



# Machine Learning

Session 2 - Supervised classification



hadrien.salem@centralelille.fr



[introduction-to-data-science](#)

# Introduction

What did we do last time?

# Course outline

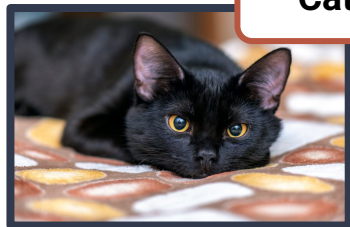
## Intro to ML course

**Session 1: Introduction to ML & Regression**

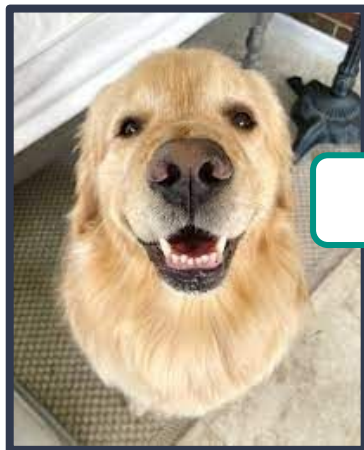
**Session 2: Supervised classification**

**Session 3: Neural networks**

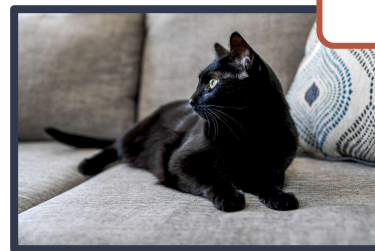
# What is classification?



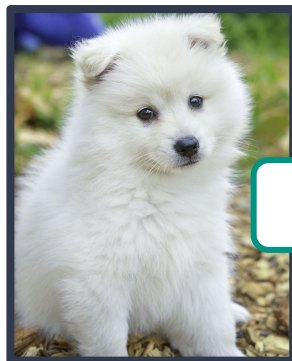
Cat



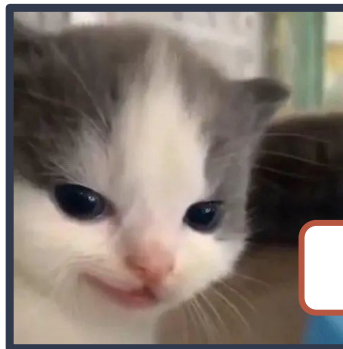
Dog



Cat



Dog



Cat



Dog

Intuitively, classification is giving objects the right label

$$f^*(x) = \arg \max_k \mathbb{P}(C_k|x)$$

Where  $f^*$  is a rule for classification,  $C_k$  are the **classes**, and  $x$  the **examples**

The goal of a classification algorithm is to find this rule

Formally, classification is finding the most probable class for an example



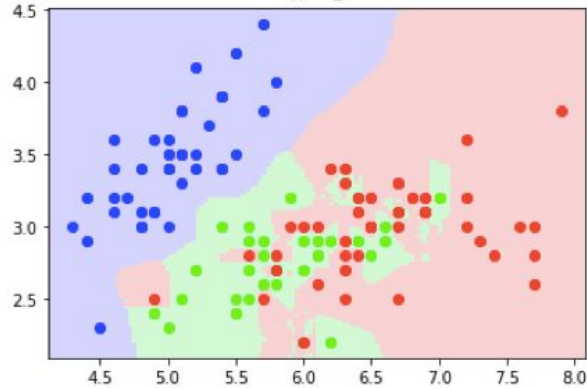
# Families of classification models

While they are all classification models, they have different purposes

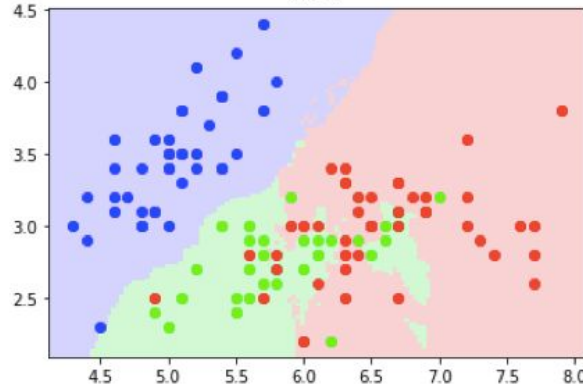
## There are three main families of classification models

- **Discriminant functions**
  - The algorithm learns a function that finds the class directly
  - *Example:* K-nearest-neighbours
- **Discriminant models**
  - The algorithm models the decision boundary
  - *Example:* Support Vector Machines
- **Generative models**
  - The algorithm models the data distribution (meaning you can generate your own data)
  - *Example:* Gaussian Mixture Model

