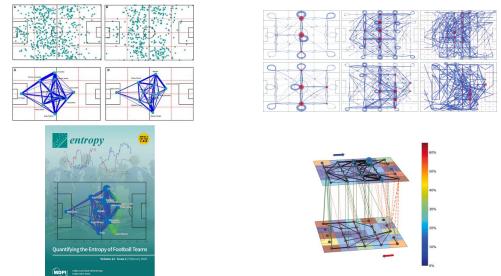


## Machine Learning for Sports Orange3



JOHANN H. MARTÍNEZ



Data, dataset and much more

[GitHub link](#)



## Skin Cancer

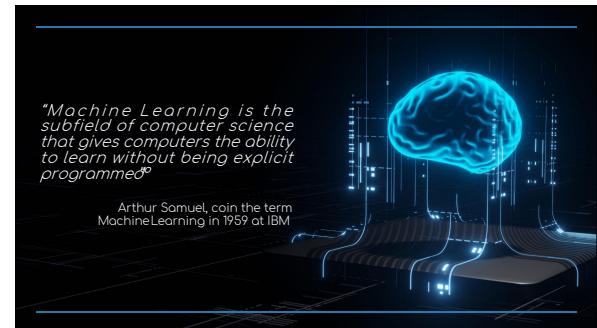
ID	Clump	UnifSize	UnifShape	MargAdh	SingEpilSize	BareNuc	BlandChrom	NormNud	Mit	Class
1000015	6	1	1	1	7	1	3	1	1	benign

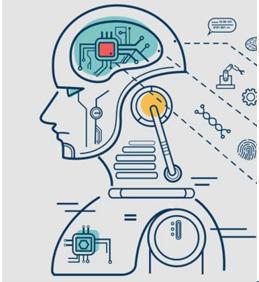
Modeling

Prediction

Accuracy = 89%

## Examples





Differences between

- Artificial Intelligence
- Machine Learning
- Deep Learning

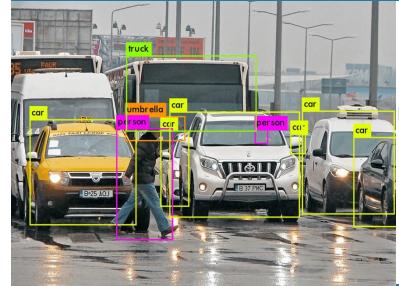


Artificial Intelligence(AI)

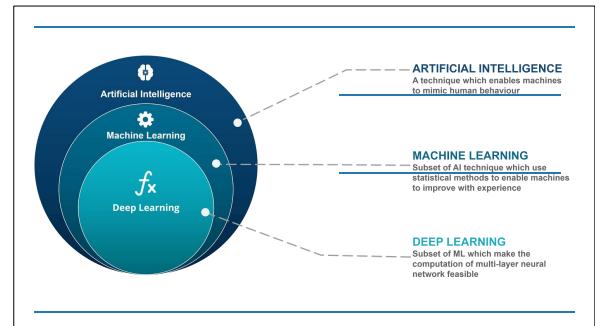
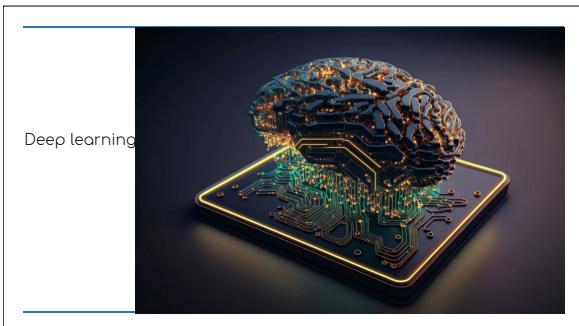
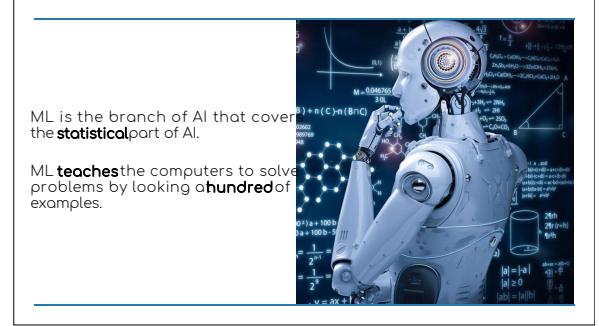
AI tries to make computers intelligent so that they mimic the cognitive functions of humans.

