

Non-parametric models

KNN and Decision Trees

Tech Lead Data Science

Master en Data Science
2022-2023

ÍNDICE

1 Parametric vs non-parametric models

2 KNN

3 Decision Trees

PARAMETRIC VS NON-PARAMETRIC

PARAMETRIC VS NON-PARAMETRIC

PARAMETRIC

The amount of parameters
does not depend on the size of
the sample.

NON-PARAMETRIC

The amount of parameters is not
fixed and **might depend on the size**
of the sample.

PARAMETRIC VS NON-PARAMETRIC

PARAMETRIC

The amount of parameters
does not depend on the size of
the sample.

Linear regression

Logistic regression

NON-PARAMETRIC

The amount of parameters is not
fixed and **might depend on the size**
of the sample.

PARAMETRIC VS NON-PARAMETRIC

PARAMETRIC

The amount of parameters
does not depend on the size of
the sample.

Linear regression

Logistic regression

NON-PARAMETRIC

The amount of parameters is not
fixed and **might depend on the size**
of the sample.

KNN (K-Nearest Neighbours)

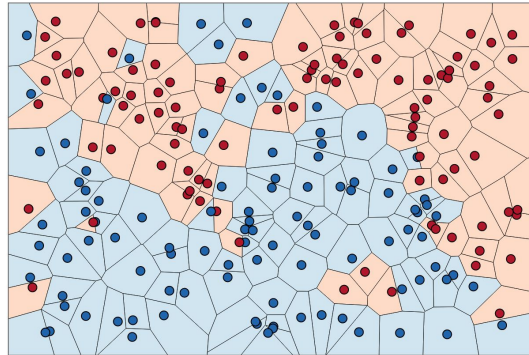
Decision Trees

K-NEAREST NEIGHBOURS

KNN MODELS

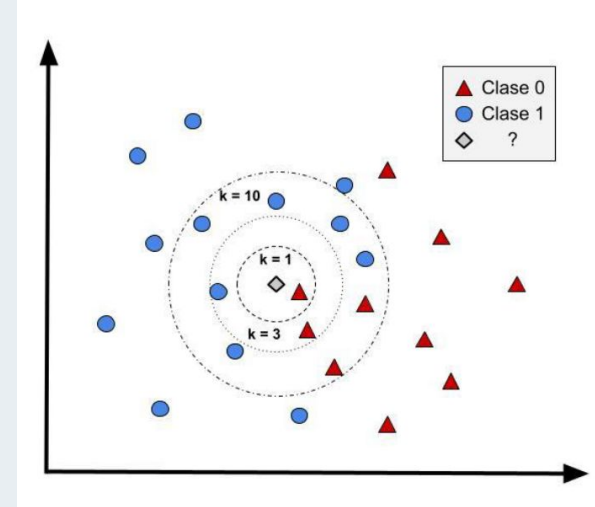
INTRODUCTION

- K-Nearest Neighbours (KNN) are **non-parametric models** that can be used for both **classification and regression** problems.
- The main idea of this models is to **use the k nearest data to predict new data**:
 - For classification, they will search for the class with highest frequency
 - For regression, they will take the mean of the k nearest data
- k will be set by the user.



KNN MODELS

HOW DOES IT WORK?

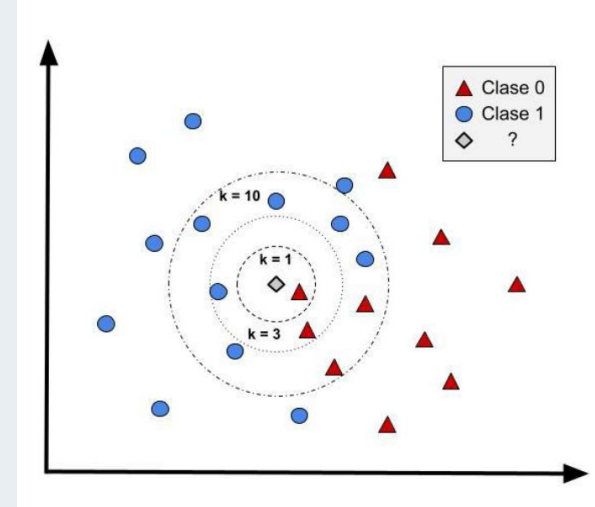


KNN MODELS

HOW DOES IT WORK?

KNN follows this structure to predict Y from X:

1. Calculates the distance from every data to Y.
2. Choose the nearest k values to Y.
3. Labeled Y using a criteria:
 - Mode for classification
 - Mean for regression

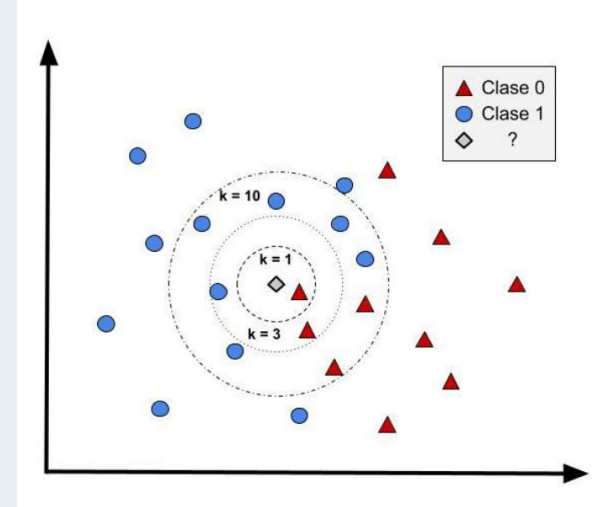


KNN MODELS

HOW DOES IT WORK?

KNN follows this structure to predict Y from X:

1. Calculates the distance from every data to Y.
2. Choose the nearest k values to Y.
3. Labeled Y using a criteria:
 - **Mode for classification**
 - Mean for regression

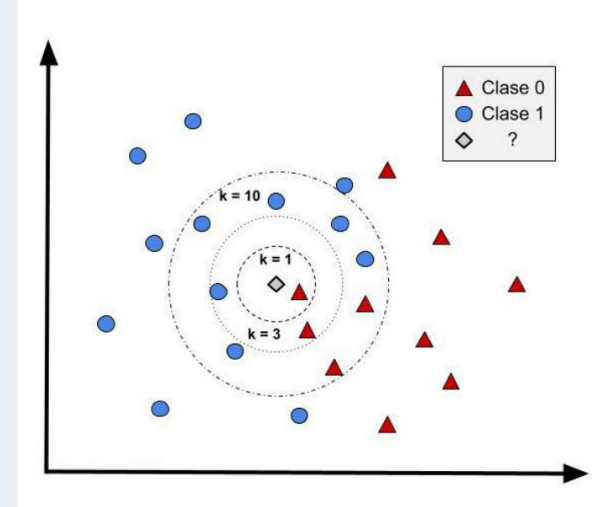


KNN MODELS

HOW DOES IT WORK?

KNN follows this structure to predict Y from X:

1. Calculates the distance from every data to Y.
2. Choose the nearest k values to Y.
3. Labeled Y using a criteria:
 - **Mode for classification**
 - Mean for regression



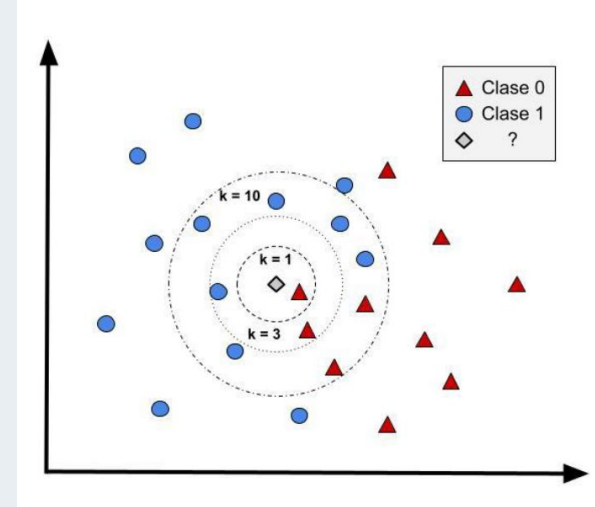
Con $k = 1$

KNN MODELS

HOW DOES IT WORK?