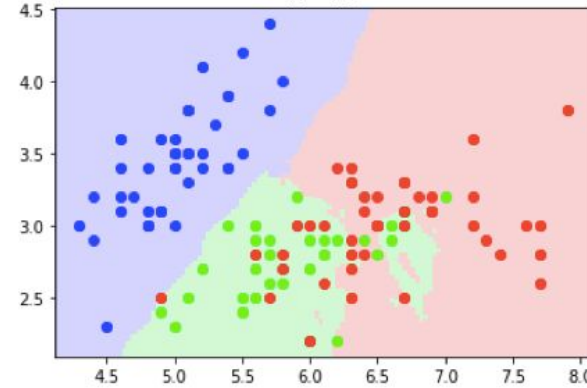
K = 1



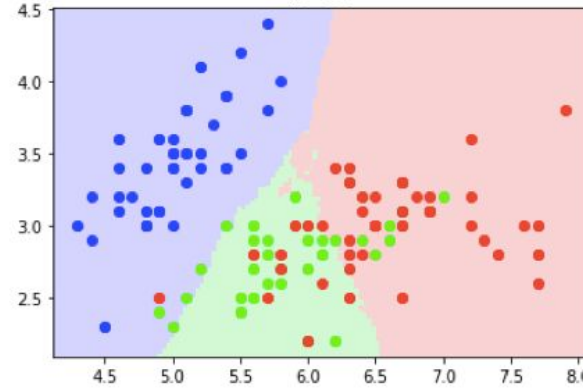
K = 5



K = 10



K = 50

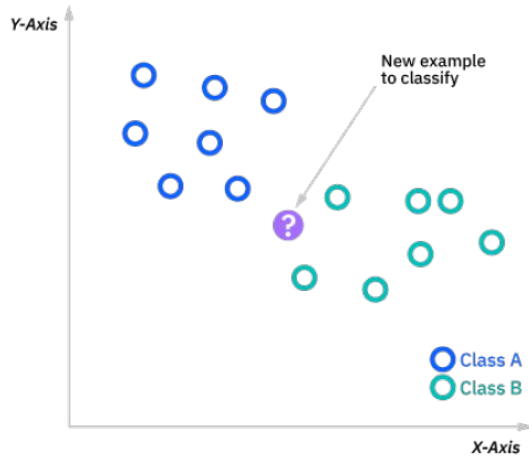


When K is small, the algorithm is very sensitive to **local variations** (risk of overfitting)

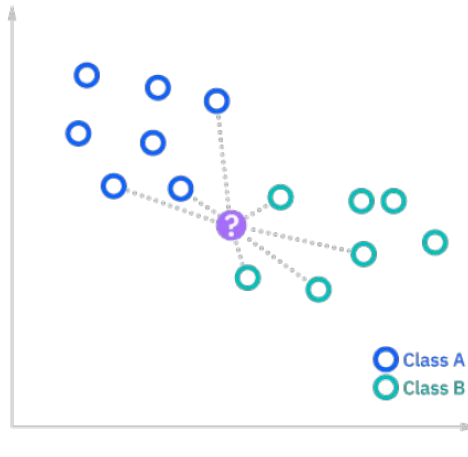
When K is large, the algorithm is **more stable**, but **does not take small variations into account** (risk of underfitting)

