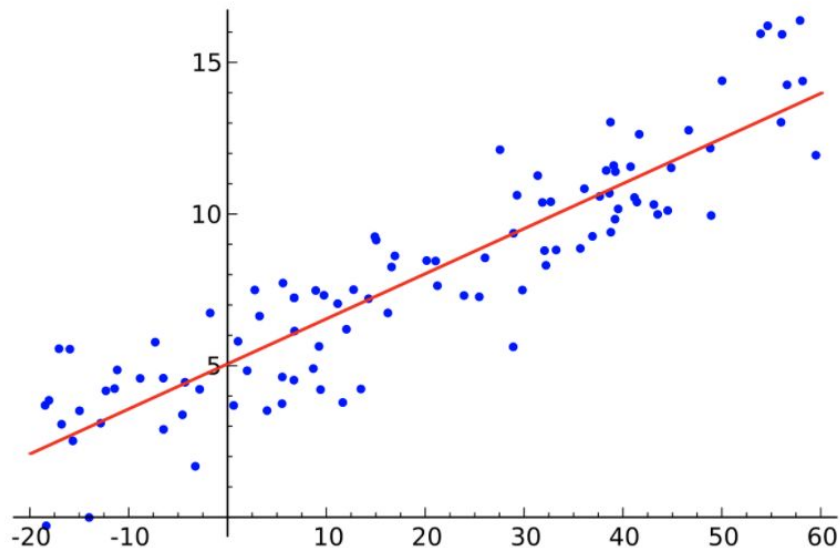
Linear / polynomial regression

Given $x \in \mathbb{R}$ and $y \in \mathbb{R}$

Find a function $f: x \rightarrow y$

How?

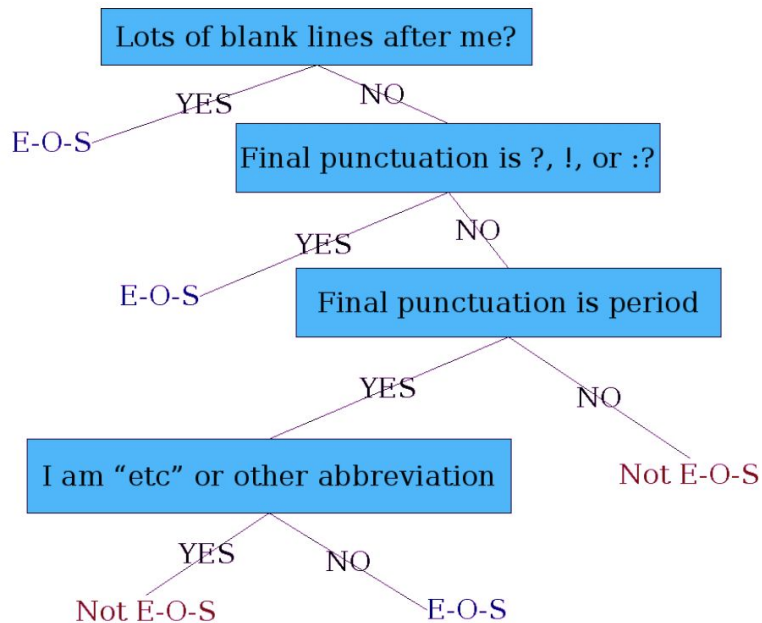
Define a **loss function** (“error”) and minimize it!



