

## Tema 3: Introducción al aprendizaje automático

M. A. Gutiérrez Naranjo

F. J. Martín Mateos

J. L. Ruiz Reina

Dpto. Ciencias de la Computación e Inteligencia Artificial  
Universidad de Sevilla

# Contenido

- Introducción
- Aprendizaje de árboles de decisión
- Aprendizaje de reglas
- Aprendizaje basado en instancias: kNN
- *Clustering*: k-medias y clustering jerárquico.

# Sección 1

## Sección 1

### Introducción

# ¿Qué es aprender?

## Definiciones de *aprendizaje*:

- Cualquier *cambio* en un sistema que le permita realizar la misma tarea de manera *más eficiente* la próxima vez  
(*H. Simon*)
- *Modificar la representación* del mundo que se está percibiendo  
(*R. Michalski*)
- Realizar cambios *útiles* en nuestras mentes  
(*M. Minsky*)
- Se dice que aprendemos de la *experiencia* a realizar alguna *tarea* si la realización de la tarea mejora con la experiencia respecto a alguna *medida de rendimiento*  
(*T. M. Mitchell*)

# ¿Qué es aprender?

## Definiciones de *aprendizaje*:

- Cualquier *cambio* en un sistema que le permita realizar la misma tarea de manera *más eficiente* la próxima vez  
(*H. Simon*)
- *Modificar la representación* del mundo que se está percibiendo  
(*R. Michalski*)
- Realizar cambios *útiles* en nuestras mentes  
(*M. Minsky*)
- Se dice que aprendemos de la *experiencia* a realizar alguna *tarea* si la realización de la tarea mejora con la experiencia respecto a alguna *medida de rendimiento*  
(*T. M. Mitchell*)

# ¿Qué es aprender?

## Definiciones de *aprendizaje*:

- Cualquier *cambio* en un sistema que le permita realizar la misma tarea de manera *más eficiente* la próxima vez  
(*H. Simon*)
- *Modificar la representación* del mundo que se está percibiendo  
(*R. Michalski*)
- Realizar cambios *útiles* en nuestras mentes  
(*M. Minsky*)
- Se dice que aprendemos de la *experiencia* a realizar alguna *tarea* si la realización de la tarea mejora con la experiencia respecto a alguna *medida de rendimiento*  
(*T. M. Mitchell*)

# ¿Qué es aprender?

## Definiciones de *aprendizaje*:

- Cualquier *cambio* en un sistema que le permita realizar la misma tarea de manera *más eficiente* la próxima vez  
(*H. Simon*)
- *Modificar la representación* del mundo que se está percibiendo  
(*R. Michalski*)
- Realizar cambios *útiles* en nuestras mentes  
(*M. Minsky*)
- Se dice que aprendemos de la *experiencia* a realizar alguna *tarea* si la realización de la tarea mejora con la experiencia respecto a alguna *medida de rendimiento*  
(*T. M. Mitchell*)

# Aprendizaje

- **Aprendizaje automático:** construir programas que mejoran automáticamente con la experiencia
- Ejemplos de tareas:
  - Construcción de bases de conocimiento a partir de la experiencia
  - Clasificación y diagnóstico
  - Minería de datos, descubrir estructuras desconocidas en grandes grupos de datos
  - Resolución de problemas, planificación y acción



# Tipos de aprendizaje y paradigmas

- Tipos de aprendizaje
  - Supervisado
  - No supervisado
  - Con refuerzo
- Paradigmas
  - Aprendizaje por memorización
  - Clasificación (*Clustering*)
  - Aprendizaje inductivo
  - Aprendizaje por analogía
  - Descubrimiento
  - Algoritmos genéticos, redes neuronales

## Ejemplo de aprendizaje

- Conjunto de entrenamiento
  - Ejemplos: días en los que es recomendable (o no) jugar al tenis
  - Representación como una lista de pares atributo–valor

EJ.	CIELO	TEMPERATURA	HUMEDAD	VIENTO	JUGAR TENIS
$D_1$	SOLEADO	ALTA	ALTA	DÉBIL	-
$D_2$	SOLEADO	ALTA	ALTA	FUERTE	-
$D_3$	NUBLADO	ALTA	ALTA	DÉBIL	+
$D_4$	LLUVIA	SUAVE	ALTA	DÉBIL	+
...					

- Objetivo: Dado el *conjunto de entrenamiento*, aprender el concepto “Días en los que se juega al tenis”
  - Se trata de *aprendizaje supervisado*
- Problema: ¿Cómo expresar lo aprendido?
  - En este tema, veremos algoritmos para aprender *árboles de decisión*, *reglas*, *modelos probabilísticos*,...