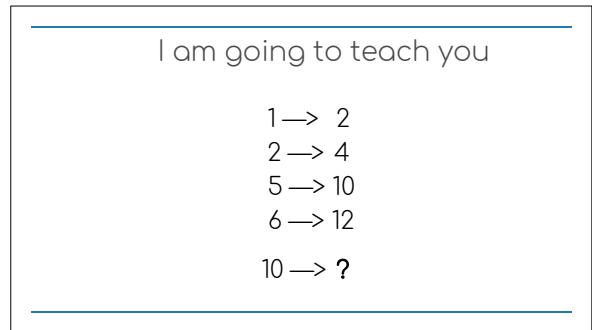
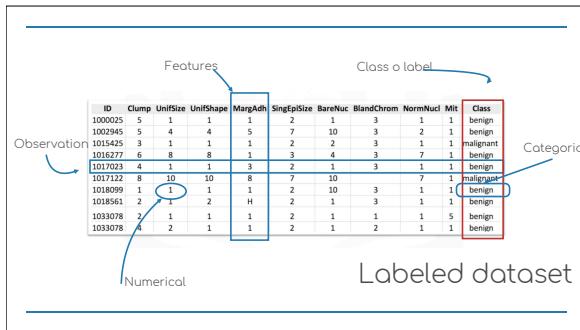
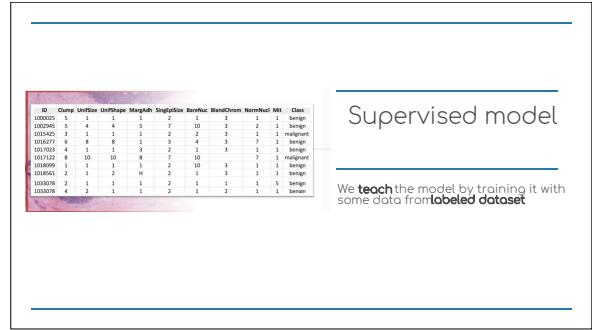
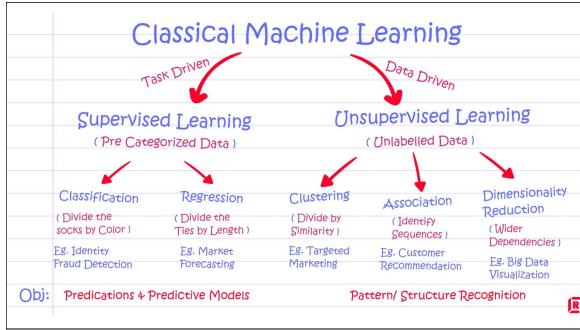


Autopilot in Tesla

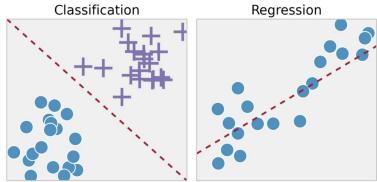


Pattern recognition  
Using computer vision





## Types of supervised techniques

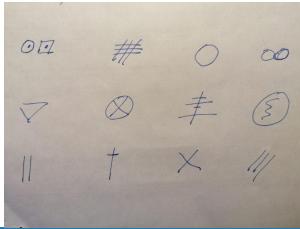


Customer ID	Age	Years Employee	Income	Card Debt	Other Debt	Address	DebtIncomratio
1	41	2	6	19.0	10.7	NAB001	6.3
2	47	1	26	100.0	4.582	8.218 NAB002	12.8
3	33	2	15	12.0	1.111	1.233 NAB003	8.3
4	29	2	4	39	0.681	5.160 NAB004	10.0
5	47	1	31	253.0	9.888	8.908 NAB005	26.0
6	40	1	23	81.0	9.998	8.908 NAB006	10.0
7	38	2	10	15.0	1.544	10.544 NAB007	1.644
8	42	3	0	64	0.279	3.945 NAB008	6.6
9	26	1	18	3.575	2.215	NAB009	1.6
10	35	2	25	15.0	1.544	10.544 NAB010	1.644
11	44	3	8	88.0	0.285	5.083 NAB010	6.1
12	34	2	9	40	0.374	2.266 NAB010	1.6

## Unsupervised model

The model **works on its own** to discover information maybe not visible to the human eye.

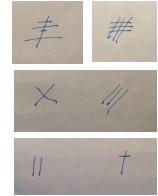
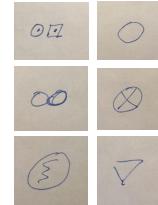
We also learn from unsupervised dataset



## Polygonal language?

## Linear language?

## Two clusters



Filled polygonal language?

Diagonal Linear language?

Four clusters?

Empty polygonal language?

Straight Linear language?

ESCUOLA DEL REAL MADRID

### First program

- File → Iris
- Data Table
- Scatter plot
- Selected data table from scatter plot!

ESCUOLA DEL REAL MADRID

### Own data & how to transform it to be read

name	gender	height	eye color	hair color
Jessica	female	5' 2"	green	grey
John	male	5' 7"	blue	black
Paul	male	5' 11"	brown	brown
Vivian	female	5' 7"	black	brown
Ava	female	5' 10"	green	brown
Sam	male	5' 10"	green	brown
Cordelia	female	5' 10" feet	grey	grey

- Construir una base de datos en Google sheet.
- Importar los datos
- Alejarse los datos
  - s: string
  - d: discrete
  - c: continuous
  - m: meta: meta data
  - Class: class
- Save data

ESCUOLA DEL REAL MADRID

### First visualizations

- Load innerNBA Players dataset
- Scatter plot
- Box plot

**ESCUOLA DEL REAL MADRID**

## First visualizations

- Download NBA Players dataset
- From xls toxlsx to be well-read
- widget Transform to select Target
- Visualize it

**Bonus** Play with Titanic's survivors

**ESCUOLA DEL REAL MADRID**

## First Predictions

- Download NBA Players dataset
- Tree → Tree viewer
- Test file
- Prediction
- Logistic Regression

**ESCUOLA DEL REAL MADRID**

ID	Gmp	Utgms	Unlrgms	Mosgms	Sndgms	Bkngms	Blckdgm	NrmNbdm	Mtr	Class
1002011	3	2	1	4	5	7	2	1	2	benign
1002941	5	4	1	4	5	7	10	3	2	benign
1003241	6	8	8	1	3	2	4	3	7	benign
1014277	6	8	8	1	3	4	4	3	7	benign
1017024	4	1	1	3	2	1	3	1	1	benign
1017334	23	20	19	7	20	17	1	1	1	benign
1018091	1	1	1	1	2	10	3	1	1	benign
1018099	1	1	2	H	2	1	3	1	1	benign
1018374	2	2	1	1	2	2	1	3	1	benign
1018378	4	2	1	1	2	1	2	1	1	benign