⇒ When choosing the parameters, there is a **compromise between the two**

# Common classification algorithms



**Introduce a new example**



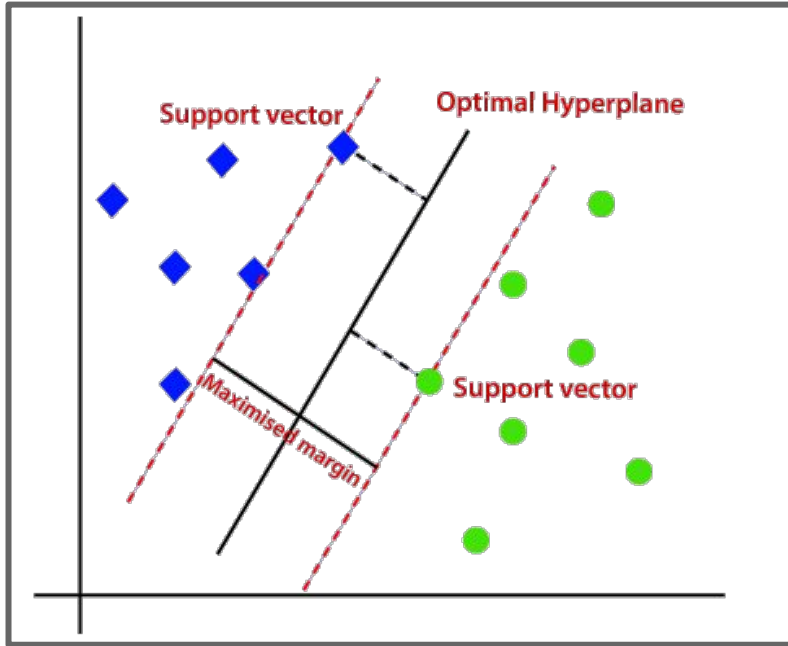
**Compute distances**



**Majority vote**

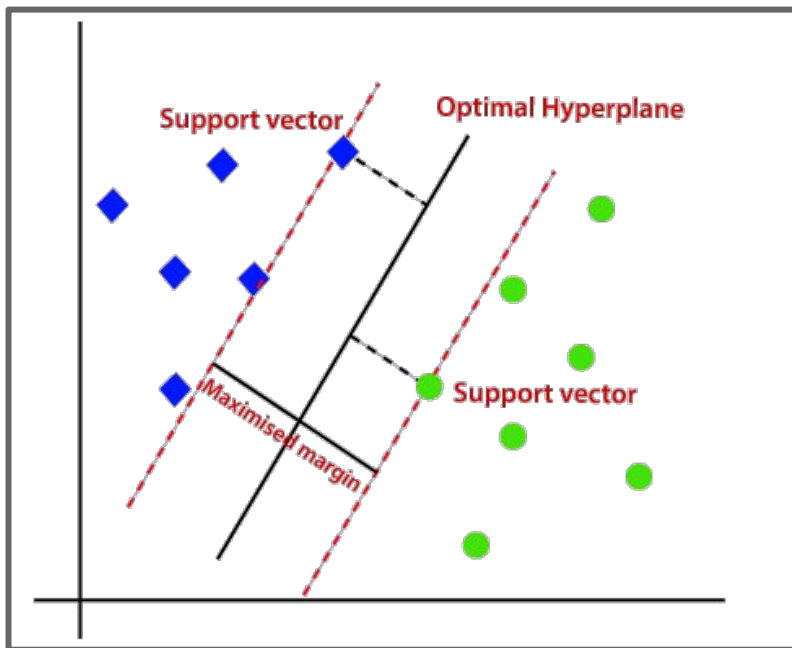
# K-nearest neighbours

<b>Decision boundary</b>	Non-linear
<b>Advantages</b>	<ul style="list-style-type: none"><li>• Easy to use and understand</li><li>• No assumptions</li></ul>
<b>Disadvantages</b>	<ul style="list-style-type: none"><li>• Slow for large datasets</li><li>• Inefficient in high dimension</li></ul>



The objective is to **find a hyperplane** such that the **margin between the two classes is maximized**.

Data can be transformed into a higher-dimensional space if it is not linearly separable in the feature space. This is achieved with **kernels** (e.g. polynomial, sigmoid, etc.).



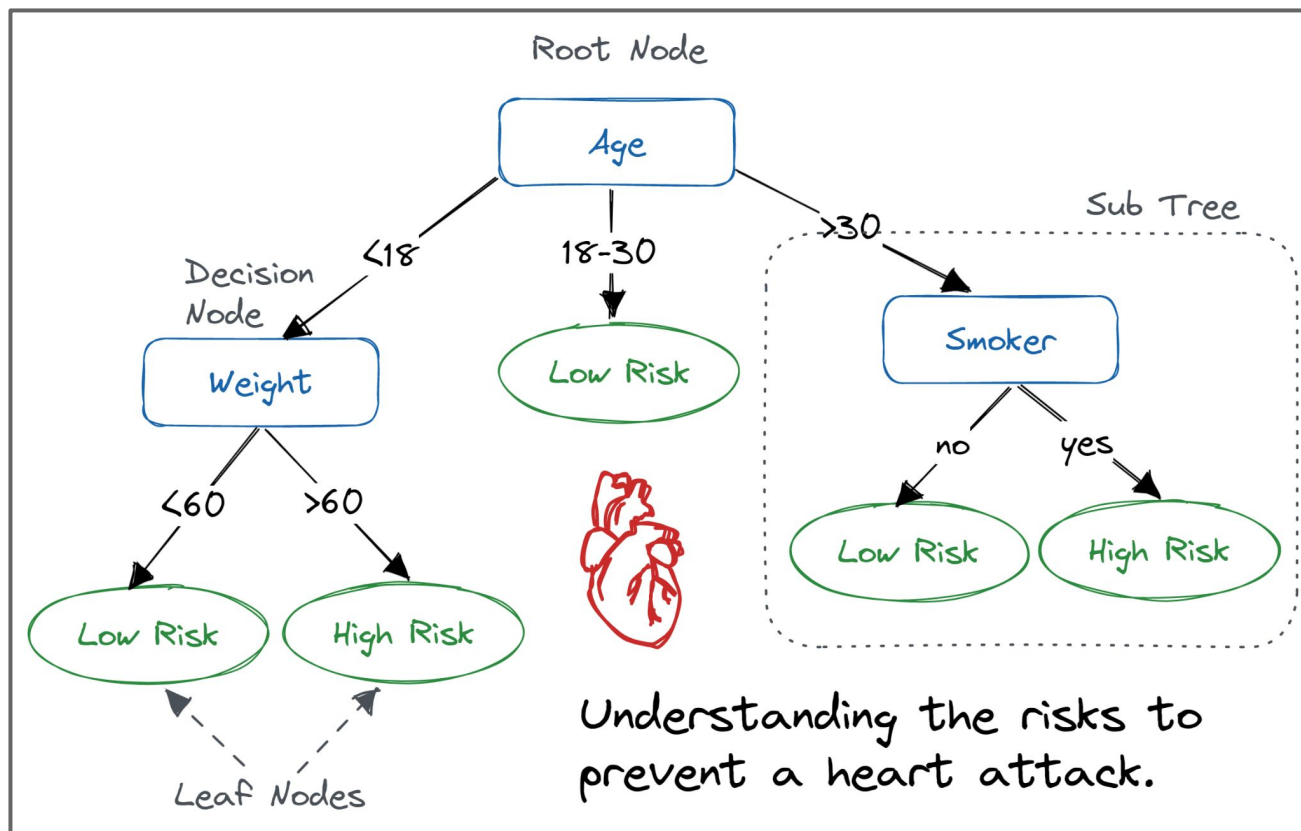
<b>Decision boundary</b>	Linear in the transformed space, can be non-linear in the feature space
<b>Advantages</b>	<ul style="list-style-type: none"> <li>• Works well in high dimension</li> <li>• Robust to outliers</li> <li>• Low memory consumption</li> </ul>
<b>Disadvantages</b>	<ul style="list-style-type: none"> <li>• Slow for large datasets</li> <li>• Choosing a kernel can be difficult</li> </ul>

