

Introduction to Machine Learning

Valentina Giunchiglia

“Machine Learning: a computer is able to learn from experience without being specifically programmed”

Self-driving cars



A smartphone is shown from a side-on perspective, its screen displaying the words "Voice Recognition" and a digital clock showing "09:21". A bright blue beam of light emanates from the bottom of the phone's screen, pointing towards a stylized profile of a person's head. The person's mouth is open as if speaking. Hexagonal icons representing various technologies, such as AI, IoT, Cloud Computing, and Big Data, are scattered in the air around the phone, suggesting a futuristic or advanced technological environment.

Speech recognition

Face recognition



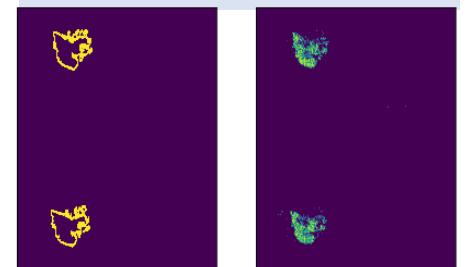
Character Recognition

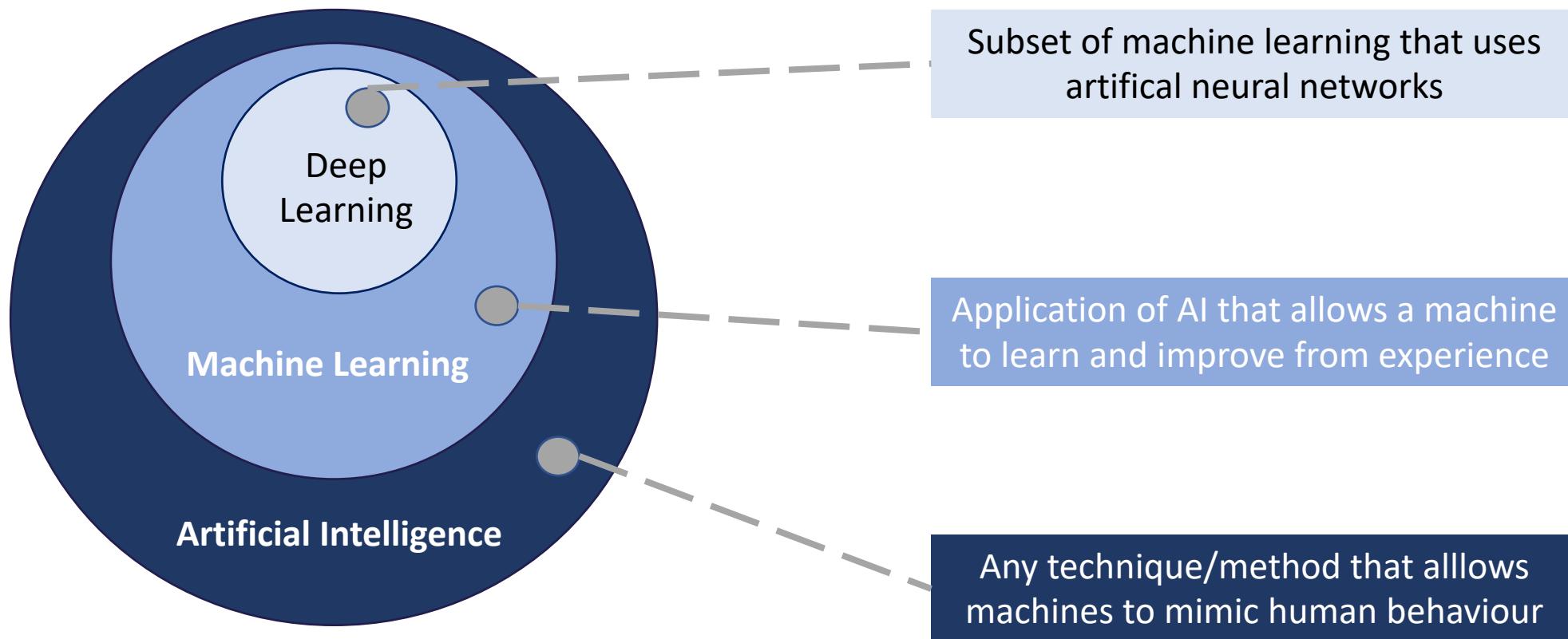
Content ads



Typing prediction

Biomedical Research

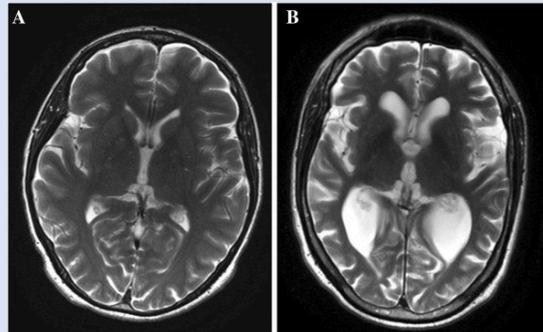




Three main **types** of Machine Learning

Supervised

ML approach that uses labeled data to learn and predict outcomes



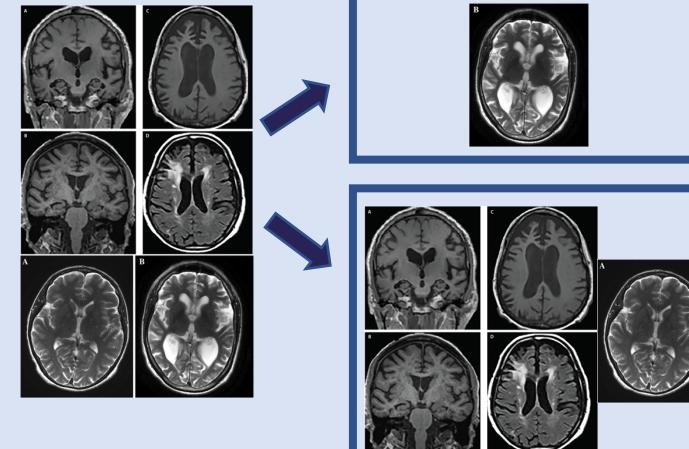
Alzheimer

Healthy

Task driven

Unsupervised

ML approach that uses unlabeled data



Reinforcement Learning

ML approach that allows an AI-driven system to learn through trial and error using feedback from its actions

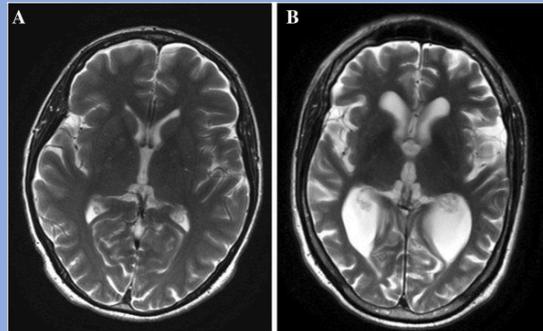
Beyond the scope of this lecture

Learning from mistakes

Three main **types** of Machine Learning

Supervised

ML approach that uses labeled data to learn and predict outcomes



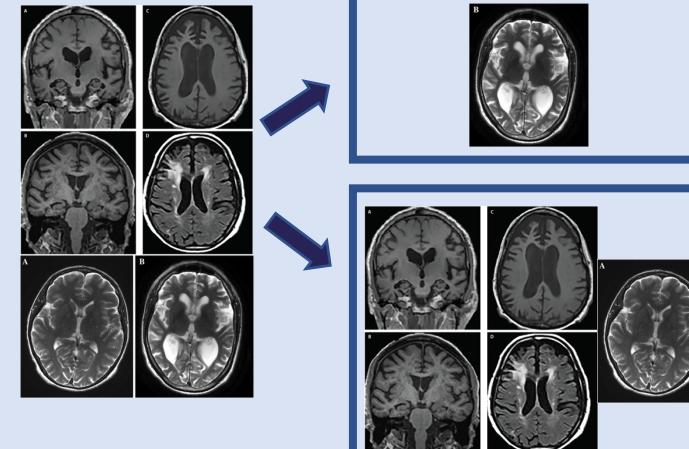
Alzheimer

Healthy

Task driven

Unsupervised

ML approach that uses unlabeled data



Data driven

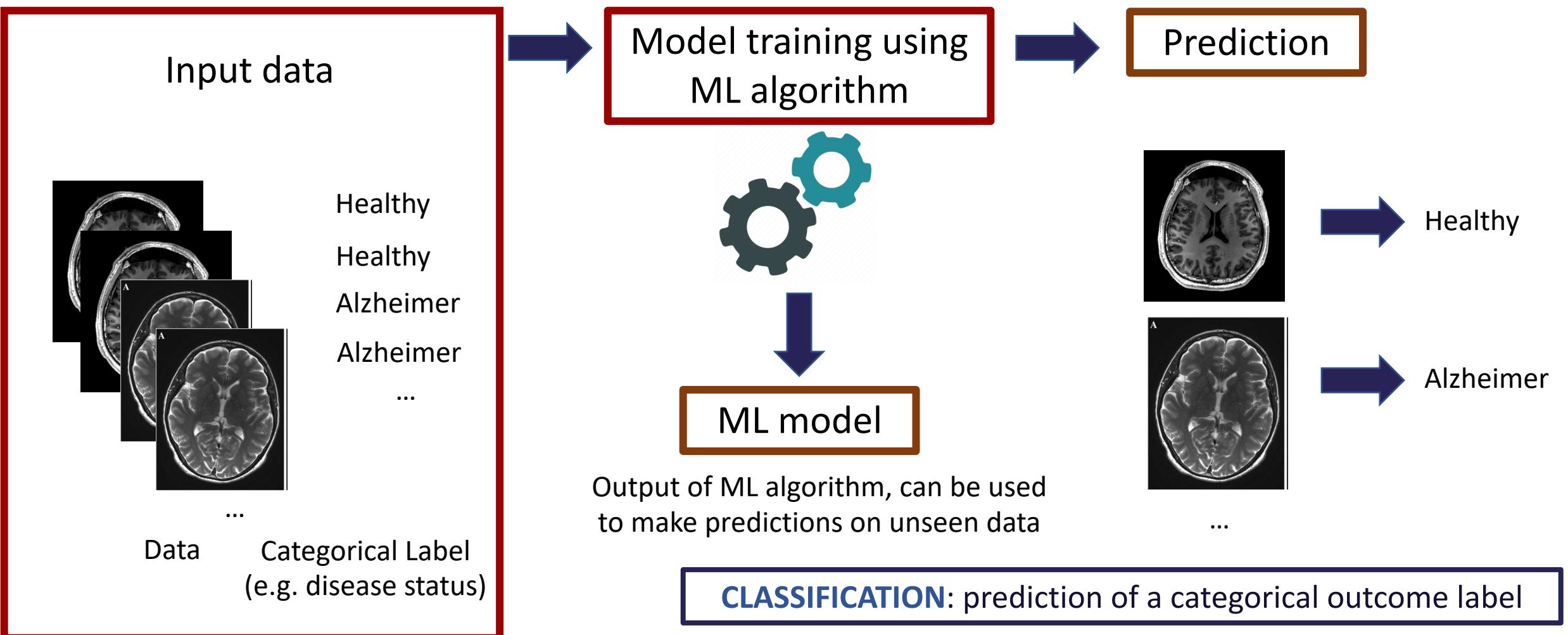
Reinforcement Learning

ML approach that allows an AI-driven system to learn through trial and error using feedback from its actions

