

Introduction to Machine Learning

Rocío Alaiz Rodríguez

Full Professor

University of León (Spain)

rocio.alaiz@unileon.es

- Machine Learning
 - Concept
 - Application fields
 - Supervised, Unsupervised & Reinforcement learning
- Approaching a problem of learning from examples
- Building Machine Learning Models
- Supervised Learning models
 - K-NN
 - Naïve Bayes
 - Neural Networks
- Evaluating classifier performance

MACHINE LEARNING: CONCEPT and APPLICATION FIELDS

Using AI, scientists find a drug that could combat drug-resistant infections

How Google Maps uses machine learning to predict bus traffic delays in real time

A.I. is transforming the job interview—and everything after

“AI-based ‘digital bridge’ enables paraplegic patient to walk”

Will AI cause or solve more problems?

Federal Oversight

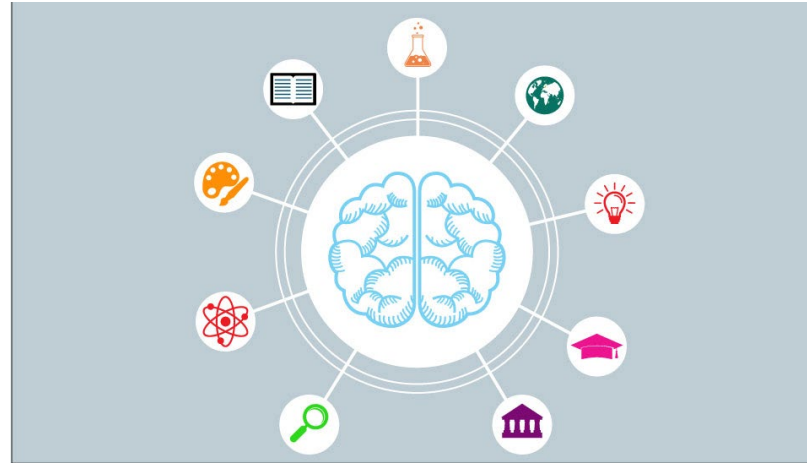
White House unveils efforts to guide federal research of AI

Artificial intelligence prominent in automotive market

Deepfakes: faces created by AI now look more real than genuine photos

What is Intelligence?

Intelligence



...the ability to reason, plan, solve problems, think abstractly, comprehend complex ideas, learn quickly and learn from experience.

Artificial
Intelligence
(AI)



AI is the ability of a machine or computer system to emulate aspects of human intelligence:

Reasoning

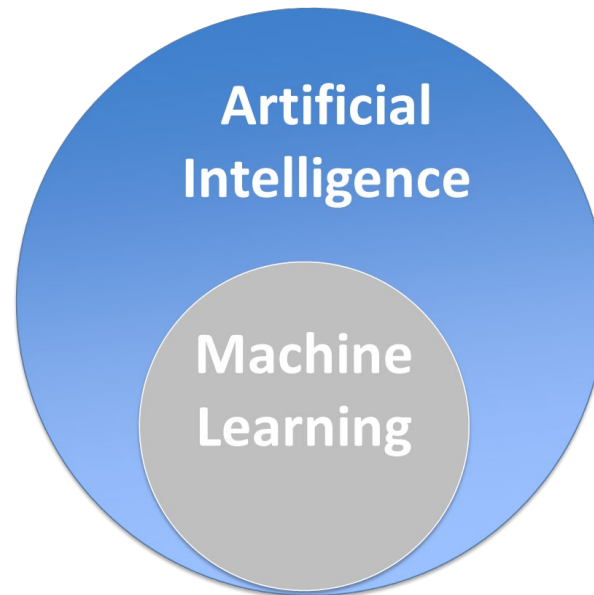
Learning

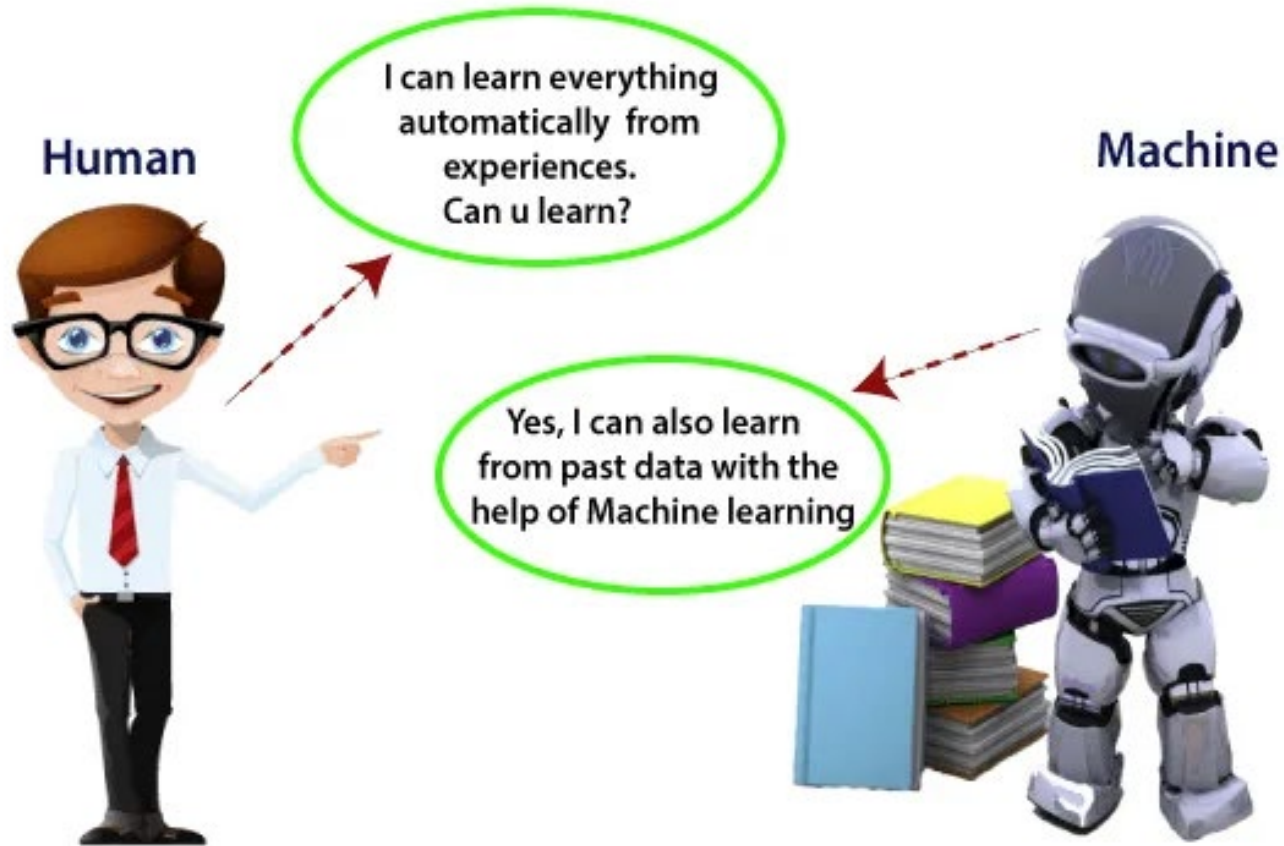
Intelligent behaviour and thought

Capable of analyzing the environment and performing action

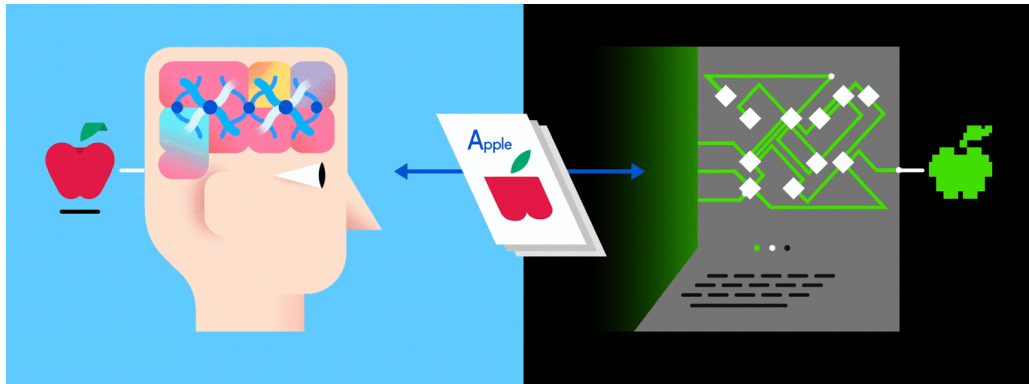
Artificial Intelligence: The general picture

- *The effort to automate intellectual tasks normally performed by humans.*
- AI is a general field that includes Machine Learning and Deep Learning, **but also other approaches that do not involve any learning at all**





Machine learning, a branch of artificial intelligence, concerns the construction and study of systems that can **learn** from data.



Source: GumGum

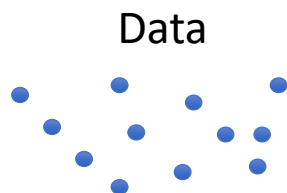
Machine Learning (Aprendizaje Automático) a branch of artificial intelligence, concerns the construction and study of systems that can **learn** from data.



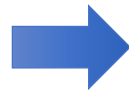
Learning from experience



Learning from **experience** → data



Data



Machine Learning
Techniques



To predict future events
Infer the causes of events

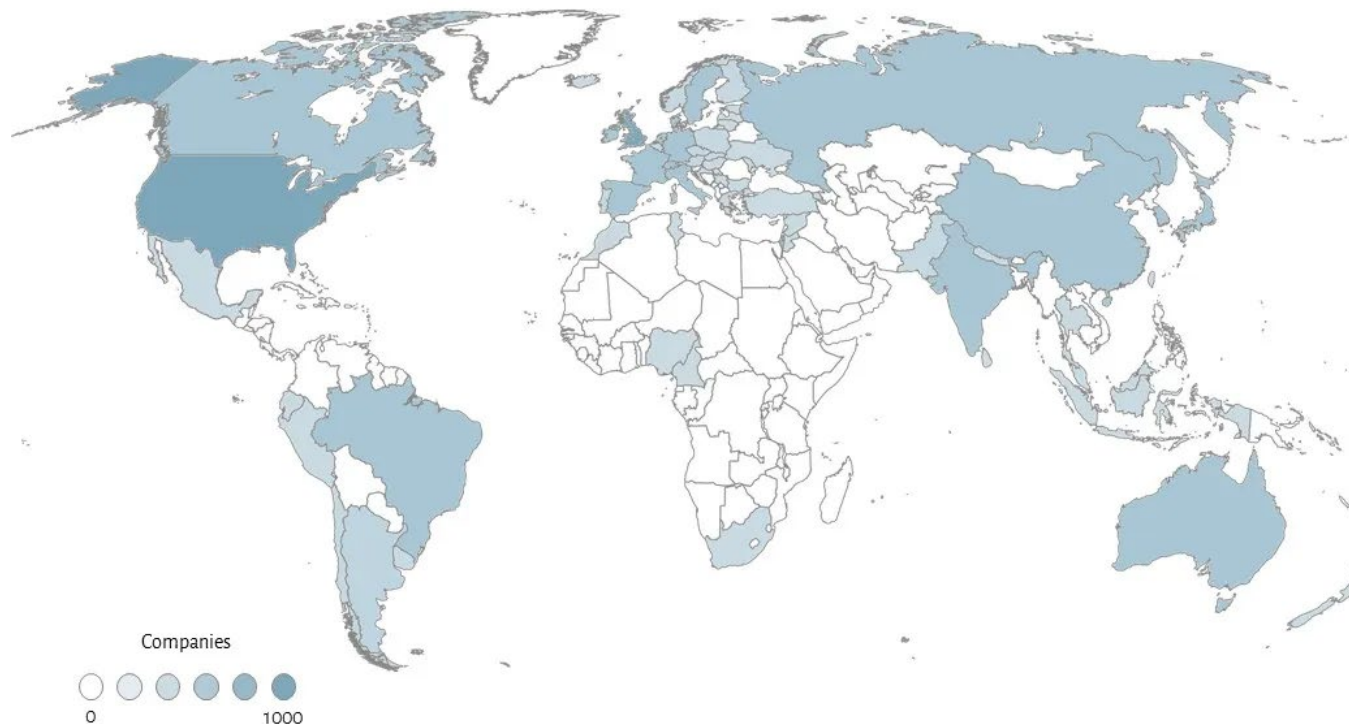
- Technology allows us to have, store and process **large amounts of data (information society)**
- **High commercial and industrial interest** in the development of techniques for extracting knowledge from data, for finding patterns. Particularly in problems:
 - where algorithms do not exist
 - not well defined
 - informally proposed
- Great progress in the **development of algorithms and models by researchers**. Development of tools:
 - Classification (Assignment to a predefined category)
 - Regression (Estimation of a numerical value)



- Why data-driven learning (machine learning)?
- Ability to mimic humans and replace them in monotonous tasks that require intelligence
 - Handwriting recognition
- Develop systems that can automatically adapt to individual users
 - Personalised news feeds, e-mail filters
- Extracting knowledge from large databases
 - Shopping basket analysis

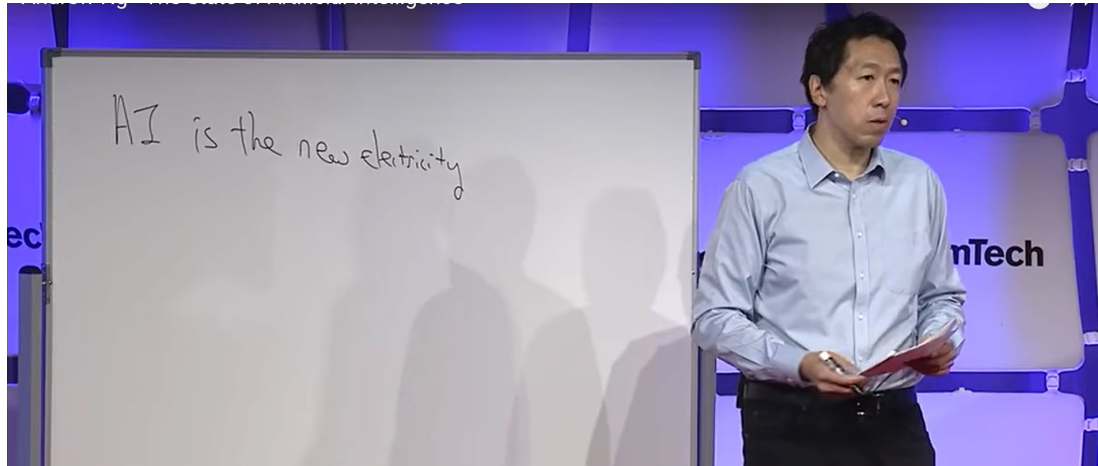
- Mundane - Ordinary Tasks (humans learn ordinary tasks since their birth but we can't describe how we do them).
Character recognition
- Expert Tasks
Quality control in manufacturing
- Problems where there are no human experts
Importance of certain genes in disease risk
- Situations where each user has his or her own target function
Newspapers with personalised news, personalised advertising
- Problems where the volume of data makes it impossible for humans to perform any analysis
Techniques capable of finding relationships within large DBs