KNN follows this structure to predict Y from X:

1. Calculates the distance from every data to Y.
2. Choose the nearest k values to Y.
3. Labeled Y using a criteria:
 - **Mode for classification**
 - Mean for regression



Con $k = 1$

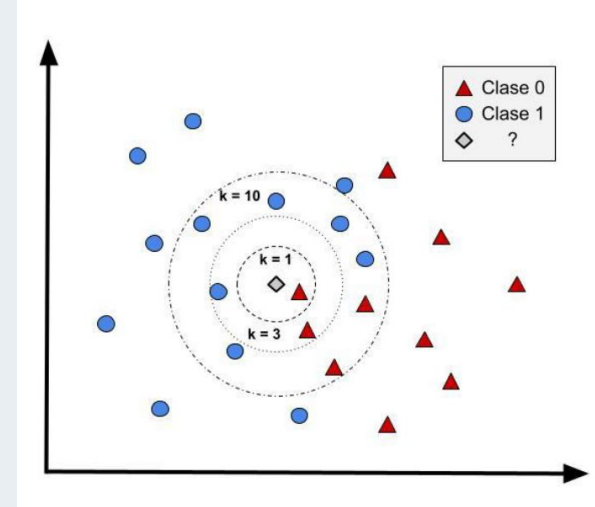


KNN MODELS

HOW DOES IT WORK?

KNN follows this structure to predict Y from X:

1. Calculates the distance from every data to Y.
2. Choose the nearest k values to Y.
3. Labeled Y using a criteria:
 - **Mode for classification**
 - Mean for regression



Con k = 1



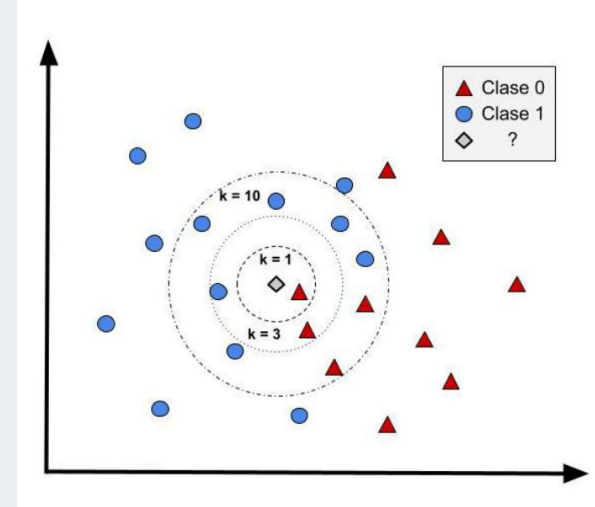
Con k = 3

KNN MODELS

HOW DOES IT WORK?

KNN follows this structure to predict Y from X:

1. Calculates the distance from every data to Y.
2. Choose the nearest k values to Y.
3. Labeled Y using a criteria:
 - **Mode for classification**
 - Mean for regression



Con $k = 1$



Con $k = 3$

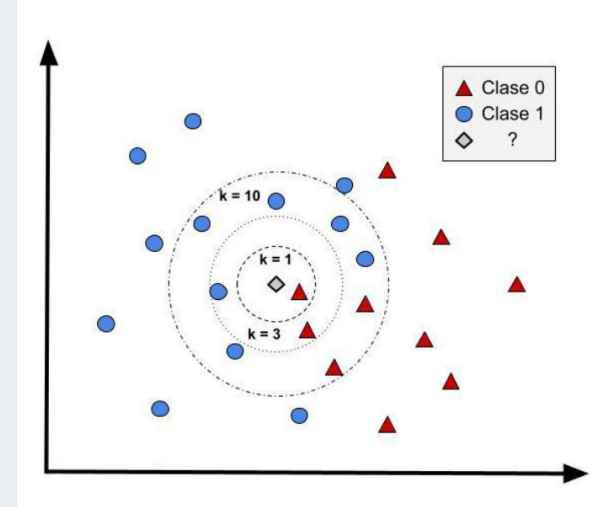


KNN MODELS

HOW DOES IT WORK?

KNN follows this structure to predict Y from X:

1. Calculates the distance from every data to Y.
2. Choose the nearest k values to Y.
3. Labeled Y using a criteria:
 - **Mode for classification**
 - Mean for regression



Con k = 1



Con k = 3



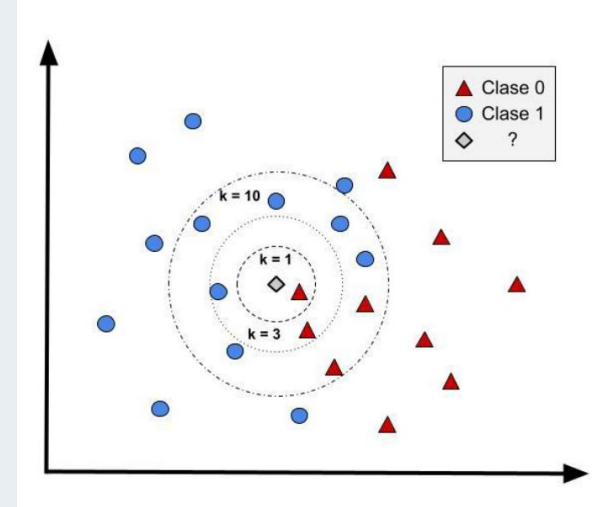
Con k = 10

KNN MODELS

HOW DOES IT WORK?

KNN follows this structure to predict Y from X:

1. Calculates the distance from every data to Y.
2. Choose the nearest k values to Y.
3. Labeled Y using a criteria:
 - **Mode for classification**
 - Mean for regression



Con k = 1



Con k = 3

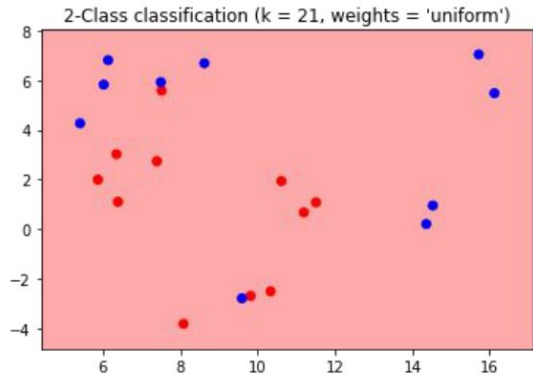
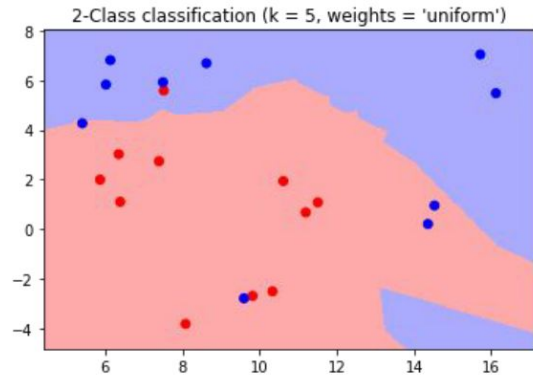
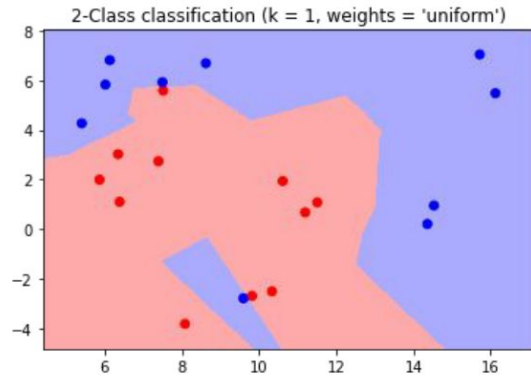


