

# Electiva Bases de datos

## Árboles de decisión

Jimmy Mateo Guerrero<sup>1</sup>

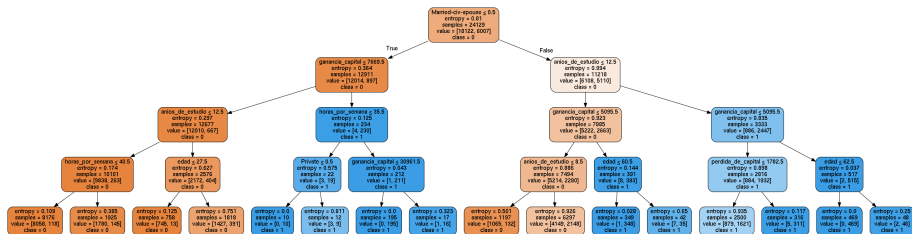
<sup>1</sup>Departamento de Sistemas  
Facultad de Ingeniería

Árboles de decisión, 2019

- 1 Induciendo árboles
- 2 Implementando árboles
  - Criterio de selección
  - Criterio de parada
  - Sobreajuste

- Probablemente el método más conocido para resolver problemas de clasificación.
- Uno de los primeros desarrollos en ML.
- Muy asociado a nuestra forma de proceder

# Árboles de decisión



# Árboles de decisión

- ID3 y sus sucesores han sido desarrollados por Ross Quinlan en la década de 1970.
- Posteriormente trabajó en Sydney Uni, Rand Corporation en California.
- Ahora dirige su propia compañía Rulequest.



## Características:

- Capaz de aprender conceptos disyuntivos.
- Tolerante al ruido en atributos y clase.
- Tolerante a valores faltantes en atributos.
- Estructura de datos recursiva definida por:
  - Nodos hoja que indican un rótulo de clase.
  - Nodos internos que contienen un test sobre el valor de un atributo.  
Para cada resultado del test hay una rama y un subárbol con la misma estructura que el árbol.
- Representación alternativa:  
conjunto de reglas: IF Condición THEN Clasificación

El problema de encontrar un árbol de decisión que concuerde con el conjunto de entrenamiento tiene una solución trivial.

- Simplemente construimos un árbol de decisión que tenga una rama por cada ejemplo.
- El problema con este árbol trivial es que solo memoriza las observaciones.
- No extrae ningún patrón a partir de los ejemplos y no puede extrapolar a ejemplos que no vio.

- Extraer un patrón significa poder describir un gran número de casos de un modo conciso. El lugar de sólo tratar de encontrar uno que concuerde con los ejemplos, deberíamos también tratar de encontrar uno conciso.
- Este es un ejemplo de un principio general del aprendizaje inductivo:
  - La hipótesis más probable es la más simple que sea consistente con todas las observaciones.



- Navaja de Occam (Occam's Razor).
- Basado en una premisa muy simple:
  - En igualdad de condiciones la solución más sencilla es probablemente la correcta.
  - No ha de presumirse de la existencia de más cosas que las absolutamente necesarias.
  - No hay que desarrollar sistemas más complejos que lo necesario.

- Desafortunadamente, encontrar el árbol de decisión óptimo (más chico) es un problema intratable, pero algunas heurísticas simples logran encontrar uno suficientemente pequeño.
- La idea básica del algoritmo es usar primero el atributo “más importante”, donde esto significa el que tiene un efecto más importante en la agrupación de los ejemplos en conjuntos que comparten la misma clasificación.

## ¿Heurística?

- Como metodología científica, la heurística es aplicable a cualquier ciencia e incluye la utilización de: medios auxiliares, principios, reglas, estrategias o algoritmos que faciliten la búsqueda de vías de solución a problemas.

## Id3(Ejemplos, Clase, Atributos)

- Si todos los ejemplos son positivos retornar un nodo positivo.
- Si todos los ejemplos son negativos retornar un nodo negativo.
- Si Atributos está vacío retornar nodo de la clase más frecuente .
- En otro caso.
  - $A = \text{ELEGIR\_MEJOR\_ATRIBUTO}(\text{Atributos}, \text{Ejemplos})$
  - $\text{Raíz} = A$
  - Para cada valor posible  $v$  de  $A$ 
    - Agregar una nueva rama debajo de  $\text{Raíz}$  con test  $A = v$
    - $\text{Ejemplos}_v = \text{elementos de Ejemplos con } A = v$
    - Si  $\text{Ejemplos}_v$  está vacío devolver un nodo de la clase más frecuente en ejemplos.
    - Si no agregar debajo de esta rama el subárbol  $\text{Id3}(\text{Ejemplos}_v, \text{Clase}, \text{Atributos} - A)$

