

Machine Learning

Session 2 - Supervised classification



hadrien.salem@centralelille.fr



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

Introduction

What did we do last time?

Course outline

Machine learning course

Session 1: Regression

Session 2: Supervised classification

Session 3: Clustering

Session 4: Decision trees and ensemble methods

Session 5: Introduction to neural networks

Session 6: Advanced neural networks

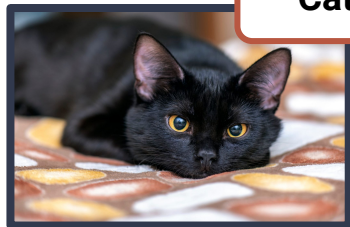
Session 7: Introduction to reinforcement learning

Session 8: Reading science papers

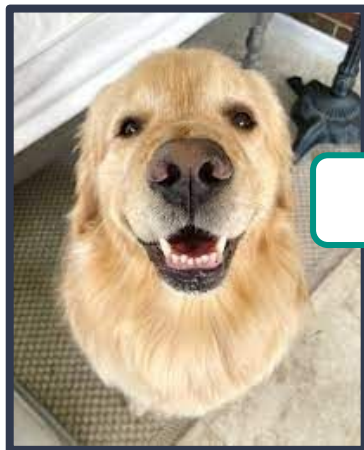


Project

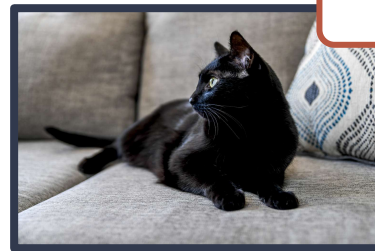
What is classification?



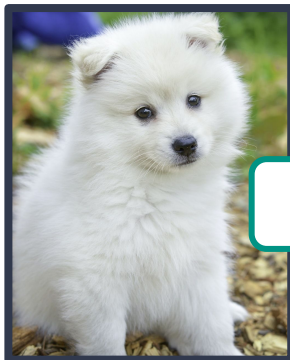
Cat



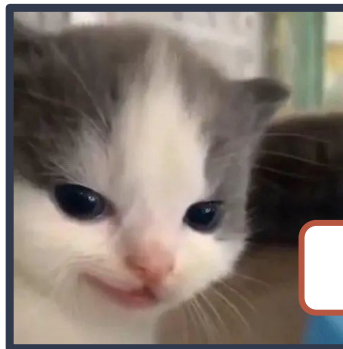
Dog



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

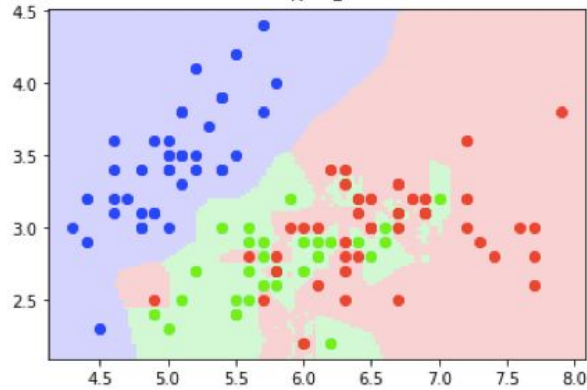
Families of classification models

While they are all classification models, they have different purposes

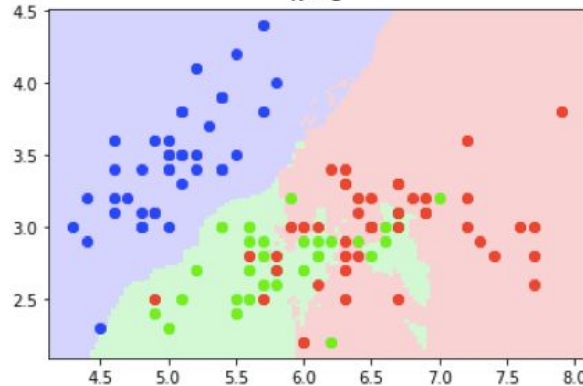
There are three main families of classification models

