# Understanding Machine Learning

## chapter1: What is machine learning?

### What is a machine learning model?

ML model is a statistical representation of a real-world process based on data, like how we organize cats, dogs...

### The 3 types of machine learning

1) Reinforcement learning

2) Supervised learning

3) Unsupervised learning

**RL:**

used for deciding sequential actions, like a robot choosing a path, or next move in a chess game

=> It is not as common as the others

=> It uses complex mathematical like game theory

**SL and UL**

Supervised learning and Unsupervised learning are the most common types:

the main difference is in the training data, one is labeled and the other one is not

**Supervised learning**

Heart disease; true and false

=> The training data is labeled (true and false )

=> The target is heart disease

=> The label is true /false

=> Observations or examples are the lines of the data ( of each patient for example ), this is the data that the model will learn from

=> age, sex, cholesterol... Those **features** are the characteristics that help us to predict

**General Notes**

ML helps us analyze many features at once, even the ones we are not sure about, and find relationships between features

we input labels and features as data to train the model

once training is done, we can give the model new input, in our case new patient

**Unsupervised learning**

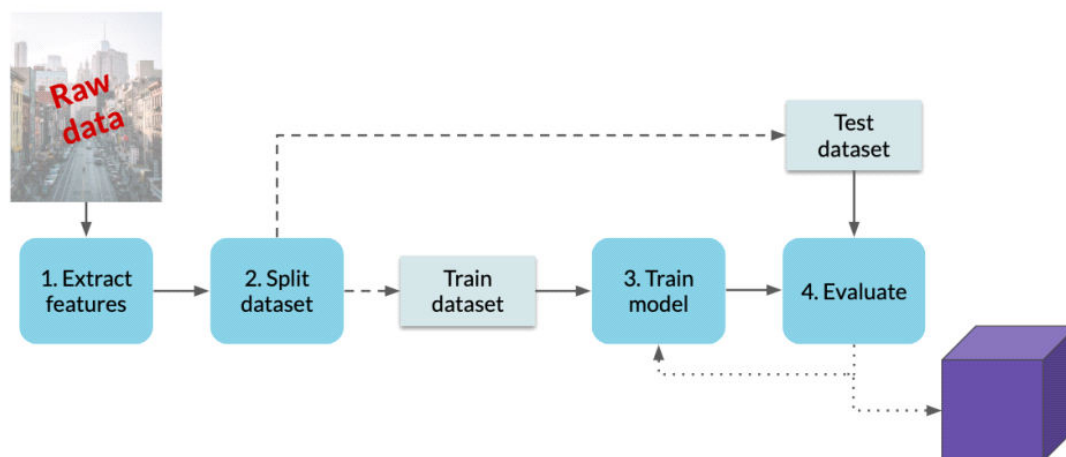We don't have labels, we only have features

=> We do anomaly detection

=> clustering

we consider just the people who have heart disease and we try to apply unsupervised learning to them to see what kind of treatment each group responds to.

**machine learning Workflow**



**step1: Extract features**

=> features that affect the target

**step2: Split data**

* train dataset

* test dataset

**step3: Train the model**

there are different types of models, from Neural Networks to Logistic Regression

**step4: Evaluation**

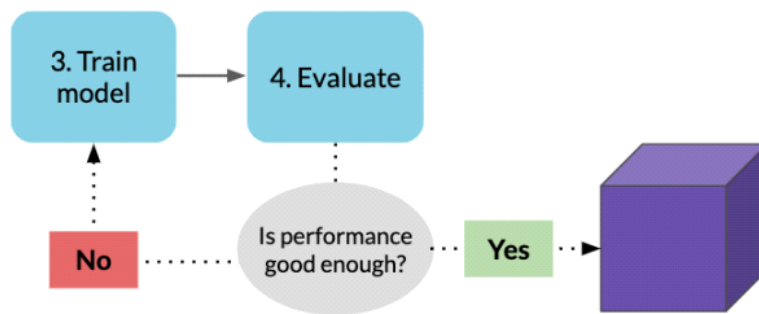*Here we use the test dataset( unseen data, never used )

*Many ways to evaluate:

What is the average error of the predictions?