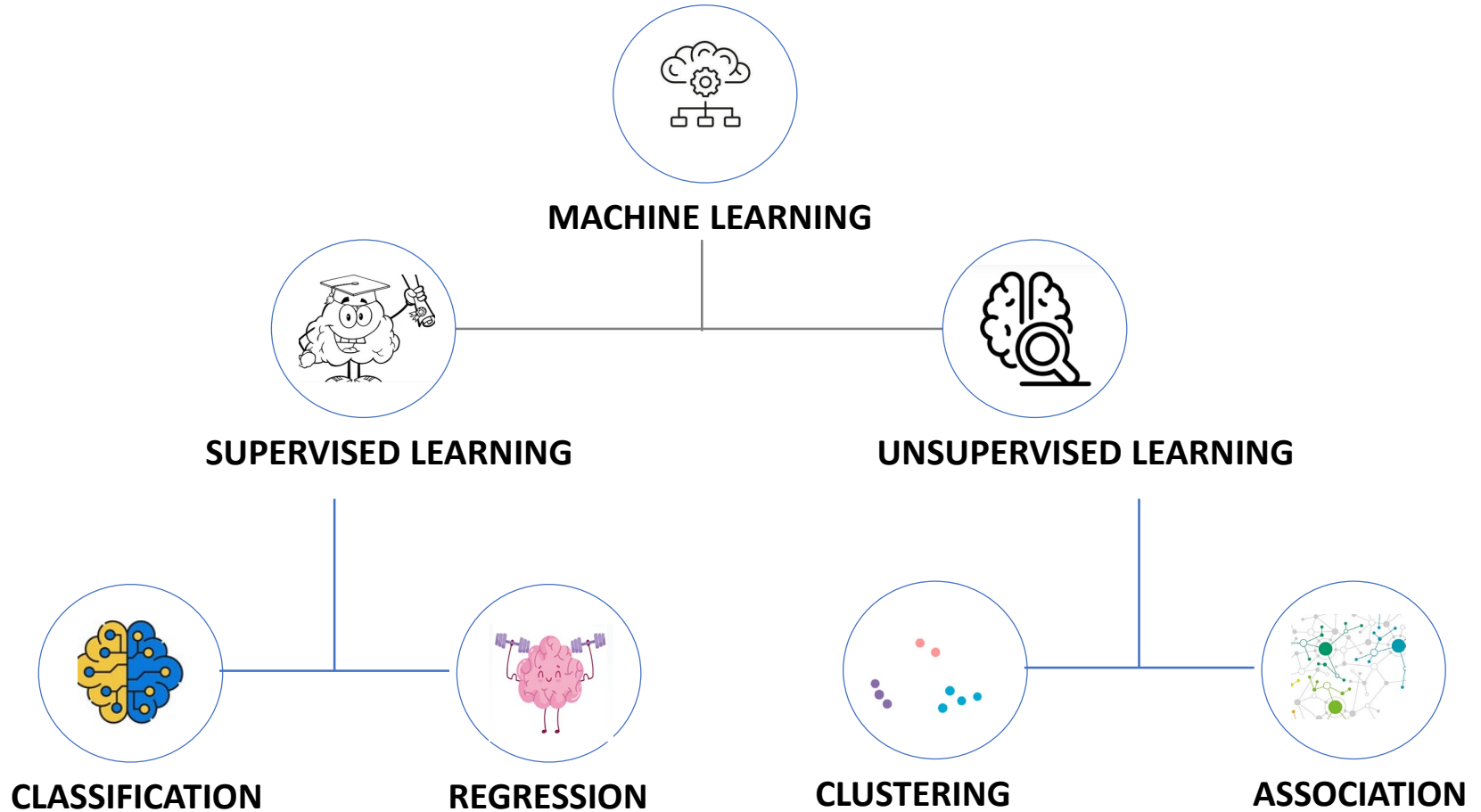
We see it practically every day, across every industry. From healthcare to agriculture, entertainment to transportation, AI applications are shaping our present and redefining our future.

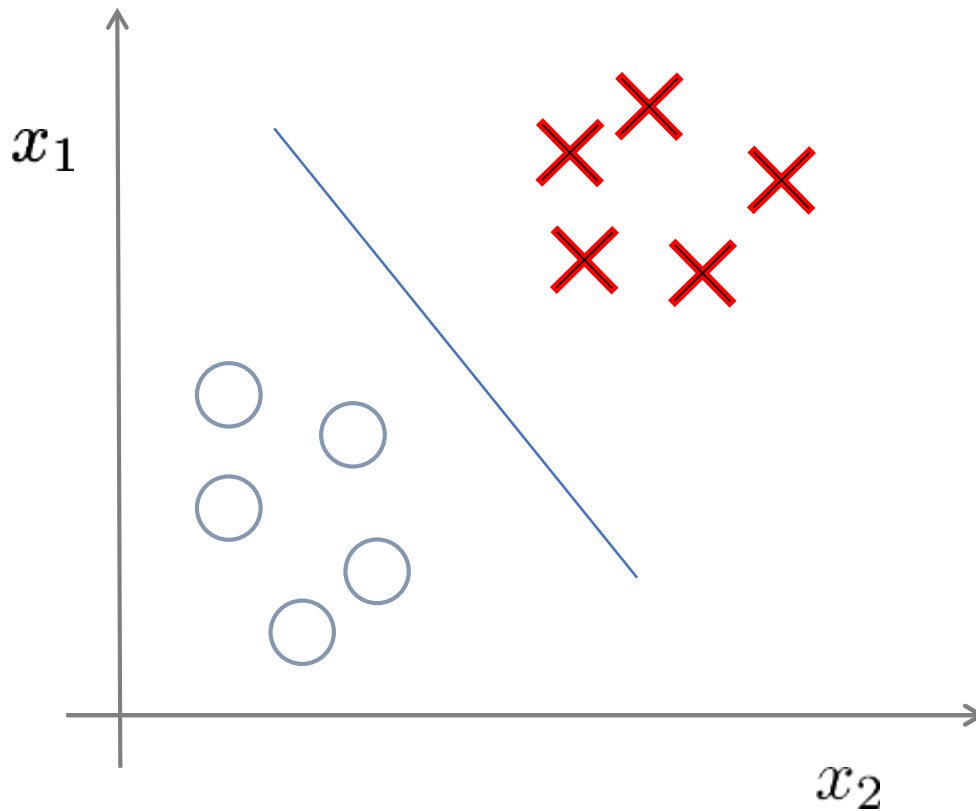


Dr. Andrew Ng:

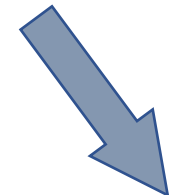
- Globally recognized leader in AI
- He is Founder of DeepLearning.AI, Founder & CEO of Landing AI, General Partner at AI Fund, Chairman and Co-Founder of Coursera
- Professor at Stanford University
- Pioneer in machine learning and online education



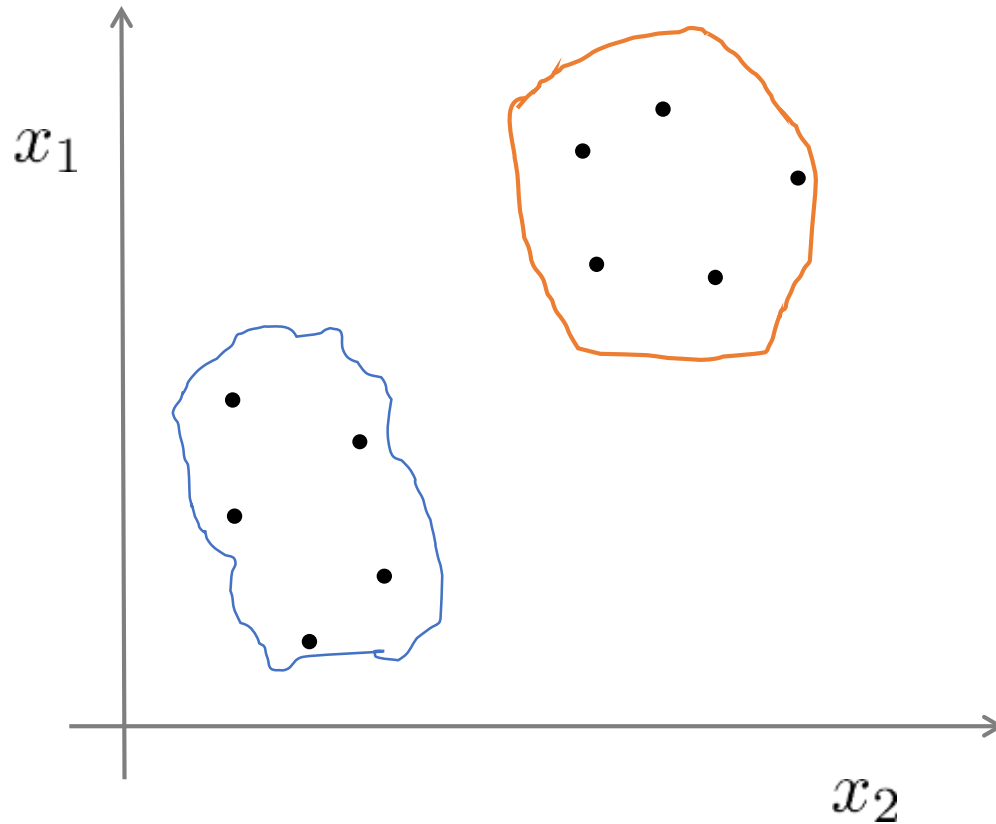




WITH labels



Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$



Training set: $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$ **WITHOUT labels**

We have a set of labelled data (examples + labels)

Classification

Feature Space \mathcal{X}

Words in a document

Label Space \mathcal{Y}

"Sports"
"News"
"Science"
...

Discrete labels (categories)

Regression









Market information up to time t

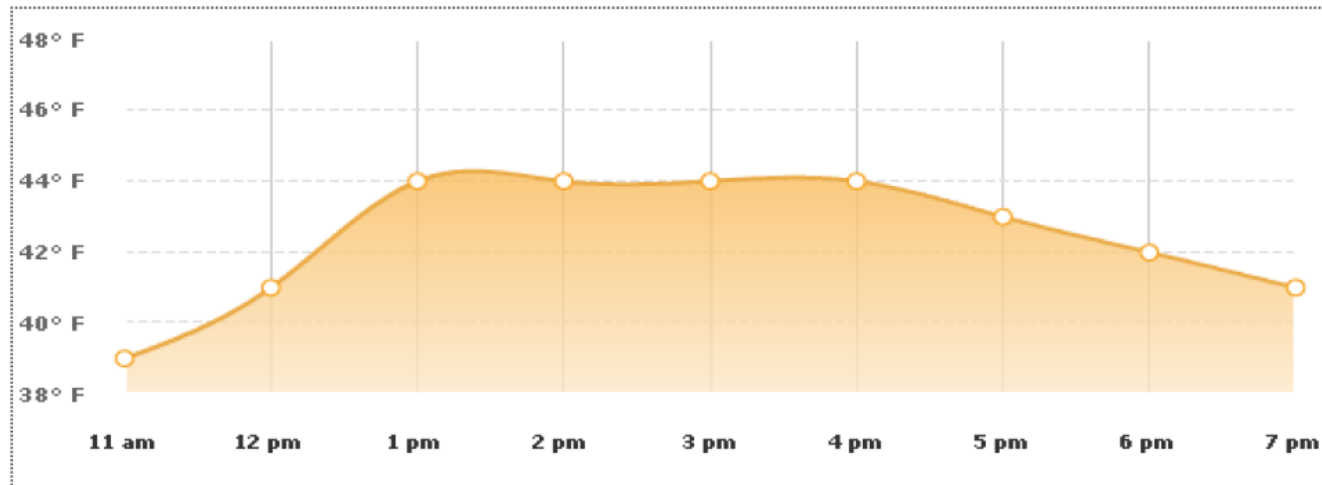
Share Price
"\$ 24.50"

Continuous numerical values

Task: Given $X \in \mathcal{X}$, predict $Y \in \mathcal{Y}$.

Regression or Classification?

11 am	12 pm	1 pm	2 pm	3 pm	4 pm	5 pm	6 pm
							
39° F	41° F	44° F	44° F	44° F	44° F	43° F	42° F
Precip: 10%	Precip: 10%	Precip: 10%	Precip: 10%	Precip: 10%	Precip: 10%	Precip: 10%	Precip: 0%



Temperature/Weather Forecast

Regression or Classification?

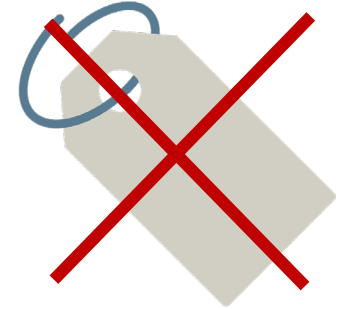
Input (x)	Output
Home features (#bedrooms, size,...)	Price
Advertisement, user info	Click on ad? (Yes/No)
Image	Object (1,2,...,1000)
Age, sex, cholesterol, #Cigarettes, blood sugar, family history	Heart disease (True/False)
Employee's attributes (seniority, income, department, distance from home,..)	How long until an employee looks for another job
Age, level of education, area, job title,..	Income

We have an unlabelled dataset (examples)
Also known as “learning without a teacher”

Feature Space \mathcal{X}

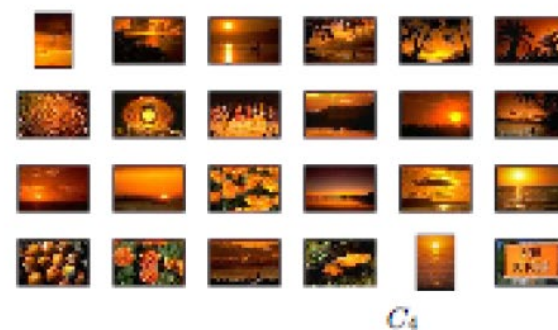


Word distribution
(Probability of a word)



Task: Given $X \in \mathcal{X}$, learn $f(X)$.

Cluster similar items, for example, images



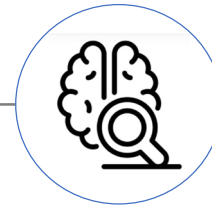
Goldberger et al.

Cluster similar items, for example, customers

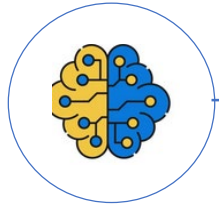




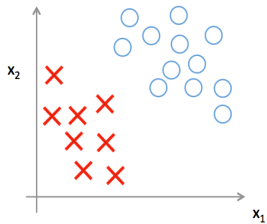
SUPERVISED LEARNING



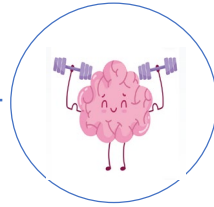
UNSUPERVISED LEARNING



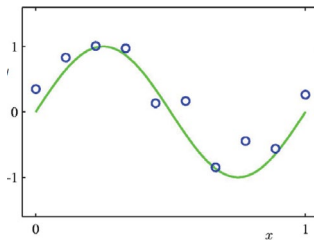
CLASSIFICATION



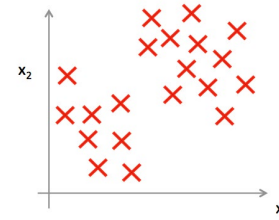
Data with labels (finite set of symbolic classes)



REGRESSION



Data with labels (numerical label)



Data

No labels



APPROACHING A PROBLEM OF LEARNING FROM EXAMPLES

Hands-on an illustrative problem



What do we need ?

- A system that verifies that the person is who he/she claims to be to enter the security room.