Con k = 10



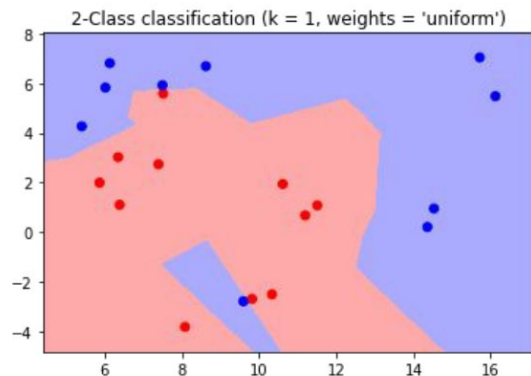
KNN MODELS

K'S IMPACT

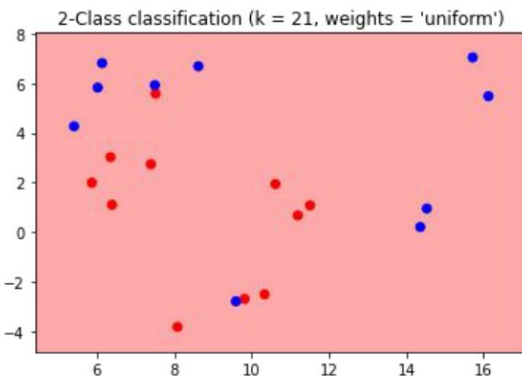
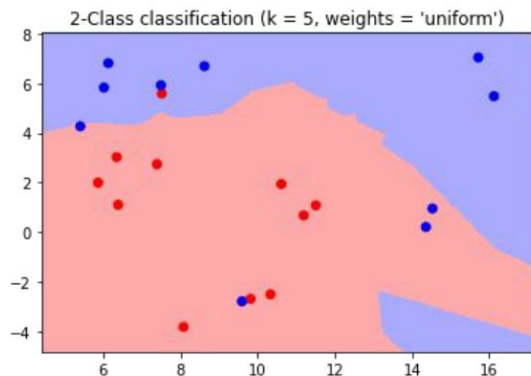


KNN MODELS

K'S IMPACT



+ Varianza
- Sesgo



- Varianza
+ Sesgo

KNN MODELS

PREDICTIVE STRATEGIES

CLASSIFICATION

REGRESSION

KNN MODELS

PREDICTIVE STRATEGIES

CLASSIFICATION

- To use the modal class:
uniform weights
- To weight the nearest
classes considering the
distances:
distance weights

$$\frac{1}{d(x, x_i)}$$

REGRESSION

KNN MODELS

PREDICTIVE STRATEGIES

CLASSIFICATION

- To use the modal class:
uniform weights
- To weight the nearest classes considering the distances:

distance weights

$$\frac{1}{d(x, x_i)}$$

REGRESSION

- **uniform weights**

$$\hat{y} = \frac{1}{k} \sum_{i \in Vec(x)} y_i$$

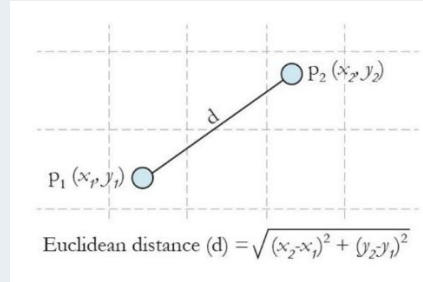
- **distance weights**

$$\hat{y} = \frac{1}{k} \sum_{i \in Vec(x)} \frac{y_i}{d(x, x_i)}$$

KNN MODELS

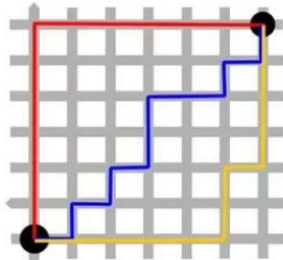
DISTANCES

The most commonly used distance is the **Euclidean distance**:



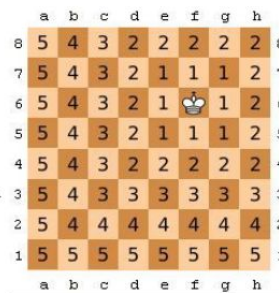
Manhattan o
distancia City-Block

$$\sum_{i=1}^m |x_i - y_i|$$



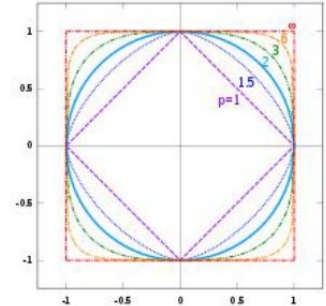
Chebyshev

$$\max_{i=1, \dots, m} \{|x_i - y_i|\}$$



Minkowski p

$$\sqrt[p]{\sum_{i=1}^m |x_i - y_i|^p}$$



KNN MODELS

PROS

- It has almost no assumptions about the data.
- It's easy to update once it is deployed to production.

CONS

KNN MODELS

PROS

- It has almost no assumptions about the data.
- It's easy to update once it is deployed to production.

CONS

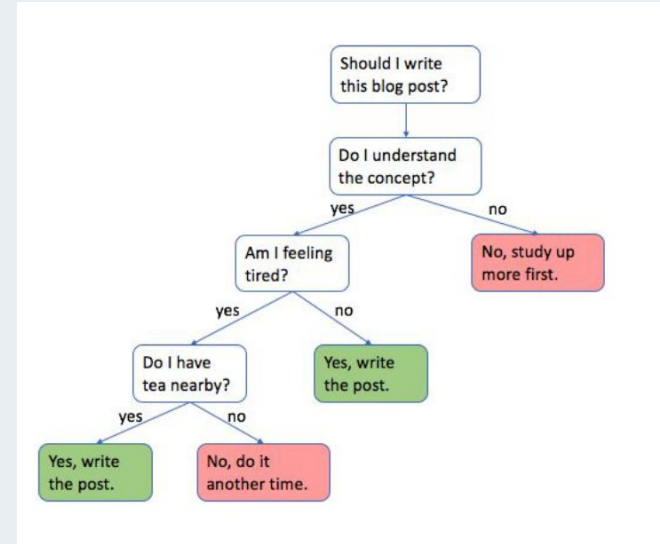
- Very sensible to outliers.
- We need to handle nas.
- It has a high computational cost:
 - Space: We need to store all the training data.
 - Time: It needs to calculate all the distances between the different data points.

DECISION TREES

DECISION TREES

INTRODUCTION

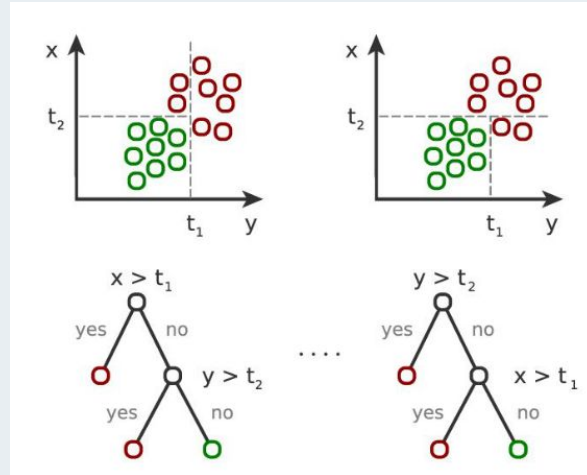
- Decision trees are **non-parametric models** that can be used for both **classification and regression** problems.
- They are highly interpretable and flexible.