What percent of apartments did the model accurately predict within a 10% margin?
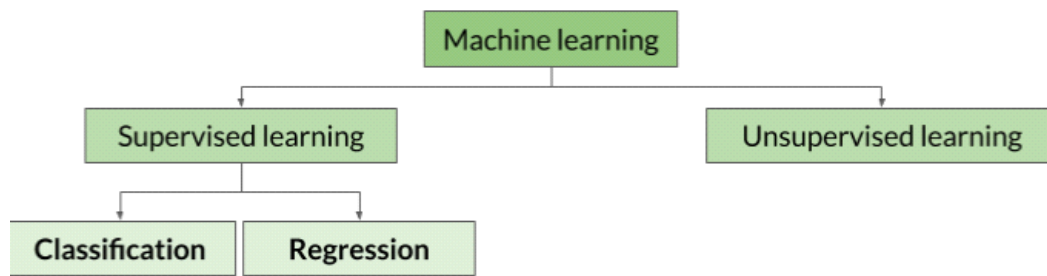
*If not, tune the model and re-train it:

e.g., change the model's options, add/remove features

# Step 4: Evaluate



## chapter2: Machine Learning Models

### 1. Supervised Learning

Machine learning

Supervised learning — Unsupervised learning

Classification — Regression

supervised learning => labeled data

## 1.1 Classification

**assigning a category**

we are predicting a discrete variable ( a variable that can only take a few different values; yes/no - red, white, orange .... )

a variable that can take only a few different values

we can use a **support vector machine** (its a line separating our points) :

* it's a linear classifier

* can be polynomial

## 1.2 Regression

**assigns a continuous variable** (a variable that can take any value like rent price, stacks ....)
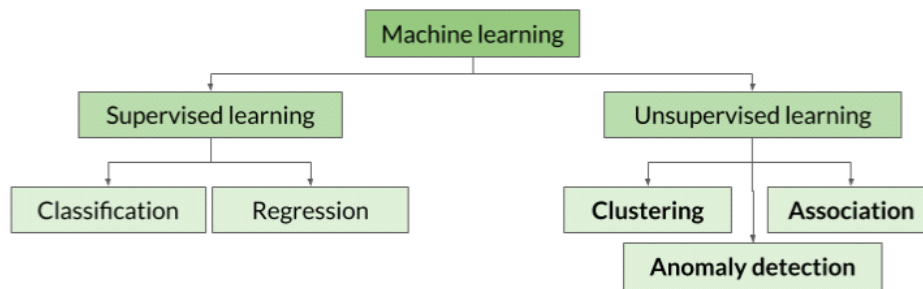
**Note:**

it's up to you if you want to work with regression or classification for example predict whether the weather is cold or warm(classification) or predict the exact temperature (37° for example)

## 2. Unsupervised Learning

no target column

checked the whole dataset to find a pattern(like self-driving cars. )

it can be divided into **clustering, association, and anomaly detection**

## 2 .1 Clustering

identifying groups in datasets that have stronger similarities than the rest of other groups

### 2.1.1 K-means :

specify the number of clusters

### 2.1.2 DBSCAN :

density-based spatial clustering of apps with noise: specify what constitutes a cluster like the minimum number of observations in one cluster

## 2.2 Association

finding relationships between observations

finding events that happen together

=> often used for market basket analysis (u watched this video, u might like to watch these as well )

## 2.3 Anomaly detection

Detecting outliers

=> observations that differ strongly from the rest (I guess cancer )

**Note:**

finding if x belongs to a group => means clustering

## 3. Evaluating performance

**for supervised learning**

the first thing to look for when evaluating is **overfitting**:

  it happens when :

* performs great on training data

* Performs poorly on testing data

* Model memorized training data and can't generalize learnings to new data

=> Use testing set to check model performance

we could measure model performance using **accuracy**

**accuracy = correctly classified observations/ all observations**
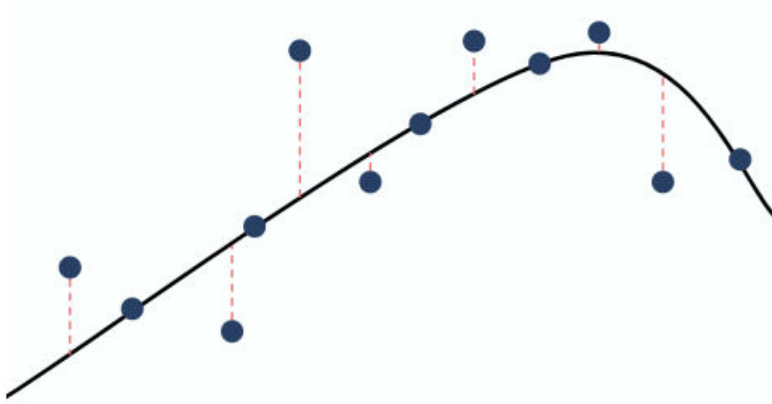
### 3.1 Evaluating classification: Confusion matrix



**sensitivity = true positives/ true positives + false negatives**

**specificity = true negatives/ true negatives + false positives**

=> **confusion matrix is    for classification problems**

### 3.2 Evaluating regression:

we want to make a difference between the actual value and the predicted line

**error = distance between the point (actual value) and line (predicted value)**

there are many ways to **calculate** this; e.g **root mean square error**

**for unsupervised learning**

this happens just for supervised learning because there are predicted values we can compare to while for unsupervised learning we can not use these methods