Other supervised classifiers

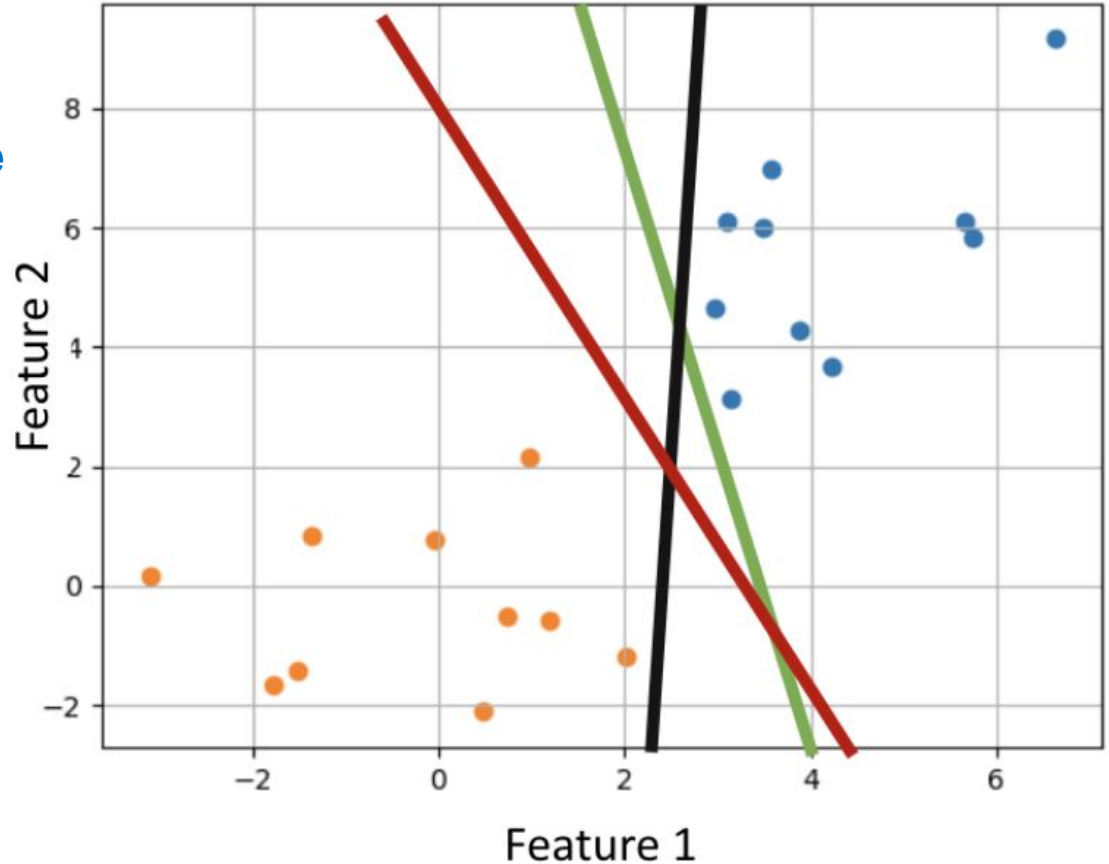
- Decision tree

Determining if a word is end-of-sentence: a Decision Tree



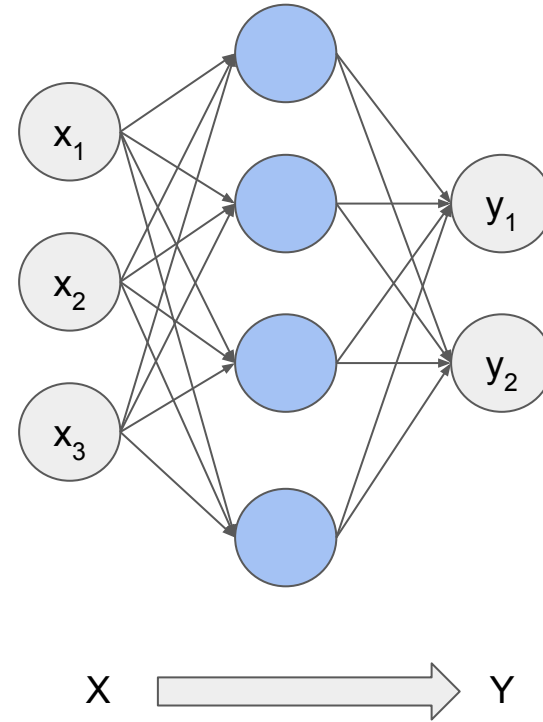
Other supervised classifiers

- **Decision tree**
- **Support Vector Machine**



Other supervised classifiers

- **Decision tree**
- **Support Vector Machine**
- **Neural Network**



How good is the model?

We define a **metric** to measure and compare accuracy.

- Precision

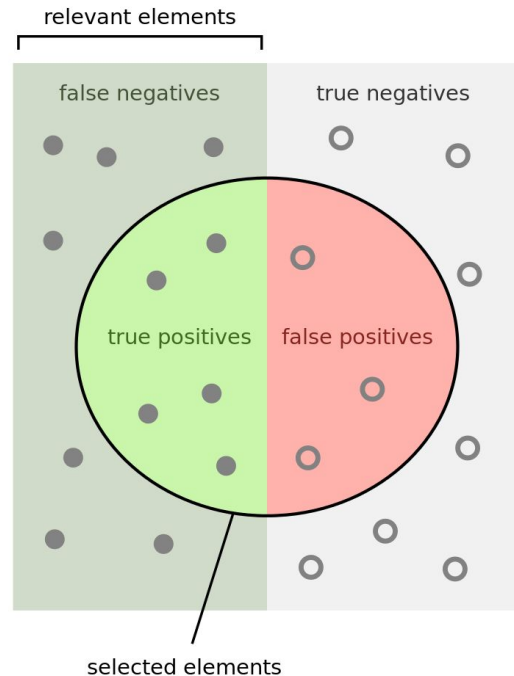
- Out of those tested positive, how many are truly positive?
- $TP / (TP + FP)$

- Recall

- Out of those truly positive, how many tested positive?
- $TP / (TP + FN)$

- F1

$$\frac{2}{\text{recall}^{-1} + \text{precision}^{-1}}$$



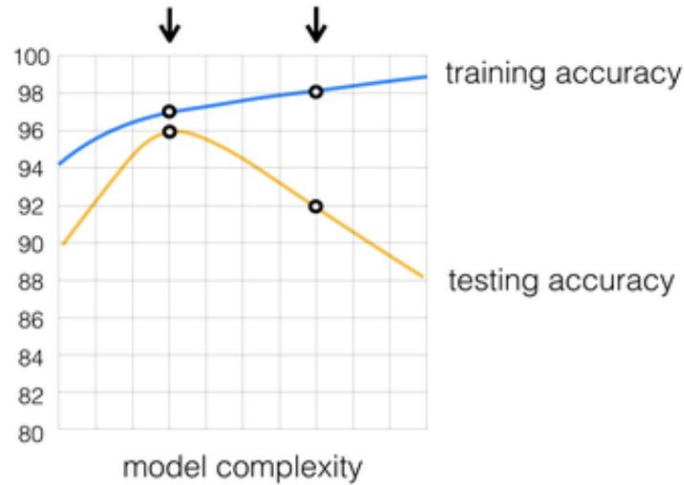
How many selected items are relevant?

$$\text{Precision} = \frac{\text{green semi-circle}}{\text{green semi-circle} + \text{red semi-circle}}$$

How many relevant items are selected?

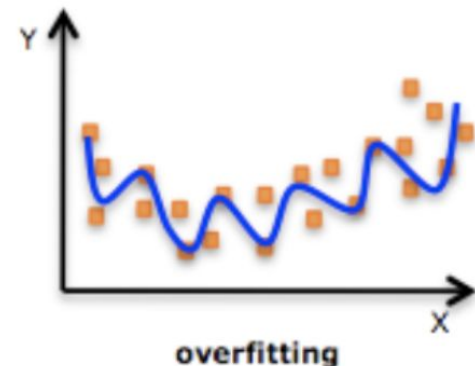
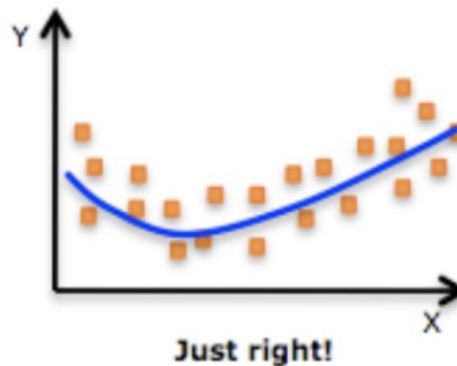
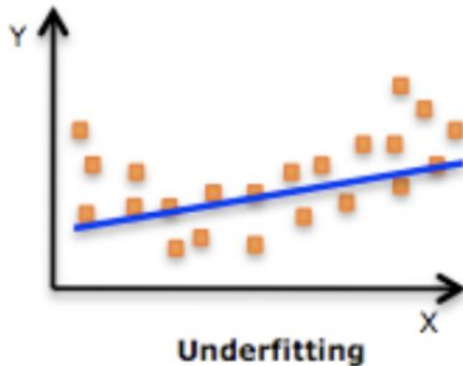
$$\text{Recall} = \frac{\text{green semi-circle}}{\text{green semi-circle} + \text{green rectangle}}$$

Overfitting



Solution: **Cross-Validation**

Split the training data into two non-overlapping sets. Train on one set, and measure performance on the other. Pick the model that does well on the data that you *didn't* train on.