DECISION TREES

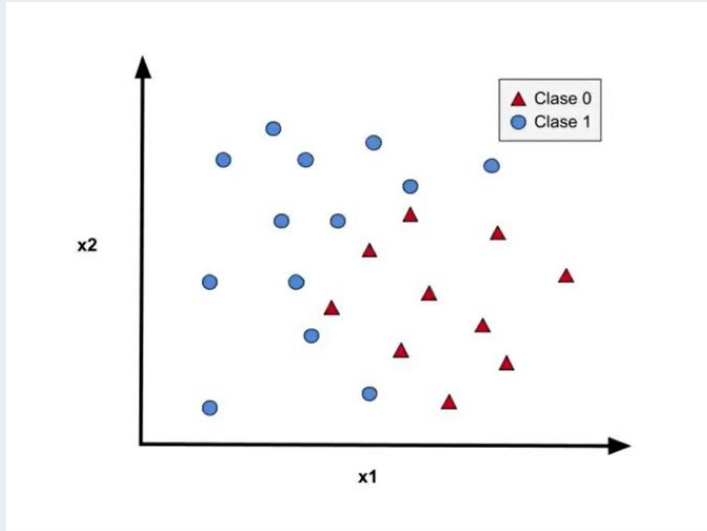
INTRODUCTION

- Decision trees are **non-parametric models** that can be used for both **classification and regression** problems.
- They are highly interpretable and flexible.



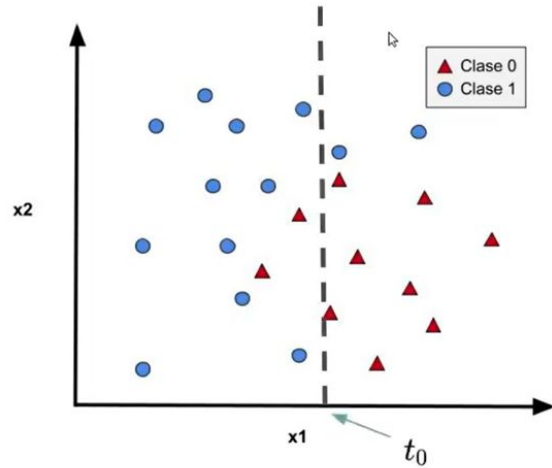
DECISION TREES

INTRODUCTION



DECISION TREES

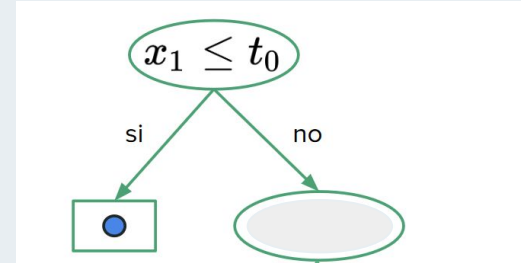
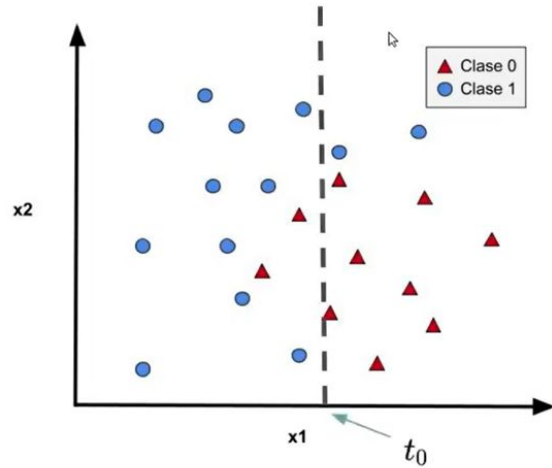
INTRODUCTION



$$x_1 \leq t_0$$

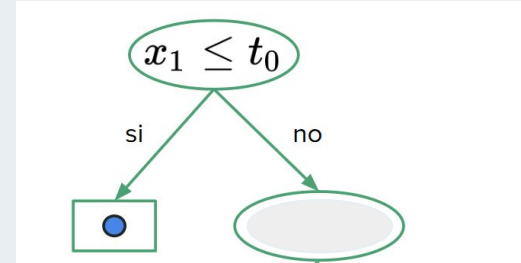
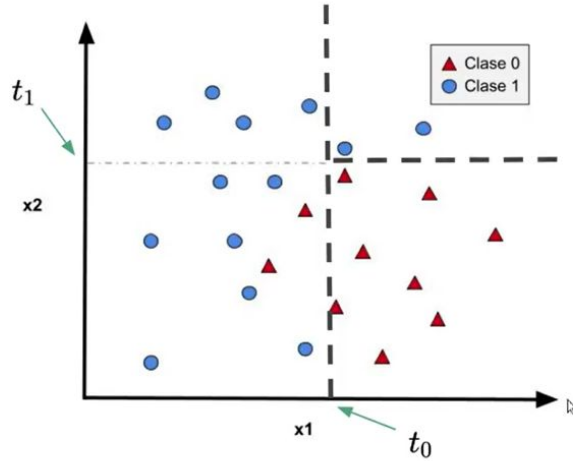
DECISION TREES

INTRODUCTION



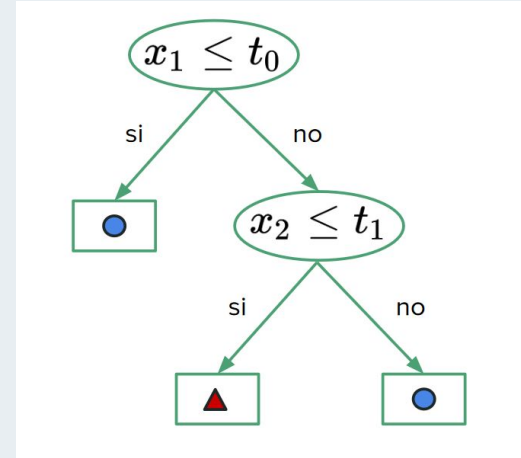
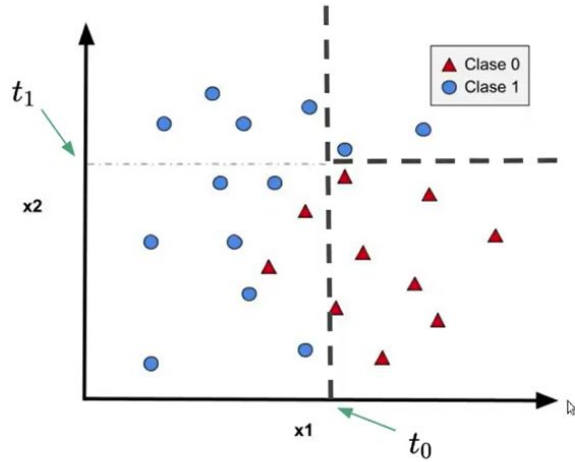
DECISION TREES

INTRODUCTION



DECISION TREES

INTRODUCTION



DECISION TREES

GOOD OR BAD SPLIT?

- In classification problems, we are going to measure the **impurity** to decide if a split was good or bad.
- The most common impurity measures are:

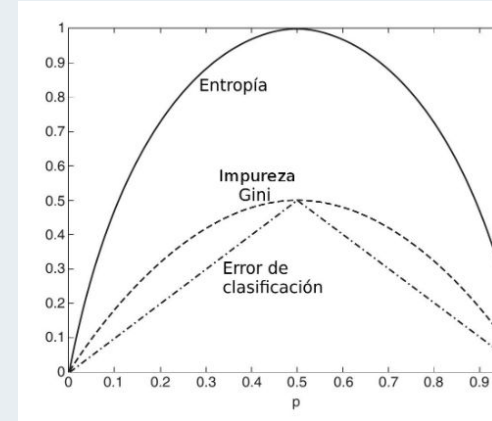
Entropía

$$Entropia(S) = \sum_{c \in Clases(S)} -p_c \cdot \log_2(p_c)$$

$$\text{donde } p_c = \frac{|\{x \in S, clase(x) = c\}|}{|S|}$$

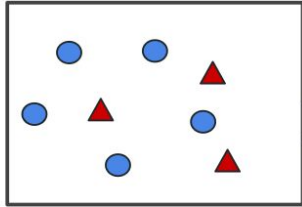
Índice de Gini

$$\begin{aligned} Gini(S) &= \sum_{c \in Clases(S)} p_c(1 - p_c) \\ &= 1 - \sum_{c \in Clases(S)} p_c^2 \end{aligned}$$



DECISION TREES

GOOD OR BAD SPLIT?