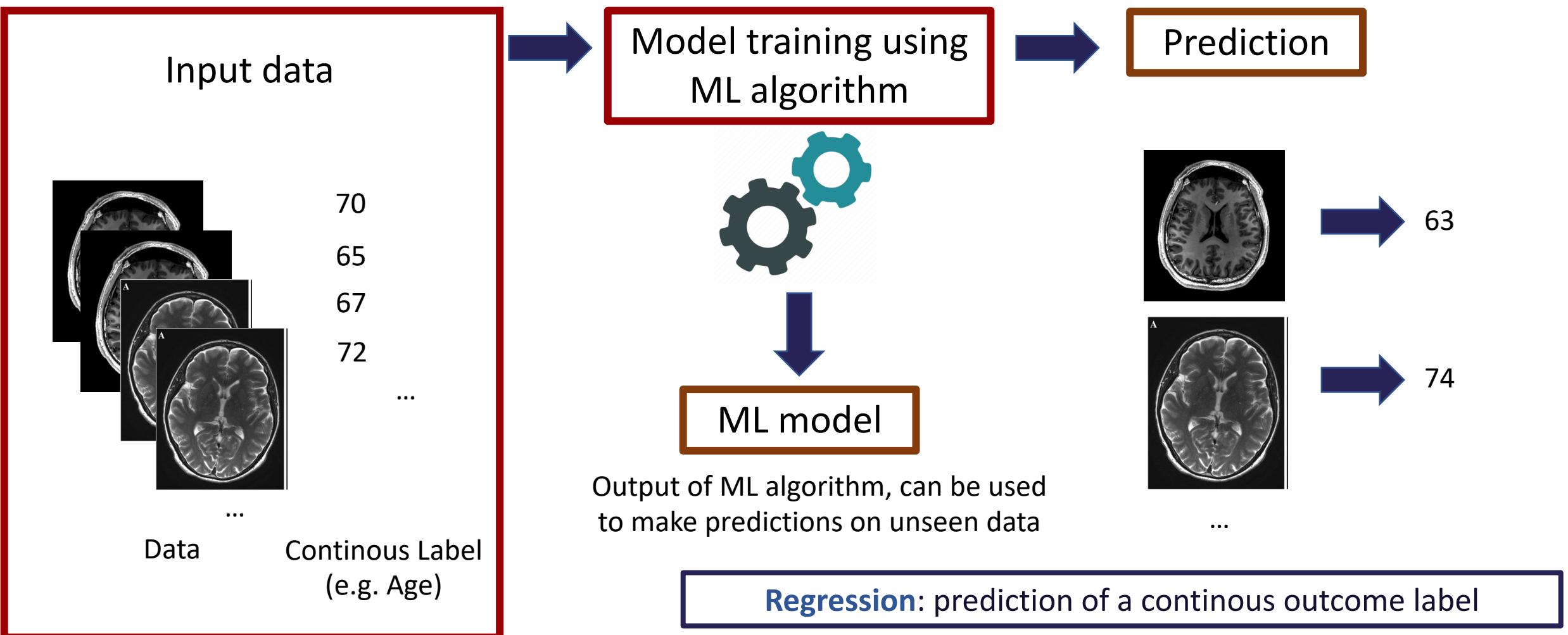
Beyond the scope of this lecture

Learning from mistakes

Supervised learning can be divided into **CLASSIFICATION** and REGRESSION

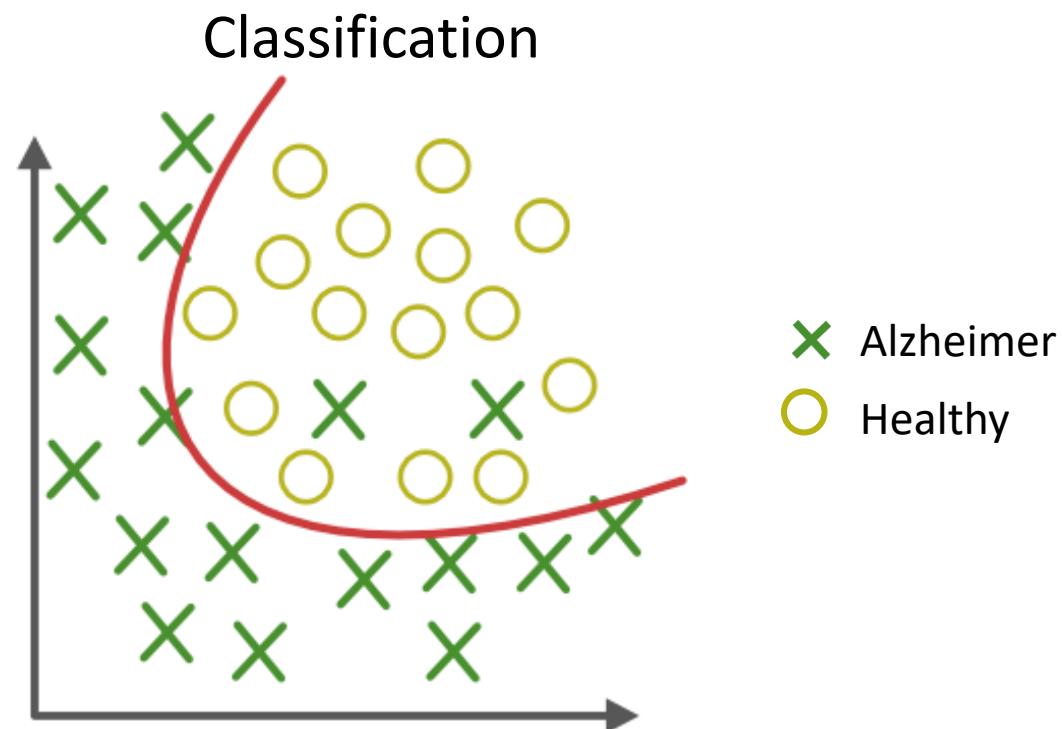


Supervised learning can be divided into CLASSIFICATION and REGRESSION

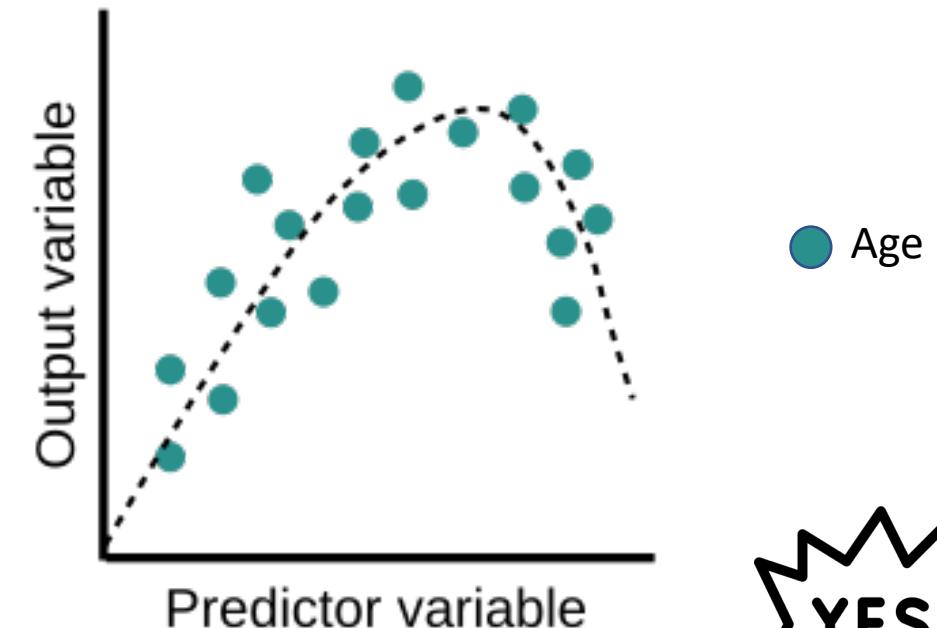


What is the ML model trying to learn? ➔

The ML model tries to learn what is the best line (in red) to separate the different classes (**CLASSIFICATION**) or to fit the data (**REGRESSION**)



Regression

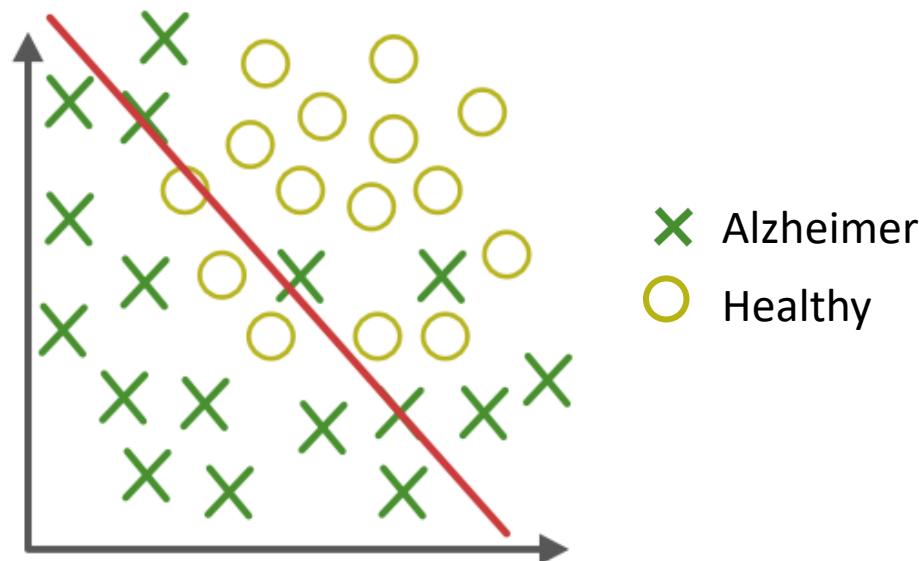


In real life there is NEVER a perfect split or fit to the data and the best line is hard to find ➔

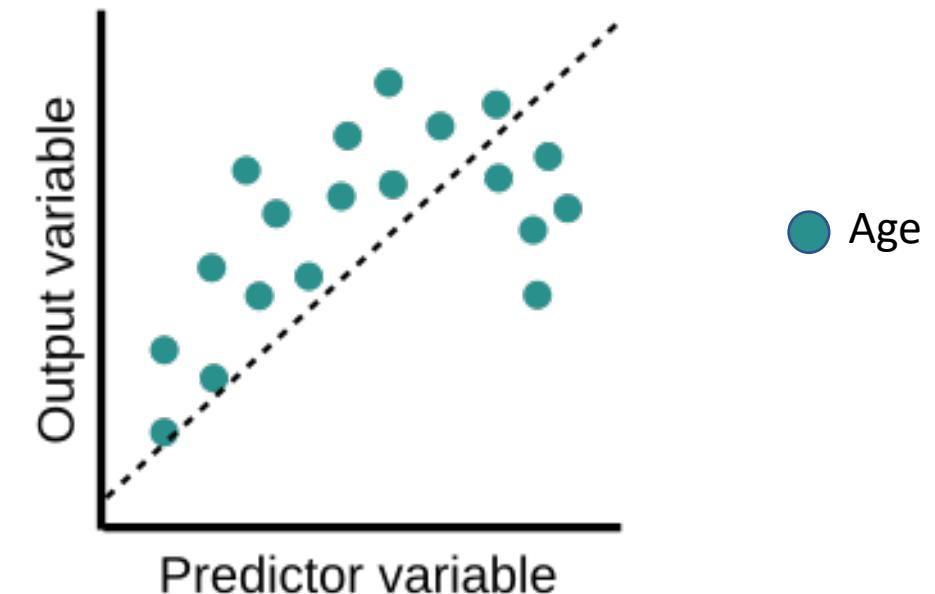
Can your model be wrong?