Supervised models

| Model | When to use it? |
|---|---|
| KNN | <ul style="list-style-type: none">➤ Little / no training time, large prediction time➤ Given small / medium dataset size |
| Linear / polynomial regression | <ul style="list-style-type: none">➤ Linear / polynomial relationship between input and output➤ Small training time➤ Given small / medium dataset size |
| Logistic regression, SVM, decision tree | <ul style="list-style-type: none">➤ Categorical output➤ Given small / medium training time |
| Neural network | <ul style="list-style-type: none">➤ Large training time, large computer➤ Given large dataset size |

Note / life tip

- **Don't re-implement it yourself!**
 - Unless you are doing research on the method itself, you are trying to learn how it works, or you are coding in an obscure language where it isn't already implemented
 - The already implemented versions are widely used and tested

Note / life tip

- **Don't re-implement it yourself!**

- Unless you are doing research on the method itself, you are trying to learn how it works, or you are coding in an obscure language where it isn't already implemented
- The already implemented versions are widely used and tested

- **Use these common tools:**

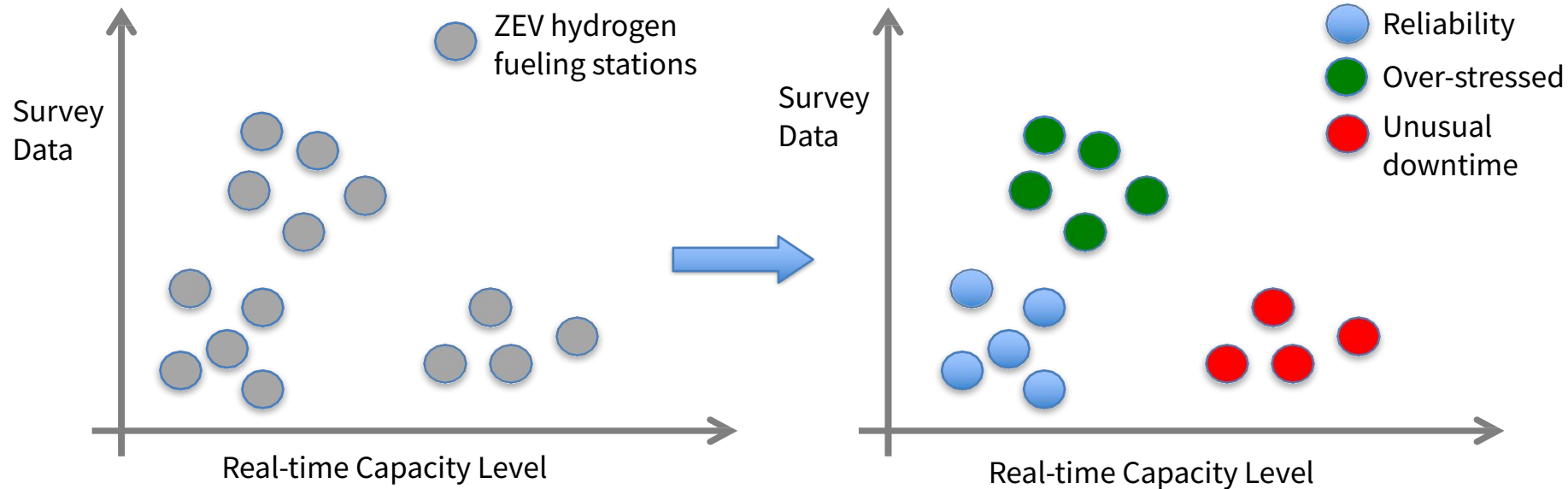
- [Scikit-learn](#) has most supervised and unsupervised methods you might need
- If you want to build a custom neural network, try using [Pytorch](#) or [Tensorflow](#)
- There are many task-specific libraries

Agenda

- What is machine learning?
- Supervised learning
- **Unsupervised learning**
- Reinforcement learning
- Generative models

Unsupervised Learning

Finding patterns in unlabeled data



Unsupervised Learning

Finding patterns in unlabeled data

Examples:

- Finding clusters
 - Customer segmentation (group customers so you can target advertising)
 - Finding user accounts that are all suspiciously similar
 - Group search results (or news / trending topics)
- Topic modeling (LDA)
- Figure out important features to use for supervised learning
- Learn vector representations for words / documents

DeSantis signs Florida bill banning offshore wind turbines >

AP

Florida Gov. Ron DeSantis signs a bill that strikes climate change from state law

9 hours ago



Callisto

Ron DeSantis signs bill scrubbing 'climate change' from Florida state laws

7 hours ago



The Washington Post

DeSantis signs bill scrubbing 'climate change' from Florida law

Yesterday · Anna Phillips



news.google.com

Clustering

1. Extract features from raw data

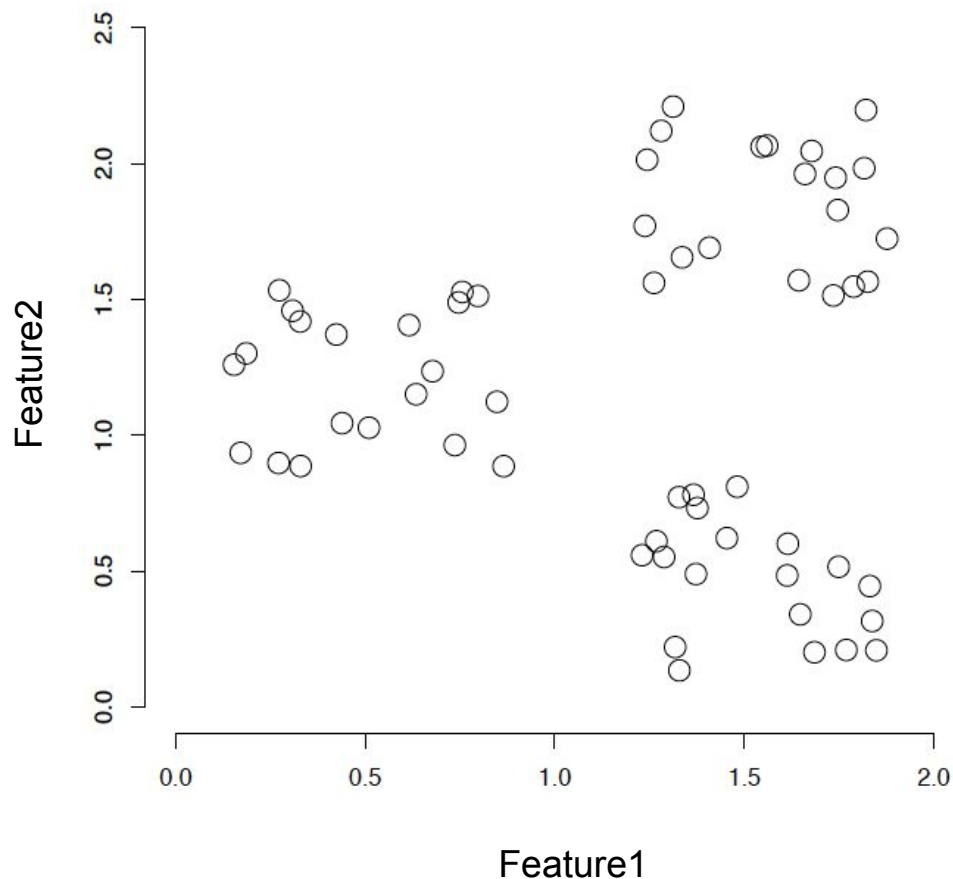
| Raw Data Item | Feature 1 | Feature 2 |
|---------------|-----------|-----------|
| Apple1 | 0.4 | 0.2 |
| Apple2 | 0.5 | 0.1 |
| Banana1 | 1.3 | 2.1 |
| . | . | . |
| . | . | . |
| . | . | . |