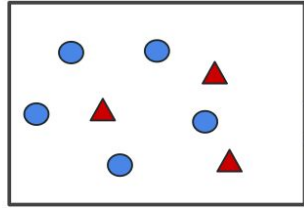


Gini: 0.468

Entropía: 0.954

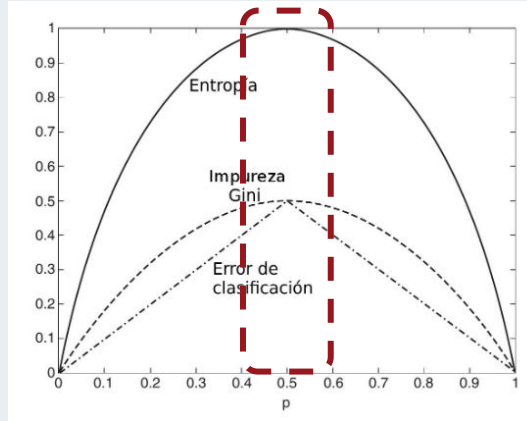
DECISION TREES

GOOD OR BAD SPLIT?



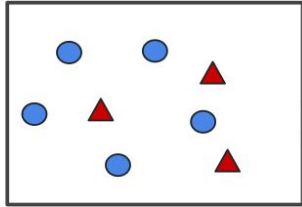
Gini: 0.468

Entropía: 0.954



DECISION TREES

GOOD OR BAD SPLIT?

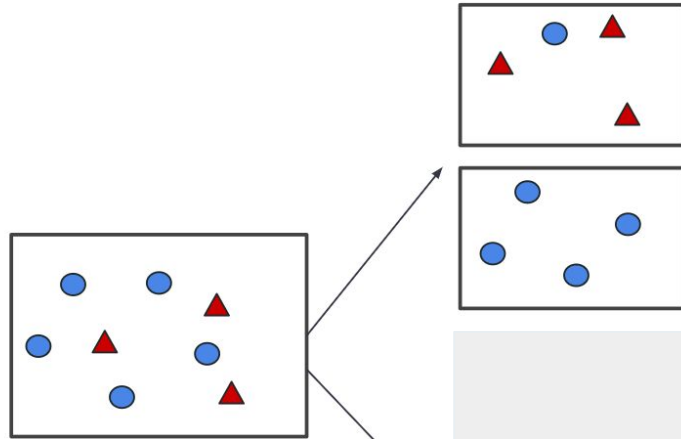


Gini: 0.468

Entropía: 0.954

DECISION TREES

GOOD OR BAD SPLIT?

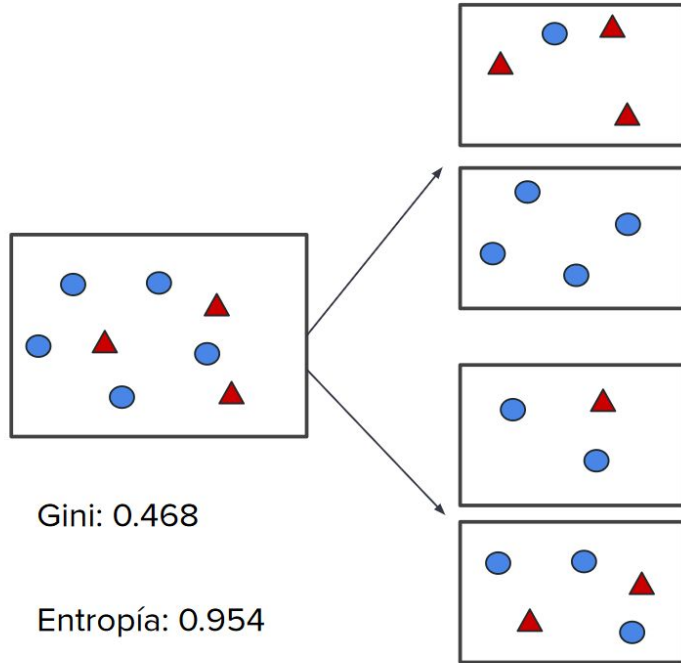


Gini: 0.468

Entropía: 0.954

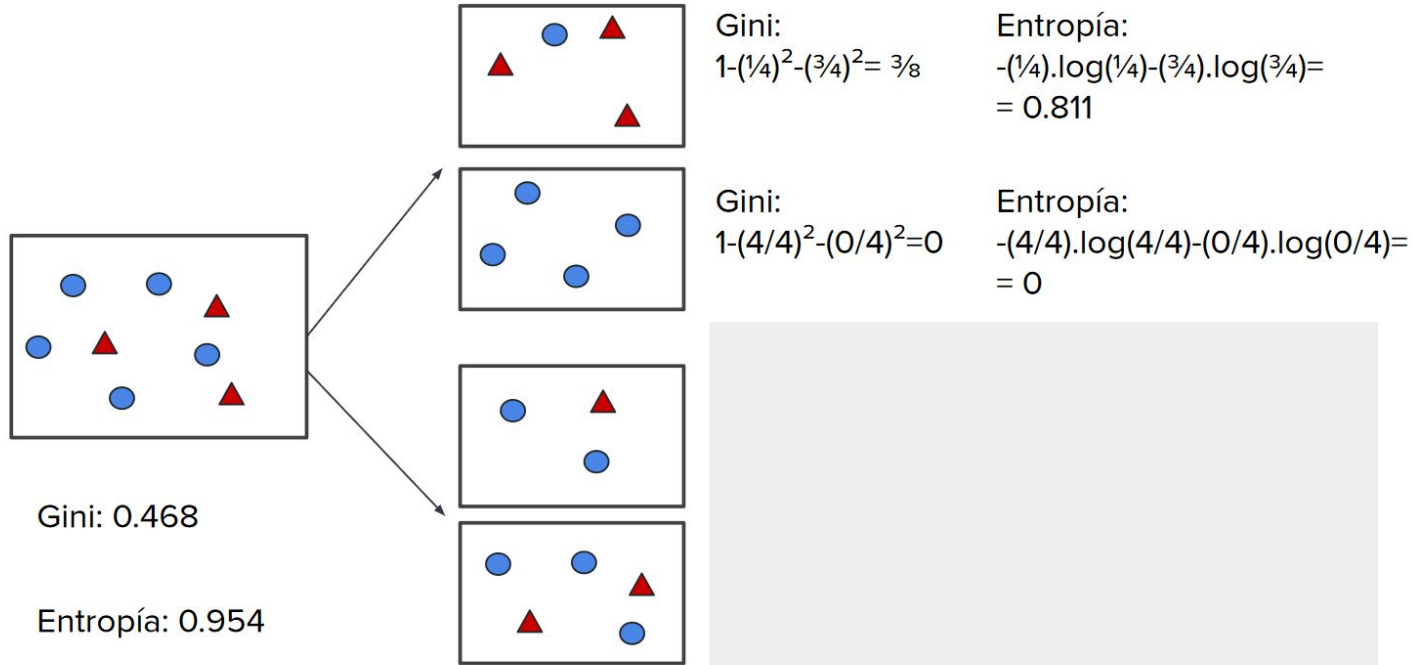
DECISION TREES

GOOD OR BAD SPLIT?