UNDERFITTING: the model cannot properly learn how to predict the categorical or continuous label and performs poorly both on the data on which it was trained and on unseen data

Classification



Regression

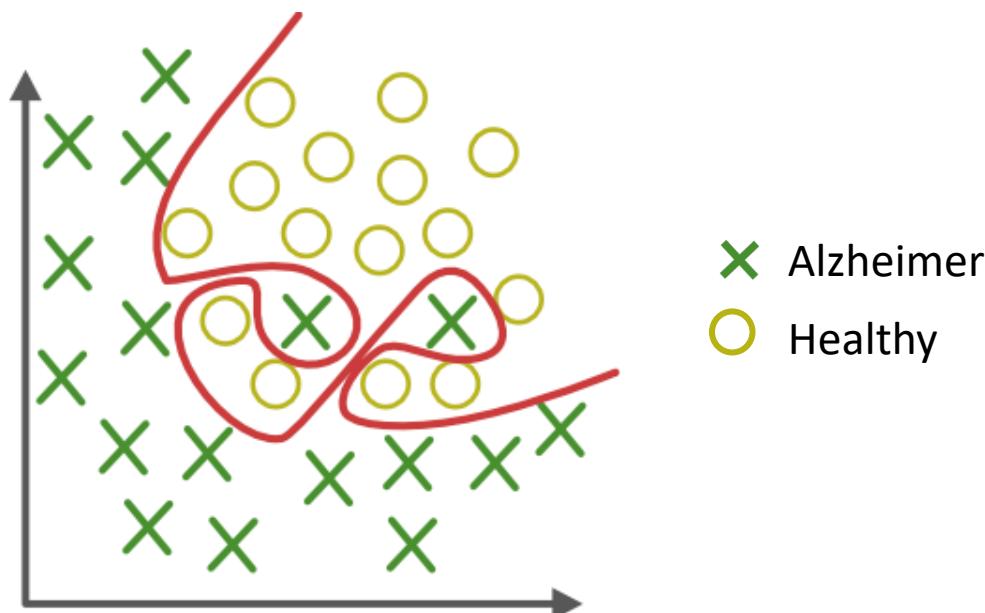


Potential **causes**:

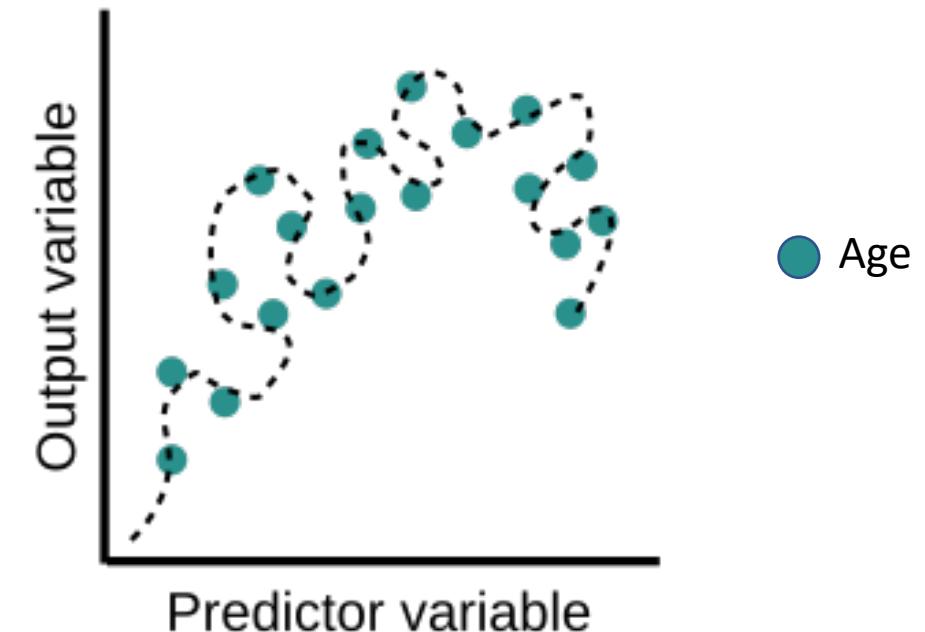
- Not enough training samples
- Features provided are not enough

OVERFITTING: the model learns the data on which it was trained *too well* and performs poorly on unseen data

Classification



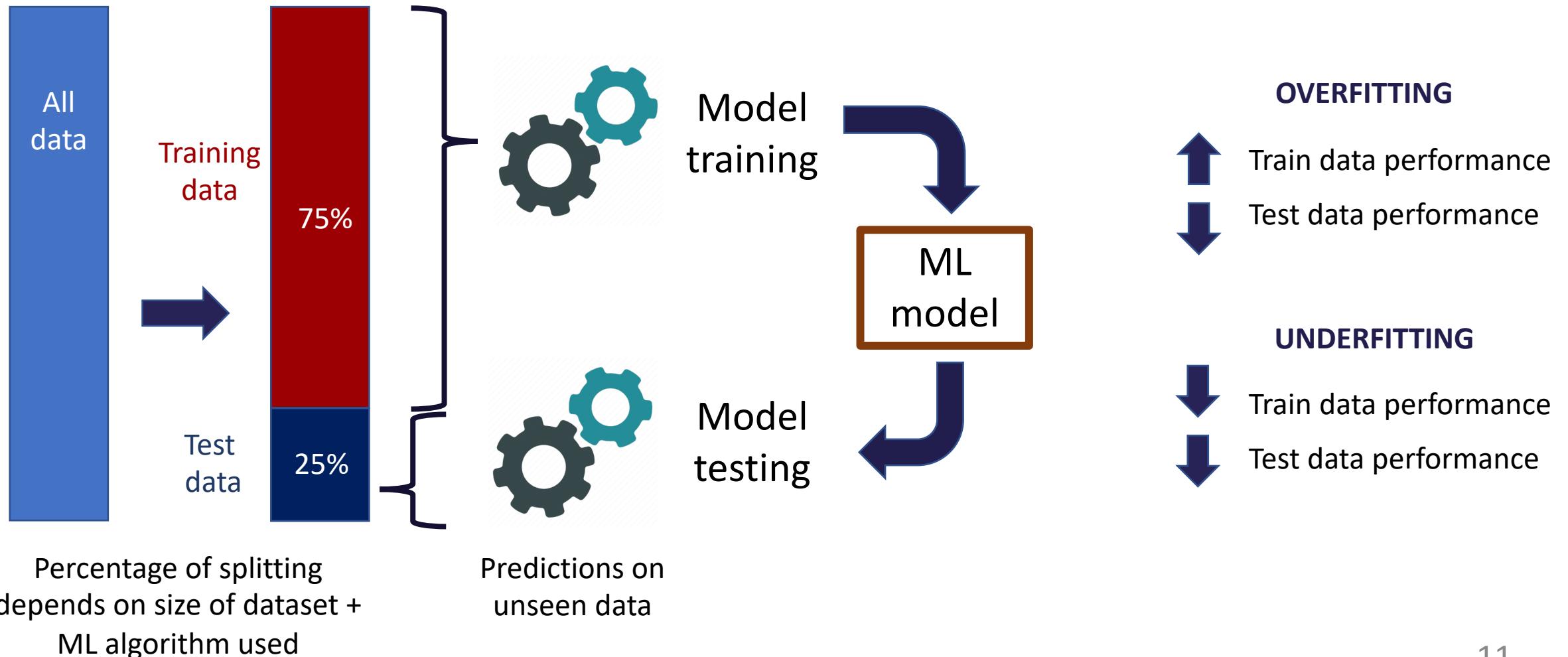
Regression



Potential **causes**:

- ✗ Wrong type of model
- ✗ Too many features in the input

The standard approach to test whether the model is underfitted/overfitted is to train it on a part of the data and test it on an hold-out dataset





Splitting in train/test is RANDOM, but based on the splitting performance can change

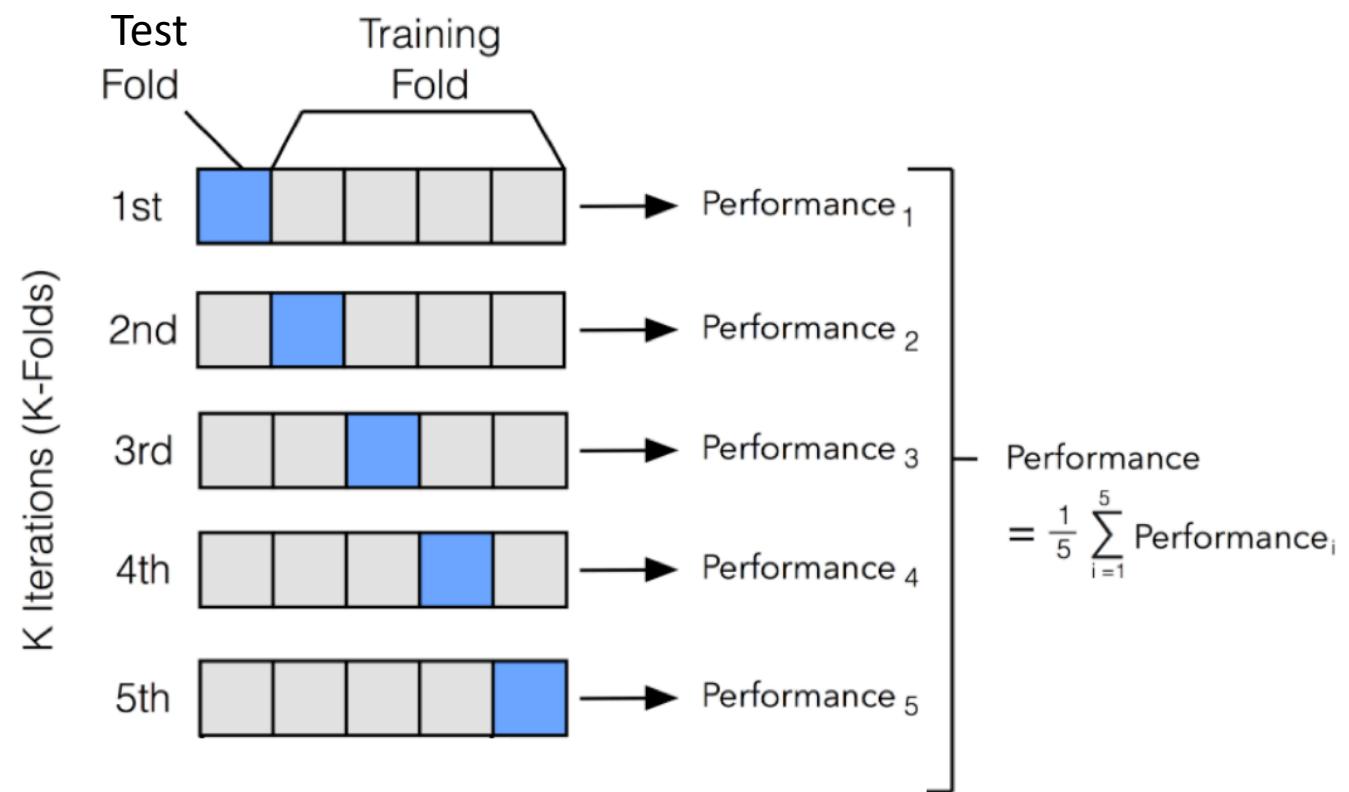


Some samples are EASIER/HARDER to predict than others



K-Fold cross validation: compute the train/test split K times, and calculate the performance on K different test sets