## 1 Induciendo árboles

## 2 Implementando árboles

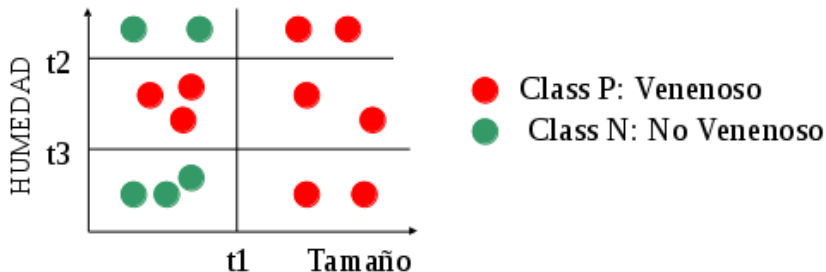
- Criterio de selección
- Criterio de parada
- Sobreajuste

- Asumamos que el Criterio de Partición se limita a seleccionar un atributo.
- ¿Cómo elegir el atributo más “importante” para dividir?
- Muchas aproximaciones:
  - Information Gain.
  - Information ratio.
  - Gini.
  - ... tantas como autores.

Existen varios aspectos prácticos a resolver a la hora de implementar árboles de decisión.

- Cuál es el objetivo del árbol?
- Cuál es un buen criterio de partición?
- Qué tan profundo debe ser el árbol?
- Qué pasa cuando hay valores faltantes?
- Cómo podemos afrontar la eficiencia computacional para grandes evidencias?

# Criterio de selección



Atributos: tamaño y humedad

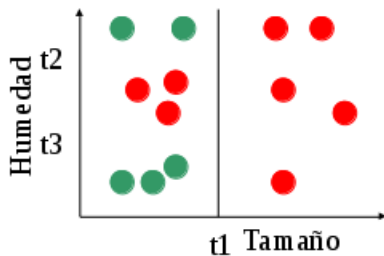
Tamaño tiene 2 valores:  $> t1$  o  $\leq t1$

Humedad tiene 3 valores:  $> t2$ , ( $> t3$  y  $\leq t2$ ),  $\leq t3$



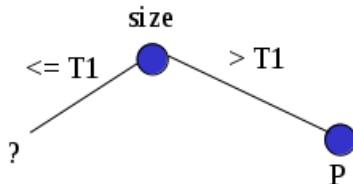
# Criterio de selección

Supongamos que elegimos el tamaño como el mejor atributo:

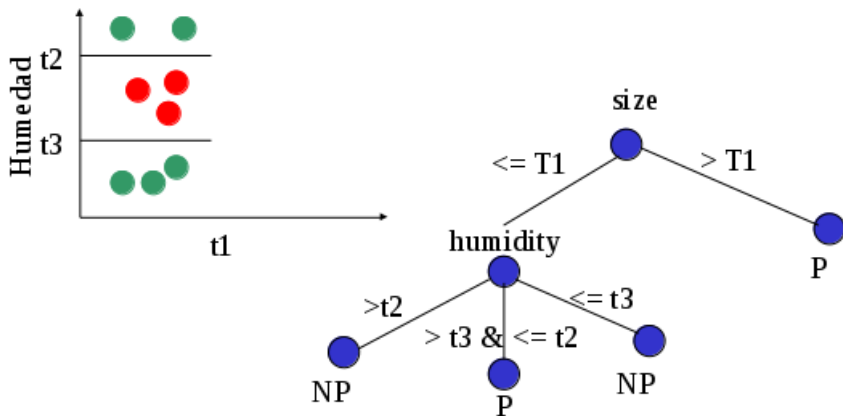


● Class P: Venenoso

● Class N: No Venenoso



Supongamos que elegimos la humedad como el siguiente mejor atributo:



Para decidir cual es el mejor atributo para abrir el árbol:

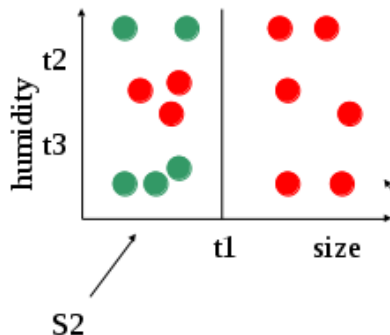
- Usamos algunas definiciones de Teoría de la Información.
- La medida de incertidumbre o entropía asociada a una variable aleatoria  $X$  se define como:

- 

$$H(X) = - \sum_{i=1} p_i \log_2(p_i)$$

- Siendo  $\log_2$  una redefinición de logaritmo en base 2 con  $\log(0)=0$ .