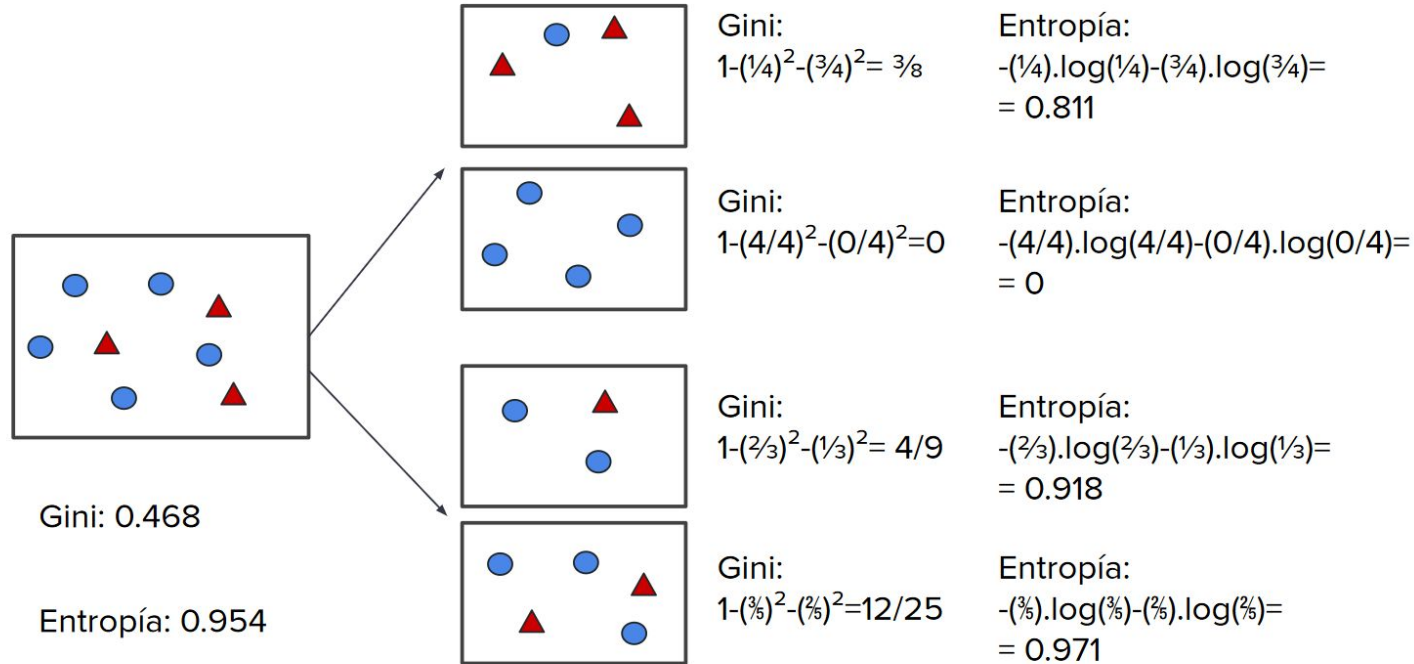
DECISION TREES

GOOD OR BAD SPLIT?



DECISION TREES

GOOD OR BAD SPLIT?



DECISION TREES

GOOD OR BAD SPLIT?

- To decide if a split is better than another one, we want to **minimize the impurity**:

$\text{Entropía}(\text{nodo padre}) - \text{Promedio pesado}(\text{Entropía}(\text{nodos hijos}))$

DECISION TREES

GOOD OR BAD SPLIT?

- To decide if a split is better than another one, we want to **minimize the impurity**:

Entropía(nodo padre) - Promedio pesado(Entropía (nodos hijos))

- **Information gain:**

Let V be the parent node with N data, X be the feature and s the split using that feature ($X \leq s$)

$$IG(X_j, s) = Entropy(V) - \sum_{V_i} \frac{N_i}{N} Entropy(V_i)$$

$$V_i = \{X \in V | X_j \leq s\} \cup \{X \in V | X_j > s\}$$

DECISION TREES

GOOD OR BAD SPLIT?

- To decide if a split is better than another one, we want to **minimize the impurity**:

Entropía(nodo padre) - Promedio pesado(Entropía (nodos hijos))

- Information gain:**

Let V be the parent node with N data, X be the feature and s the split using that feature ($X < s$)

$$IG(X_j, s) = Entropy(V) - \sum_{V_i} \frac{N_i}{N} Entropy(V_i)$$

$$V_i = \{X \in V | X_j \leq s\} \cup \{X \in V | X_j > s\}$$

THE HIGHER THE INFORMATION GAIN, THE BETTER THE ENTROPY REDUCTION

DECISION TREES

GOOD OR BAD SPLIT?

- To decide if a split is better than another one, we want to **minimize the impurity**:

$$\text{Entropía(nodo padre)} - \text{Promedio pesado(Entropía (nodos hijos))}$$

- Information gain:**

Let V be the parent node with N data, X be the feature and s the split using that feature ($X < s$)

$$IG(X_j, s) = Entropy(V) - \sum_{V_i} \frac{N_i}{N} Entropy(V_i)$$

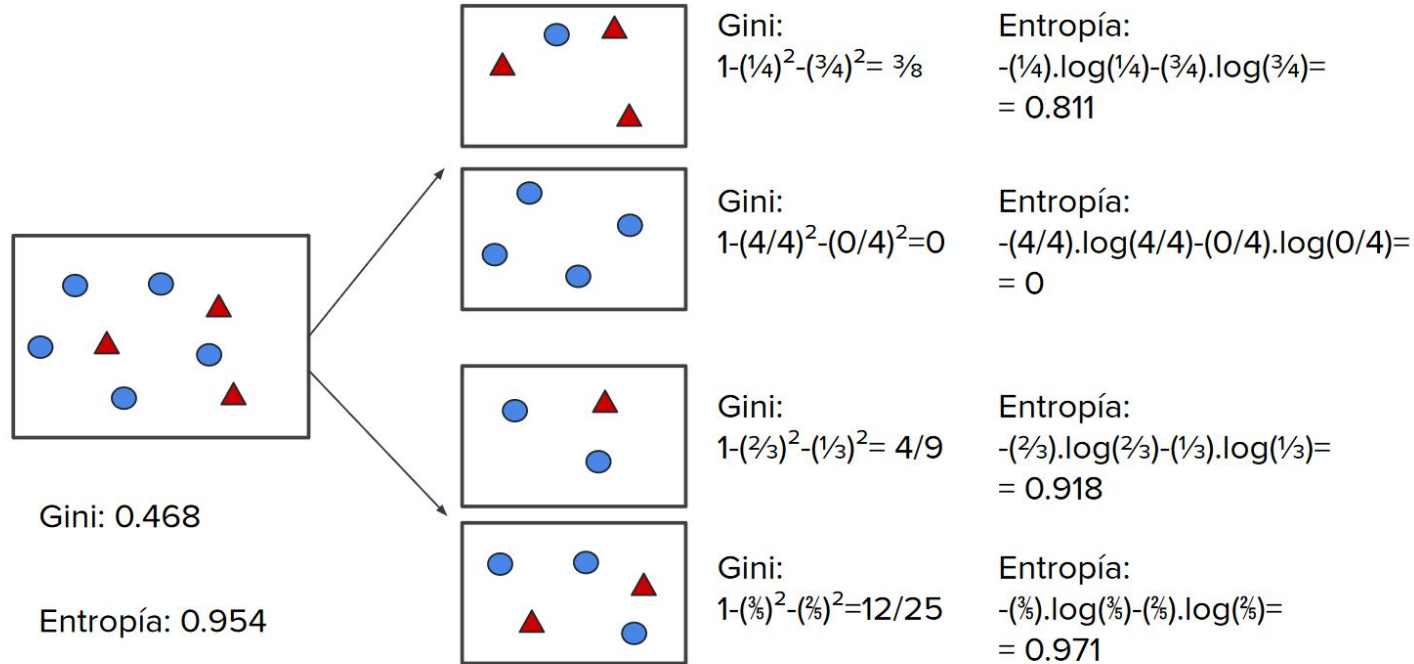
$$V_i = \{X \in V | X_j \leq s\} \cup \{X \in V | X_j > s\}$$

THE HIGHER THE INFORMATION GAIN, THE BETTER THE ENTROPY REDUCTION

If we use Gini index instead of entropy, it is called Gini Gain.