→ Final performance = Average



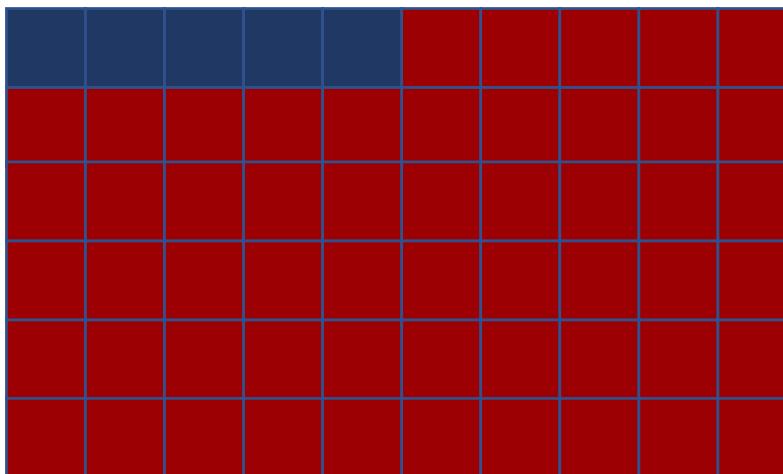
CLASSIFICATION

$$\text{Accuracy} = \frac{\text{Number of corrected predictions}}{\text{Total number of predictions}}$$

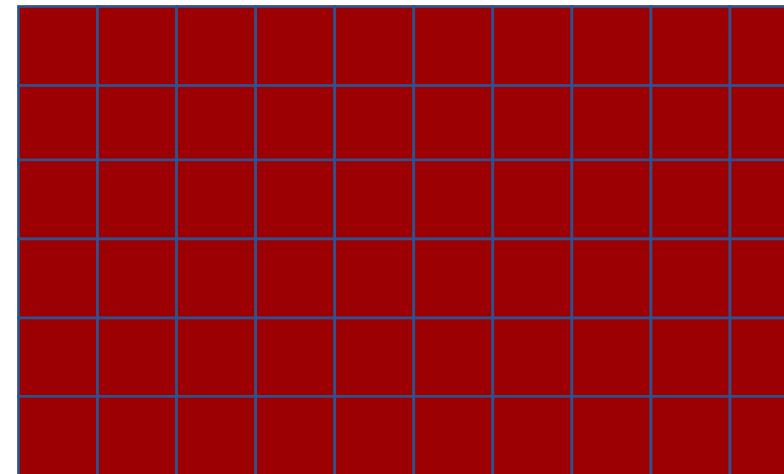


Not robust to **CLASS IMBALANCE** (i.e.
unequal number of samples in each class)

Actual Labels



Predicted Labels



Accuracy:
 $55/60 = 0.92$



High accuracy **AND**
bad model



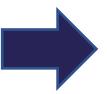
Healthy



Alzheimer

CLASSIFICATION

CONFUSION
MATRIX



Robust to CLASS
IMBALANCE

TP: 0	}	$TPR = \frac{TP}{Actual\ Positive} = \frac{TP}{TP + FN}$
FP: 0		$FNR = \frac{FN}{Actual\ Positive} = \frac{FN}{TP + FN}$
TN: 55		$TNR = \frac{TN}{Actual\ Negative} = \frac{TN}{TN + FP}$
FN: 5		$FPR = \frac{FP}{Actual\ Negative} = \frac{FP}{TN + FP}$

Predicted Label

Actual Label

		Alzheimer	Healthy
Alzheimer	Alzheimer	True positive	False positive
	Healthy	False negative	True negative

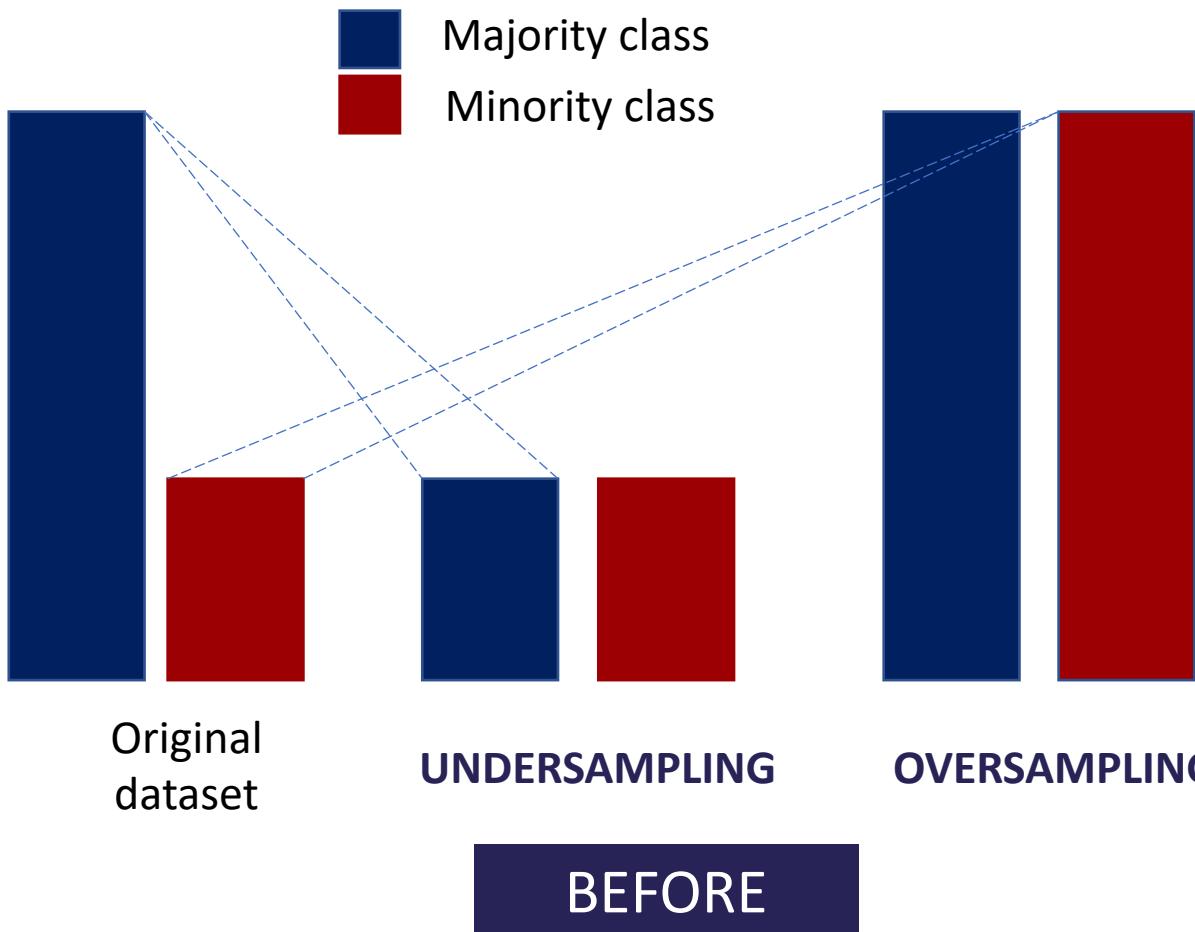
Class imbalance



Class imbalance should be addressed
DURING or BEFORE model training



Model will learn to predict everything as
MAJORITY CLASS (i.e. class with most samples)

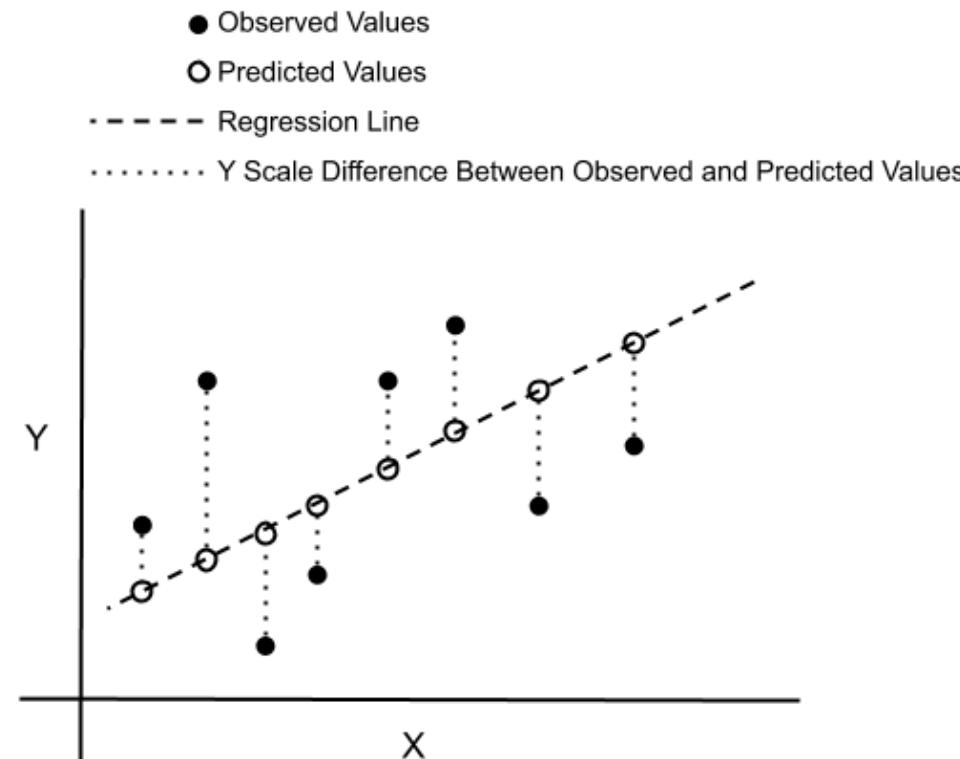


Model PENALISATION



Influence the training of the model
so that the model makes fewer
mistakes with the minority class

REGRESSION



We cannot calculate accuracy in case of regression models

MEAN SQUARED ERROR: measures how close a fitted regression line is to the actual data