## Classification

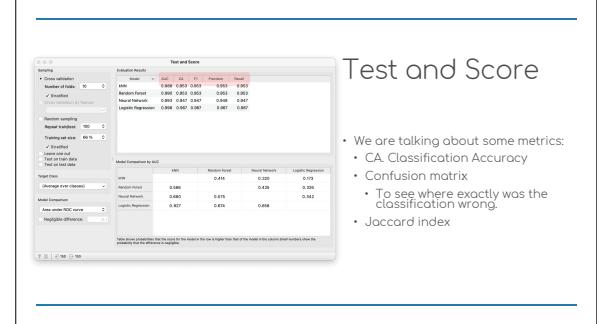
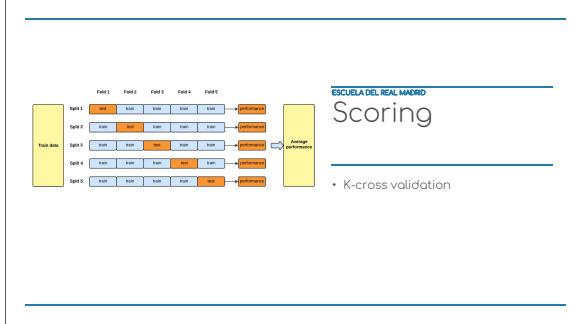
The goal is to predict *class label*

- Binary Classification
- Multiclass Classification

**ESCUOLA DEL REAL MADRID**

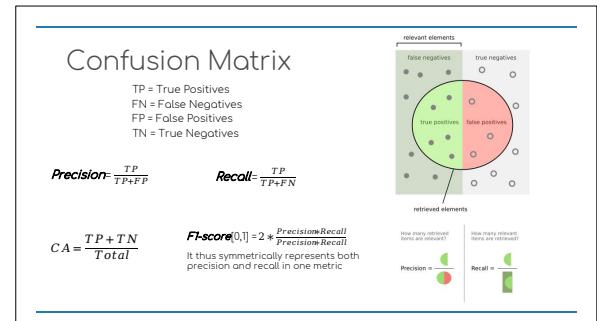
## Scoring

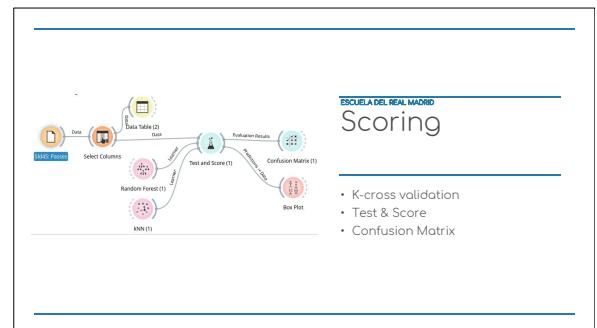
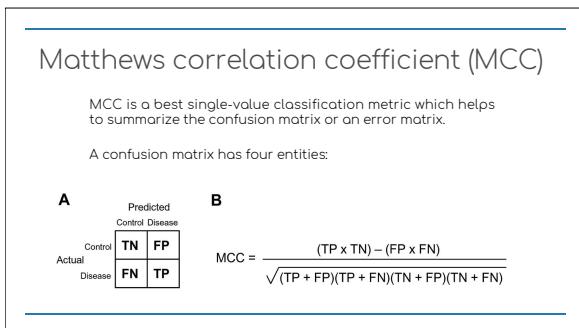
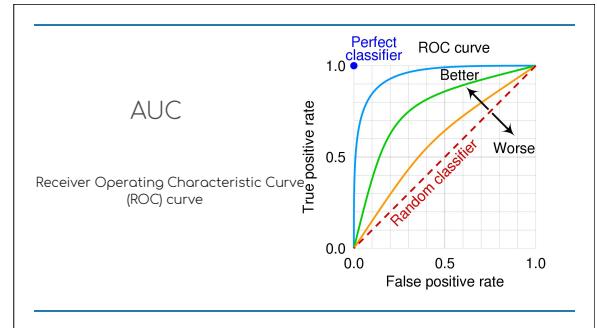
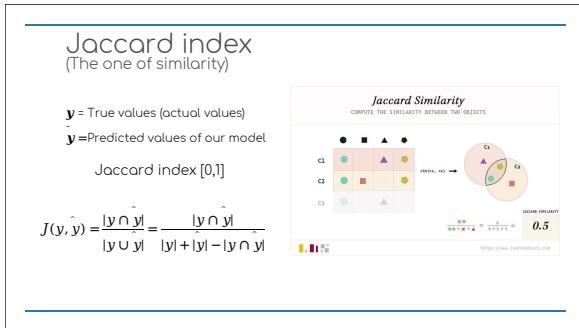
- K-cross validation
- Test & Score
- Confusion Matrix