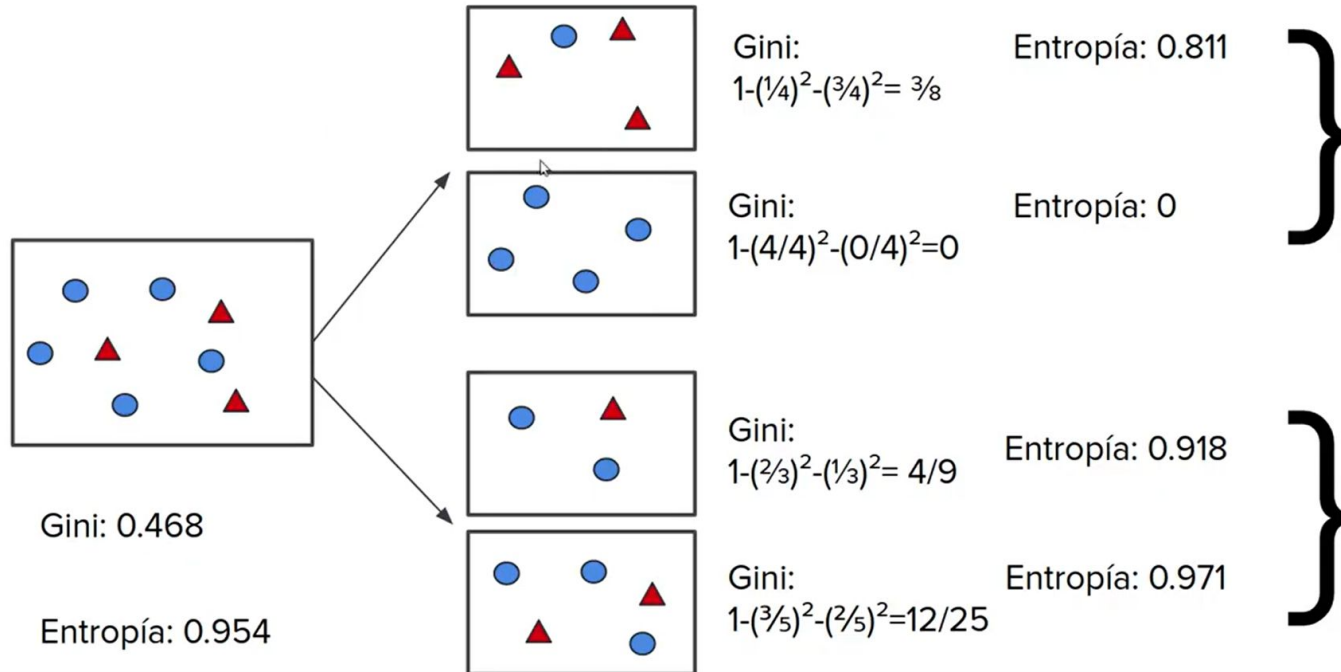
DECISION TREES

GOOD OR BAD SPLIT?



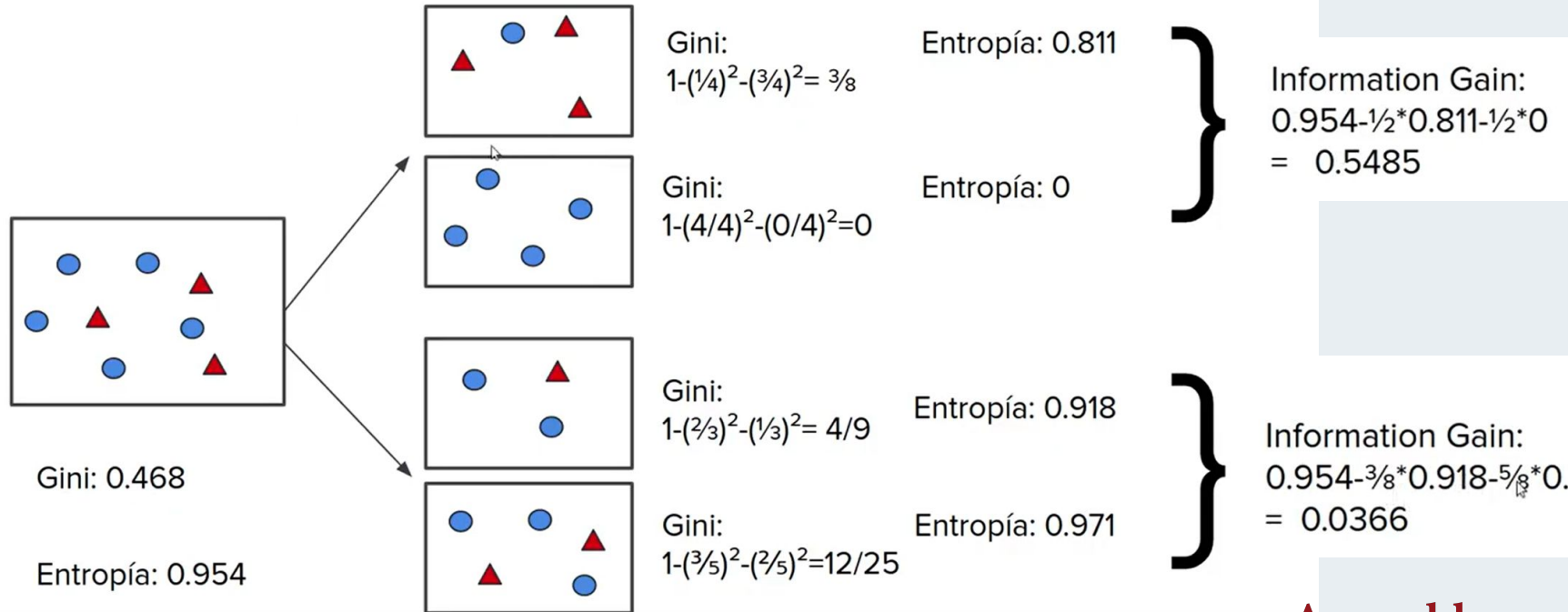
DECISION TREES

GOOD OR BAD SPLIT?