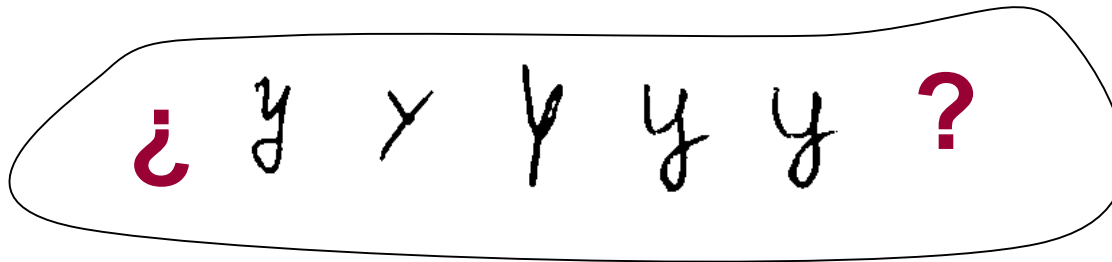


¿user?



Identity verification with biometric features

Handwriting



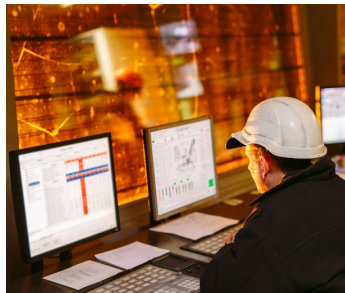
Classification

- It verifies that the person is who he/she claims to be.
- It is based on a pattern recognition system.

y ?
=



x ?
=



Basic steps to create the system

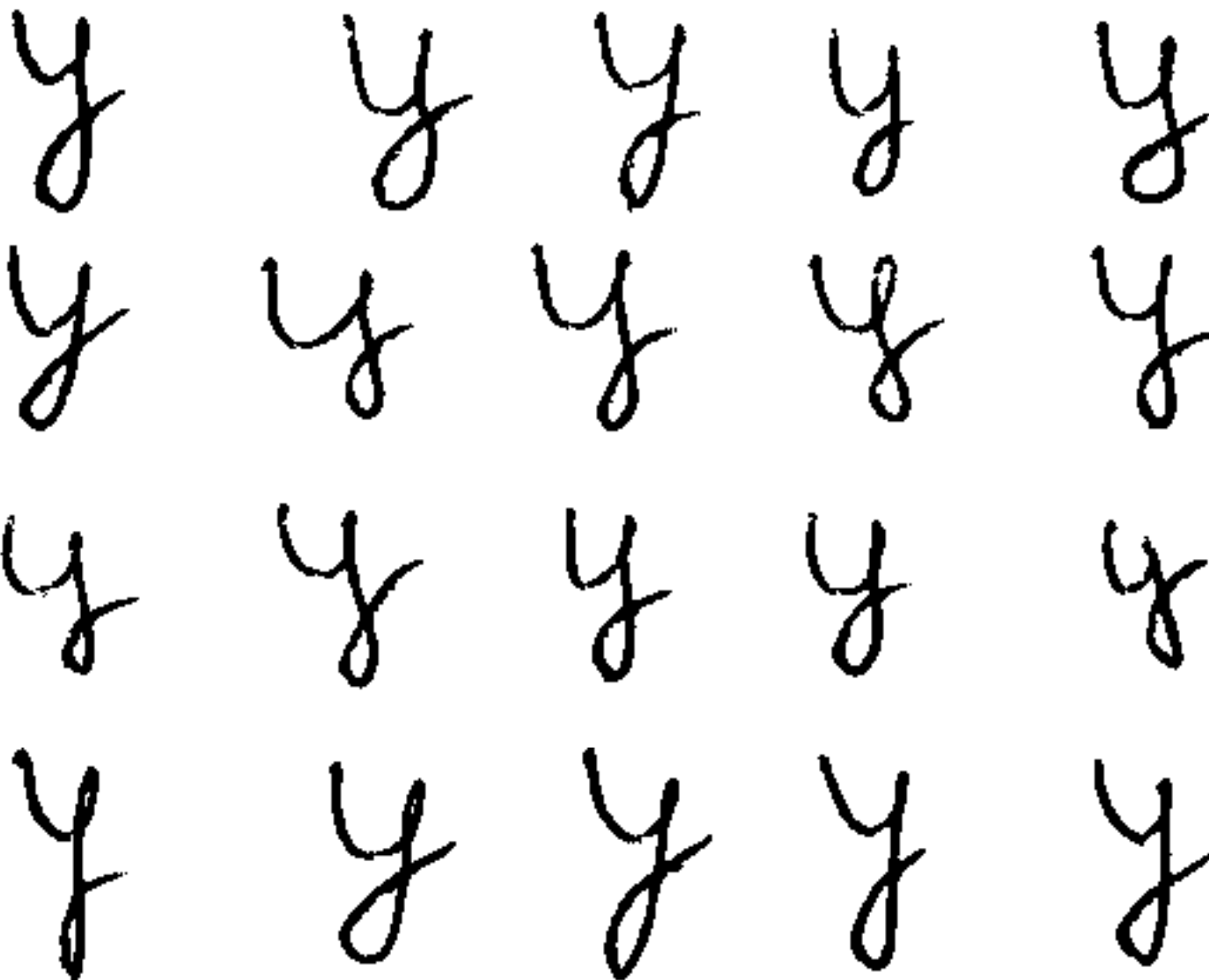
- Dataset
- Training the model
- Model test
- Model evaluation

Basic steps to create the system

- Dataset
- Training the model
- Model test
- Model evaluation

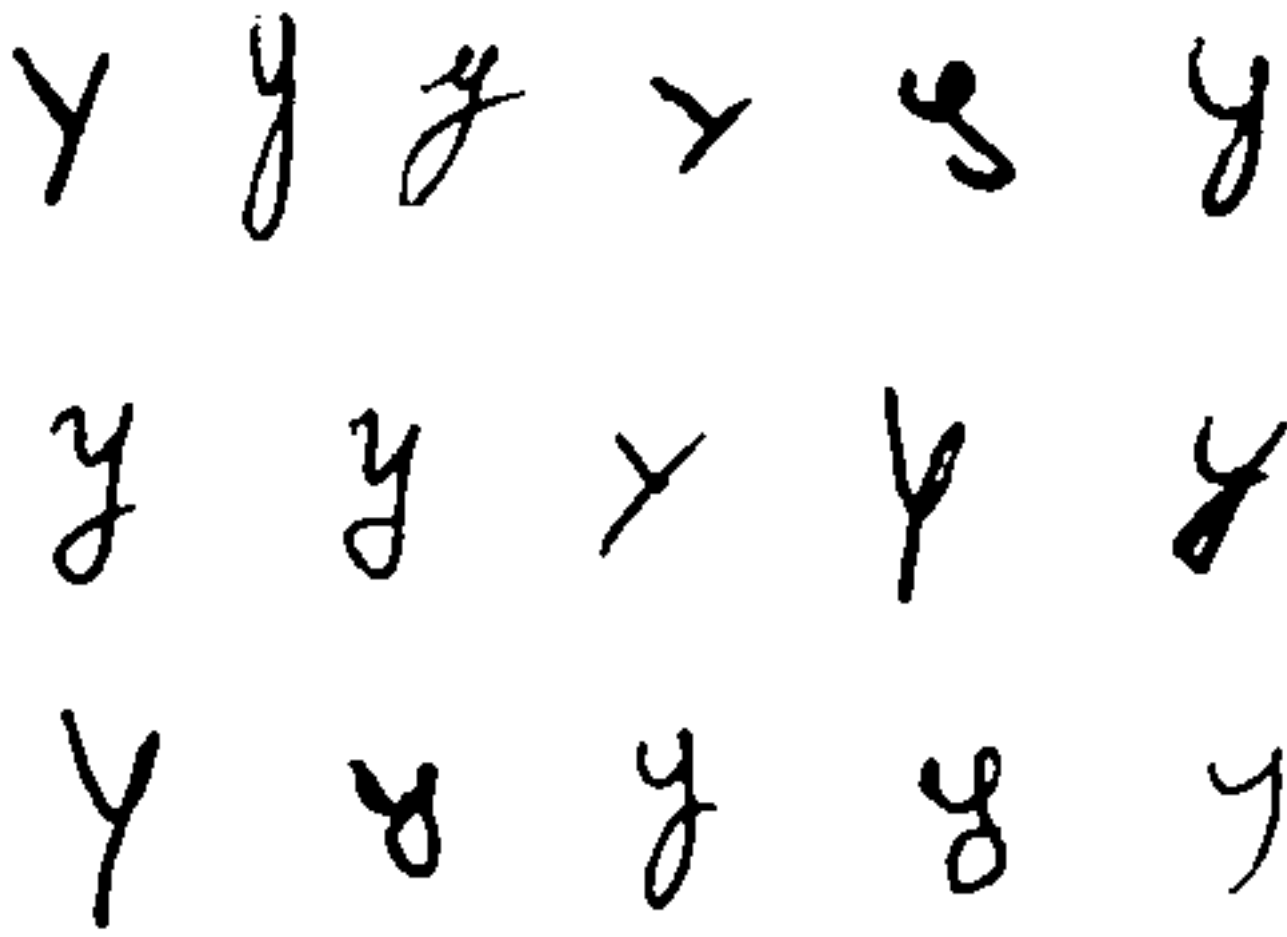
User X

Handwriting "y"

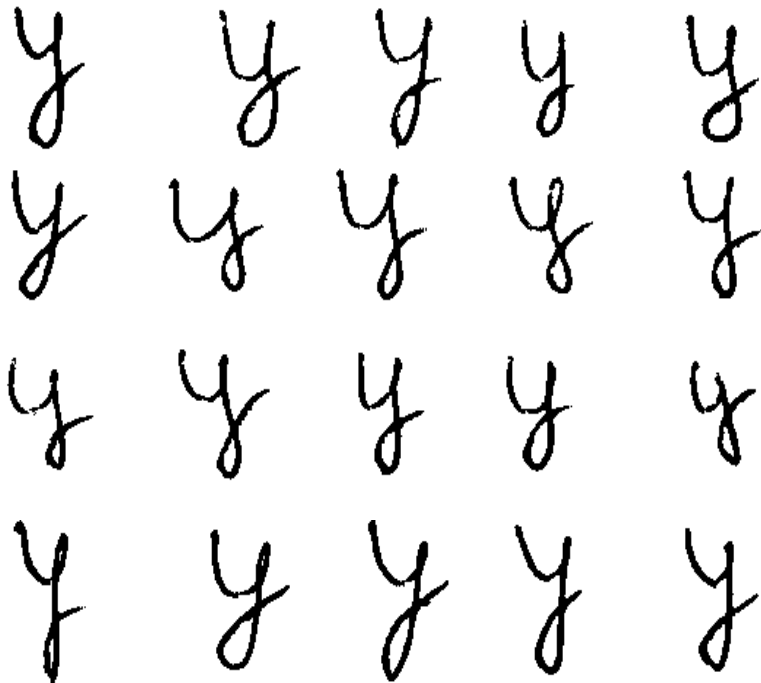
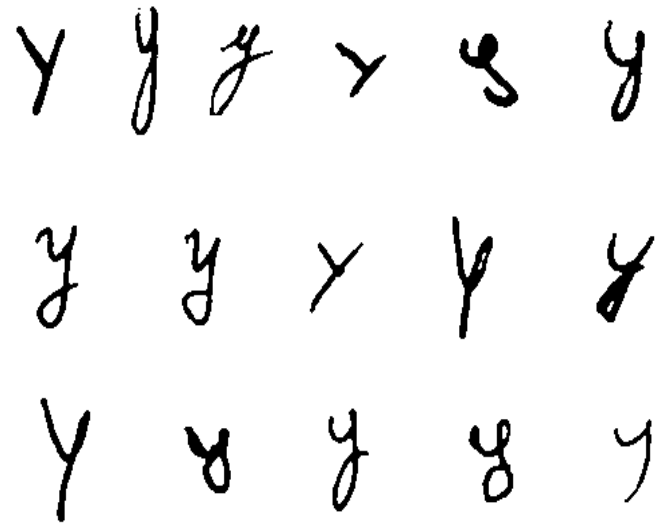


Other users

Handwriting "y"



Available dataset

 $d = 1$ **User X** $d = 0$ **Other users**

Basic steps to create the system

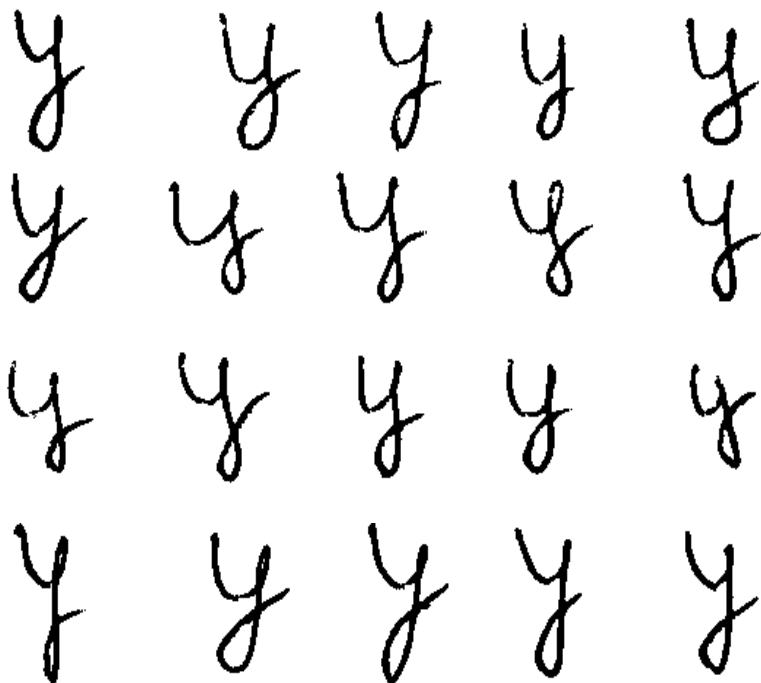
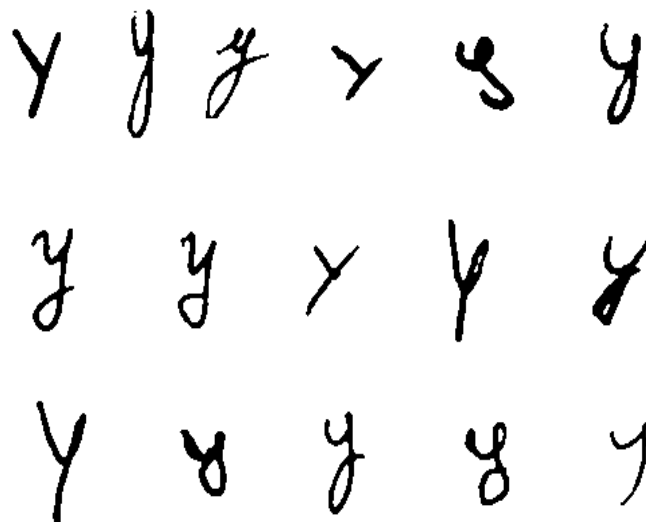
- Dataset
- Training the model
- Model test
- Model evaluation

What can YOU do?