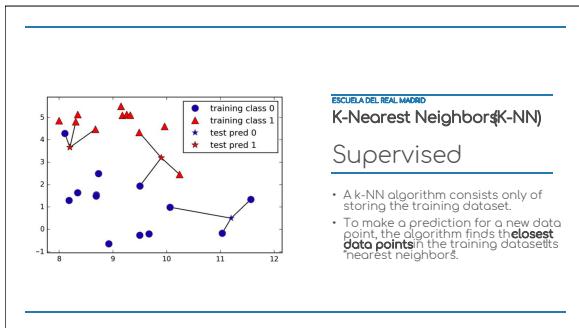


Matriz de confusión		Estimado por el modelo			
		Negativo (N)	Positivo (P)		
Real	Negativo	a: (TN)	b: (FP)	Precisión ("precision") Porcentaje de predicciones positivas correctas	d/(b+d)
	Positivo	c: (FN)	d: (TP)		
		Sensibilidad, exhaustividad o "Recall" Porcentaje casos positivos detectados	Especificidad ("specificity") Porcentaje casos negativos detectados	Exactitud ("accuracy") Porcentaje de predicciones correctas (no se incluyen por equivocados)	(a+d)/(a+b+c+d)
		d/(d+c)	a/(a+b)		

Main diagonal: Correct estimated values (TP + TN)  
Secondary diagonals:  
• Upper means False Positives (FP)  
• Bottom means False Negatives (FN)







**K-nearest Neighbors**

- Iris and Poses
- KNN → Test Score → Confusion Matrix
- Different values of K and distances.
- Scatter Plot (KNN) + Box plot

**K-nearest Neighbors**

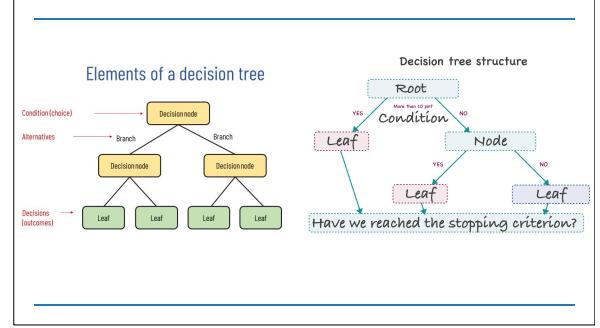
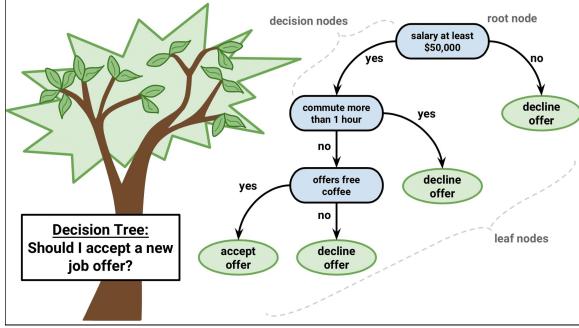
- Iris and Poses
- KNN → Test Score → Confusion Matrix
- Different values of K and distances.
- Scatter Plot (KNN) + Box plot

**K-nearest Neighbors in Python**

- Iris
- KNN → Test Score → Confusion Matrix

**K-nearest Neighbors in R**

- Iris
- KNN → Test Score → Confusion Matrix



Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	High	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A
p15	Middle-age	F	Low	Normal	?

**ESCUADA DEL REAL MADRID**

## Decision Tree

- What exactly is a decision tree?
- We're doctors and we want to predict the best drug for a kind of patients

1. Choose an attribute (Age):

2. Split it into its categories (Y, M-A, S)

- It seems that M-A arrive to the solution to drug B
- ... but Y, and S?
- Split the remaining categories ageing based on the other attributes:

For Y, now split it into Sex

- Oh, look! It seems we arrive to solution A (for Y-female), and B (for Y-male)

For S, split it into Cholesterol

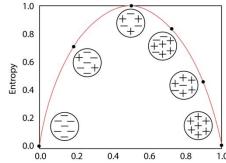
- Oh, look! It seems we arrive to solution A (for S-Chol-H), and B (for S-Chol-N)

**How to build a decision tree?**

```

graph TD
    Age[Age] -- Young --> SexF[Sex F]
    Age -- Middle-age --> SexM[Sex M]
    Age -- Senior --> Cholesterol[Cholesterol]
    SexF --> DrugA1([Drug A])
    SexM --> DrugB1([Drug B])
    Cholesterol -- High --> DrugA2([Drug A])
    Cholesterol -- Normal --> DrugB2([Drug B])
    
```

## How to compute the significance?



In every step, we want to dismiss the **impurity** of the nodes.

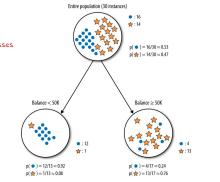
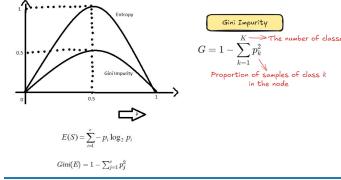
**Impurity** is computed by **Entropy** of data in the node

$$E = - \sum_i p_i \log_2(p_i)$$

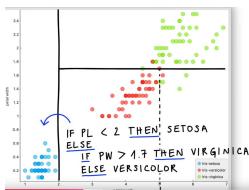
**Entropy** is the amount of information disorder. If the samples are completely homogeneous the entropy is zero.

## How to compute the significance?

**Gini impurity**: a measure of how often a randomly chosen element from the set would be incorrectly labeled if it were randomly labeled according to the distribution of labels in the branch.



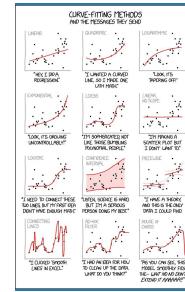
## Which attributes are the best?