Clustering

1. Extract features from raw data

2. Find natural groupings

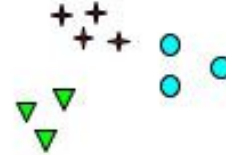
| Raw Data Item | Feature 1 | Feature 2 |
|---------------|-----------|-----------|
| Apple1 | 0.4 | 0.2 |
| Apple2 | 0.5 | 0.1 |
| Banana1 | 1.3 | 2.1 |
| . | . | . |
| . | . | . |
| . | . | . |



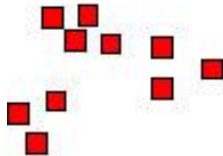
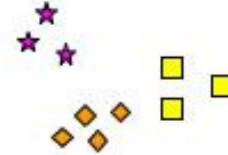
Clusters are ambiguous



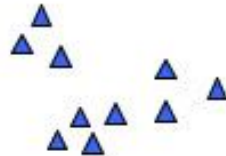
How many clusters?



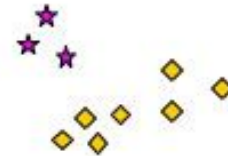
Six Clusters



Two Clusters



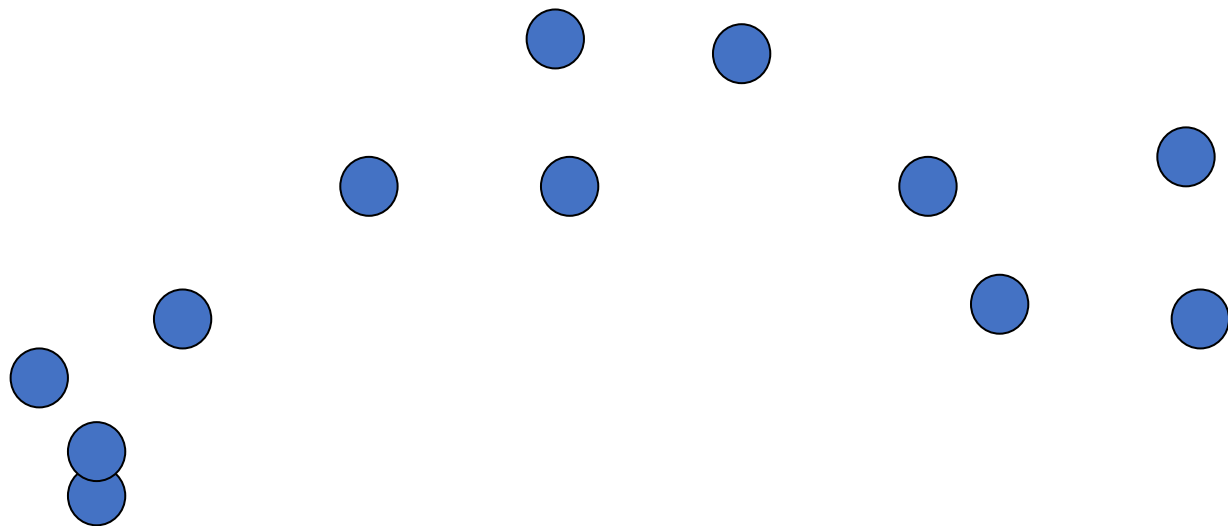
Four Clusters



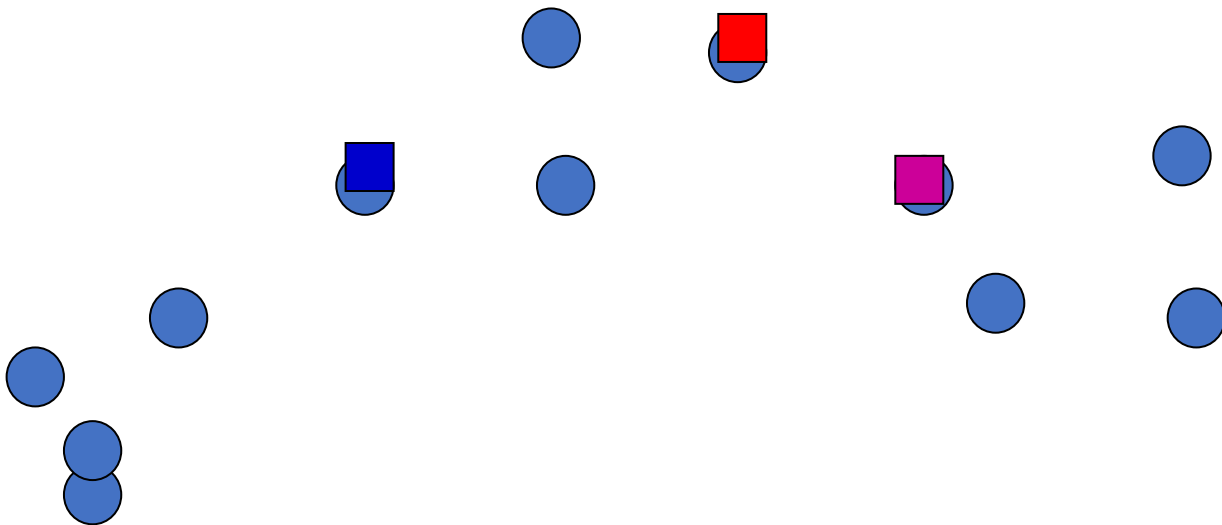
K-means

- Most well-known popular clustering algorithm
- Usually a baseline
- The algorithm:
 - Iterate until the clusters stop changing:
 - Assign / cluster each example to the closest center
 - Recalculate the centers as the mean of the points in their cluster