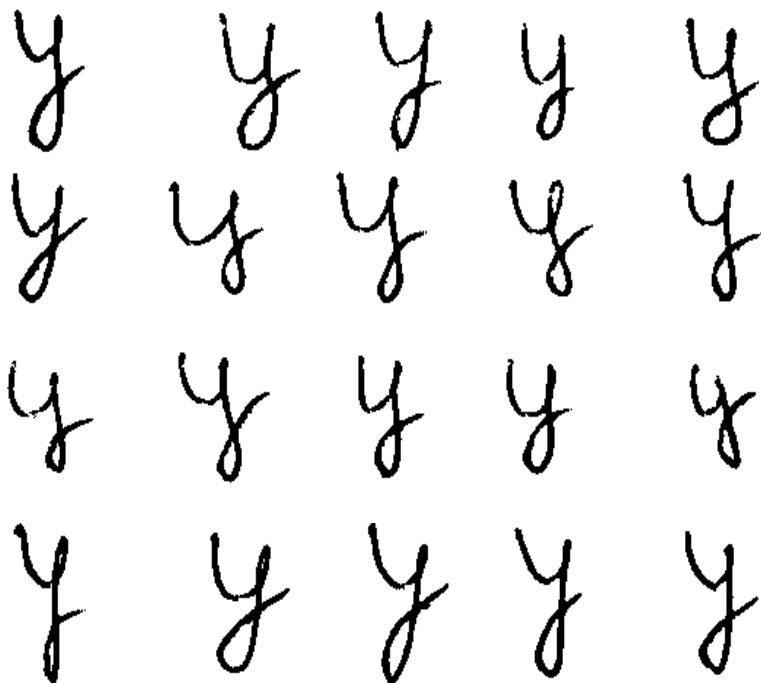
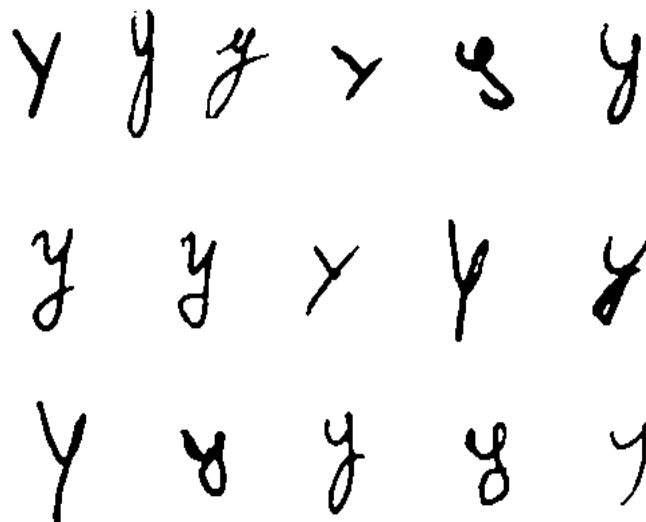


TRY TO LEARN

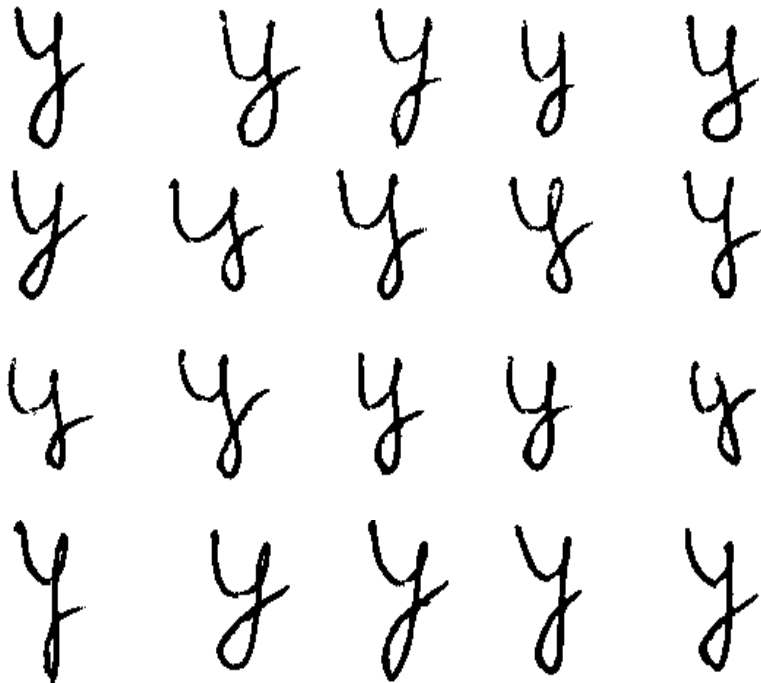
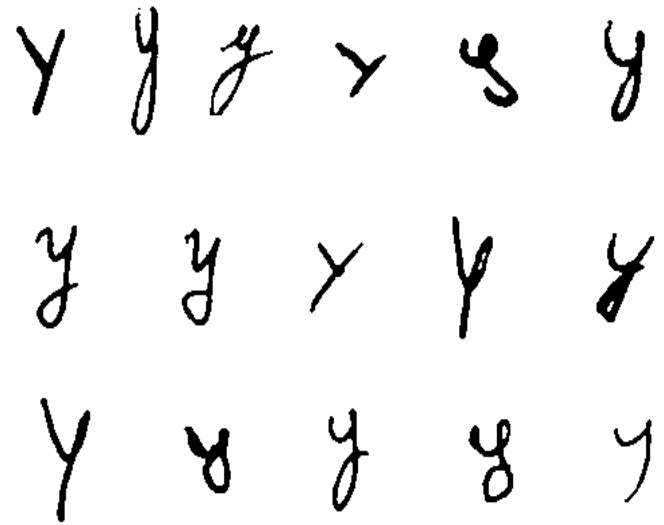
Model training

 $d = 1$ **User X** $d = 0$ **Other users**

Model training

 $d = 1$ **User X** $d = 0$ **Other users**

Available dataset

 $d = 1$ **User X** $d = 0$ **Other users**

Basic steps to create the system

- Dataset
- Training the model
- Model test
- Model evaluation

$$\hat{d} = 1$$

A handwritten digit '4' in black ink, written in a cursive style.

User X

$$d = 1$$

$$\hat{d} = 0$$

y

User X

$$d = 1$$

False rejection

Basic steps to create the system

- Dataset
- Training the model
- Model test
- Model evaluation

TEST

¿ $\hat{d} = 0$ or $\hat{d} = 1$?



$d =$ $\hat{d} =$



$d =$ $\hat{d} =$



$d =$ $\hat{d} =$



$d =$ $\hat{d} =$



$d =$ $\hat{d} =$



$d =$ $\hat{d} =$



$d =$ $\hat{d} =$



$d =$ $\hat{d} =$



$d =$ $\hat{d} =$



$d =$ $\hat{d} =$

¿ $\hat{d} = 0$ or $\hat{d} = 1$?



$\hat{d} = 1$



$\hat{d} = 0$



$\hat{d} = 1$



$\hat{d} = 0$



$\hat{d} = 0$



$\hat{d} = 1$



$\hat{d} = 0$



$\hat{d} = 0$




$\hat{d} = 1$




$\hat{d} = 0$


¿ $\hat{d} = 0$ or $\hat{d} = 1$?




$d = 1$ $\hat{d} = 1$




$d = 0$ $\hat{d} = 0$




$d = 0$ $\hat{d} = 1$




$d = 1$ $\hat{d} = 0$




$d = 0$ $\hat{d} = 0$




$d = 1$ $\hat{d} = 1$




$d = 1$ $\hat{d} = 0$



$d = 0$ $\hat{d} = 0$



$d = 1$ $\hat{d} = 1$



$d = 0$ $\hat{d} = 0$

Evaluation metrics

Error rate?

False Rejection
Rate?

False
Acceptance
Rate?



SUPERVISED LEARNING MODELS

Supervised Learning Models

K-NN

Naïve
Bayes

SVM

Decision
Trees

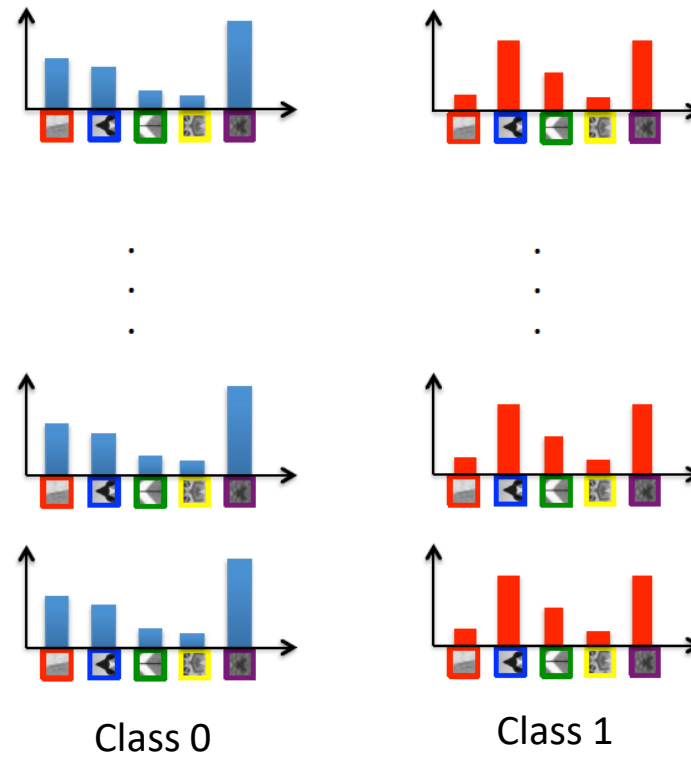
Random
Forest

Neural
Networks

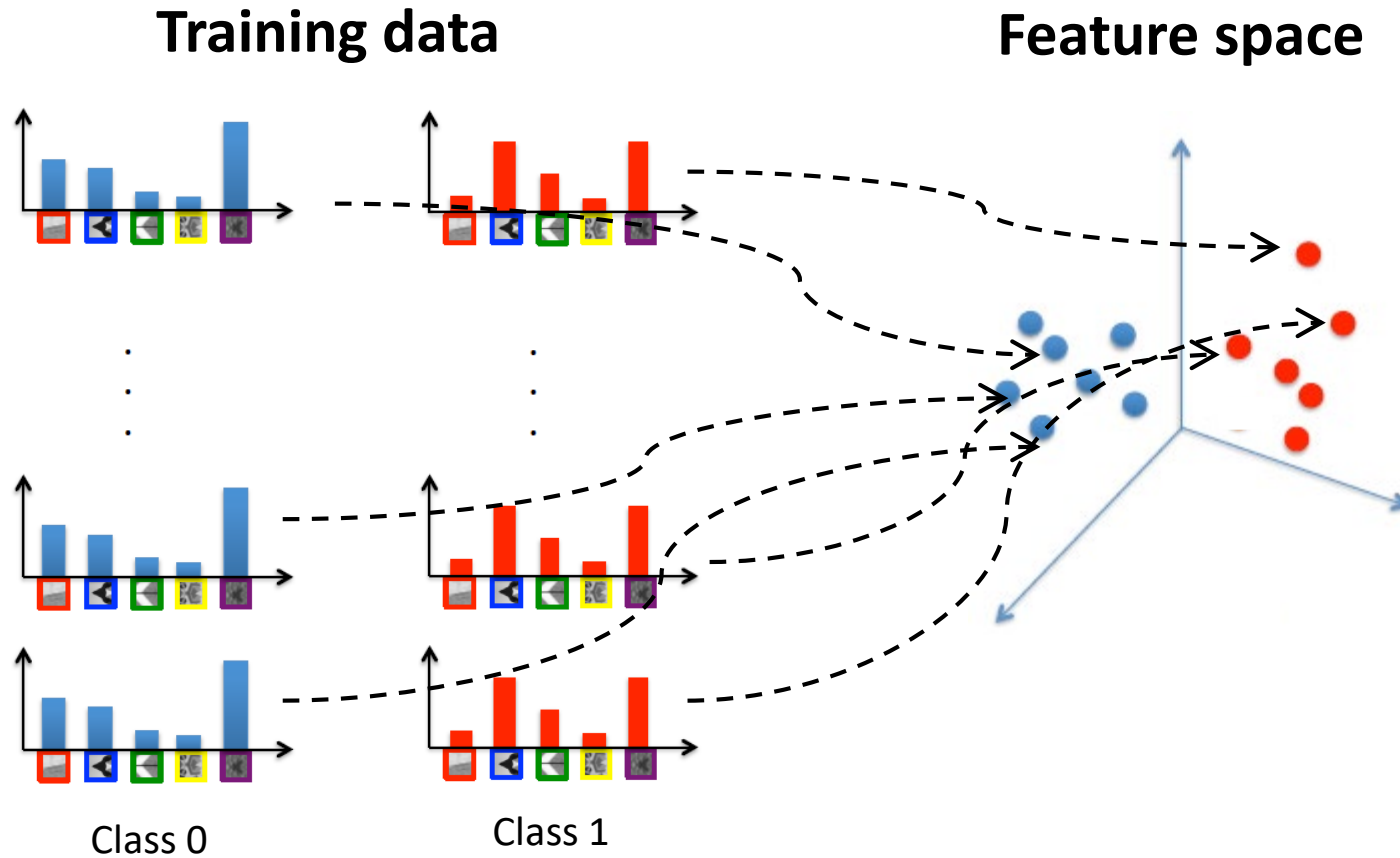
SUPERVISED LEARNING: K-NEAREST NEIGHBOURS (K-NN)

Use a set of **training data** to adjust the model

kNN Assumes that the samples from the same class are **close** in the feature space

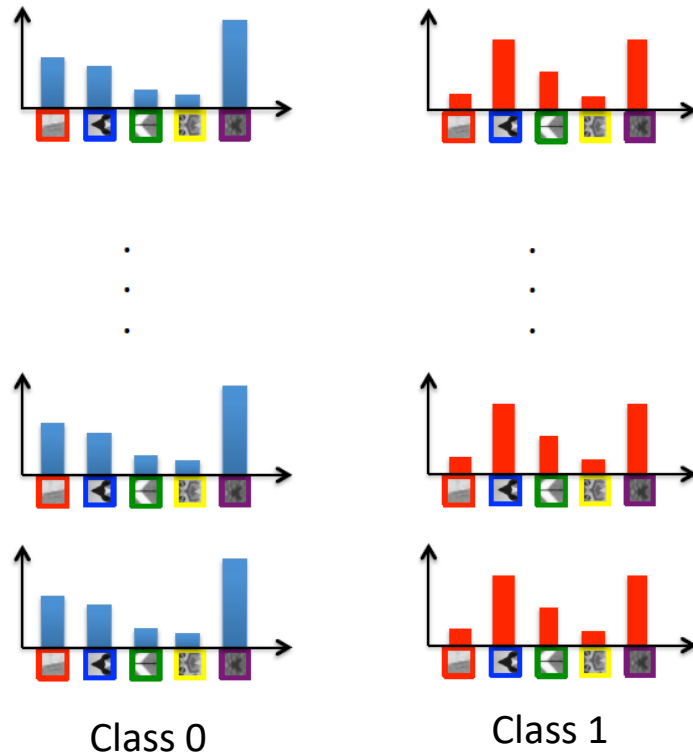


K-nearest Neighbours (kNN)