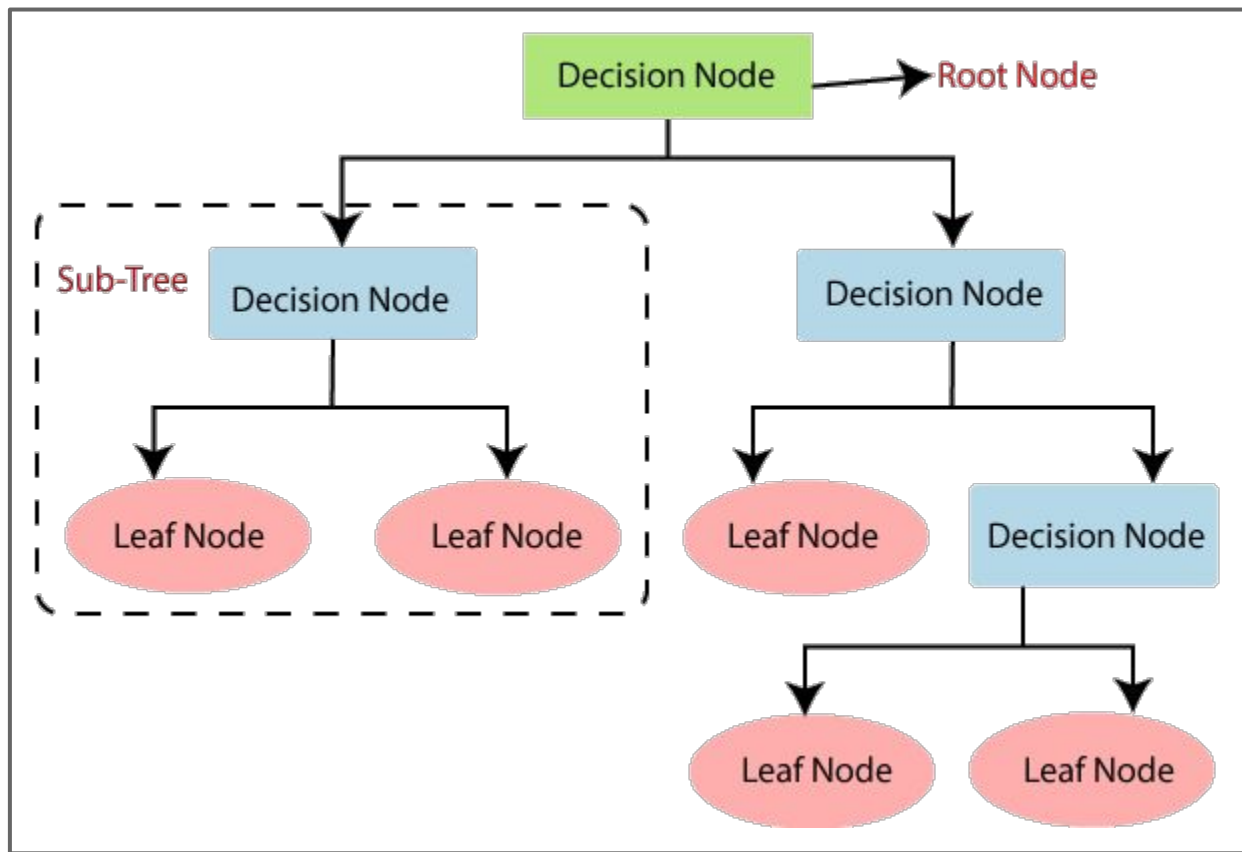
# Other methods for supervised classification

- **Logistic Regression**
  - THIS IS A CLASSIFICATION METHOD
  - Only works when data is linearly separable
  - Easy to use, good baseline method
- **Naive Bayes**
  - Assumes that features are independent
  - Non-linear decision boundary (computes class membership probabilities)
  - Low-cost, also good baseline
- **Linear / Quadratic discriminant analysis**
  - Assumes that data follows a normal distribution
  - Limited to linear / quadratic decision boundaries
  - Good baseline
- **And other algorithms we will study later**
  - Decision trees / Random Forests
  - Ensemble methods
  - Neural networks



# Decision trees





# Definition

## Purity of a node

**A node is 100% pure when all of its data belongs to a single class.**

**It is 100% impure when it contains the same proportion of each class.**

**(e.g. 50/50 for binary classification)**

Several functions can be used to compute the impurity of a node:

- **Gini Index**
- **Cross-entropy**
- **Misclassification error**

$$\Delta\phi(s, m) =$$

Frequency of a class  
at node  $m$

Proportion of data in  
the left subtree

Proportion of data in  
the right subtree

$$\underbrace{\phi(p_m)}$$

$$- \underbrace{(\pi_L \phi(p_{m_L}) + \pi_R \phi(p_{m_R}))}$$

Purity before splitting

Purity after splitting

Quality of split  $s$  at node  $m$

**Different splits are tested recursively to find the best partitioning**

# Strength and weaknesses of decision trees

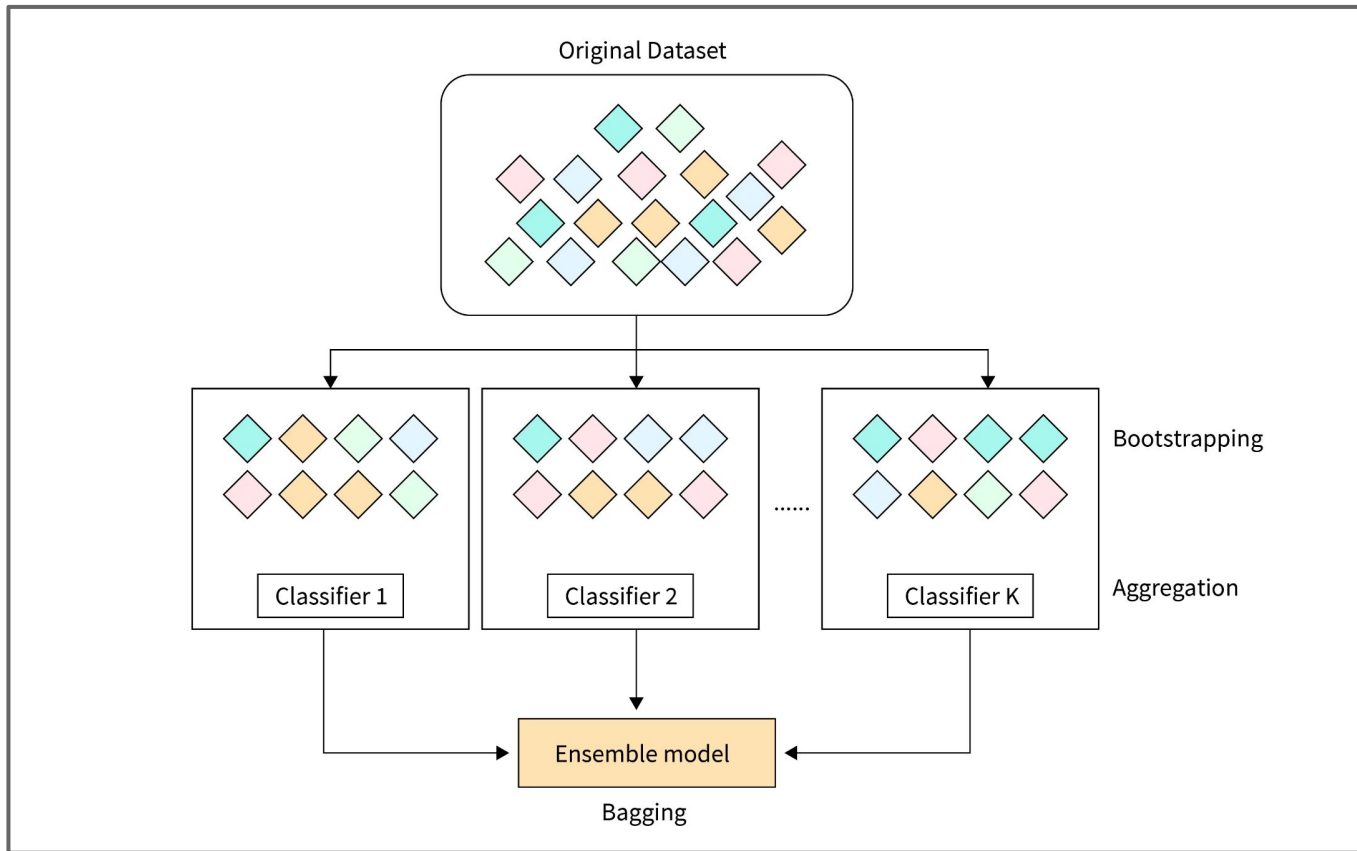
## Strengths

- **Flexible** (few hypotheses)
  - They can also apply to regression!
- Easy to **interpret** (explicit rules)
- **Non-linear** (complex decision boundaries)

## Weaknesses

- Prone to **overfitting**
- Unstable to **noise**
- **Expensive** on large datasets