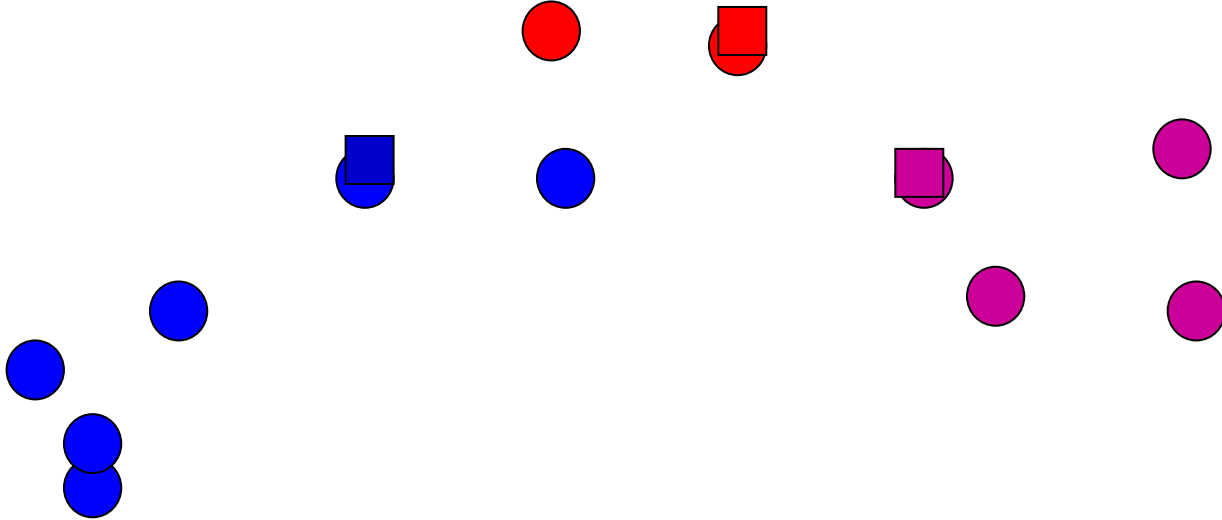
K-means example



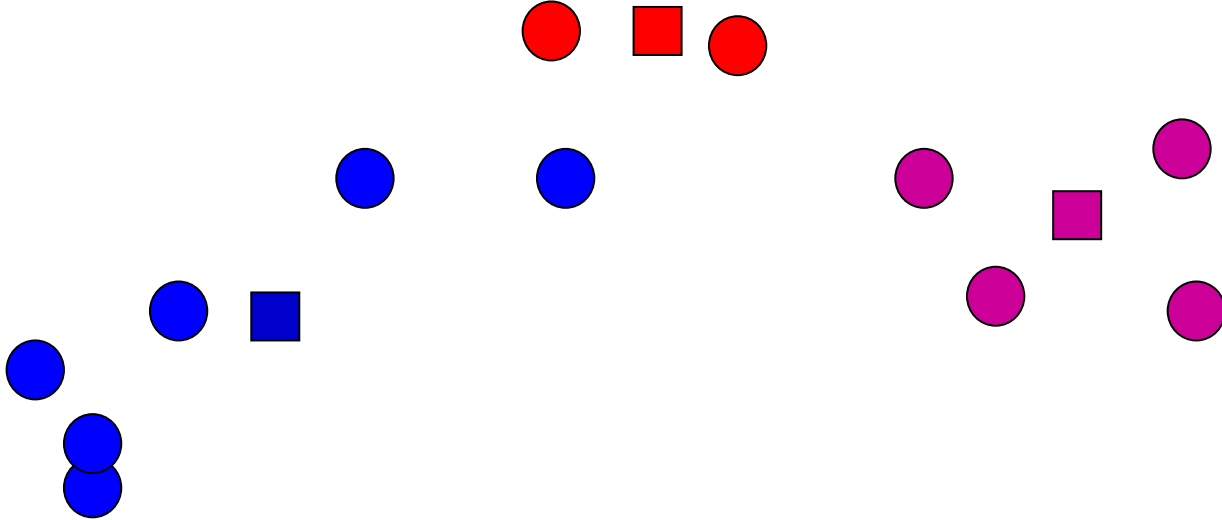
K-means example: initialize centers randomly



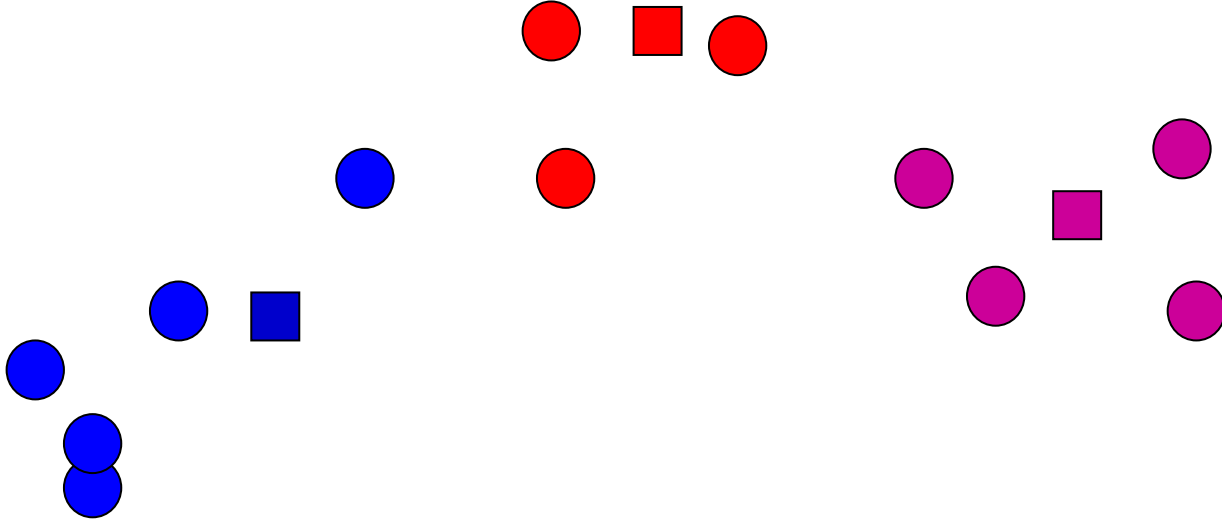
K-means example: assign points to nearest center