- **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

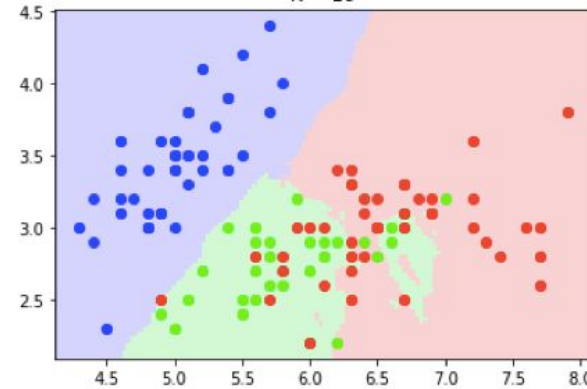
K = 1



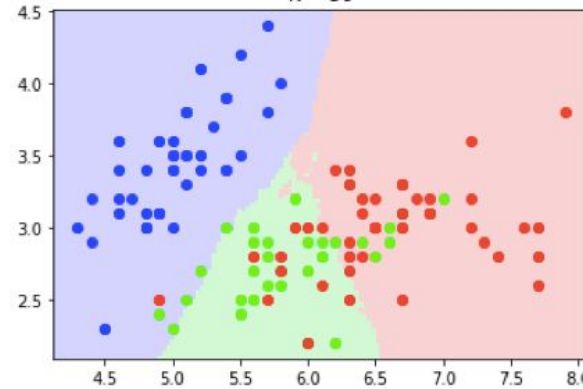
K = 5



K = 10



K = 50

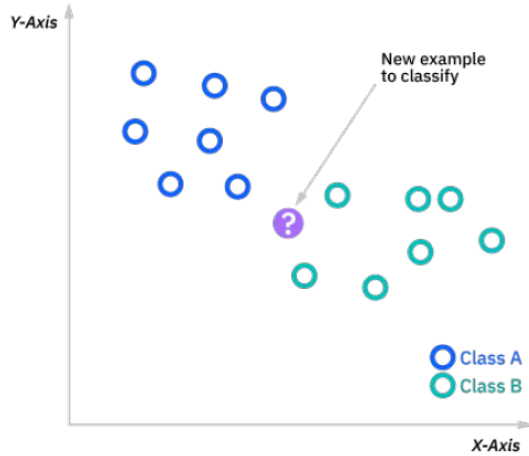


When K is small, the algorithm is very sensitive to **local variations** (risk of overfitting)

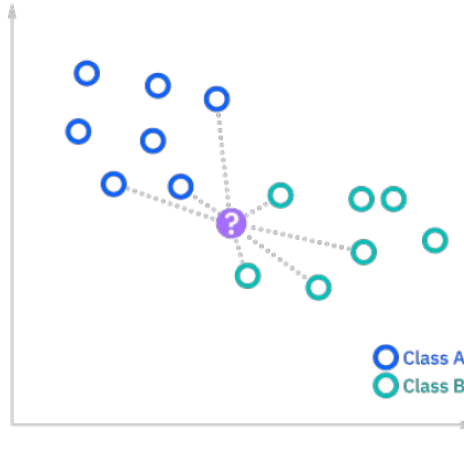
When K is large, the algorithm is **more stable**, but **does not take small variations into account** (risk of underfitting)