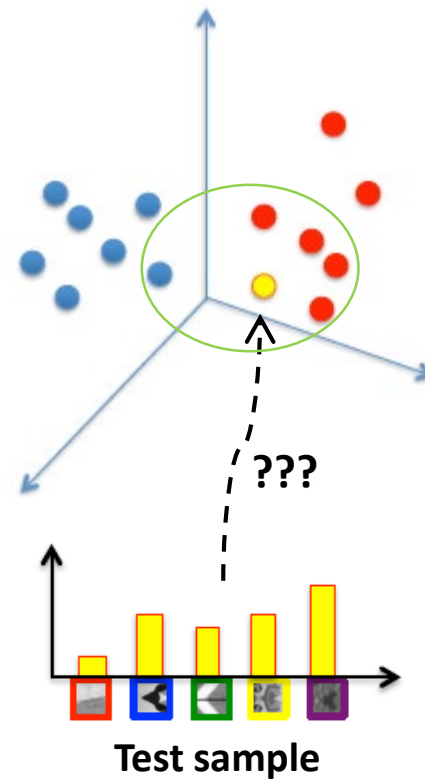


K-nearest Neighbours (kNN)

Training data

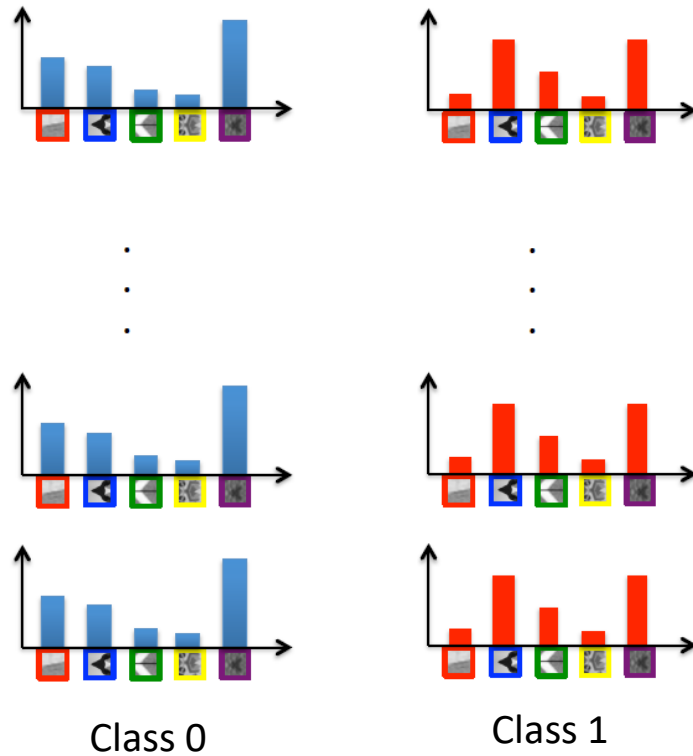


Feature space

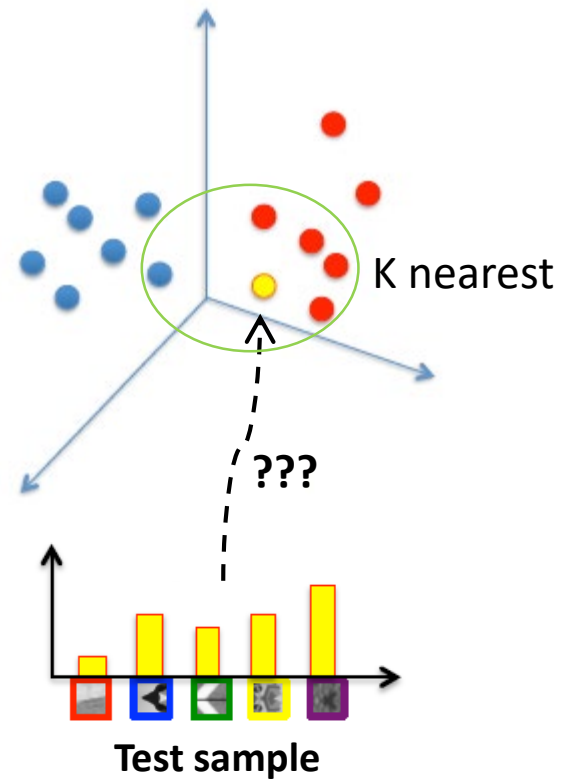


K-nearest Neighbours (kNN)

Training data

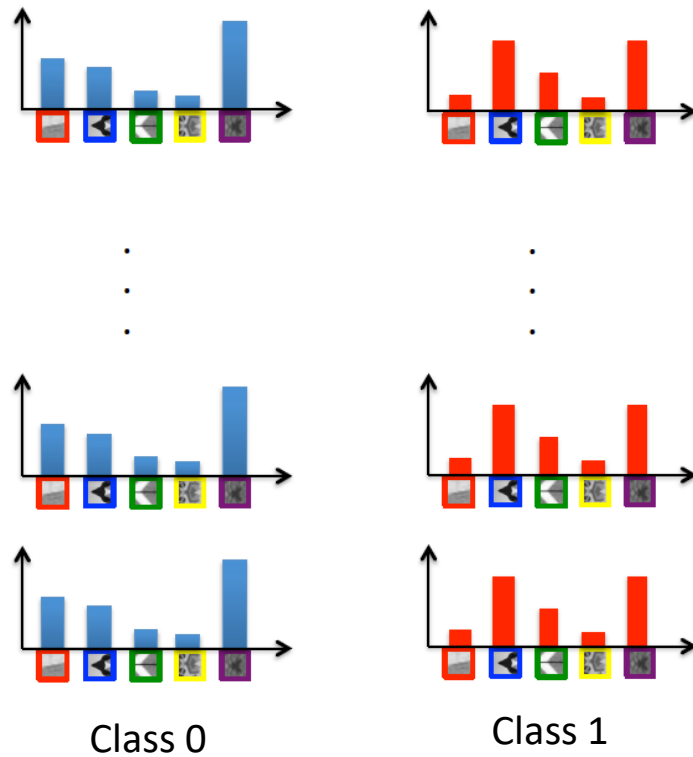


Feature space

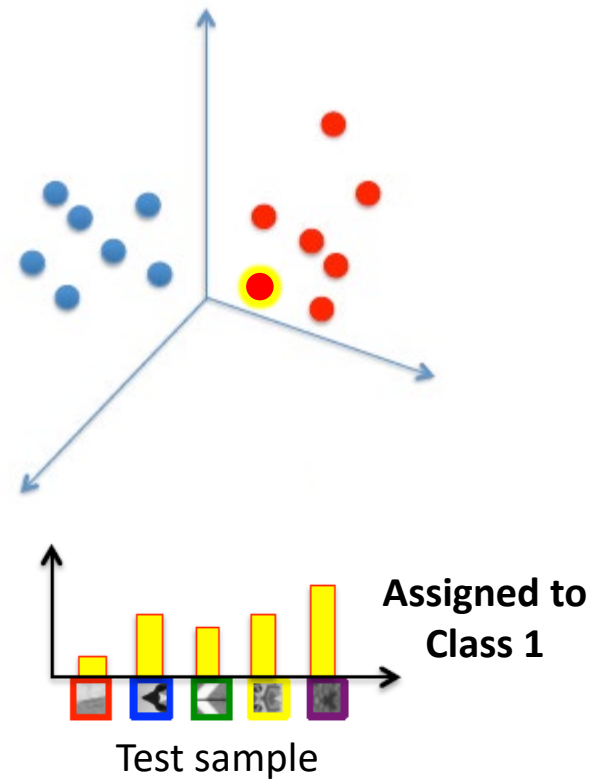


K-nearest Neighbours (kNN)

Training data



Feature space



Euclidean distance:
$$d(\mathbf{x}^{(r)}, \mathbf{x}^{(s)}) = \sqrt{\sum_{i=1}^n (x_i^{(r)} - x_i^{(s)})^2}$$

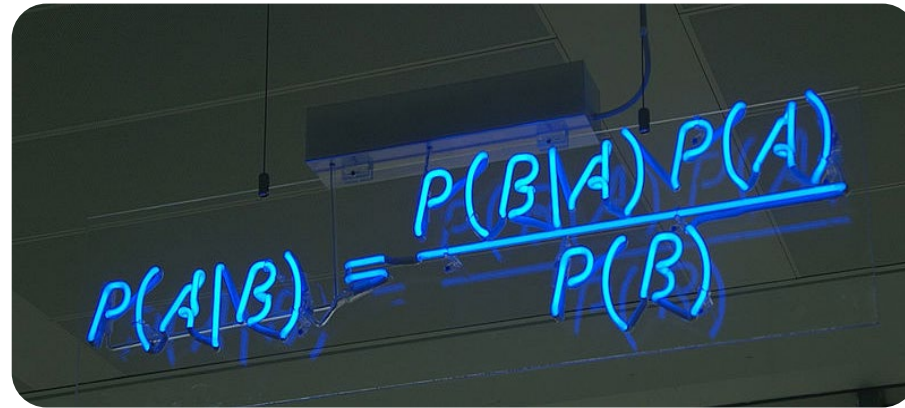
Manhattan distance:
$$d(\mathbf{x}^{(r)}, \mathbf{x}^{(s)}) = \sum_{i=1}^n |x_i^{(r)} - x_i^{(s)}|$$

Chebyshev distance:
$$d(\mathbf{x}^{(r)}, \mathbf{x}^{(s)}) = \max_{i=1,2,\dots,n} |x_i^{(r)} - x_i^{(s)}|$$

Cosine distance:
$$d(\mathbf{x}^{(r)}, \mathbf{x}^{(s)}) = \arccos \left(\frac{\sum_{i=1}^n x_i^{(r)} x_i^{(s)}}{\sqrt{\sum_{i=1}^n x_i^{(r)}} \cdot \sqrt{\sum_{i=1}^n x_i^{(s)}}} \right)$$

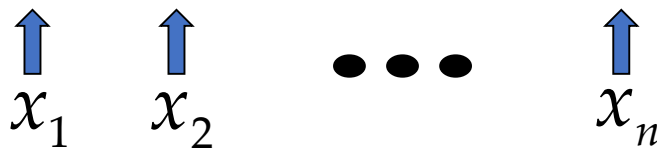
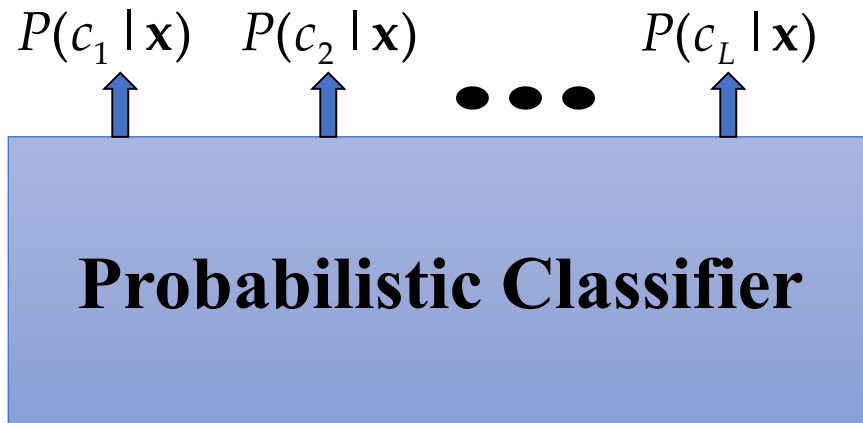
SUPERVISED LEARNING: NAIVE BAYES

Naive Bayes Classifier

A photograph of a whiteboard with a blue marker formula. The formula is $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$. The whiteboard has a dark background and the text is written in bright blue. There is a reflection of the whiteboard on the surface below it.
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- Establishing a probabilistic model for classification

$$P(c|\mathbf{x}) \quad c = c_1, \dots, c_L, \quad \mathbf{x} = (x_1, \dots, x_n)$$



$$\mathbf{x} = (x_1, x_2, \dots, x_n)$$

Prior, conditional and joint probability for random variables

- Prior probability: $P(c)$
- Conditional probability: $P(x_1 | x_2), P(x_2 | x_1)$
- Joint probability: $\mathbf{x} = (x_1, x_2), P(\mathbf{x}) = P(x_1, x_2)$
- Relationship: $P(x_1, x_2) = P(x_2 | x_1)P(x_1) = P(x_1 | x_2)P(x_2)$
- Independence: $P(x_2 | x_1) = P(x_2)$
 $P(x_1 | x_2) = P(x_1)$ $P(x_1, x_2) = P(x_1)P(x_2)$

Bayesian Rule

$$P(c | \mathbf{x}) = \frac{P(\mathbf{x} | c)P(c)}{P(\mathbf{x})}$$

$$Posterior = \frac{Likelihood \times Prior}{Evidence}$$

- **M**aximum **A** Posterior (**MAP**) classification rule
 - For an input \mathbf{x} , find the largest one from L probabilities output by a probabilistic classifier $P(c_1 | \mathbf{x}), \dots, P(c_L | \mathbf{x})$.
 - Assign \mathbf{x} to label c^* if $P(c^* | \mathbf{x})$ is the largest.
- Classification with the MAP rule
 - Apply Bayesian rule to get posterior probabilities

$$P(c_i | \mathbf{x}) = \frac{P(\mathbf{x} | c_i)P(c_i)}{P(\mathbf{x})} \propto P(\mathbf{x} | c_i)P(c_i)$$

for $i = 1, 2, \dots, L$

Common factor for all L probabilities

Bayes classification

$$P(c | \mathbf{x}) \propto P(\mathbf{x} | c)P(c) = P(x_1, \dots, x_n | c)P(c) \text{ for } c = c_1, \dots, c_L.$$

Difficulty: learning the joint probability $P(x_1, \dots, x_n | c)$

Naïve Bayes classification

- Assume **all input features are class conditionally independent!**

$$\begin{aligned} P(x_1, x_2, \dots, x_n | c) &= \frac{P(x_1 | x_2, \dots, x_n, c)P(x_2, \dots, x_n | c)}{=} P(x_1 | c)P(x_2, \dots, x_n | c) \\ &= P(x_1 | c)P(x_2 | c) \cdots P(x_n | c) \end{aligned}$$

Applying the independence assumption

- Apply the MAP classification rule: assign $\mathbf{x}' = (a_1, a_2, \dots, a_n)$ to c^* if

$$\underbrace{[P(a_1 | c^*) \cdots P(a_n | c^*)]P(c^*)}_{\text{estimate of } P(a_1, \dots, a_n | c^*)} > \underbrace{[P(a_1 | c) \cdots P(a_n | c)]P(c)}_{\text{estimate of } P(a_1, \dots, a_n | c)}, \quad c \neq c^*, c = c_1, \dots, c_L$$

Learning Phase: Given a training set \mathbf{S} , with L classes and n features

For each target value of c_i ($c_i = c_1, \dots, c_L$)

$\hat{P}(c_i) \leftarrow$ estimate $P(c_i)$ with examples in \mathbf{S} ;

For every attribute value x_{jk} of each attribute x_j ($j = 1, \dots, n; k = 1, \dots, N_j$)

$\hat{P}(x_{jk} | c_i) \leftarrow$ estimate $P(x_{jk} | c_i)$ with examples in \mathbf{S} ;

Output: conditional probability tables; for x_j , $N_j \times L$ elements

Test Phase: Given an unknown instance $\mathbf{x}' = (a'_1, \dots, a'_n)$

Look up tables to assign the label c^* to \mathbf{X}' if

$$[\hat{P}(a'_1 | c^*) \cdots \hat{P}(a'_n | c^*)] \hat{P}(c^*) > [\hat{P}(a'_1 | c) \cdots \hat{P}(a'_n | c)] \hat{P}(c), \quad c \neq c^*, c = c_1, \dots, c_L$$

PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Learning Phase: Given a training set \mathbf{S} , with L classes and n features

For each target value of c_i ($c_i = c_1, \dots, c_L$)

$\hat{P}(c_i) \leftarrow$ estimate $P(c_i)$ with examples in \mathbf{S} ;

For every attribute value x_{jk} of each attribute x_j ($j = 1, \dots, n; k = 1, \dots, N_j$)

$\hat{P}(x_{jk} | c_i) \leftarrow$ estimate $P(x_{jk} | c_i)$ with examples in \mathbf{S} ;

Output: conditional probability tables; for x_j , $N_j \times L$ elements

L= ?

n= ?