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

n = number of samples

Y_i = observed values

\hat{Y}_i = predicted values

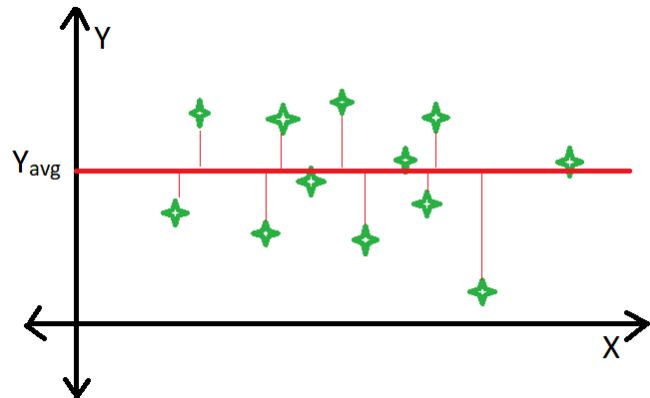
MSE = mean squared error



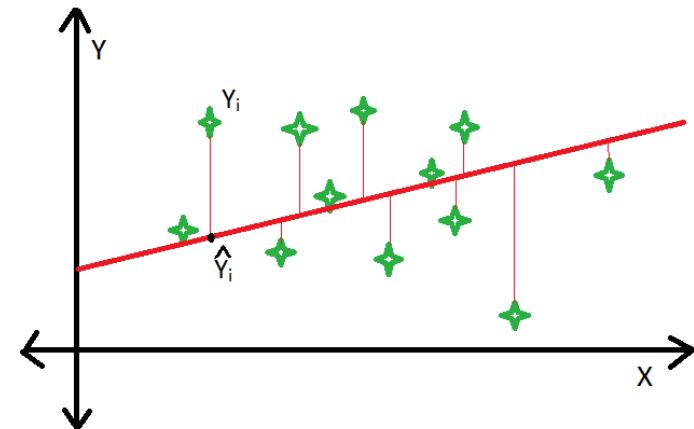
Hard to interpret on its own

REGRESSION

$$SS_{total} = \sum_{i=1}^n (Y_i - Y_{mean})^2$$



$$SS_{res} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$



R2 (or Coefficient of determination): the R2 represents how better your model is compared to the mean model

→ Describes how much of the total variation in your data is explained by the independent variables

$$R^2 = 1 - \frac{SS_{res}}{SS_{total}}$$

R2 = R squared

SS_{res} = residual sum of squares

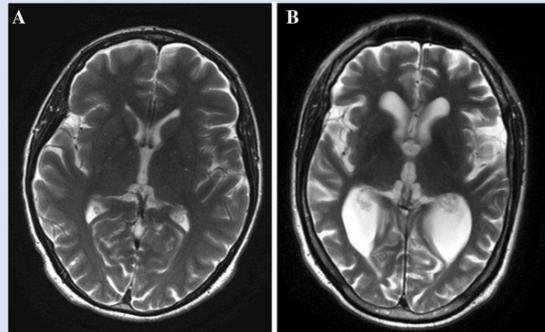
SS_{tot} = total sum of squares

R1 = 1 → Perfect fit

Three main **types** of Machine Learning

Supervised

ML approach that uses labeled data to learn and predict outcomes



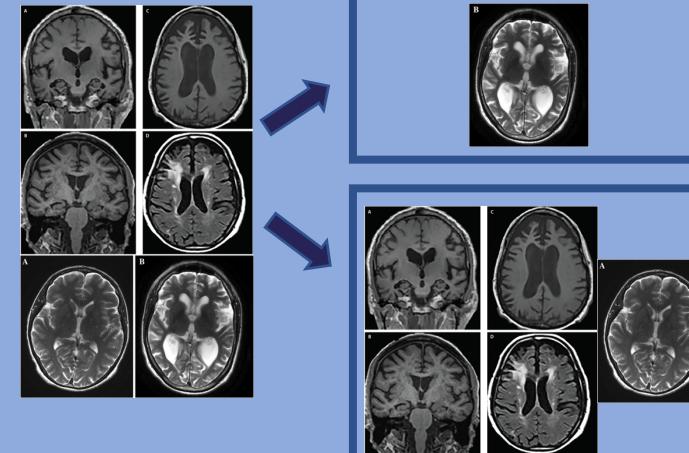
Alzheimer

Healthy

Task driven

Unsupervised

ML approach that uses unlabeled data to discover patterns



Data driven

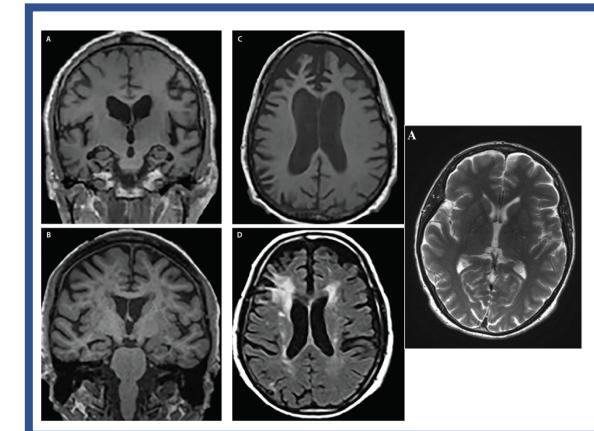
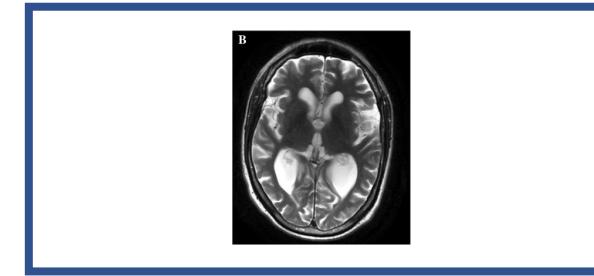
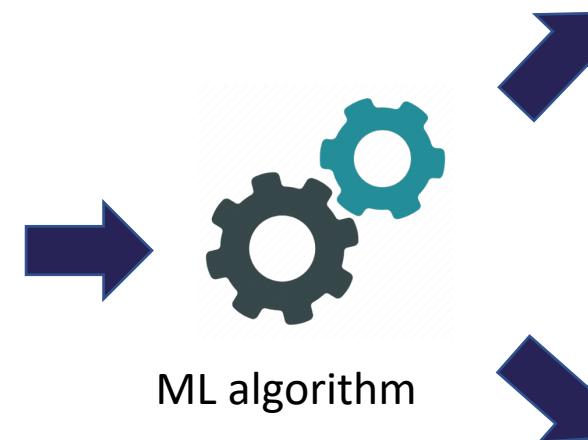
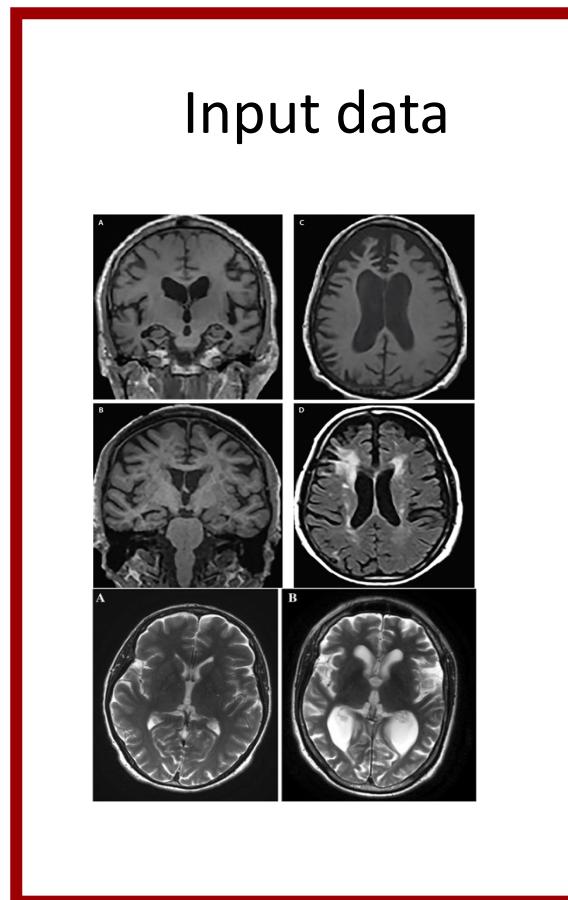
Reinforcement Learning

ML approach that allows an AI-driven system to learn through trial and error using feedback from its actions

Beyond the scope of this lecture

Learning from mistakes

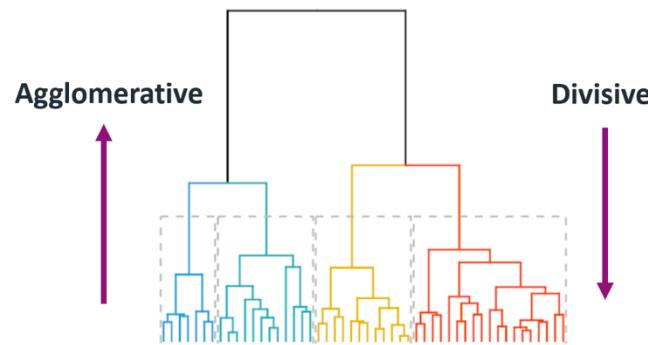
Unsupervised learning is commonly used for **CLUSTERING** and **DIMENSIONALITY REDUCTION**



Clustering: grouping of unlabeled data – **similar** data / data points are grouped together

CONNECTIVITY-BASED

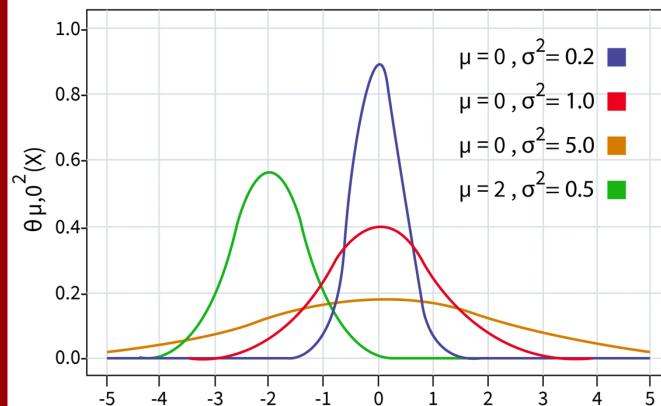
Clusters are created by building a hierarchical tree-type structure



HIERARCHICAL CLUSTERING

DISTRIBUTION-BASED

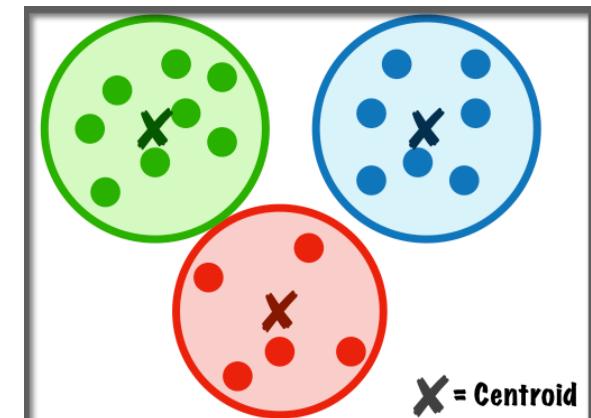
Data points are grouped together based on their probability of belonging to the same probability distribution



GAUSSIAN MIXTURE MODELS

CENTROID-BASED

Similarity is defined based on the distance between the data points and the centroid/center of the cluster

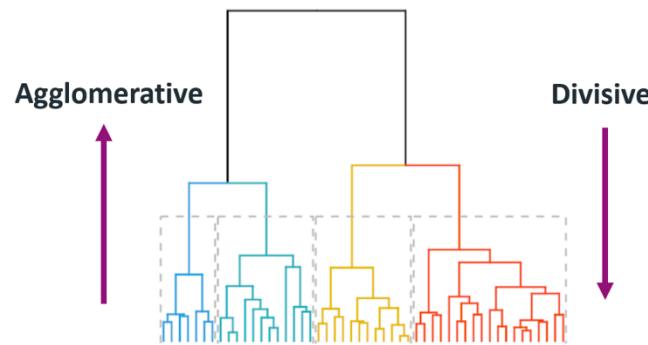


K-MEANS CLUSTERING

This list includes only some of the most common types, but not all of them!