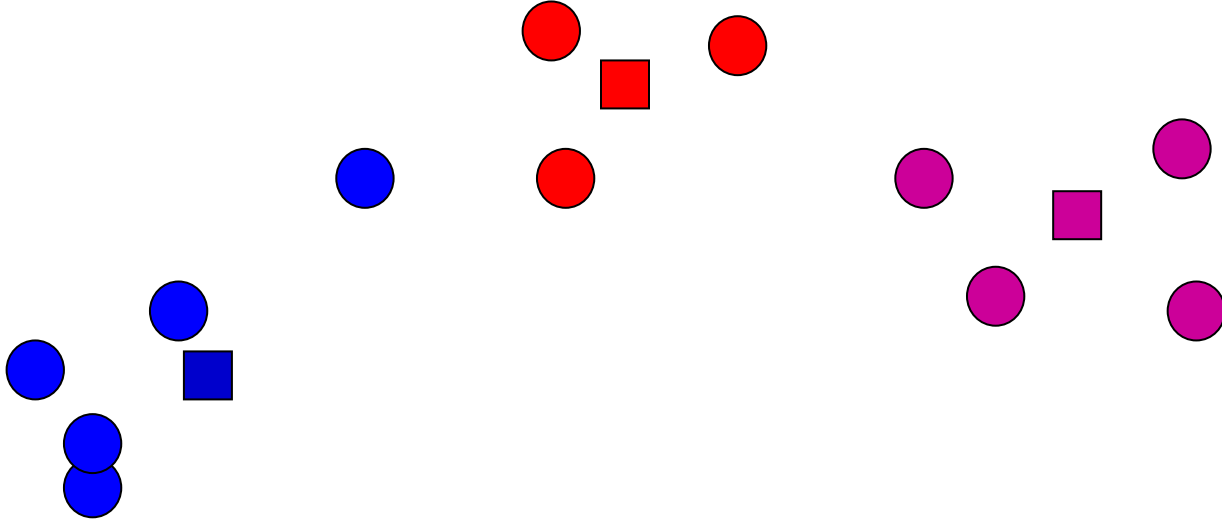
K-means example: recalculate centers



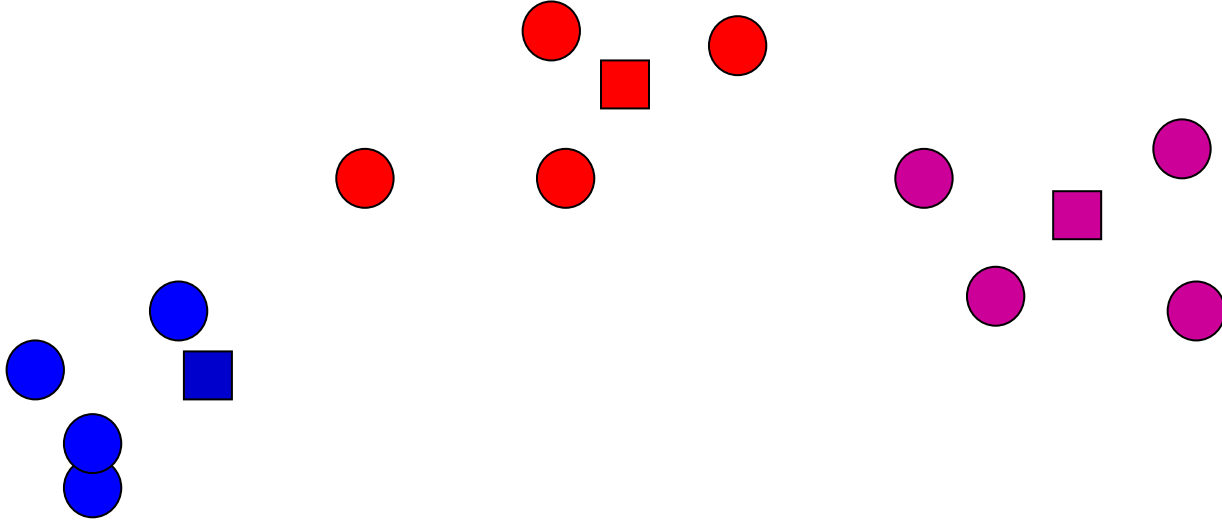
K-means example: assign points to nearest center