Learning Phase: Given a training set \mathbf{S} , with L classes and n features

For each target value of c_i ($c_i = c_1, \dots, c_L$)

$\hat{P}(c_i) \leftarrow$ estimate $P(c_i)$ with examples in \mathbf{S} ;

For every attribute value x_{jk} of each attribute x_j ($j = 1, \dots, n; k = 1, \dots, N_j$)

$\hat{P}(x_{jk} | c_i) \leftarrow$ estimate $P(x_{jk} | c_i)$ with examples in \mathbf{S} ;

Output: conditional probability tables; for x_j , $N_j \times L$ elements

$$P(\text{Play}=\text{Yes}) = 9/14 \quad P(\text{Play}=\text{No}) = 5/14$$

Learning Phase: Given a training set \mathbf{S} , with L classes and n features

For each target value of c_i ($c_i = c_1, \dots, c_L$)

$\hat{P}(c_i) \leftarrow$ estimate $P(c_i)$ with examples in \mathbf{S} ;

For every attribute value x_{jk} of each attribute x_j ($j = 1, \dots, n; k = 1, \dots, N_j$)

$\hat{P}(x_{jk} | c_i) \leftarrow$ estimate $P(x_{jk} | c_i)$ with examples in \mathbf{S} ;

Output: conditional probability tables; for x_j , $N_j \times L$ elements

Outlook	Play=Yes	Play=No
Sunny	2/9	3/5
Overcast	4/9	0/5
Rain	3/9	2/5

Example: Play Tennis

Learning Phase:

Outlook	Play=Yes	Play=No
Sunny	2/9	3/5
Overcast	4/9	0/5
Rain	3/9	2/5

Temperature	Play=Yes	Play=No
Hot	2/9	2/5
Mild	4/9	2/5
Cool	3/9	1/5

Humidity	Play=Yes	Play=No
High	3/9	4/5
Normal	6/9	1/5

Wind	Play=Yes	Play=No
Strong	3/9	3/5
Weak	6/9	2/5

$$P(\text{Play=Yes}) = 9/14$$

$$P(\text{Play=No}) = 5/14$$

Test Phase: Given a new instance, predict its label

- $\mathbf{x}' = (\text{Outlook}=\textit{Sunny}, \text{Temperature}=\textit{Cool}, \text{Humidity}=\textit{High}, \text{Wind}=\textit{Strong})$
- Look up tables achieved in the learning phrase

$$P(\text{Outlook}=\textit{Sunny} \mid \text{Play}=\textit{Yes}) = 2/9$$

$$P(\text{Temperature}=\textit{Cool} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Humidity}=\textit{High} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Wind}=\textit{Strong} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Play}=\textit{Yes}) = 9/14$$

$$P(\text{Outlook}=\textit{Sunny} \mid \text{Play}=\textit{No}) = 3/5$$

$$P(\text{Temperature}=\textit{Cool} \mid \text{Play}=\textit{No}) = 1/5$$

$$P(\text{Humidity}=\textit{High} \mid \text{Play}=\textit{No}) = 4/5$$

$$P(\text{Wind}=\textit{Strong} \mid \text{Play}=\textit{No}) = 3/5$$

$$P(\text{Play}=\textit{No}) = 5/14$$

- Decision making with the MAP rule

$$P(\text{Yes} \mid \mathbf{x}') \approx [P(\textit{Sunny} \mid \textit{Yes})P(\textit{Cool} \mid \textit{Yes})P(\textit{High} \mid \textit{Yes})P(\textit{Strong} \mid \textit{Yes})]P(\text{Play}=\textit{Yes}) = 0.0053$$

$$P(\text{No} \mid \mathbf{x}') \approx [P(\textit{Sunny} \mid \textit{No})P(\textit{Cool} \mid \textit{No})P(\textit{High} \mid \textit{No})P(\textit{Strong} \mid \textit{No})]P(\text{Play}=\textit{No}) = 0.0206$$

Given the fact $P(\text{Yes} \mid \mathbf{x}') < P(\text{No} \mid \mathbf{x}')$, we label \mathbf{x}' to be “No”.

Pros and cons of Naive Bayes

Advantages

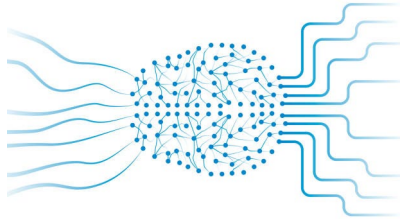
- It's relatively simple to understand and build
- It's easily trained, even with a small dataset
- It's fast!
- It's not sensitive to irrelevant features
- Test is straightforward; just looking up tables

Disadvantages

- It assumes every feature is independent, which isn't always the case

SUPERVISED LEARNING: NEURAL NETWORKS

Inspired by the biological processes: Nervous system



- Simplified mathematical models
- Try to mimic neurons in a very basic setting
- Without claiming to faithfully reflect the real behaviour of the nervous system

- Capable of handling uncertainty
- Robust Solutions
- Not an algorithm
- Non- linear models

S. Ramón y Cajal → Neuron discovery

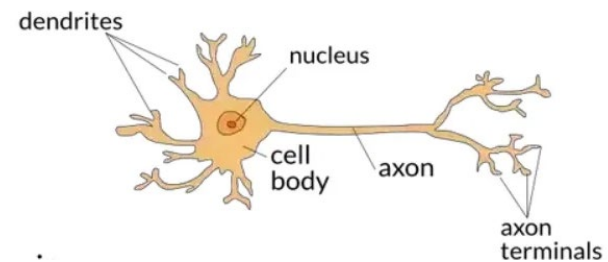
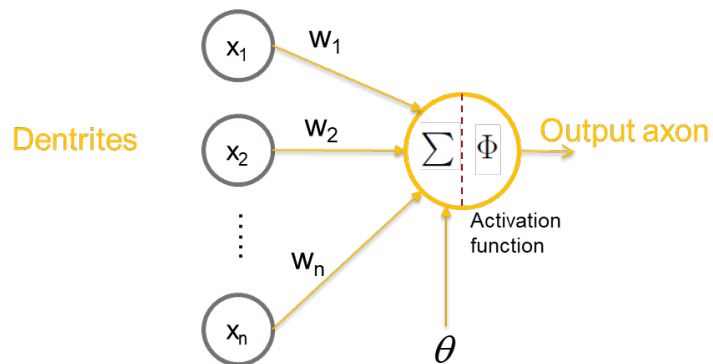


1906 Nobel Prize in Medicine

Artificial Neuron

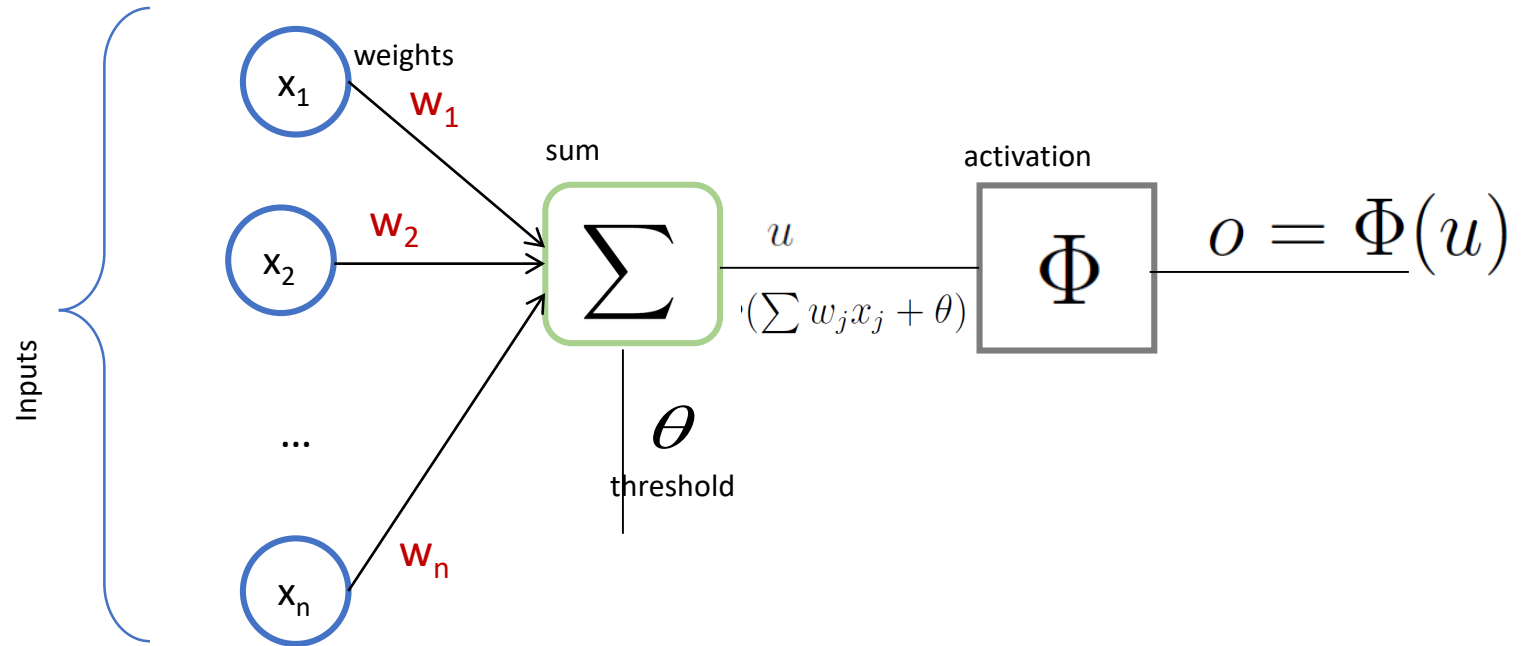
vs.

Biological Neuron



Perceptron (by Rosenblatt)

Architecture

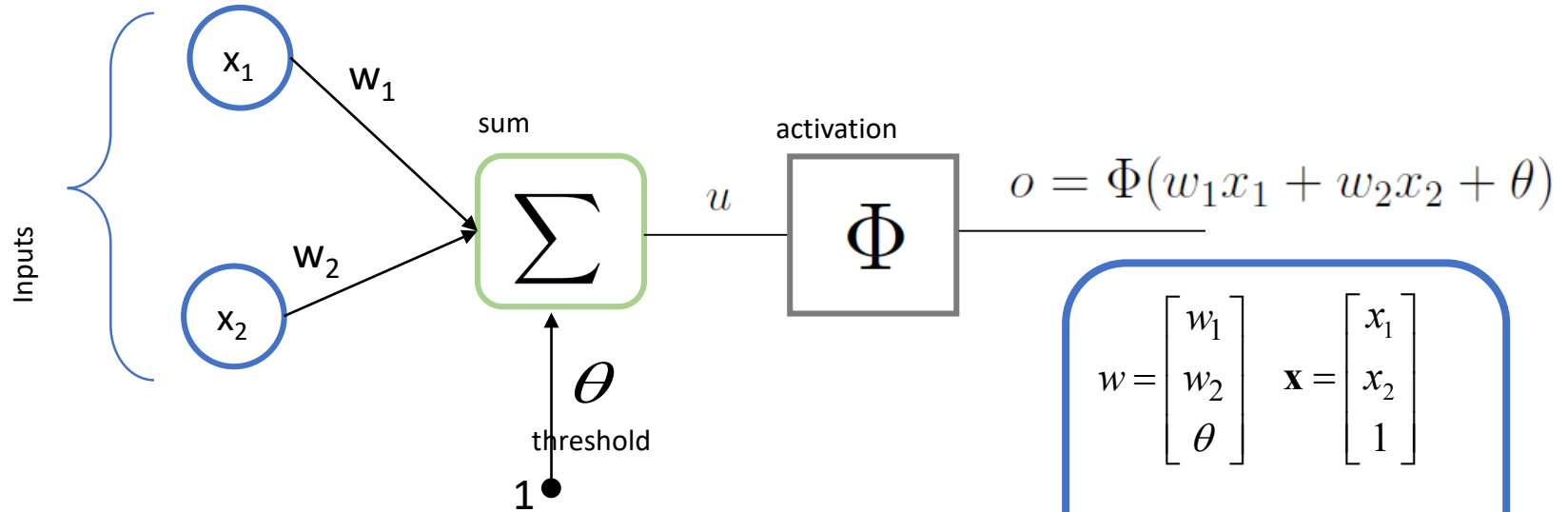


$$o = \Phi(u) = \Phi\left(\sum_j w_j x_j + \theta\right)$$

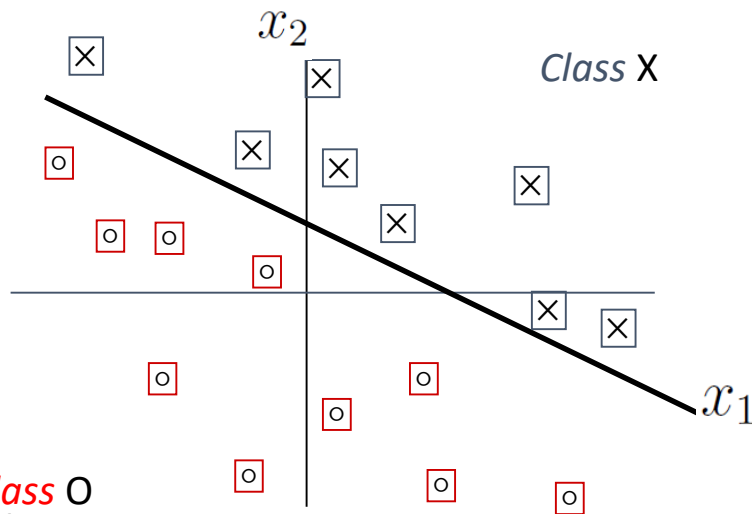
$$\Phi(u) = \text{sgn}(u) = \begin{cases} 1, & u \geq 0 \\ -1, & u < 0 \end{cases}$$

Perceptron (by Rosenblatt)

Architecture



$$\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \theta \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$
$$o = \Phi(\mathbf{w}^T \mathbf{x})$$