CONNECTIVITY-BASED

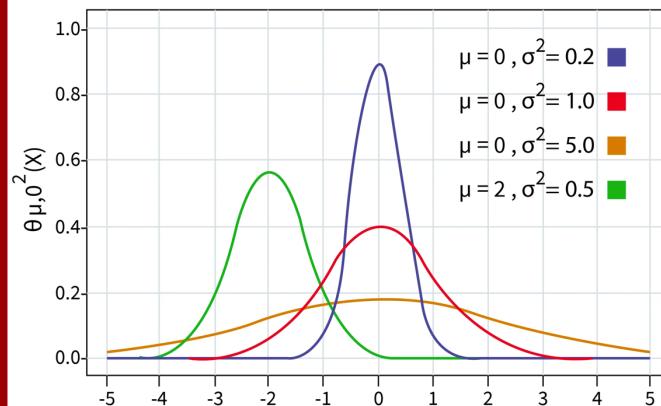
Clusters are created by building a hierarchical tree-type structure



HIERARCHICAL CLUSTERING

DISTRIBUTION-BASED

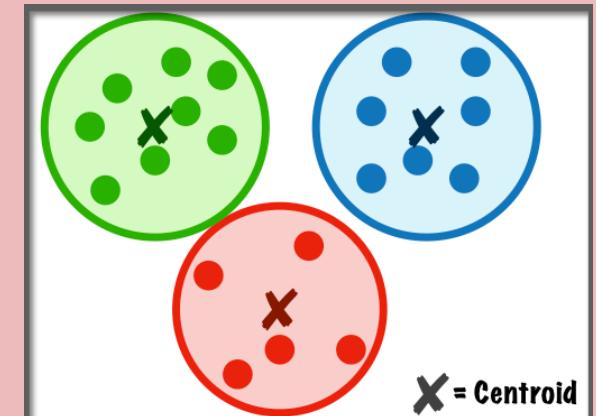
Data points are grouped together based on their probability of belonging to the same probability distribution



GAUSSIAN MIXTURE MODELS

CENTROID-BASED

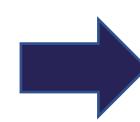
Similarity is defined based on the distance between the data points and the centroid/center of the cluster



K-MEANS CLUSTERING

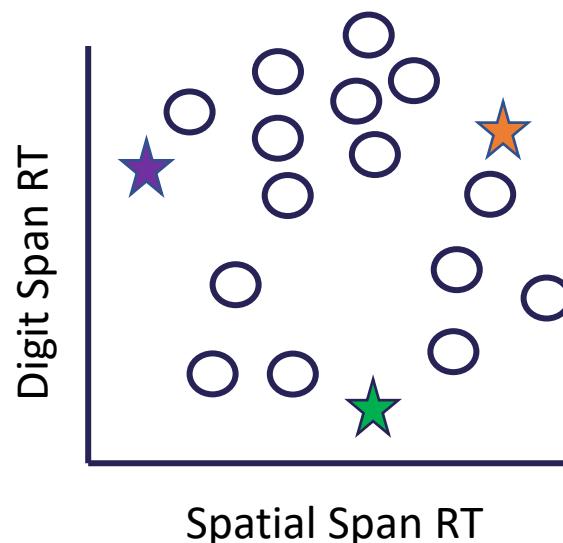
This list includes only some of the most common types, but not all of them!

K-means clustering aims to minimise the total intra-cluster variation (or within cluster sum of squares)



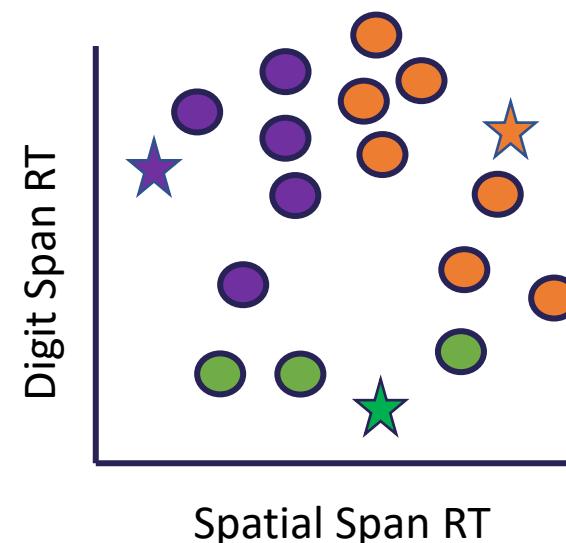
! The number of clusters is pre-defined

Initial centroids



Spatial Span RT

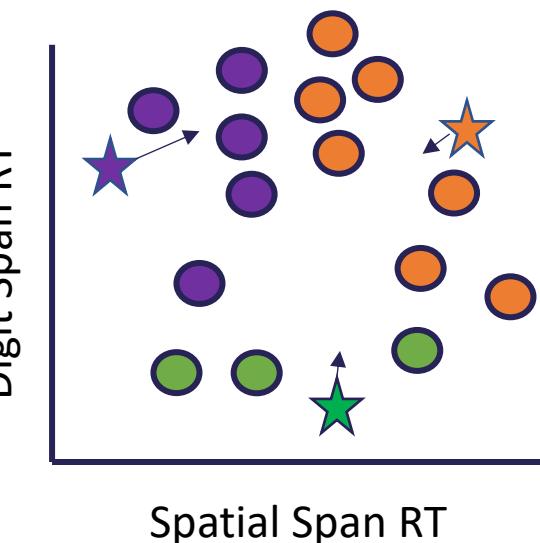
Centroids are randomly selected



Spatial Span RT

Data points are grouped with the cluster of the closest centroid

Iteration 1



Spatial Span RT

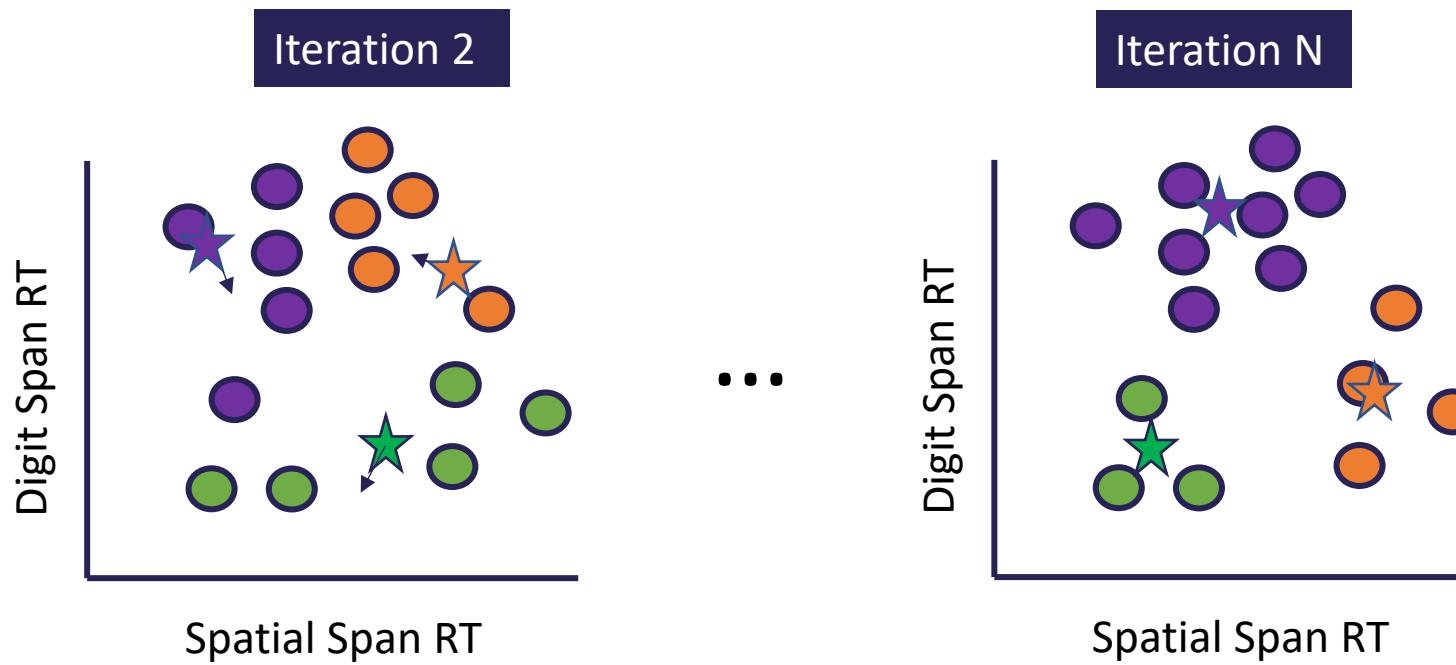
New centroids are located at the cluster means



K-means clustering aims to minimise the total intra-cluster variation (or within cluster sum of squares)

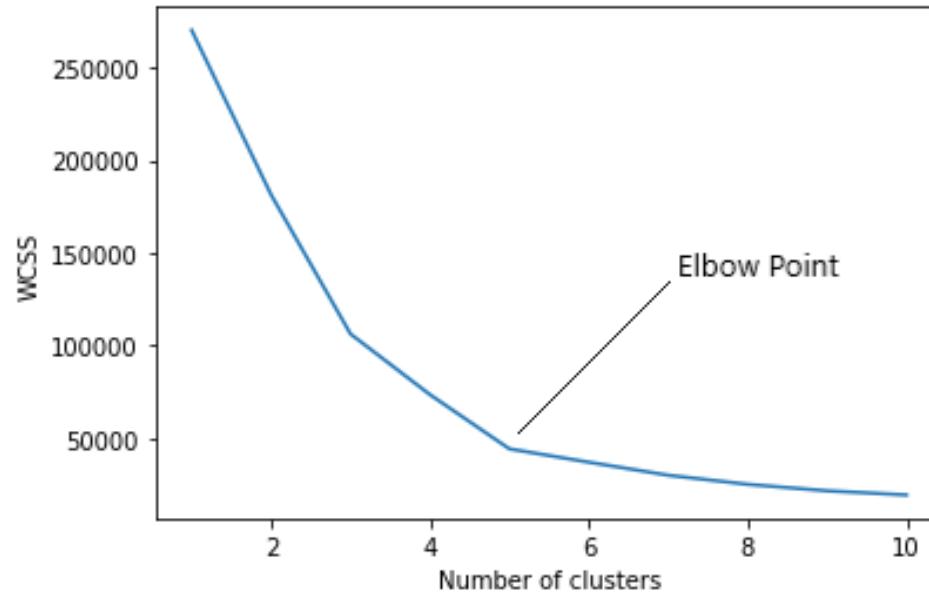


! The number of clusters is pre-defined



How do select the optimal number of clusters?