## Decision Tree

- Iris dataset
- Test and Score → confusion Matrix
- Tree → Tree Viewer

## Linear Models for Classification



## Linear Models for Classification

### Linear Models for classification



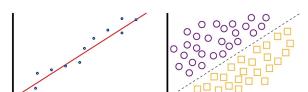
- They are a class of models that is **widely used in practice.**



- They make a prediction **using a linear function** of the input features

### Linear Models for classification

Regression	Classification
The task of predicting a continuous quantity is known as regression.	The classification process involves anticipating a discrete class label.
The task is to map the input value(s) to the continuous output variable ( $y$ ).	The task is to map the input value(s) to the discrete output variable ( $y$ ).
A regression model can predict a discrete value, but only in the form of an integer quantity.	A classification model can predict a continuous value but it is in the form of a probability for a class label.
The output variable in regression must be continuous.	The output variable in Classification must be discrete.



## Linear Binary classification

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

Dependent Variable  
 Population Y intercept  
 Population Slope Coefficient  
 Independent Variable  
 Random Error term  
 Linear component  
 Random Error component

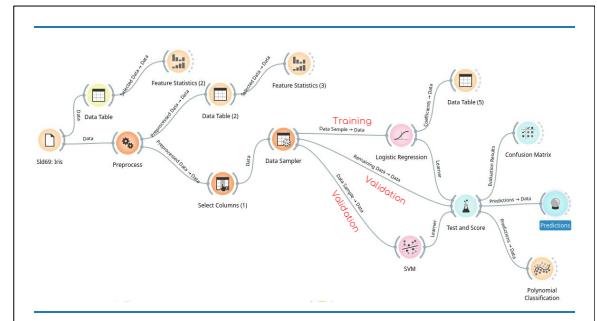
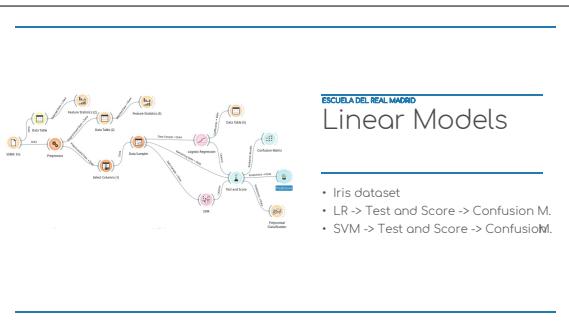
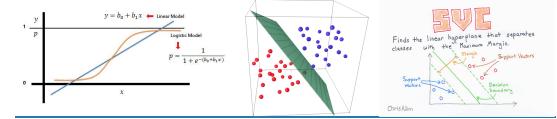
$X_i [0..j]$  Features (columns)  
 $\beta [0...i]$  Parameters of the model  
 $Y_i$  Prediction of the model

## Linear Models for classification

A linear classifier separates classes using a

- line
- plane
- hyperplane

**High-dimensional** problems: Where the number of features  $p$  is much larger than the number of observations  $N$  often written  $p \gg N$



$P(I'M\ NEAR\ |\ I\PICKED\ UP) = P(I\PICKED\ UP\ |\ I'M\ NEAR)\ P(I'M\ NEAR)$

$P(I\PICKED\ UP\ |\ A\ SEASHELL) \cdot P(I'M\ NEAR\ |\ THE\ OCEAN)$

STATISTICALLY SPEAKING, IF YOU PICK UP A SEASHELL AND DON'T HOLD IT TO YOUR EAR, YOU CAN PROBABLY HEAR THE OCEAN.

### Naive Bayes Classifiers

$P(A) = \frac{\text{# } A}{\text{# } U}$ ,  $P(B|A) = \frac{\text{# } A \cap B}{\text{# } A}$

$P(B) = \frac{\text{# } B}{\text{# } U}$ ,  $P(A|B) = \frac{\text{# } A \cap B}{\text{# } B}$

$P(A) \cdot P(B|A) = \frac{\text{# } A}{\text{# } U} \times \frac{\text{# } A \cap B}{\text{# } A} = \frac{\text{# } A \cap B}{\text{# } U}$

$P(B) \cdot P(A|B) = \frac{\text{# } B}{\text{# } U} \times \frac{\text{# } A \cap B}{\text{# } B} = \frac{\text{# } A \cap B}{\text{# } U}$

$= P(A) \cdot P(B|A)$ , i.e.

$$P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)}$$

$$P(B|A) = \frac{P(B) \cdot P(A|B)}{P(A)}$$

**Bayes' theorem.**  
Knowing the probability **having a headache** given that one **has the flu** one could know if one **has any additional data**: the probability **of having the flu** if one **has a headache**

**Naive** means independency between features.

**Bayes** Based on Bayes' theorem.

**Naive Bayes classifier**  
Assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature

e.g., A fruit can be considered an apple if it is

- red
- round
- about 7 cm in diameter.

Even if these features depend on each other or upon the existence of the other features

**BAYES THEOREM**  
Posterior probability  
 $p(A|B) = \frac{p(B|A)p(A)}{p(B)}$

**BAE'S THEOREM**  
Lamontoon law probability  
 $p(A|B) = \frac{p(A)p(B|A)}{p(A)p(B) + p(\neg A)p(\neg B|A)}$

(1). Given a dataset

	Outlook	Play
0	Rainy	No
1	Sunny	No
2	Overcast	No
3	Rainy	Yes
4	Sunny	No
5	Rainy	Yes
6	Sunny	Yes
7	Overcast	Yes
8	Rainy	No
9	Sunny	No
10	Sunny	Yes
11	Rainy	No
12	Overcast	No
13	Overcast	Yes

(2) Convert it into frequency tables.