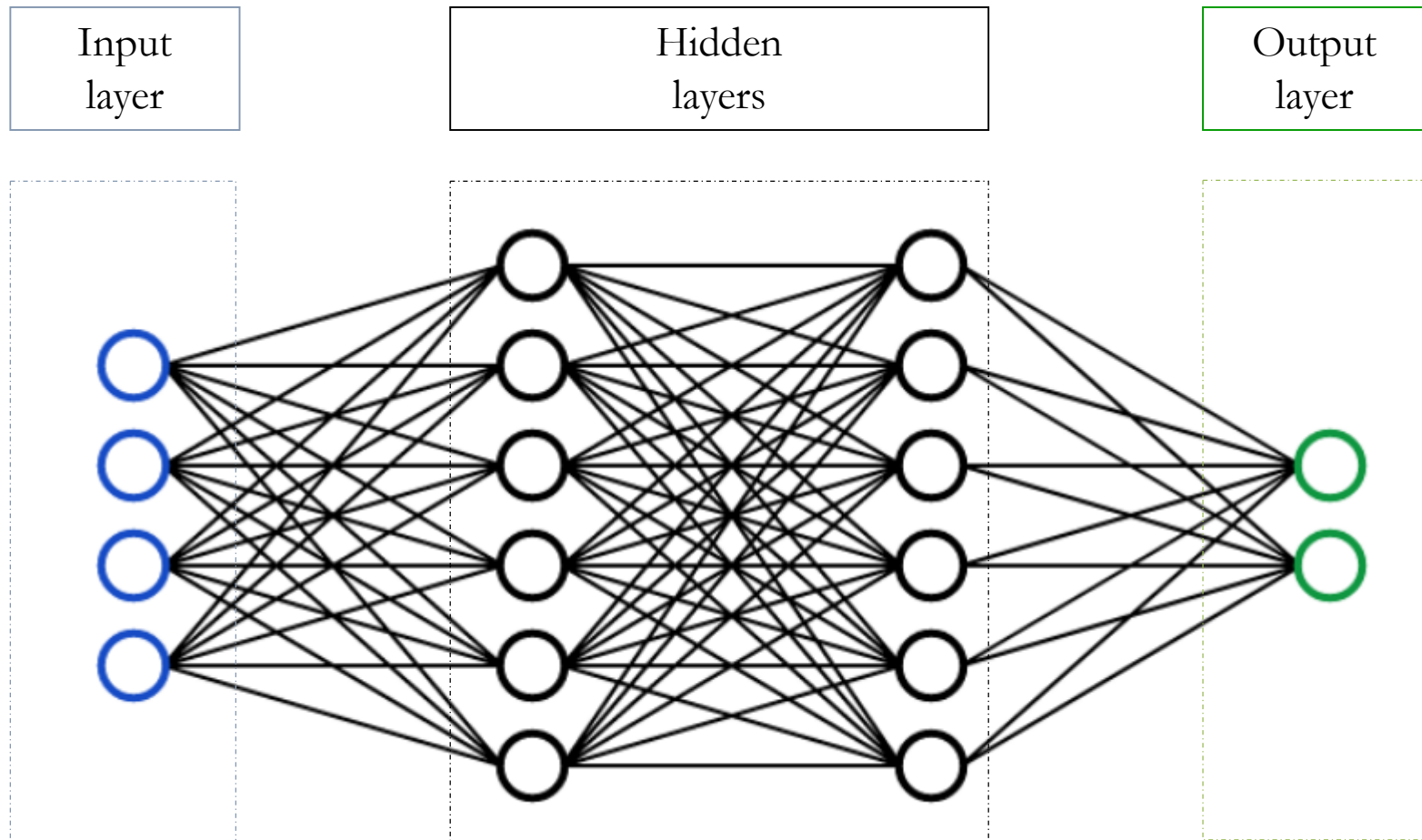


Decision boundary

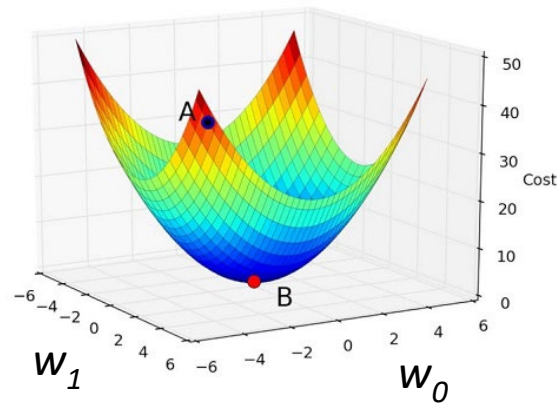
$$(\mathbf{w}^T \mathbf{x}) = 0$$
$$w_1x_1 + w_2x_2 + \theta = 0$$

MultiLayer Perceptron(MLP)

Architecture



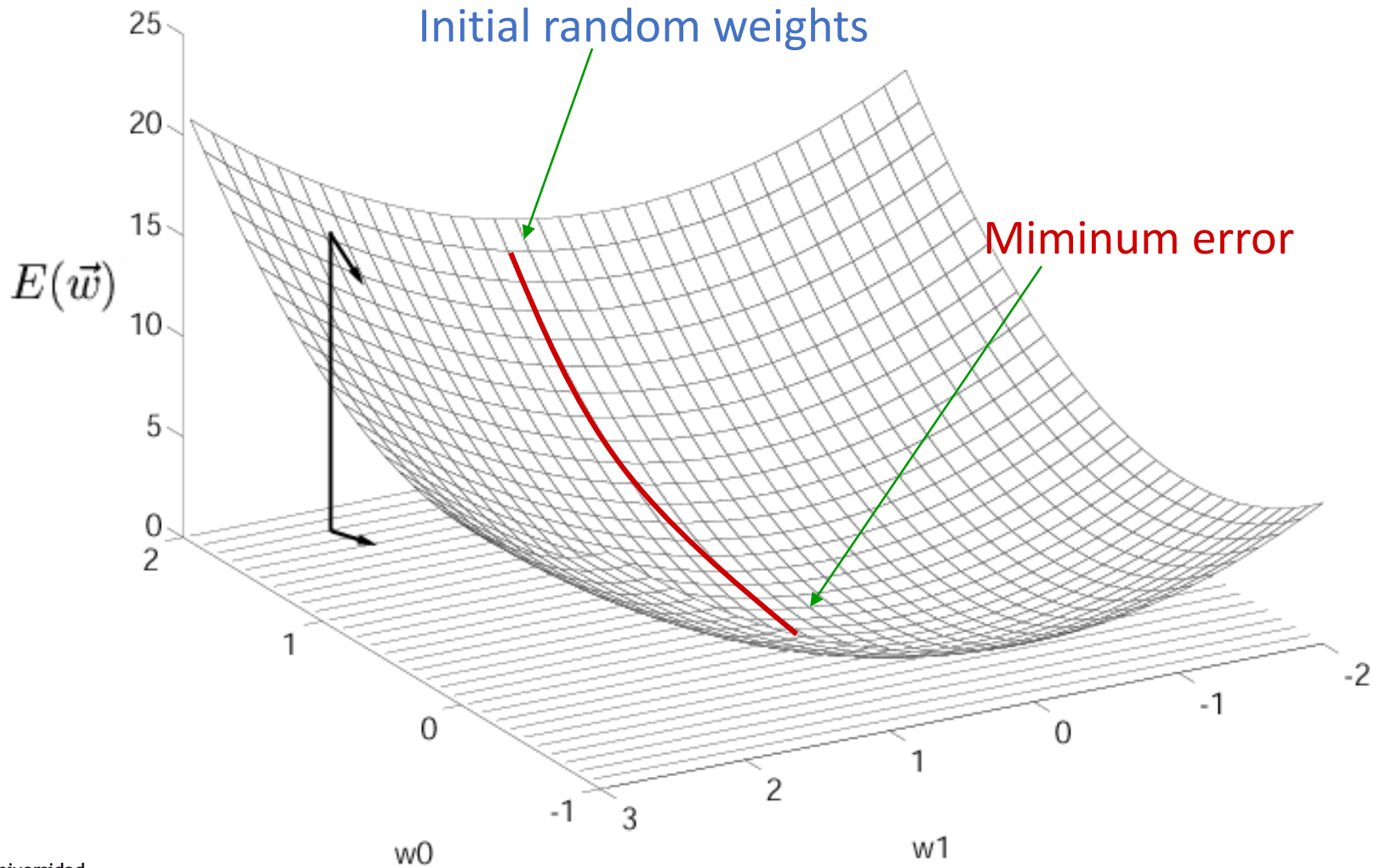
Cost function



- The key idea is to use gradient descent to search the hypothesis space of possible weights to find the network weights that best fit the training examples.

... Training Neural Network

Cost function



Building Machine Learning Models

Pipeline





I-Problem Definition and Data Collection

Clearly define the problem statement and objective of the machine learning model.

Identify the type of machine learning task (classification, regression, clustering, etc.).

Collect and preprocess the relevant data required for training and evaluation.





II-Data Exploration and Preparation

- Perform exploratory data analysis (EDA) to gain insights into the dataset.
- Handle missing values, outliers, and inconsistencies in the data.
- Feature selection, feature extraction.
- Split the data into training, validation, and test sets.



III-Model Selection

Understand different types of machine learning algorithms (decision trees, neural networks, etc.).

Consider factors such as model complexity, interpretability, and scalability.

Choose an appropriate model based on the problem requirements and dataset characteristics.



IV-Model training

- Prepare the data for model training (normalization, encoding categorical variables, etc.).
- Optimize the model parameters using a suitable optimization algorithm (e.g., gradient descent).
- Evaluate the model's performance on the validation set to monitor its progress.



V-Model Fine-Tuning

Adjust hyperparameters (e.g., learning rate, regularization) to improve model performance.



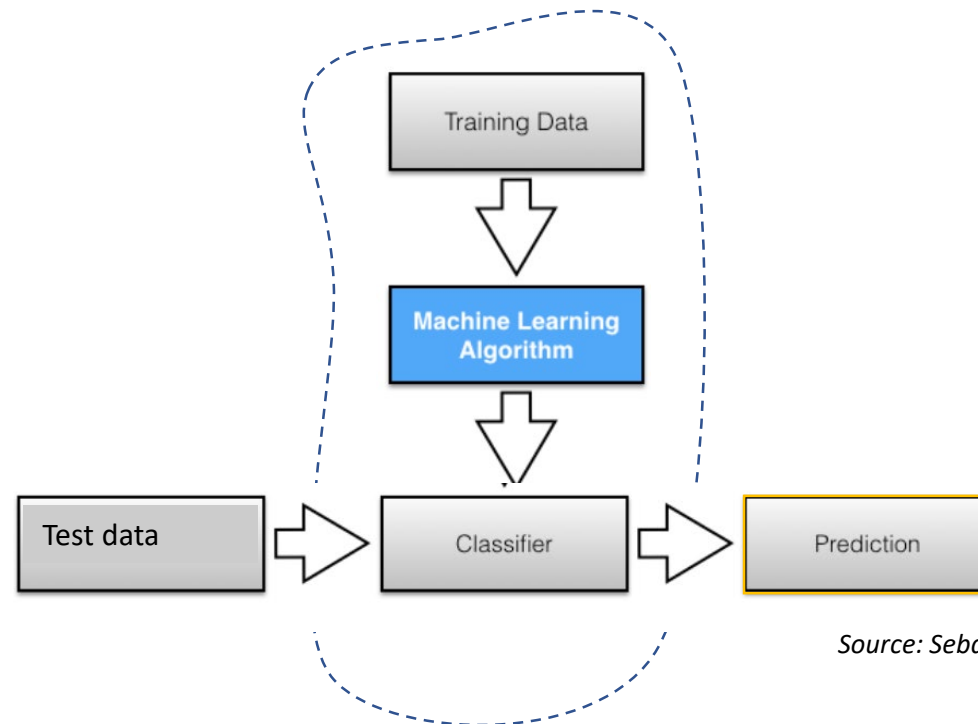
Use techniques like cross-validation or grid search to find the optimal hyperparameter values.

Measure the model's performance using appropriate evaluation metrics (accuracy, precision, recall, etc.).

Analyze the model's strengths, weaknesses, and potential sources of error.



VI-Model testing



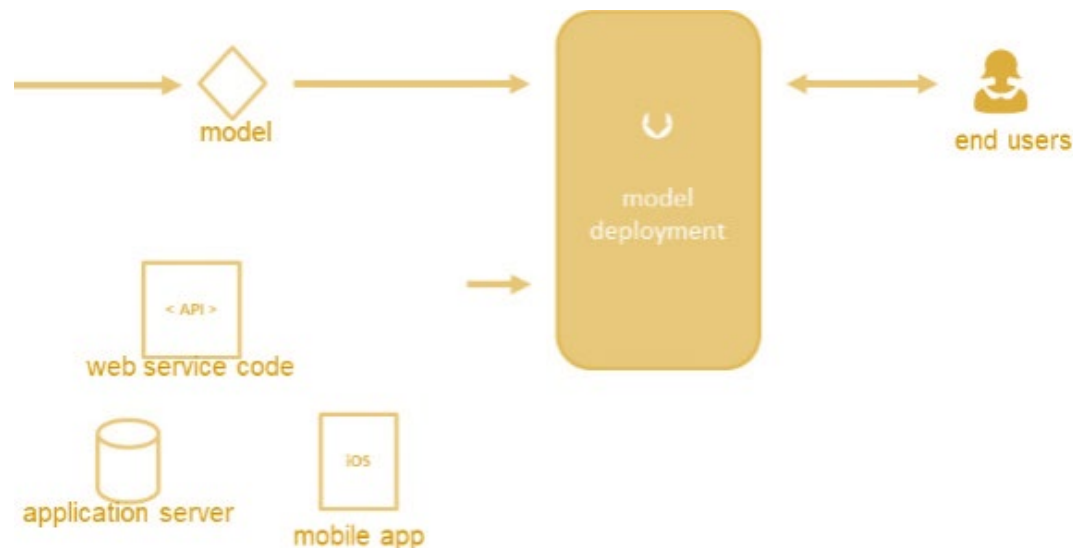
Source: Sebastian Raschka

- Assess the final model's performance on the test set.



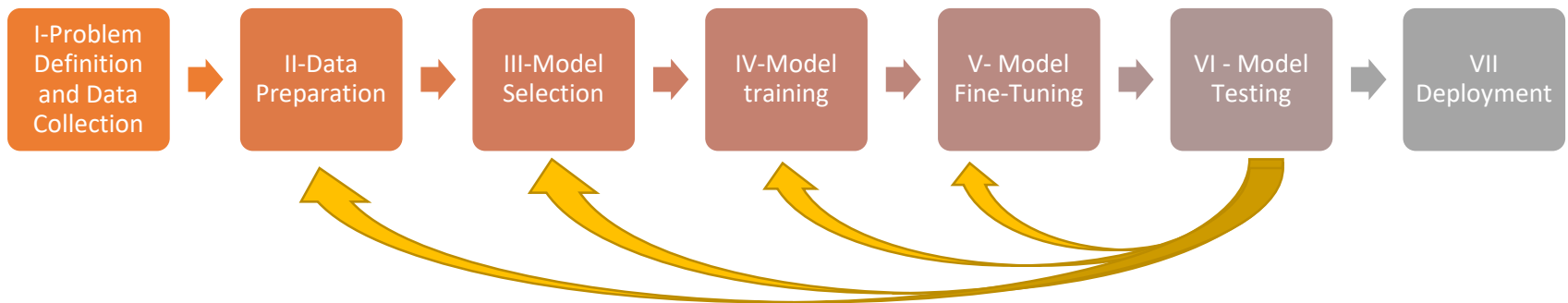
VII-Model deployment

- Implement the model in a production environment.
- Monitor the model's performance in production and iterate on improvements as needed.