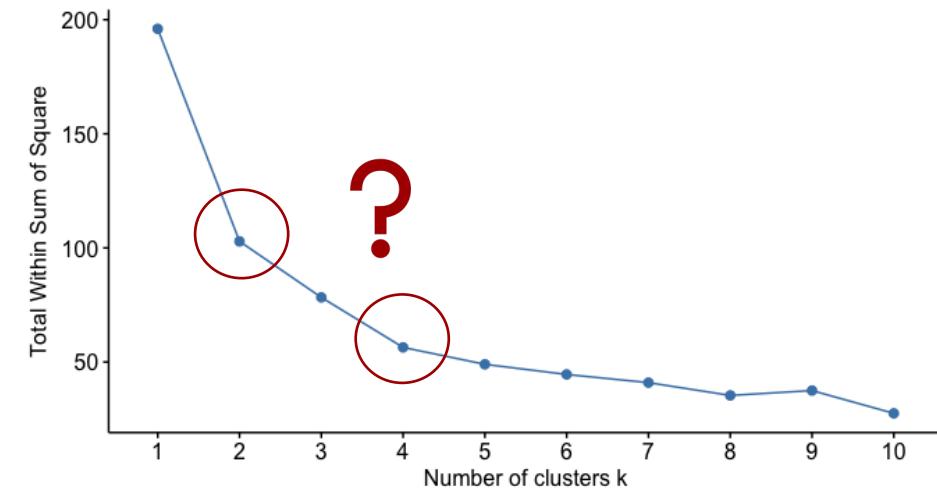
The iterations are continued until convergence
→ The data points don't change cluster anymore



The elbow is very subjective and often hard to identify

WCSS = within-cluster sum of squares

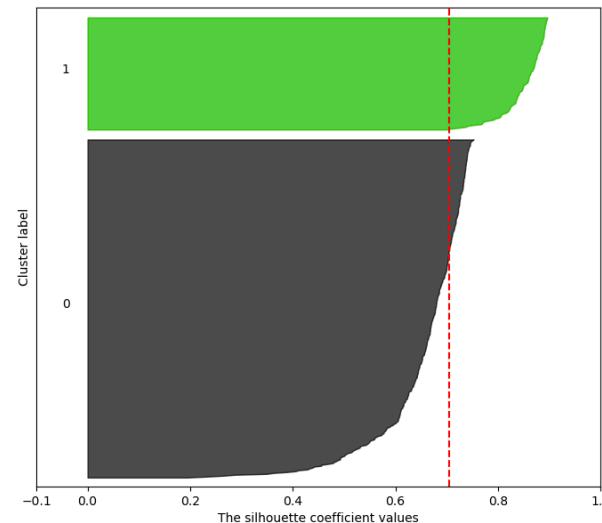
Elbow point = point in the curve where the knee is located → WCSS falls suddenly



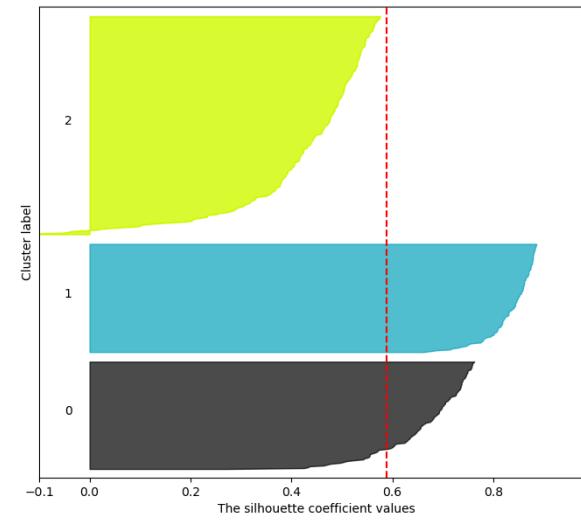
Silhouette score: measure of how close each data point in one cluster is to the data points of the neighboring cluster

- $S(i) = 1$ Point i is far from neighboring cluster
- $S(i) = 0$ Point i is at the border
- $S(i) = -1$ Point i is in the wrong cluster

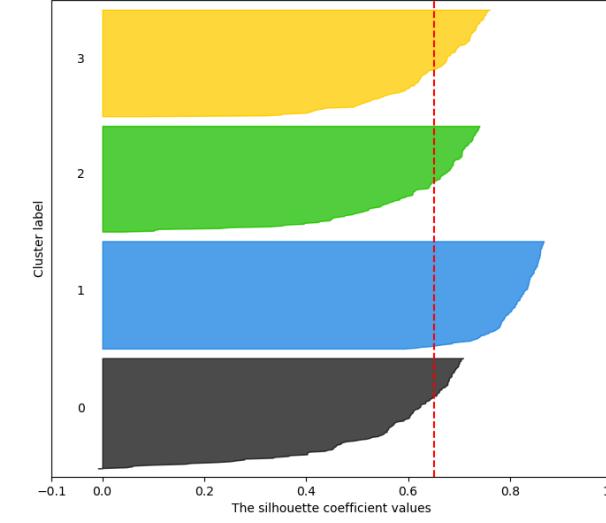
$S(\text{average}) = 0.70$



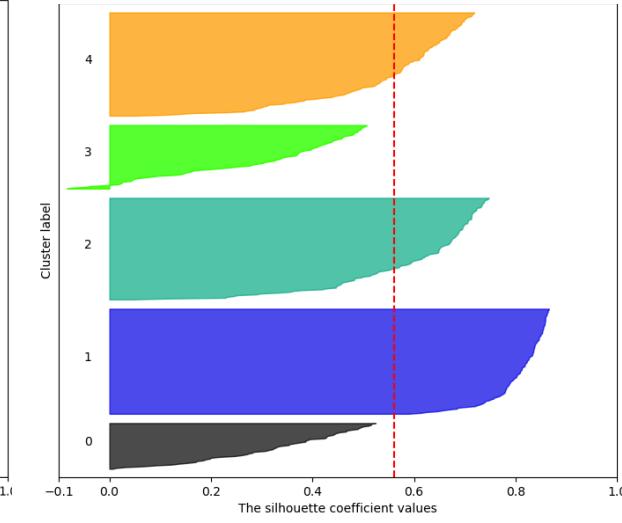
$S(\text{average}) = 0.59$



$S(\text{average}) = 0.65$



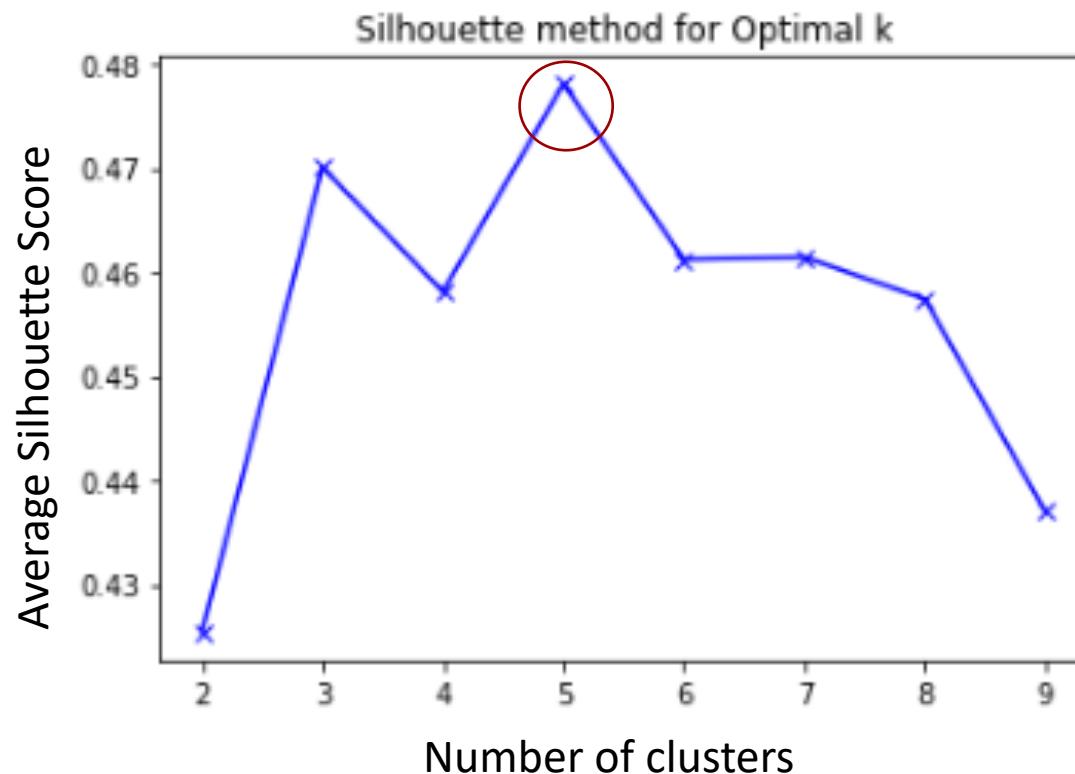
$S(\text{average}) = 0.56$



Silhouette scores

Silhouette score: measure of how close each data point in one cluster is to the data points of the neighboring cluster

- $S(i) = 1$ Point i is far from neighboring cluster
- $S(i) = 0$ Point i is at the border
- $S(i) = -1$ Point i is in the wrong cluster



The optimal K has the highest average silhouette score



The silhouette score is more objective!



What if the data are
categorical?