# Criterio de selección



El tamaño divide la muestra en dos.

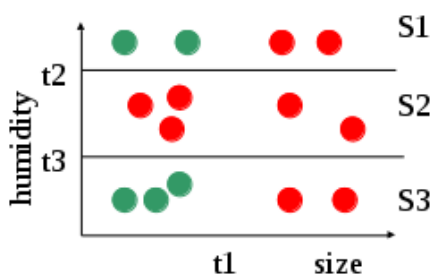
$S1 = \{ 6P, 0NP \}$  ( $size > t1$ )

$S2 = \{ 3P, 5NP \}$  ( $size < t1$ )

$$H(S1) = -\left(\frac{6}{6}\right)\log_2\left(\frac{6}{6}\right) + (0/6)\log_2(0/6) = 0$$

$$H(S2) = -\left(\frac{3}{8}\right)\log_2\left(\frac{3}{8}\right) - \left(\frac{5}{8}\right)\log_2\left(\frac{5}{8}\right) = 0.95$$

# Criterio de selección



La humedad divide la muestra en tres.

$S_1 = \{2P, 2NP\}$   $humidity > t_2$

$S_2 = \{5P, 0NP\}$   $t_3 < humidity < t_2$

$S_3 = \{2P, 3NP\}$   $humidity < t_3$

$$\begin{aligned} H(S_1) &= -(2/4)\log_2(2/4) + \\ &\quad -(2/4)\log_2(2/4) \\ &= 1 \end{aligned}$$

$$H(S_2) = 0$$

$$\begin{aligned} H(S_3) &= -(2/5)\log_2(2/5) + \\ &\quad -(3/5)\log_2(3/5) = 0.97 \end{aligned}$$

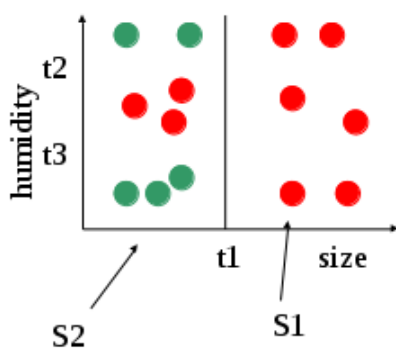
- Definimos la función de Ganancia de Información como:

- 

$$IG(A) = H(S) - \sum_v (S_v/S) H(S_v)$$

- Donde:
  - $H(S)$  es la entropía del conjunto de todos los ejemplos
  - $H(S_v)$  es la entropía del subconjunto de los ejemplos con valor  $v$  para el atributo  $A$
- Elegiremos al atributo que maximice la  $IG$ , por lo tanto debemos calcular  $IG$  para todos los atributos:

# Criterio de selección



$$H(S1) = 0$$

$$H(S2) = 0.95$$

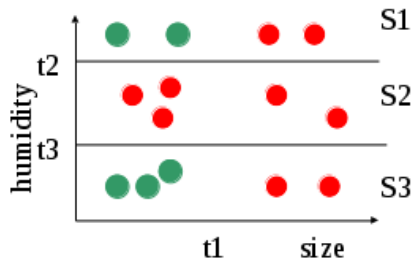
$$H(S) = -(9/14)\log_2(9/14) + \\ -(5/14)\log_2(5/14) = 0.94$$

$$|S1|/|S| = 6/14$$

$$|S2|/|S| = 8/14$$

$$IG(\text{size}) = H(S) - ( |S1|/|S| * H(S1) + |S2|/|S| * H(S2) ) \\ = 0.94 - ( 6/14 * 0 + 8/14 * 0.95 ) = 0.40$$

# Criterio de selección



$$H(S_1) = 1$$

$$H(S_2) = 0$$

$$H(S_3) = 0.97$$

$$H(S) = 0.94$$

$$|S_1|/|S| = 4/14$$

$$|S_2|/|S| = 5/14$$

$$|S_3|/|S| = 5/14$$

$$\begin{aligned} IG(\text{humidity}) &= H(S) - ( |S_1|/|S| * H(S_1) + |S_2|/|S| * H(S_2) + |S_3|/|S| * H(S_3) ) \\ &= 0.94 - ( 4/14 * 1 + 5/14 * 0 + 5/14 * 0.97 ) \\ &= \mathbf{0.31} \end{aligned}$$



## 1 Induciendo árboles

## 2 Implementando árboles

- Criterio de selección
- Criterio de parada
- Sobreajuste

Hasta cuando dividimos?