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

Common classification algorithms



Introduce a new example



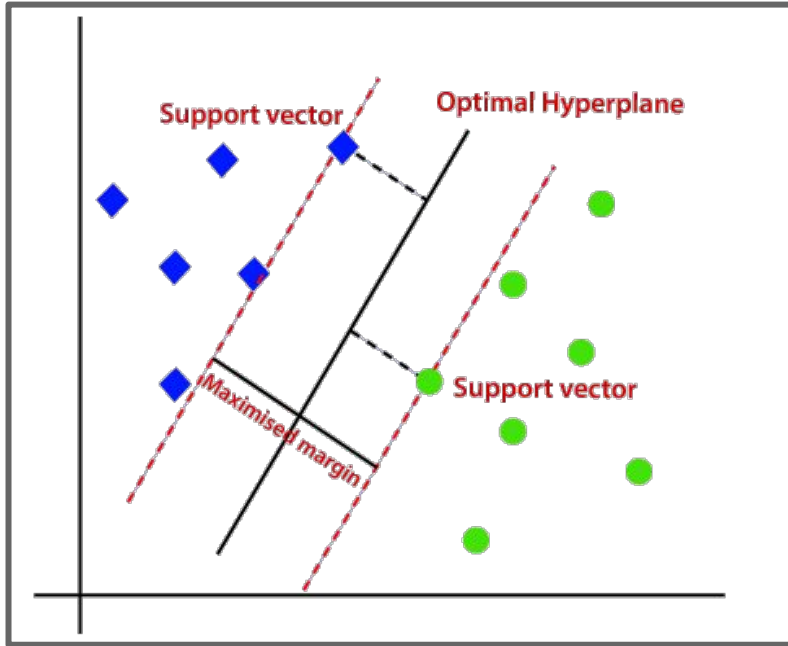
Compute distances



Majority vote

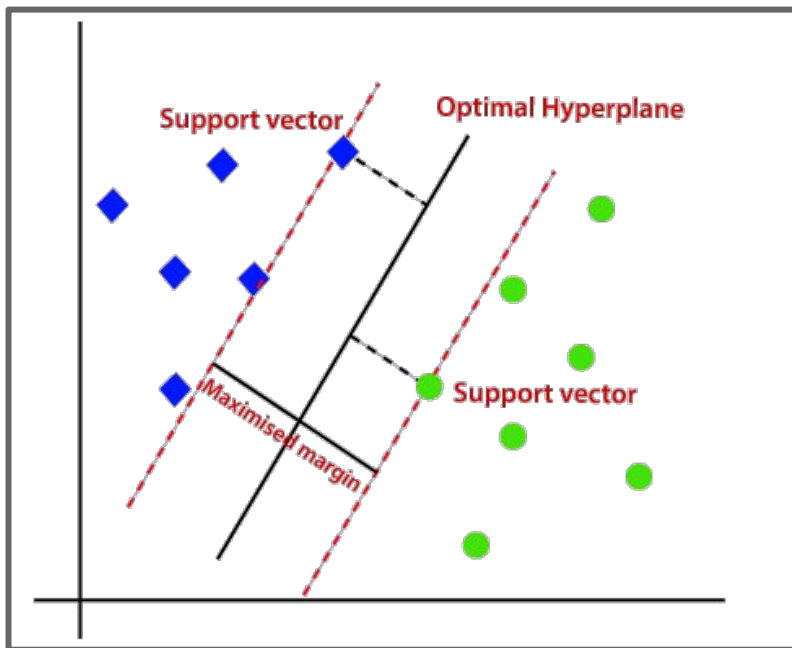
K-nearest neighbours

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



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.).