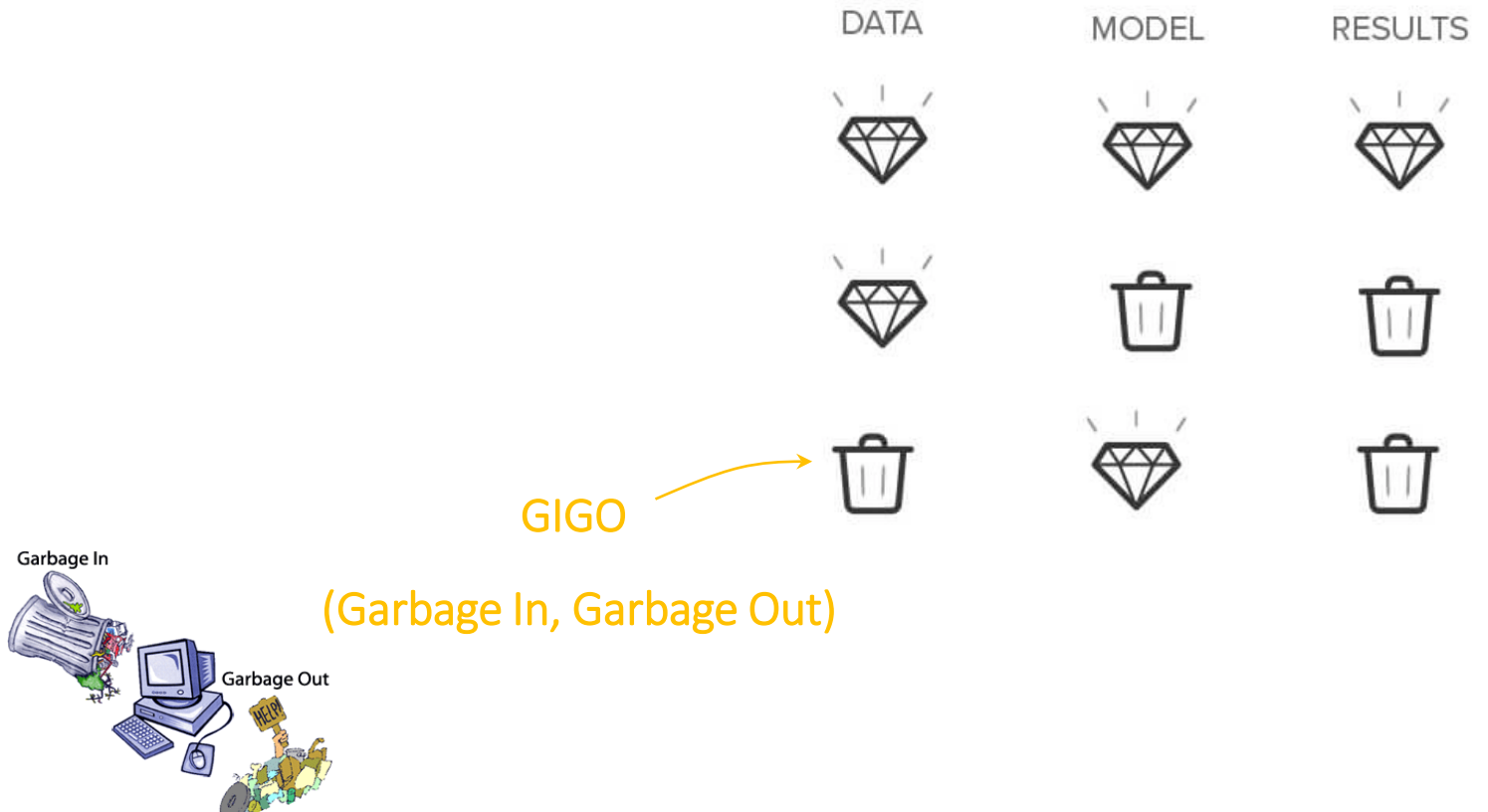
Machine Learning — An iterative process

The development of Machine Learning-based applications is a highly empirical process.



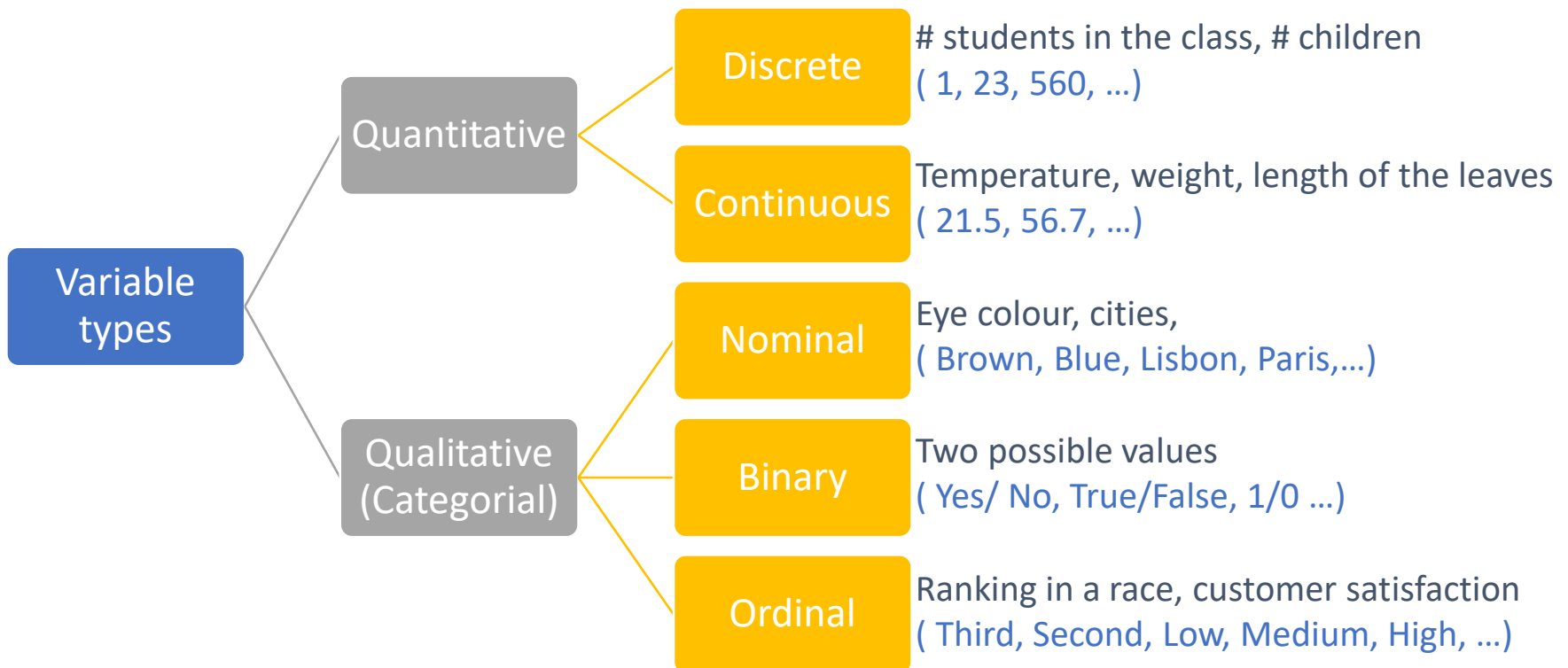
Remarks: Data quality

The quality of the results strongly depends on the quality of the training data.



Remarks: Variable (Feature, Attribute) Types

- Objective function: It is the true function f that we aim to learn.



MODEL EVALUATION

What if I use the training set to assess the quality of the model?

In that case, we would be rewarding models that MEMORIZE TRAINING DATA



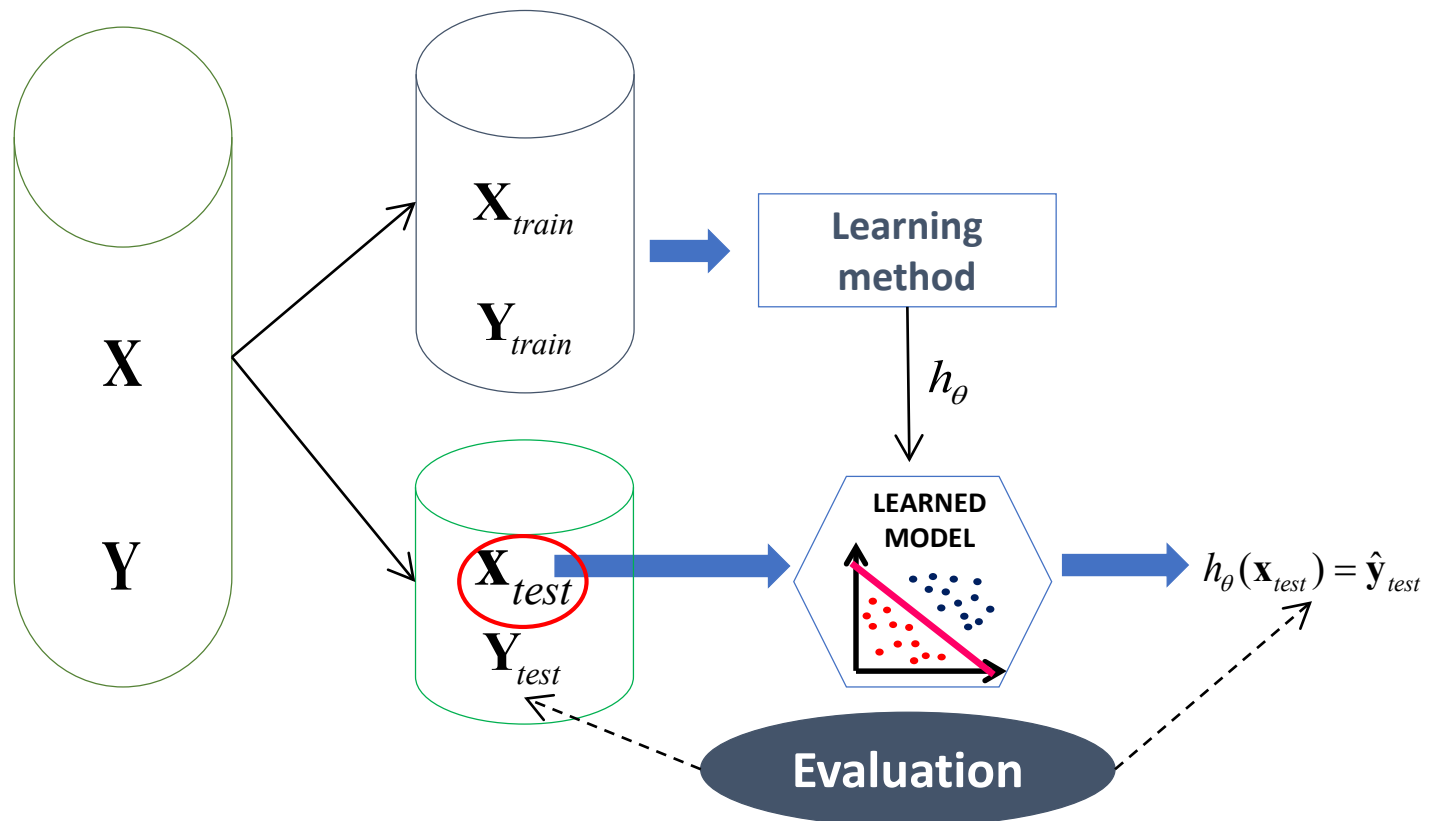
PROBLEM!!!! We want our classifier to generalize well for new (unknown) cases.

SOLUTION Evaluate the models on a dataset different from the training dataset.

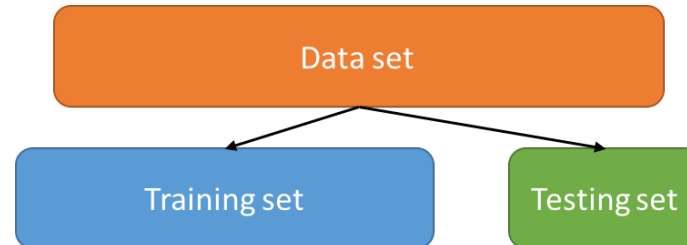
How can we get an **unbiased estimate** of the **accuracy** of a **learned model**?

Train-Test split

Splitting the dataset into train and test set



Train-Test split



The available data set D is divided into two disjoint subsets,

- the **training set** X_{train} (for learning a model) (70%)
- the **test set** X_{test} (for testing the model) (30%)

This method is mainly used when the data set X is large.

IMPORTANT!!!!

The training set should not be used in testing and the test set should not be used in learning.


Unseen test set provides a unbiased estimate of accuracy.

Confusion matrix


It is used to **describe the performance of a classification model** (or "classifier") on a set of test data for which the true values are known.

		Real Class	
		1	0
Predicted Class	1	TP	FP
	0	FN	TN


¿ $\hat{d} = 0$ o $\hat{d} = 1$?




$d = 1$ $\hat{d} = 1$




$d = 0$ $\hat{d} = 0$




$d = 0$ $\hat{d} = 1$




$d = 1$ $\hat{d} = 0$




$d = 0$ $\hat{d} = 0$




$d = 1$ $\hat{d} = 1$




$d = 1$ $\hat{d} = 0$



$d = 0$ $\hat{d} = 0$



$d = 1$ $\hat{d} = 1$



$d = 0$ $\hat{d} = 0$

Confusion matrix

		Real Class	
		1	0
Predicted Class	1	TP	FP
	0	FN	TN

FN, FP, TN, TP?

When someone performs a test to check if s/he is an authorized user:

False negative: When the model says you are not the user but you actually are.

False positive: When the model says you are the user but you are not.

True negative: When the model says you are not the user and you are not.

True positive: When the model says you are the user and you are.

Confusion matrix

		Real Class	
		1	0
Predicted Class	1	TP	FP
	0	FN	TN

CONFUSION MATRIX

		Real Class	
		1	0
Predicted Class	1	3	1
	0	2	4

Confusion matrix

Accuracy is the number of correctly classified examples

		Real Class	
		1	0
Predicted Class	1	TP	FP
	0	FN	TN

CONFUSION MATRIX

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Confusion matrix

Precision is the number of correctly classified positive examples divided by the total number of examples that are classified as positive.

$$\textit{Precision} = \frac{TP}{TP + FP}$$

Recall is the number of correctly classified positive examples divided by the total number of actual positive examples in the test set.

$$\textit{Recall} = \frac{TP}{TP + FN}$$

$$\textit{FScore} = \frac{2 \cdot \textit{Precision} \cdot \textit{Recall}}{\textit{Precision} + \textit{Recall}}$$

		Real Class	
		1	0
Predicted Class	1	TP	FP
	0	FN	TN

CONFUSION MATRIX

Confusion matrix

Accuracy is the number of correctly classified examples

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad 0,70$$

		Real Class	
		1	0
Predicted Class	1	TP	FP
	0	FN	TN

CONFUSION MATRIX

		Real Class	
		1	0
Predicted Class	1	3	1
	0	2	4

Confusion matrix

Precision is the number of correctly classified positive examples divided by the total number of examples that are classified as positive.

$$Precision = \frac{TP}{TP + FP}$$

$$Precision = \frac{TP}{TP + FP} \quad 0,75$$

		Real Class	
		1	0
Predicted Class	1	TP	FP
	0	FN	TN

CONFUSION MATRIX

		Real Class	
		1	0
Predicted Class	1	3	1
	0	2	4

Confusion matrix

Recall is the number of **correctly classified positive examples** divided by the total number of actual positive examples in the test set.

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad 0,60$$

		Real Class	
		1	0
Predicted Class	1	TP	FP
	0	FN	TN

CONFUSION MATRIX

		Real Class	
		1	0
Predicted Class	1	3	1
	0	2	4

Confusion matrix

F-score

$$Precision = \frac{TP}{TP + FP} \quad 0,75$$

$$Recall = \frac{TP}{TP + FN} \quad 0,60$$

$$FScore = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall} \quad 0,66$$

- Machine Learning
 - Concept
 - Application fields
 - Supervised, Unsupervised & Reinforcement learning
- Approaching a problem of learning from examples
- Building Machine Learning Models
- Supervised Learning models
 - K-NN
 - Naïve Bayes
 - Neural Networks
- Evaluating classifier performance

Introduction to Machine Learning

Thank you!

Rocío Alaiz Rodríguez

rocio.alaiz@unileon.es