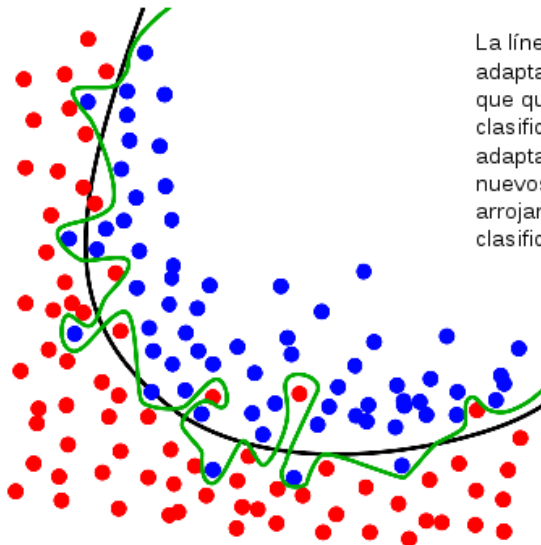
- Hasta que tenga clases puras en todos los nodos
- Hasta que no tenga más variables disponibles
- O hay algo mejor?

- La versión canónica dice que el árbol se sigue desarrollando hasta que se llega a una clase uniforme. Esto en la práctica no suele ser efectivo.
- Existen muchas opciones de criterio de parada:
  - No hay más atributos.
  - Clase uniforme.
  - Se supera un nivel de profundidad máximo.
  - El número de casos es menor a un mínimo.
  - El mejor valor del criterio de partición es menor de un umbral.
  - El número de casos de los nodos hijos es menor que un mínimo.

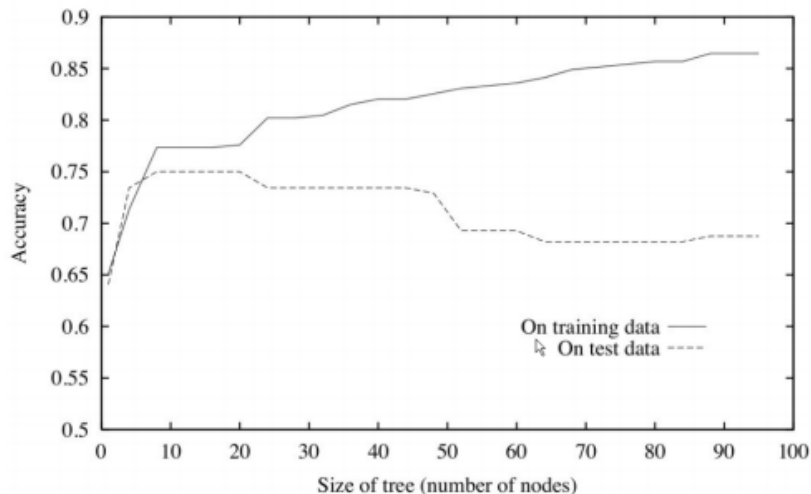
## 1 Induciendo árboles

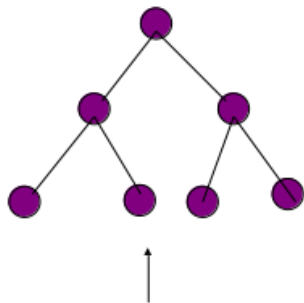
## 2 Implementando árboles

- Criterio de selección
- Criterio de parada
- Sobreajuste



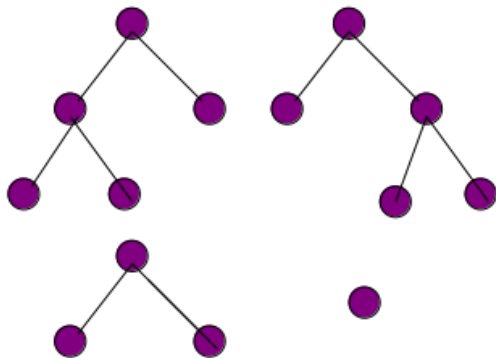
La línea verde como clasificador se adapta mejor a los datos con los que se ha entrenado al clasificador, pero está demasiado adaptada a ellos, de forma que ante nuevos datos probablemente arrojará más errores que la clasificación usando la línea negra.





Árbol original

Posibles árboles después de la poda:





Aurélien Géron

*Hands-On Machine Learning with Scikit-Learn and TensorFlow:  
Concepts, Tools, and Techniques to Build Intelligent Systems.*

Publisher: O'Reilly Media, Year: 2017



Jake VanderPlas.

Python Data Science Handbook: Essential Tools for Working with  
Data

*O'Reilly Media, Year: 2016*



Brett Lantz.

Machine learning with R

*Packt Publishing, Year: 2013*