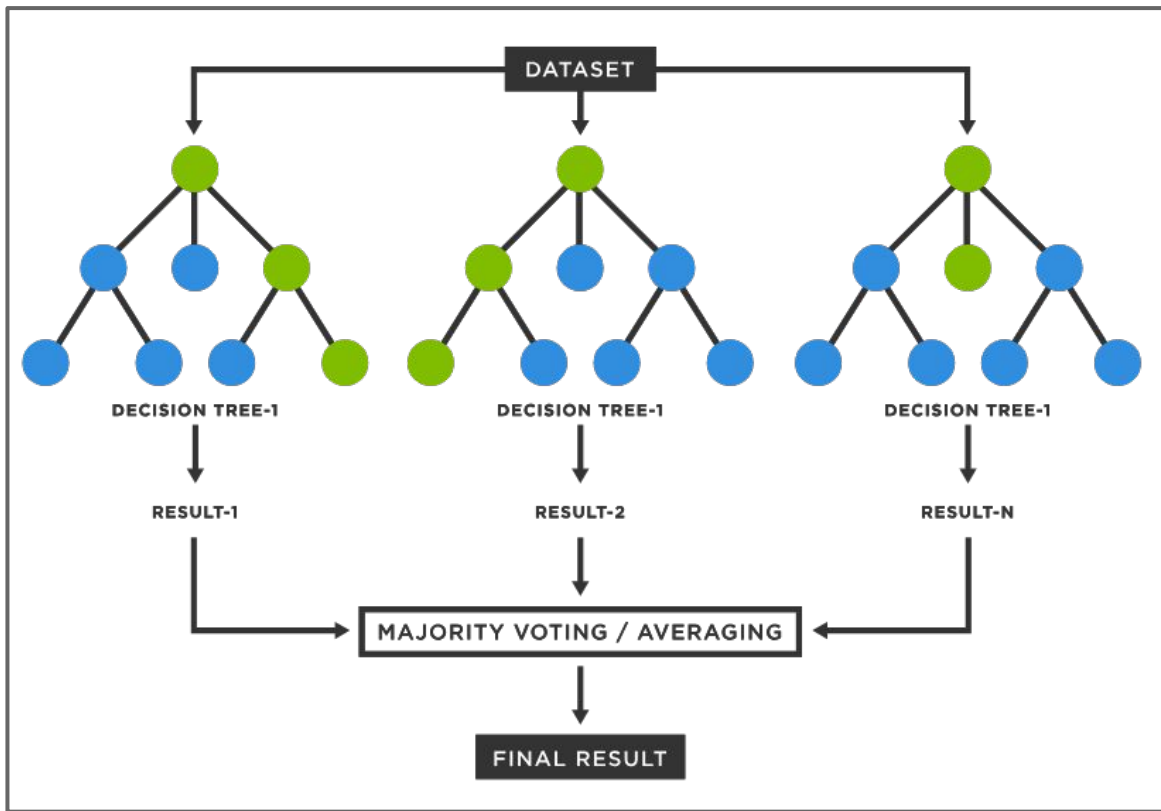
# Ensemble methods



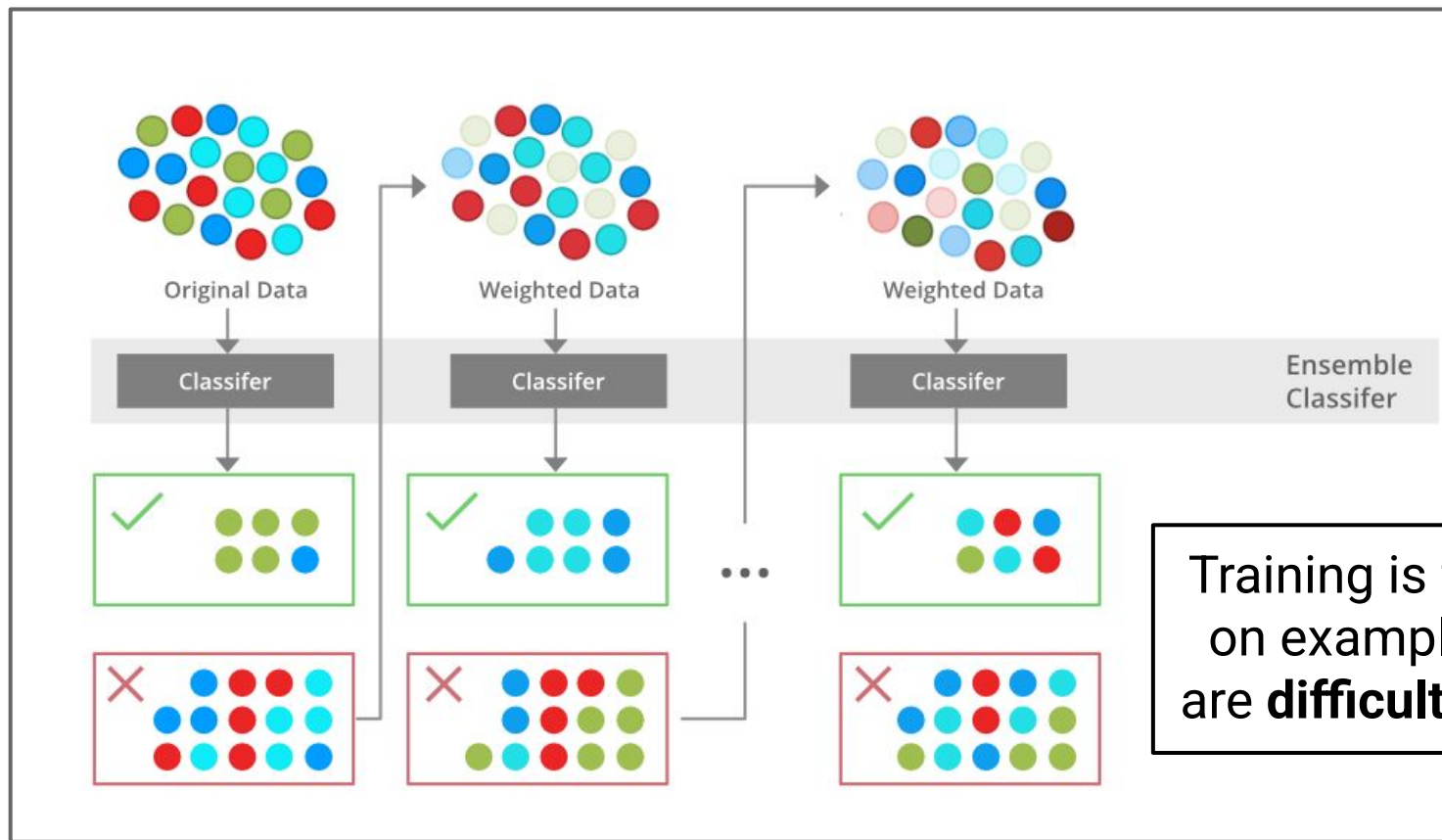
**Bootstrapping**  
Recombining  
existing data to  
create datasets