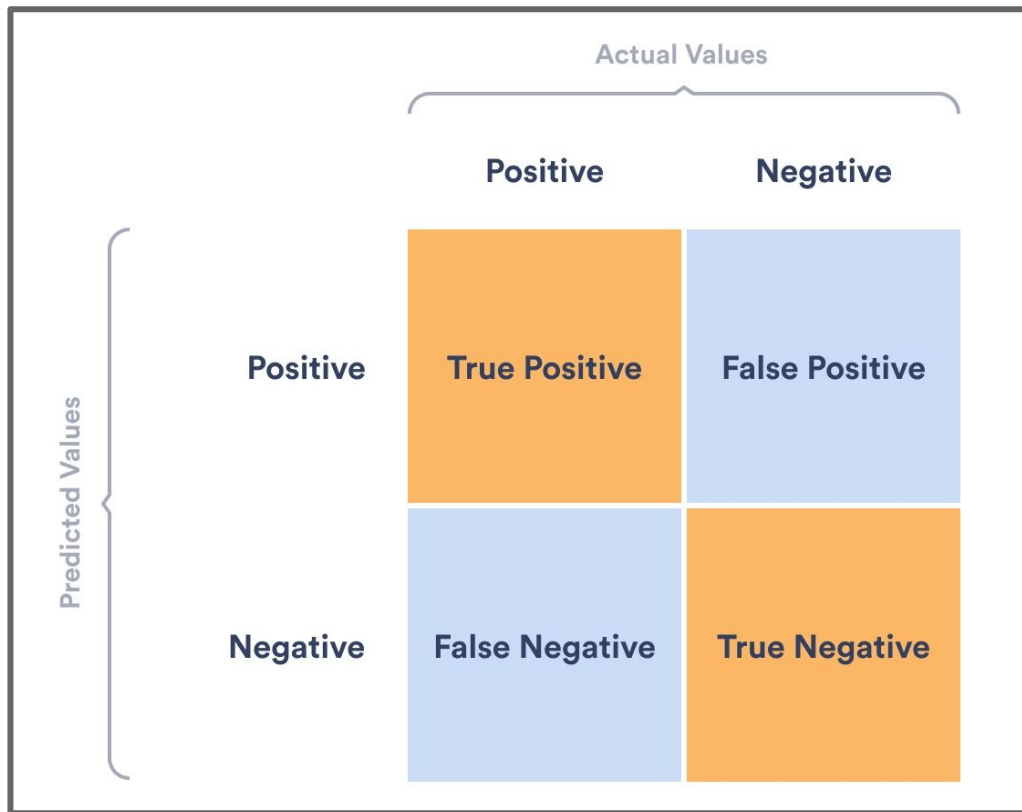


Decision boundary	Linear in the transformed space, can be non-linear in the feature space
Advantages	<ul style="list-style-type: none"> • Works well in high dimension • Robust to outliers • Low memory consumption
Disadvantages	<ul style="list-style-type: none"> • Slow for large datasets • Choosing a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

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

Evaluating a classification algorithm



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

[Image source](#)

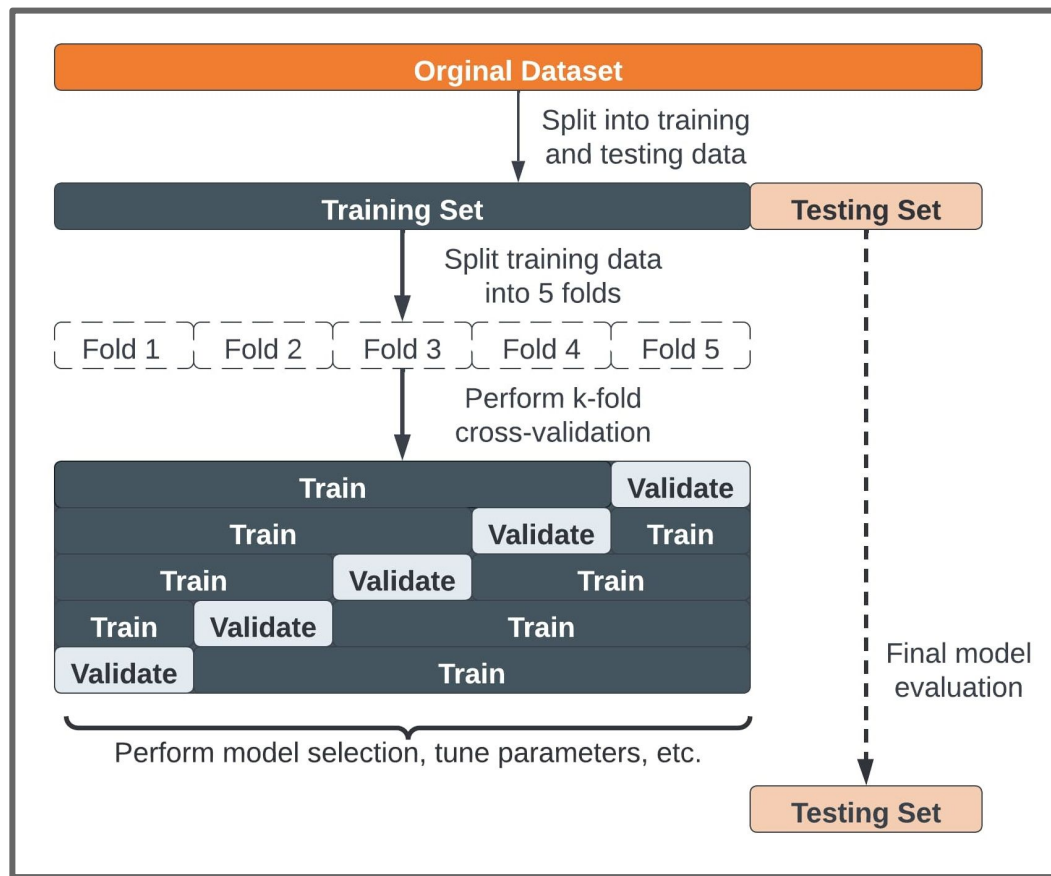
		Real (Actual, Observed)		
		Real Negatives TN+FP	Real Positives TP+FN	
Predicted	Predicted Negatives TN+FN	↑ true negatives (TN)	↑ false negatives (FN)	
	Predicted Positives TP+FP	← false positives (FP)	→ true positives (TP)	Precision = true positives/PREdiCted positives $TP/(TP+FP)$
		Specificity SPIN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