Weather	No	Yes
Overcast	0	5
Rainy	2	2
Sunny	2	3
All	4	10

(3). Then, into probability one.

Weather	No	Yes
Overcast	0	5
Rainy	2	2
Sunny	2	3
All	4/14=0.29	10/14=0.71

(3) Apply the algorithm

Weather	No	Yes
Overcast	0	5
Rainy	2	2
Sunny	2	3
All	4/14=0.29	10/14=0.71

$P(\text{Sunny}|\text{No}) = 2/4 = 0.5$

$P(\text{Sunny}|\text{Yes}) = P(\text{Sunny}|\text{Yes}) * P(\text{Yes}) / P(\text{Sunny})$   
 $P(\text{Sunny}|\text{Yes}) = 2/4 = 0.5$   
 $P(\text{Sunny}) = 0.35$   
 $P(\text{No}) = 0.29$   
 $P(\text{No}|\text{Sunny}) = 0.5 * 0.29 / 0.35 = 0.41$

$P(\text{Yes}|\text{Sunny}) = P(\text{Sunny}|\text{Yes}) * P(\text{Yes}) / P(\text{Sunny})$   
 $P(\text{Sunny}|\text{Yes}) = 3/10 = 0.3$   
 $P(\text{Sunny}) = 0.35$   
 $P(\text{Yes}) = 0.71$   
 $P(\text{Yes}|\text{Sunny}) = 0.3 * 0.71 / 0.35 = 0.6$

**Slang**

**GAUSSIAN NAIVE BAYES CLASSIFIER**

"Gaussian" because this is a normal distribution

$P(\text{class}|\text{data}) = \frac{P(\text{data}|\text{class}) * P(\text{class})}{P(\text{data})}$

This is our prior belief

We don't believe this is more logical than other classifiers

**Jargon**

The Posterior The Evidence The Prior

The probability of getting this evidence if this hypothesis were true

The probability of H being true, before gathering evidence

$$P(H|E) = \frac{P(H|E)P(H)}{P(E)}$$

The probability that the hypothesis (H) is true given the evidence (E)

The marginal probability of the evidence (Prob of E over all possibilities)

**Naive Bayes models**

- Naive Bayes classifiers are **family** that are quite similar to the linear models.
- They tend to be **faster** in training.
- They often provide **generalizations** slightly worse than that of linear models.

**ESCUOLA DEL REAL MADRID**

## Naive Bayes models

- There are 3 kinds of Naive Bayes classifiers:
  - GaussianNB
  - Continuous data
  - BernoulliNB
  - Binary data
  - MultinomialNB
  - Count data
- Their greatest baseline models often used **very large datasets** where training a line can take too long.

## Decision Trees

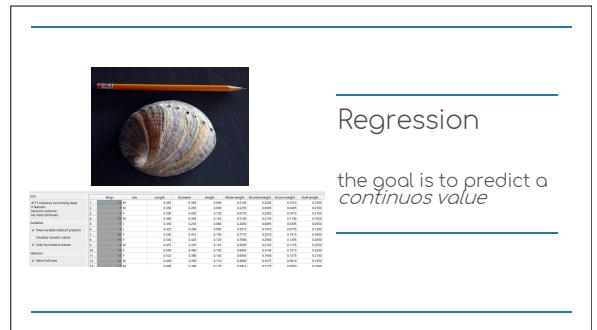
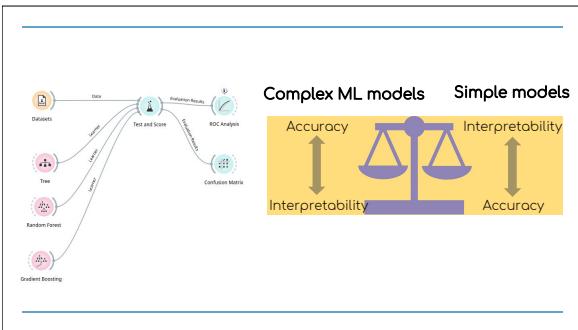
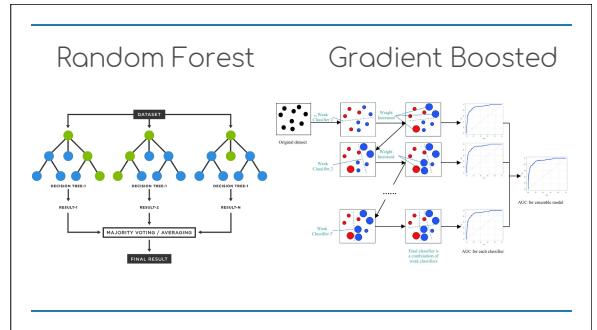
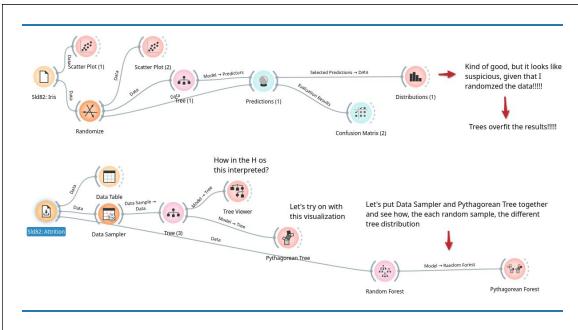
**ESCUOLA DEL REAL MADRID**