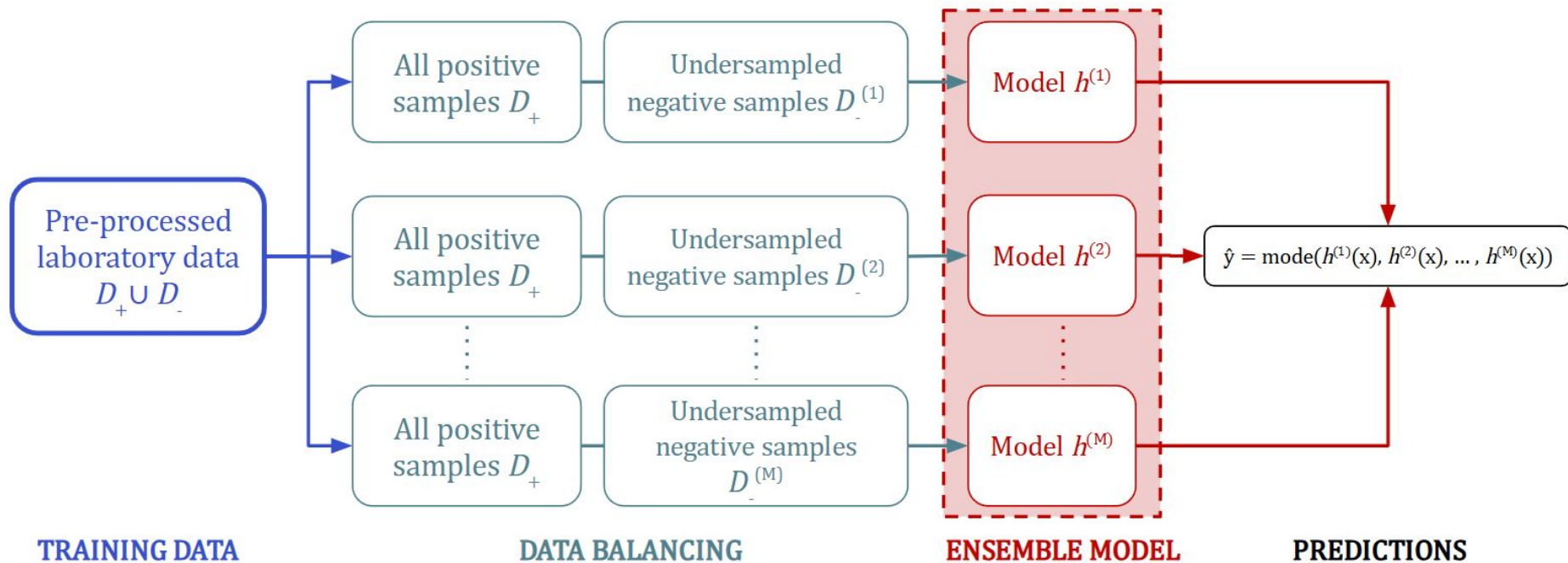
**Aggregating**  
Training an  
algorithm for  
each dataset





The principle is similar to bagging, except **trees are built upon random subsets of features**





# Strength and weaknesses of ensemble methods

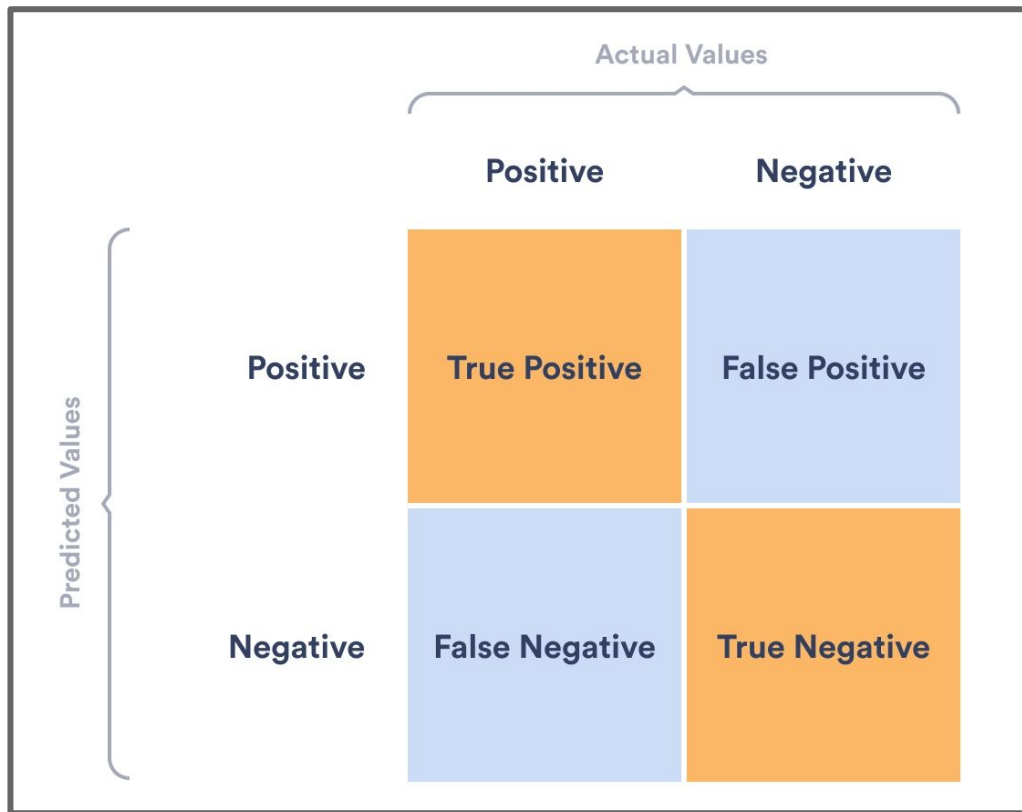
## Strengths

- Tends to increase accuracy
- Robust to noise
- Helps reduce overfitting

## Weaknesses

- Requires more resources
- Makes interpretation more difficult

# Evaluating a classification algorithm



Confusion matrices allow you to analyse how well each class is handled

[Image source](#)

