

What is machine learning?

Let us consider the following set of (x, y) pairs:

$(1, 1)$, $(2, 4.06)$, $(2.5, 6.1)$, $(3, \text{csfs})$, $(3, 9)$, $(4, _)$

Now, our task is given a new x to find an estimate of the corresponding y . So if $x = 20$, how much is y ?

What is machine learning (cont'd)

- Machine learning aims to learn useful things from data collected from the world.
- Useful things:
 - Classifiers (return one of several discrete labels)
 - Regressors (return a continuous value)
 - Policies / behaviors / skills (return a recommendation for a next action)
 - Predictors (return a likely next item in a series)
 - Clustering (group items based on similarity)
 - Generators (create an item with specified properties)
 - ... etc

Relationship to artificial intelligence

- In this class, we will see machine learning as a subset of artificial intelligence $ML \subset AI$
- Not everything in AI needs ML, some things can be calculated without data
 - Theorem proving
 - Path planning
 - Game playing etc.
- Sometimes ML can help even in AI problems that technically don't need it
 - Eg. learning from previous chess games

Relationship to statistics

- Many statisticians believe $ML \subset STAT$
 - They have a case!
 - But statistics had been around for a long time and the explosion of ML/AI is recent.
- Without angering statisticians, we can say that **classical statistics** dealt with situations of scarce data, preferring simple, linear and understandable models, careful proofs of significance etc.
- Machine learning often deals with large amount of low quality data, and we learned to love huge, complex, nonlinear, difficult-to-understand models.
 - Because they work.

Unsupervised learning

- Let us say we have some data $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots\}$
 - We can collect this data from just observing the world!
 - There are no labels: **unsupervised**
- \mathbf{x} is **bold**, because it can be a vector
 - e.g. $\mathbf{x} = [x_1, x_2, x_3]$
 - But also: list of features, text, picture, video etc.

Unsupervised learning (cont'd)

- Can we learn something from this?
 - Group the data into similar clusters: **clustering**
 - Assume that all the data is drawn from a probability distribution $\mathbf{x} \sim P(x)$.
Try to learn what P is?
 - For a given \mathbf{x}_{test} check if it is extremely unlikely in P : **anomaly detection**
 - Generate new data by sampling from P : **generative adversarial networks, diffusion models, large language models** etc

NOTE: We will discuss some unsupervised learning in this class.

Supervised learning

- Let us say we have some data $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_3, y_3), \dots\}$
 - There are the labels y which needs to come from somewhere, for instance from a supervising human: **supervised**
 - Labeling is major pain in the neck, and we try many tricks to avoid it
- What can we learn from this?
 - Learn a function $f(\mathbf{x}) \rightarrow \hat{y}$ such that \hat{y} is a good approximation of the real y (which we don't know)

NOTE: Most of this class is about supervised learning.

Reinforcement learning

- The setting
 - Consider an **agent** A that is in a **world** in **state** s .
 - The agent takes an **action** a , which changes the world into a state s'
 - At the same time the agent receives a **reward** r
- The problem:
 - How should the agent behave / what actions should it take to maximize its rewards over time?

NOTE: We will **not** discuss some RL in this class.