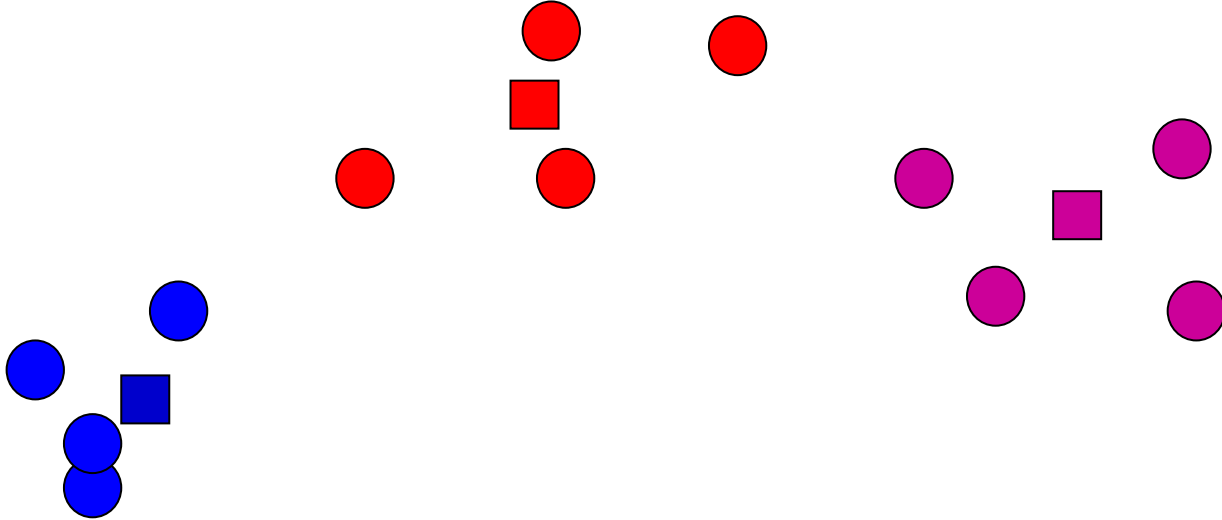
K-means example: recalculate centers



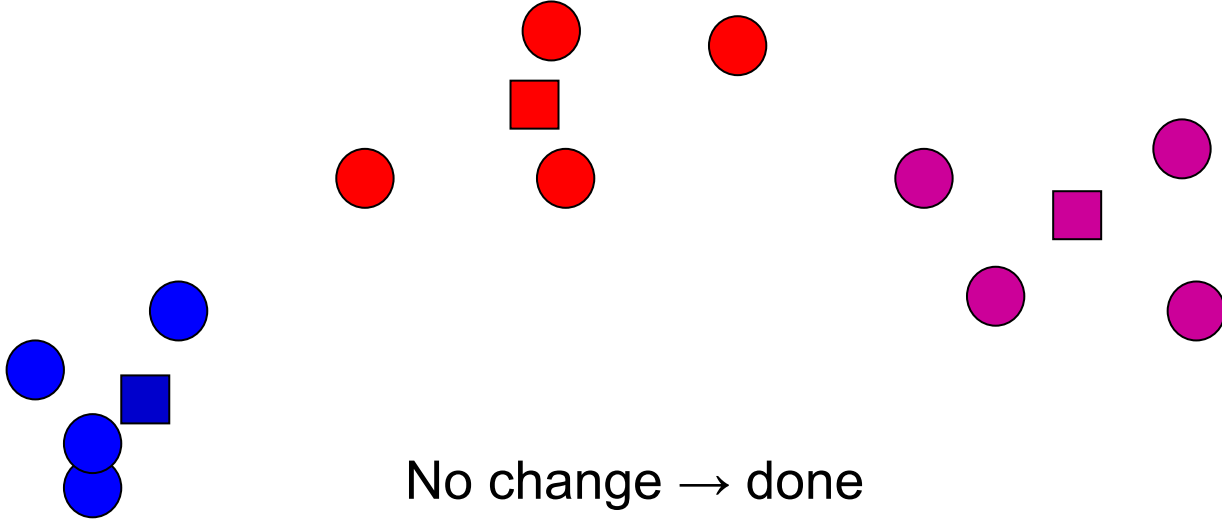
K-means example: assign points to nearest center



K-means example: recalculate centers



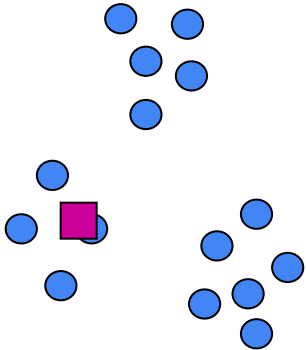
K-means example: assign points to nearest center





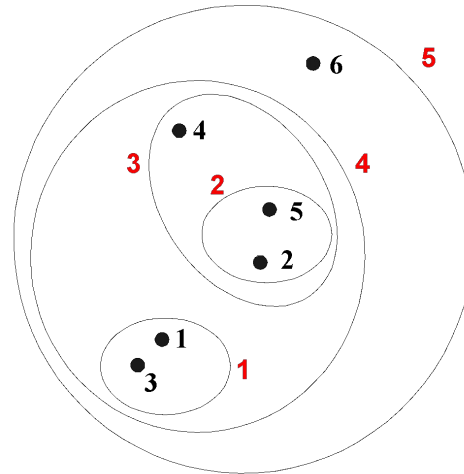
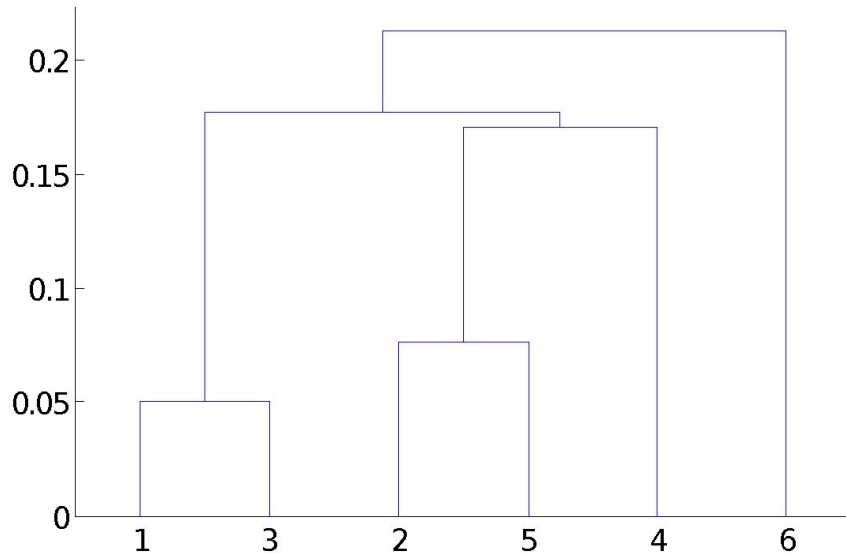
A Problem with K-Means: **Outliers**

- Centroid has to move all the way to the outlier
- Each outlier takes up an entire cluster



Hierarchical Clustering

- Produces a set of nested clusters organized as a hierarchical tree
- Can be visualized as a **dendrogram**
 - A tree like diagram that records the sequences of merges or splits



Clustering in Scikit-Learn

| Method name | Parameters | Scalability | Usecase | Geometry (metric used) |
|------------------------------|--|---|---|--|
| K-Means | number of clusters | Very large <code>n_samples</code> , medium <code>n_clusters</code> with <code>MiniBatch</code> code | General-purpose, even cluster size, flat geometry, not too many clusters | Distances between points |
| Affinity propagation | damping, sample preference | Not scalable with <code>n_samples</code> | Many clusters, uneven cluster size, non-flat geometry | Graph distance (e.g. nearest-neighbor graph) |
| Mean-shift | bandwidth | Not scalable with <code>n_samples</code> | Many clusters, uneven cluster size, non-flat geometry | Distances between points |
| Spectral clustering | number of clusters | Medium <code>n_samples</code> , small <code>n_clusters</code> | Few clusters, even cluster size, non-flat geometry | Graph distance (e.g. nearest-neighbor graph) |
| Ward hierarchical clustering | number of clusters or distance threshold | Large <code>n_samples</code> and <code>n_clusters</code> | Many clusters, possibly connectivity constraints | Distances between points |
| Agglomerative clustering | number of clusters or distance threshold, linkage type, distance | Large <code>n_samples</code> and <code>n_clusters</code> | Many clusters, possibly connectivity constraints, non Euclidean distances | Any pairwise distance |
| DBSCAN | neighborhood size | Very large <code>n_samples</code> , medium <code>n_clusters</code> | Non-flat geometry, uneven cluster sizes | Distances between nearest points |
| OPTICS | minimum cluster membership | Very large <code>n_samples</code> , large <code>n_clusters</code> | Non-flat geometry, uneven cluster sizes, variable cluster density | Distances between points |
| Gaussian mixtures | many | Not scalable | Flat geometry, good for density estimation | Mahalanobis distances to centers |
| Birch | branching factor, threshold, optional global clusterer. | Large <code>n_clusters</code> and <code>n_samples</code> | Large dataset, outlier removal, data reduction. | Euclidean distance between points |