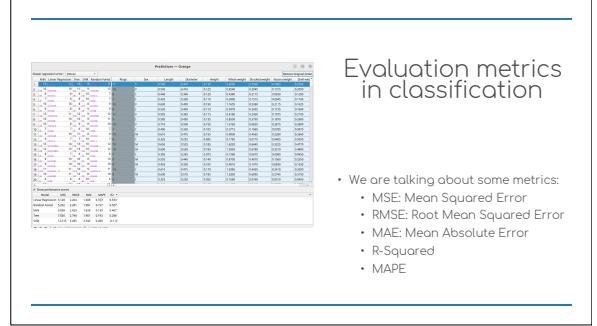
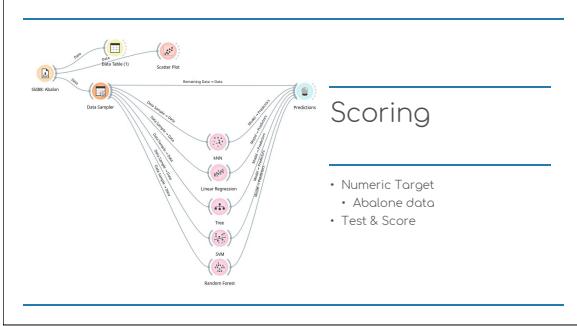
## Ensembles of Decision Trees

- Ensembles** are combinations of multiple ML models to create powerful ones.
- There are two ensemble models that have proven to be effective, both of which use decision trees:
  - Random forests(RF)
  - Gradient boosted decision trees.

## Random Forest

- The main drawback of decision trees is that they **tend to overfit** the training data. RF are one way to address this problem.
- A RF is essentially a **collection of decision trees** where each tree is slightly different from the others. The idea behind random forests is that each tree might do a relatively good job of predicting, but will likely overfit on part of the data. If we build many trees, all of which work well and overfit in different ways, we can reduce the amount of overfitting by averaging their results.





## Error coefficients

## Mean Squared Error

Average squared difference between the estimated values and the actual value

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

## Root Mean Squared Error

-Quadratic mean of the differences between the observed values and predicted ones

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

## Error coefficients

### Mean Absolute Percentage Error

Average magnitude of the absolute errors between the predicted and actual values.

$$MAPE = 100 \frac{1}{n} \sum_{i=1}^n \left| \frac{A_i - F_i}{A_i} \right|$$

*A<sub>i</sub>: Actual value  
F<sub>i</sub>: Forecast value*

## Coefficient of determination

Proportion of the variation in the dependent variable that is predictable from the independent variable.

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2}$$

**Scoring**

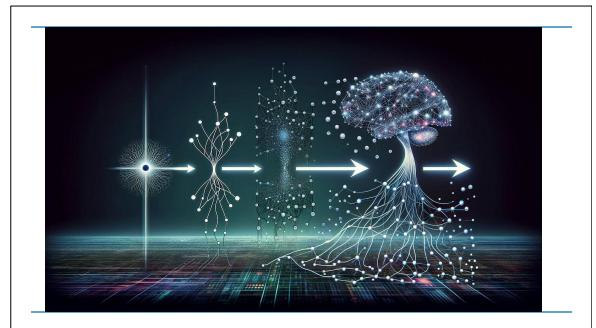
- Numeric Target
  - Number of goles (La Liga)
- Test & Score

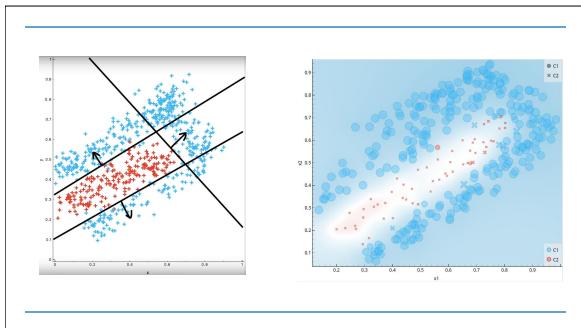
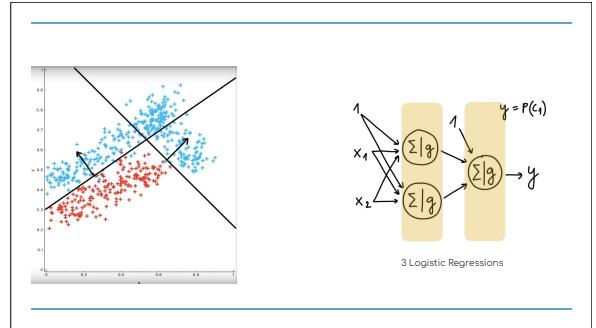
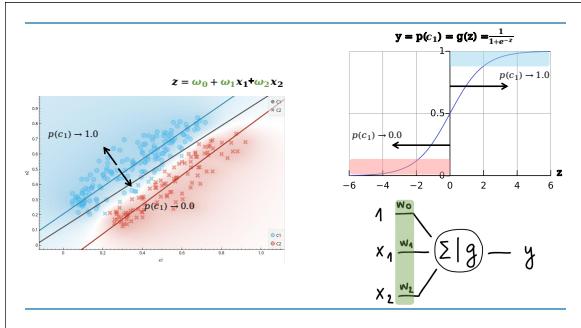
**Scoring**

- Numeric Target
  - Goal Keepers
- Test & Score

**Regression Algorithms**

- Linear Regression
- Polynomial Regression
- Support Vector Regression (SVR)
- Decision Tree Regression
- Random Forest Regression





## Challenges in Unsupervised Learning

Usually applied to unlabeled data

### Challenges

- We don't know what the right output should be.
- It is tough to say whether a model "did well".
- It is tough to say if what model learned "is useful".