		Real (Actual, Observed)		
		Real Negatives TN+FP	Real Positives TP+FN	
Predicted	Predicted Negatives TN+FN	↑ true negatives (TN)	↑ false negatives (FN)	
	Predicted Positives TP+FP	← false positives (FP)	true positives (TP)	<b>Precision</b> = true positives/PREdiCted positives $TP/(TP+FP)$
		<b>Specificity</b> <b>SPIN (SPecificity Is Negative)</b> true negatives/real negatives $TN/(TN+FP)$	<b>Sensitivity</b> <b>SNIP (SeNsitivity Is Positive)</b> true positives/real positives $TP/(TP+FN)$	<b>Accuracy</b> true predictions/all predictions $(TP+TN)/(TP+TN+FP+FN)$
			<b>Recall</b> true positives/REAL positives $TP/(TP+FN)$ Recall = Sensitivity	

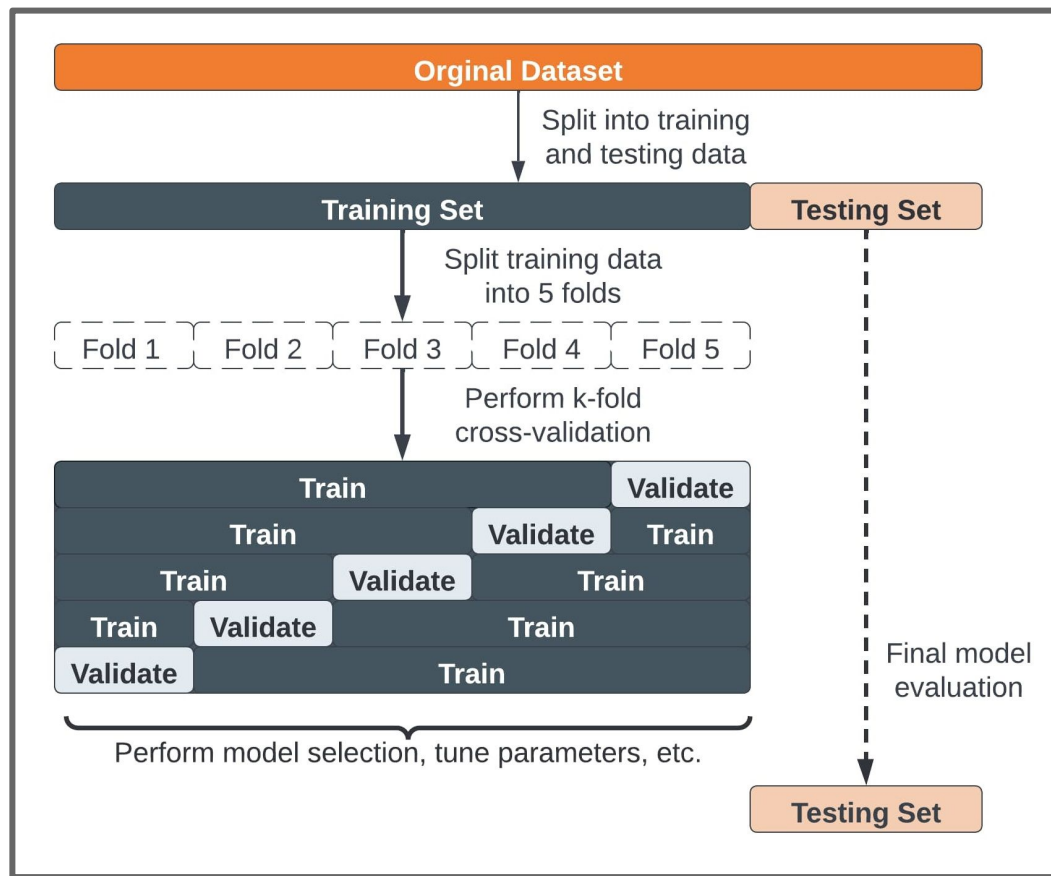
Several performance indicators can be computed from the confusion matrix

[\*Image source\*](#)

$$\begin{aligned}\text{F1 Score} &= \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} \\ &= \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}\end{aligned}$$

F1-Score is a synthesis of the previous metrics and can be more meaningful than accuracy





Cross-validation can allows for a better estimation of the model's performance

[\*Image source\*](#)

# Practical work

The notebook contains all the necessary instructions

# Debrief

# Debrief

**What did we learn today?**

**What could we have done better?**

**What are we doing next time?**

# Machine Learning

Session 2 - Supervised classification



hadrien.salem@centralelille.fr



[introduction-to-data-science](#)