**4. Improving performance**

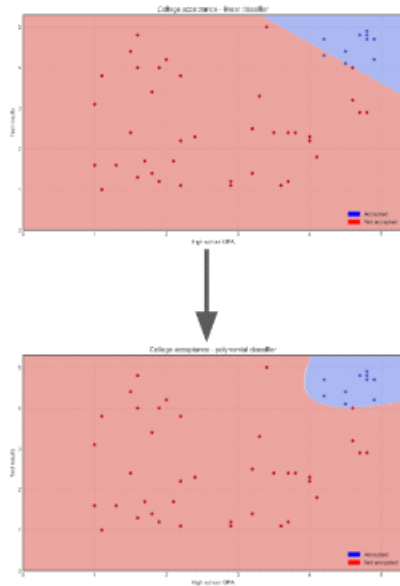In order to improve our model if we are not satisfied with its performance we can use one of these :

**4.1 Dimensionality reduction:**

reducing the number of features for those possible reasons;

  * some of the features might be irrelevant

* some of the features might be highly correlated, if you use one, it's like you used another

*we can collapse multiple features in just one( like weight and height to a BMI)

**4.2 Hyperparameter tuning:**

like changing the kernel from linear to poly

## 4.3 Ensemble methods:

A technique that combines several models in order to produce one optimal model

in the case of classification; we pick the majority

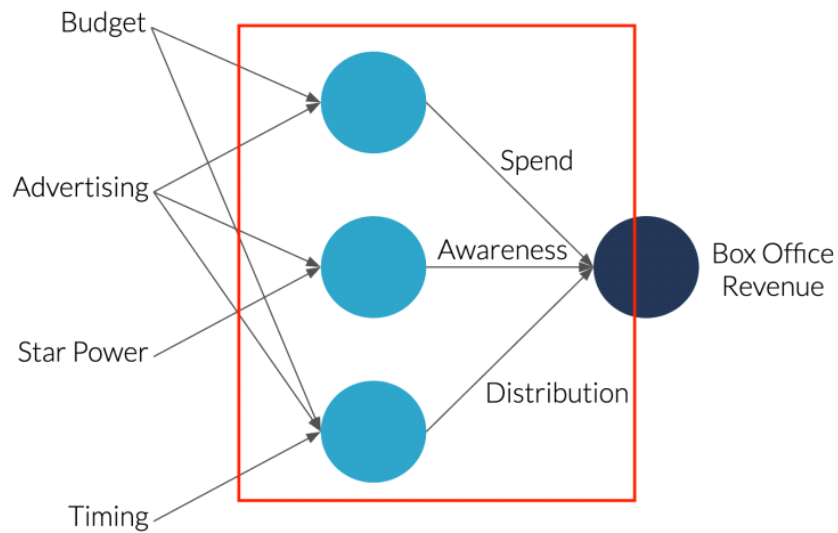in case of regression, we pick the average

# chapter3: Deep Learning

## 1. what is deep learning

* AKA: Neural Networks

* Basic unit: neurons (nodes)

* A special area of Machine Learning

* Requires more data

* Best when inputs are not structured like images or text
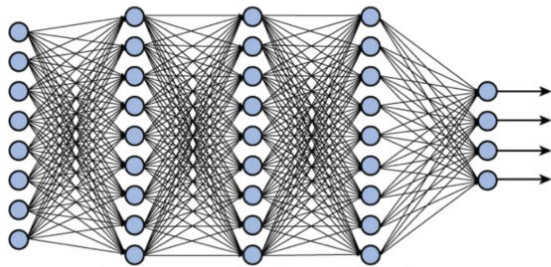
* we use it when we don't have much knowledge

**example:**

# Predicting box office revenue



# Deep learning



- Neural networks are much larger
- Deep learning: neural network with many neurons
- Can solve complex problems

**2. what is computer vision**

helps computers see and understand the content of digital images

**steps that happen inside the neural network:**

* images are transformed into numbers

* pixels are fed into NN

* neurons will learn to detect edges

* neurons will learn more complex objects like wheels, doors, and windows

* image is classified as a car or a truck

## 3. what is natural language processing

the ability for computers to understand the meaning of human language

n-gram: sequence of words

Applications of NLP :

  Language translation, Chatbots, Personal assistants, Sentiment analysis...

**=> There are 2 types of deep learning; computer vision and natural language processing**

## 4. limits of machine learning

**data quality:**

the output quality depends on the input quality(Garbage in, garbage out)

high-quality data requires:

* data analysis

* review of outliers

* domain expertise

* documentation

**explainability**:

often machine learning models are considered black boxes