K-MODES
CLUSTERING

Centroids are random points



Centroids are random data points among those available

Measure of **distance**



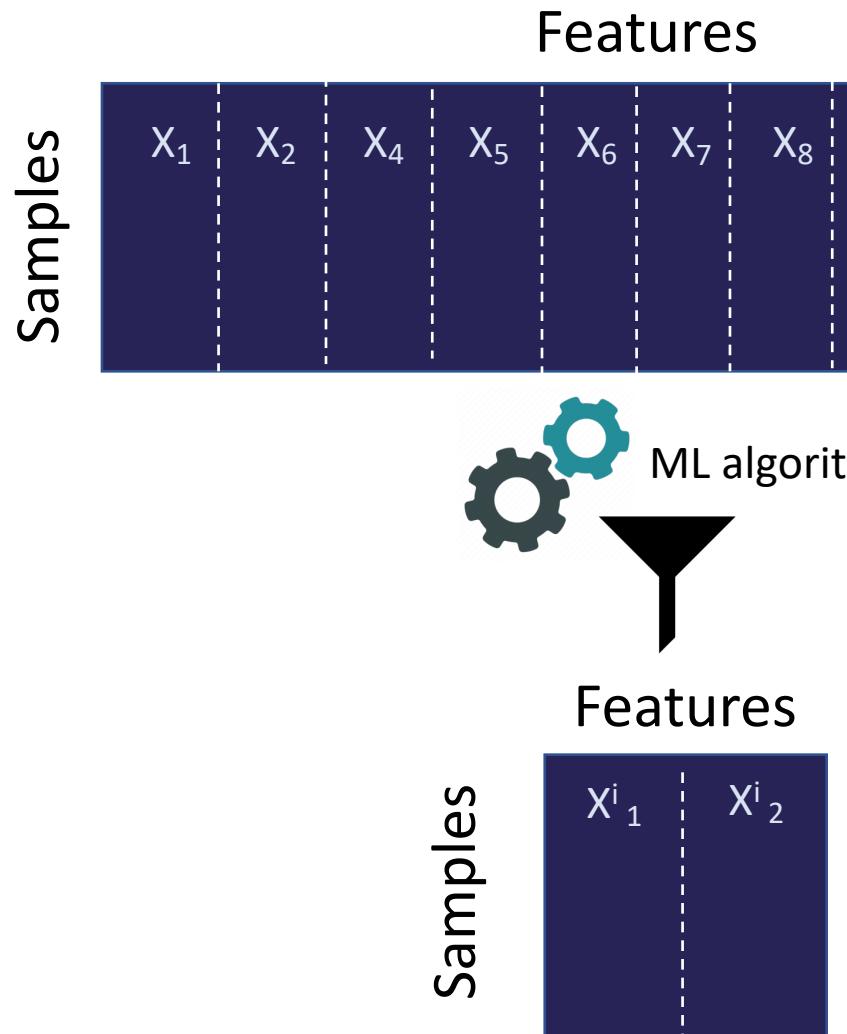
Measure of **dissimilarity** (categories in common with the centroids)

At each iteration, the centroids are redefined based on the **mean**



At each iteration, the centroids are redefined based on the **mode**

Unsupervised learning is commonly used for CLUSTERING and **DIMENSIONALITY REDUCTION**



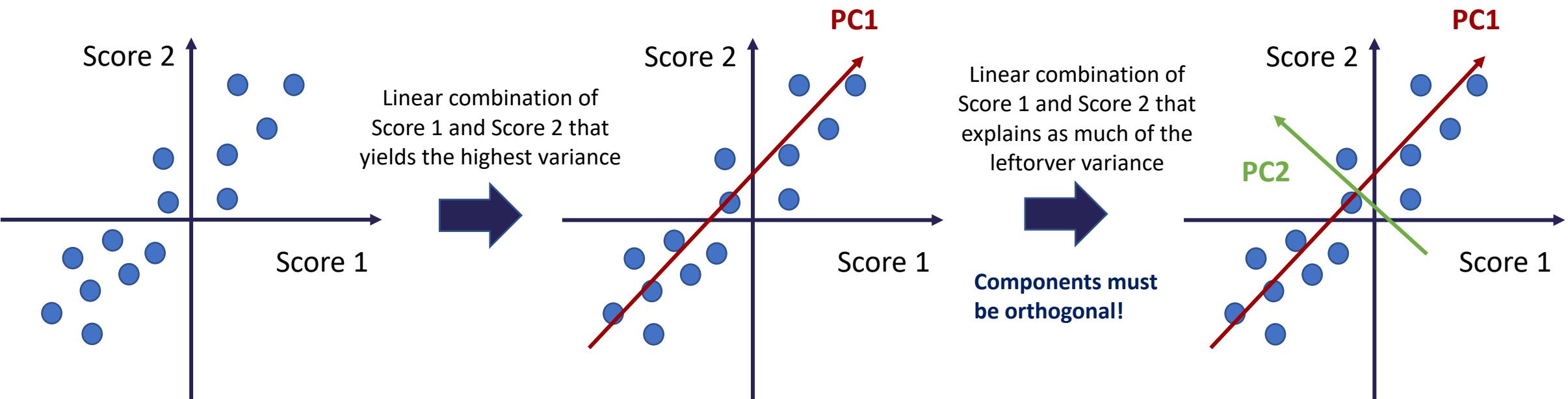
WHY?

1. Visualisation
2. Decrease model complexity and prevent overfitting
3. Feature selection

! Feature selection differs from DR because it simply selects a subset of features without changing them

Dimensionality Reduction: transformation of data from high to low dimensional space

Principal Component Analysis (PCA) is a common unsupervised ML technique for dimensionality reduction that transforms a dataset by creating new dimensions that try to explain as much variance in the data as possible.



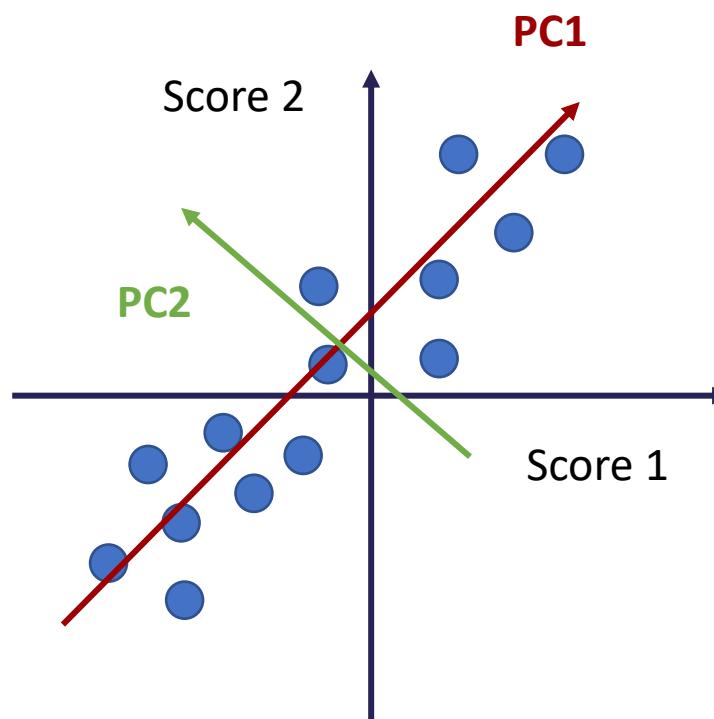
Let's imagine data with only 2 features
(e.g. scores in 2 cognitive tests)

PC1 and PC2 are the 2 new lower dimensional features

You can have as many components as features

How do you interpret the principal components?

NOT STRAIGHTFORWARD



Principal components are **linear combinations** of the original features

$$PC1 = \boxed{1.5} * Score1 + \boxed{2} * Score2$$

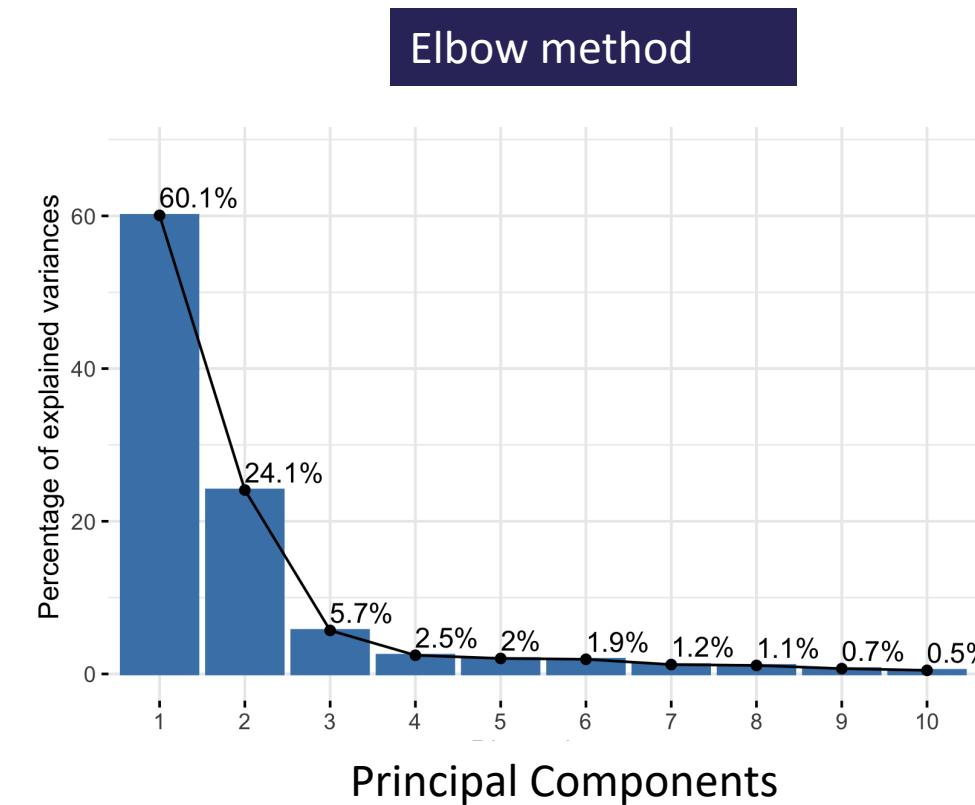
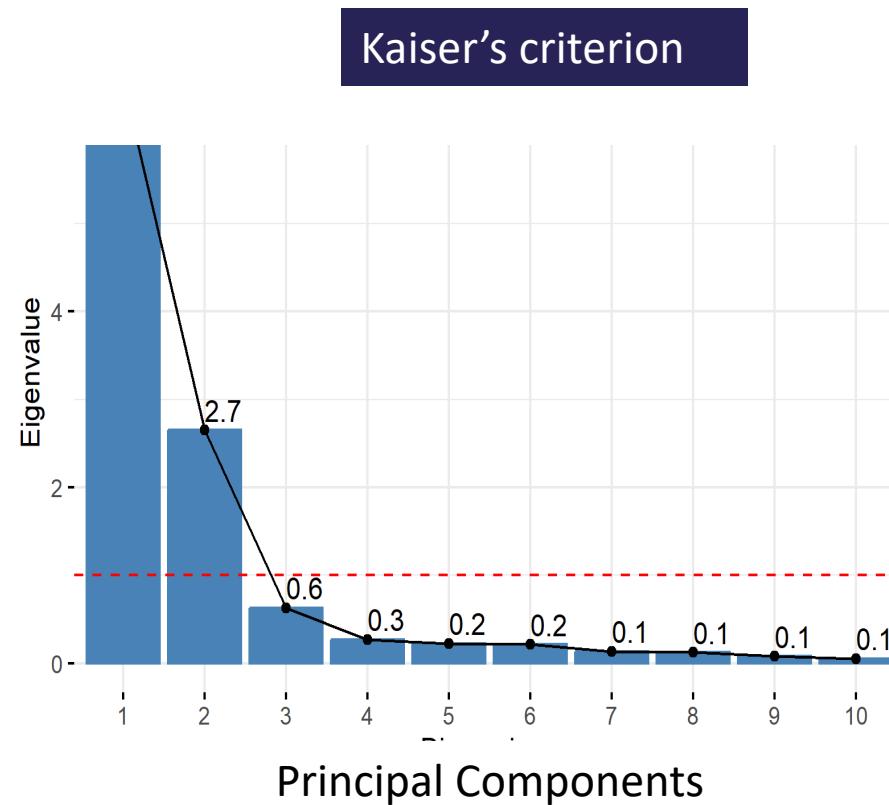
↓ ↓
Loadings Loadings

Features with higher loadings contributed the most to define the PC

→ Can be potentially used for **FEATURE SELECTION**



How do you select the optimal number of components?



Questions?