Example: climate policy documents

<https://www.climatechange.ai/papers/neurips2022/59>

Given: Many companies' climate policy documents

Want to know: What is in these documents? Understand vague general categories

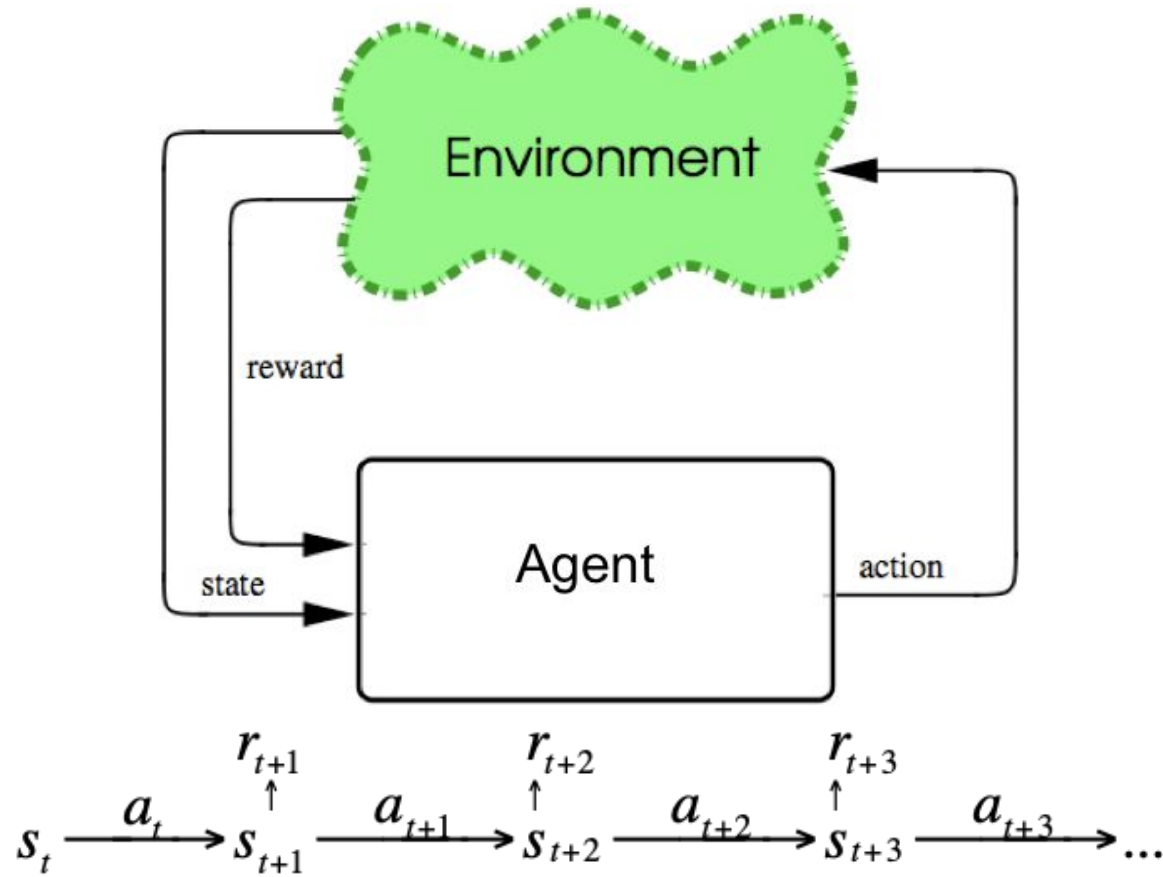
It's not labeled. **Which unsupervised algorithm might work best?**

- A. Flat clustering (KMeans)
- B. Hierarchical clustering
- C. Topic modeling (LDA)

Agenda

- What is machine learning?
- Supervised learning
- Unsupervised learning
- **Reinforcement learning**
- Generative models

Reinforcement Learning



Example rewards: PacMan

- One example:
 - 1 if you eat a pill
 - -10 if you get caught by a ghost
 - 2 if you eat a power pill or eat a ghost
 - 0 otherwise
- Another example:
 - -1 at every time step
 - 1,000,000 if you win the level



Exploration vs. Exploitation

- **Exploitation:** take good actions in each state already taken before to maximize reward
- **Exploration:** take a chance on actions that may have lower value in order to learn more, and maybe find true best action to later exploit

Need to balance the two!

Example: battery charging / discharging policy

<https://www.climatechange.ai/papers/iclr2024/16>

Given: batteries which can charge and discharge in a complex power grid

Want to know: an optimal charging / discharging policy

Which type(s) of learning could we use?

- A. Supervised learning
- B. Unsupervised learning
- C. Reinforcement learning

Agenda

- What is machine learning?
- Supervised learning
- Unsupervised learning
- Reinforcement learning
- **Generative models**

Generative language models predict the next word.

Please turn your homework...

- A. in
- B. over
- C. the
- D. agriculture

Large language models (LLMs)

- **Huge** and trained on **large amounts of text** (the internet)
 - Large emissions to train (<https://www.jmlr.org/papers/v24/23-0069.html>)
- Can “hallucinate” facts
- Reproduces the (social) bias from its training set (the internet)

When do we use generative models like GPT?

For each: yes or no?

- ?? To make decisions which we use without a human in the loop
- ?? To look up facts and information which we verify afterwards
- ?? To write code which we run without verifying
- ?? To help write text or code faster while verifying everything it writes

Key take-aways

- Use ML for tasks which:
 - Ask the computer to find and use patterns in data (with errors)
 - Are more cost-effective with ML
 - Have appropriate data
- **Supervised** vs. **unsupervised** vs. **reinforcement** learning
 - vs. **generative models**
- **Categorical** vs. **continuous** data
 - Images: each pixel is 3 continuous features (RGB)
 - Text: each word is a categorical feature
- Most things can be done in a couple lines of code using [Scikit-learn](#)
 - Make use of their [code examples](#)