## Unsupervised learning algorithms

### Clustering

K-means  
Polynomial  
Hierarchical  
FuzzyC-means

Identify clusters

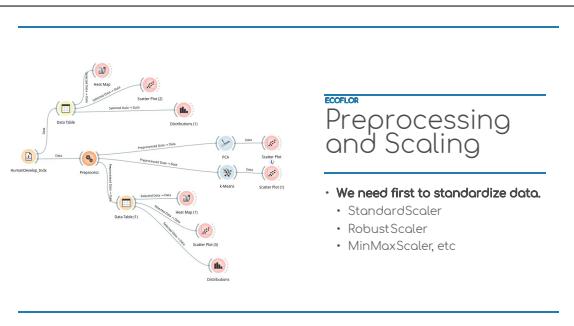
### Dimensionality reduction

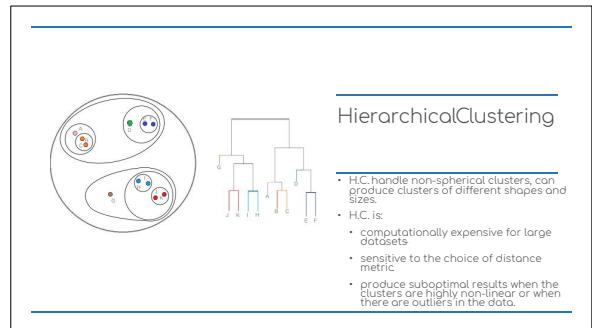
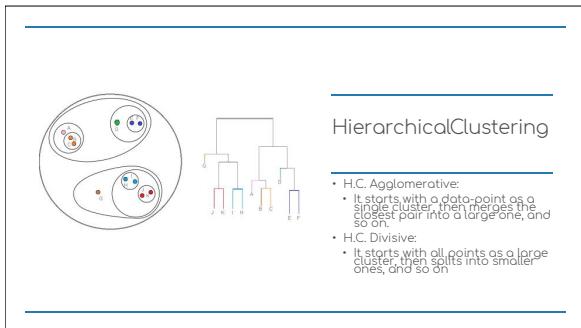
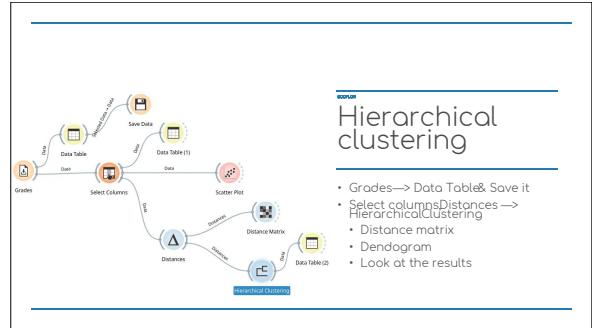
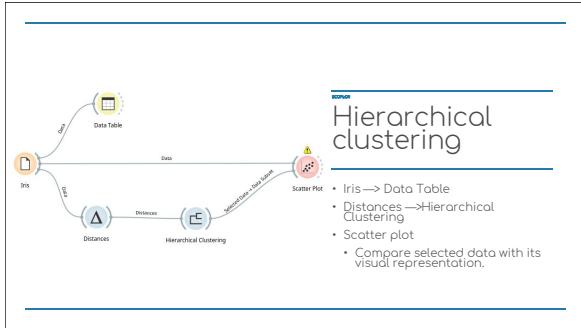
PCA  
Kernel-PA

Transform data to better understanding

Association  
Apriori Algorithm  
Eclat Algorithm  
FP-growth Algorithm

Recommendation list





**K-Means**  
Just distances among points and emergent centroids

- H.G. isn't the best for large datasets
  - Could you imagine a H.C. for genetic data?
  - Matrix distance of millions...
  - Almost 8Tb in memory...

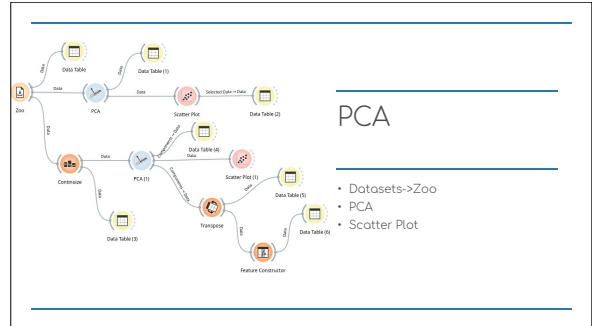
**K-means**

- Download Educational widget
- Point 3 data clusters
- InteractiveK-M compared with Data
- Fats Algorithm
  - Look! Each point is linked to a centroid
  - Then, it recompute centroids...
  - You can add centroids, and re-run.

**K-means**

- Iris
- Grades





**Data Mining  
Fruitful and Fun**  
Open source machine learning and data mining  
johannheinrichmuller.com  
[Download ZIP](#)

## Bibliography

---

- [@JohannH\\_M](#)
- <https://johannheinrichmuller.com/neurocomplexity>
- Orange Data Mining YouTube Channel
- The Creativity Code, Marcus du Sautoy
- Introduction to Machine Learning with Python, A. Müller and S. Guido