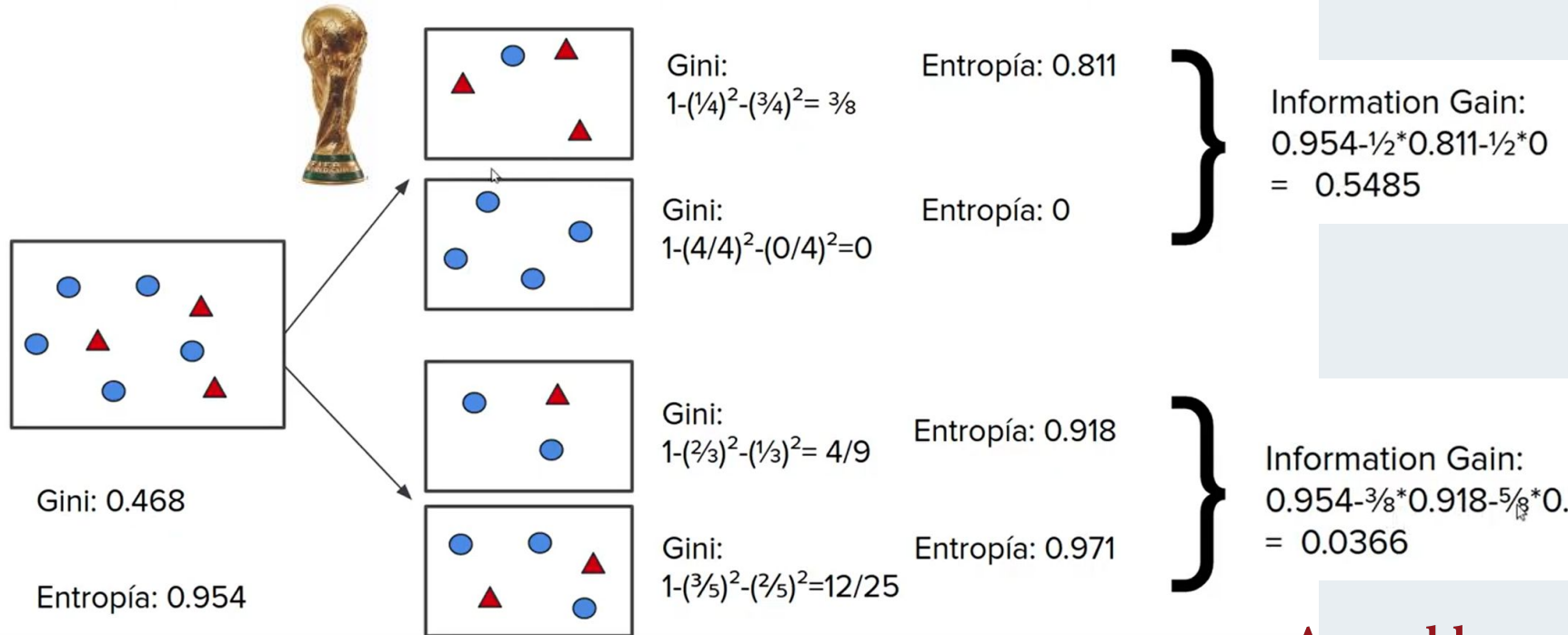
DECISION TREES

GOOD OR BAD SPLIT?



DECISION TREES

GOOD OR BAD SPLIT?



DECISION TREES

TRAINING

- For each feature, the model will try different splits.
- The best split is selected according to the information gain value.
- It checks the stop constraints:
 - Maximum depth
 - Minimum amount of data to split
 - Minimum amount of data in a leaf to consider a split

DECISION TREES

TRAINING

- For each feature, the model will try different splits.
- The best split is selected according to the information gain value.
- It checks the stop constraints:
 - Maximum depth
 - Minimum amount of data to split
 - Minimum amount of data in a leaf to consider a split

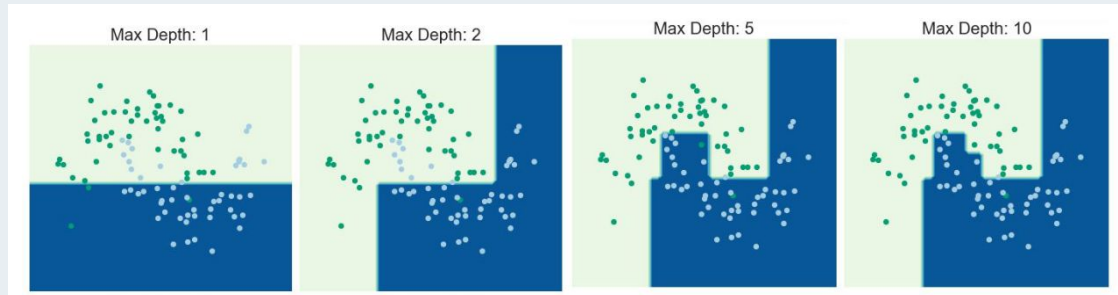
**Without them, the tree keeps growing
until it finds a “perfect” model**

DECISION TREES

TRAINING

- For each feature, the model will try different splits.
- The best split is selected according to the information gain value.
- It checks the stop constraints:
 - Maximum depth
 - Minimum amount of data to split
 - Minimum amount of data in a leaf to consider a split

**Without them, the tree keeps growing
until it finds a “perfect” model**

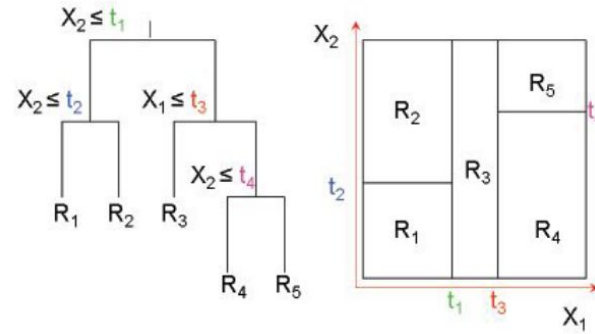


DECISION TREES

REGRESSION

- We'll try to minimize the Mean Squared Error

$$ECM = \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

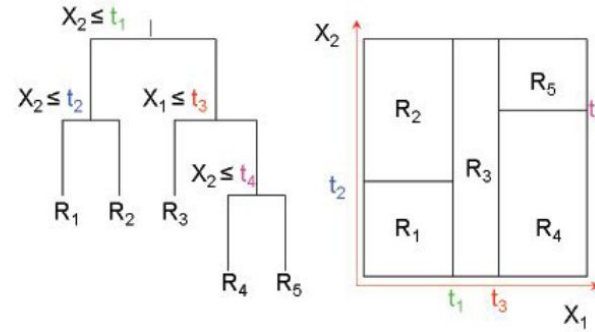


DECISION TREES

REGRESSION

- We'll try to minimize the Mean Squared Error

$$ECM = \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$



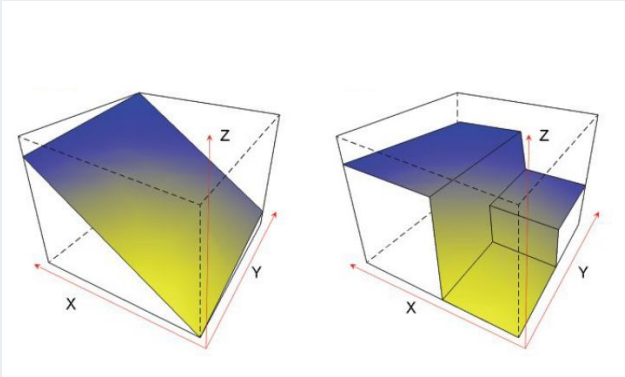
$$\sum_{i: x_i \in R_1(j,s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: x_i \in R_2(j,s)} (y_i - \hat{y}_{R_2})^2$$

$$R_1(j, s) = \{X | X_j < s\} \text{ and } R_2(j, s) = \{X | X_j \geq s\}$$

$$\Delta = ECM(\text{padre}) - \sum_{j \in \text{hijos}} \frac{N_j}{N} ECM(\text{hijo}_j)$$

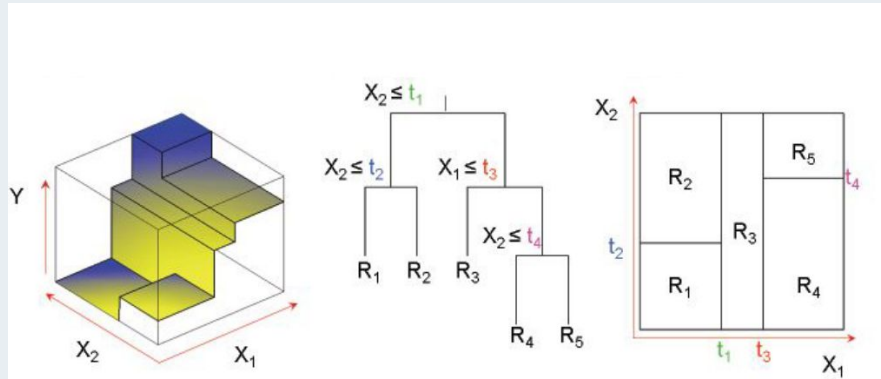
DECISION TREES

REGRESSION



DECISION TREES

REGRESSION



DECISION TREES

PROS

- Easy to interpret
- Fast training
- They can handle continuous and discrete data and missing values

CONS

- They tend to overfit
- Their performance is usually worse than the one for other classical models