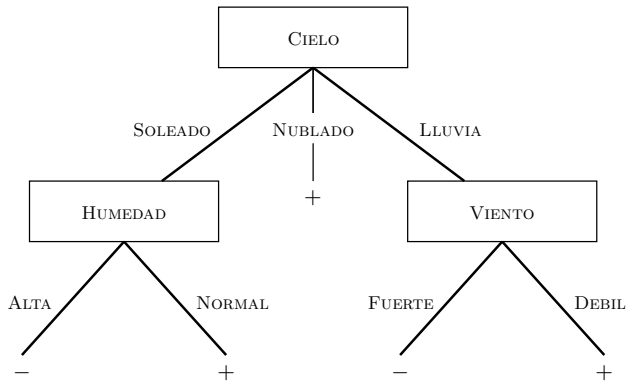
# Sección 2

## Sección 2

### Aprendizaje de árboles de decisión

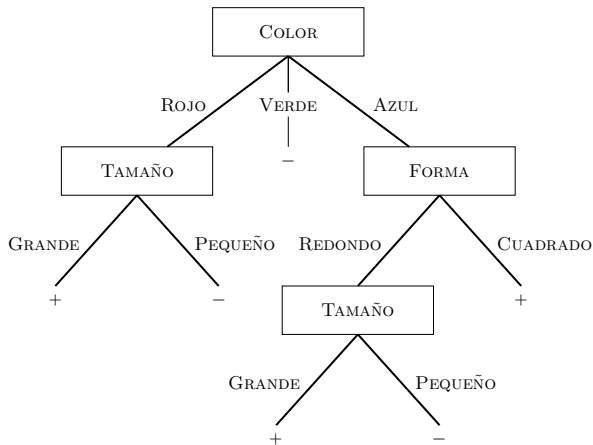
# Árboles de decisión

- Un **árbol de decisión** es un grafo etiquetado que representa un concepto.
- Ejemplos de árboles de decisión



# Árboles de decisión

- Ejemplos de árboles de decisión



# Árboles de decisión

- Árboles de decisión
  - Nodos interiores: atributos
  - Arcos: posibles valores del nodo origen
  - Hojas: valor de clasificación (usualmente + ó -, aunque podría ser cualquier conjunto de valores, no necesariamente binario)
  - Representación de una función objetivo
- Disyunción de reglas proposicionales:
$$\begin{aligned} &(\text{CIELO}=\text{SOLEADO} \wedge \text{HUMEDAD}=\text{ALTA} \rightarrow \text{JUGAR TENIS}=-) \\ \vee &(\text{CIELO}=\text{SOLEADO} \wedge \text{HUMEDAD}=\text{NORMAL} \rightarrow \text{JUGAR TENIS}=+) \\ \vee &(\text{CIELO}=\text{NUBLADO} \rightarrow \text{JUGAR TENIS}=+) \\ \vee &(\text{CIELO}=\text{LLUVIOSO} \wedge \text{VIENTO}=\text{FUERTE} \rightarrow \text{JUGAR TENIS}=-) \\ \vee &(\text{CIELO}=\text{LLUVIOSO} \wedge \text{VIENTO}=\text{DEBIL} \rightarrow \text{JUGAR TENIS}=+) \end{aligned}$$
- Capaz de representar cualquier subconjunto de instancias

# Aprendizaje de árboles de decisión

- Objetivo: aprender un árbol de decisión consistente con los ejemplos, para posteriormente clasificar ejemplos nuevos
- Ejemplo de conjunto de entrenamiento:

EJ.	CIELO	TEMPERATURA	HUMEDAD	VIENTO	JUGAR TENIS
$D_1$	SOLEADO	ALTA	ALTA	DÉBIL	-
$D_2$	SOLEADO	ALTA	ALTA	FUERTE	-
$D_3$	NUBLADO	ALTA	ALTA	DÉBIL	+
$D_4$	LLUVIA	SUAVE	ALTA	DÉBIL	+
...					

# Árboles de decisión

$$D = [34^+, 27^-]$$



# Árboles de decisión

$$D = [34^+, 27^-]$$

COLOR

# Árboles de decisión

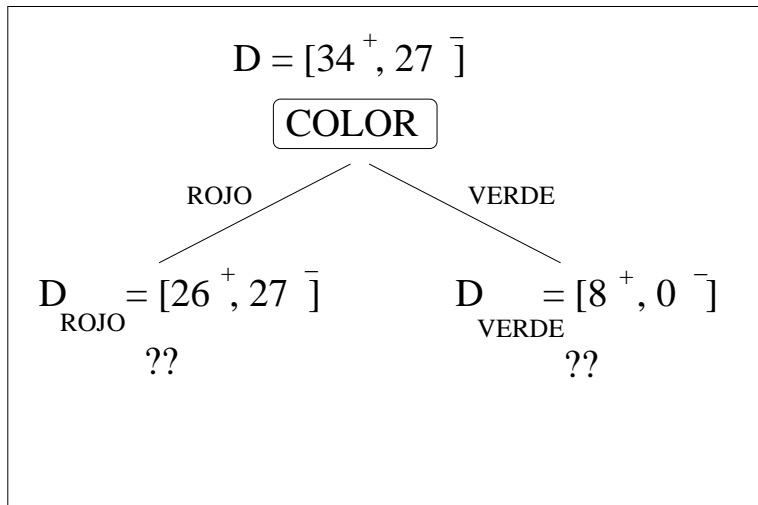
$D = [34^+, 27^-]$

COLOR

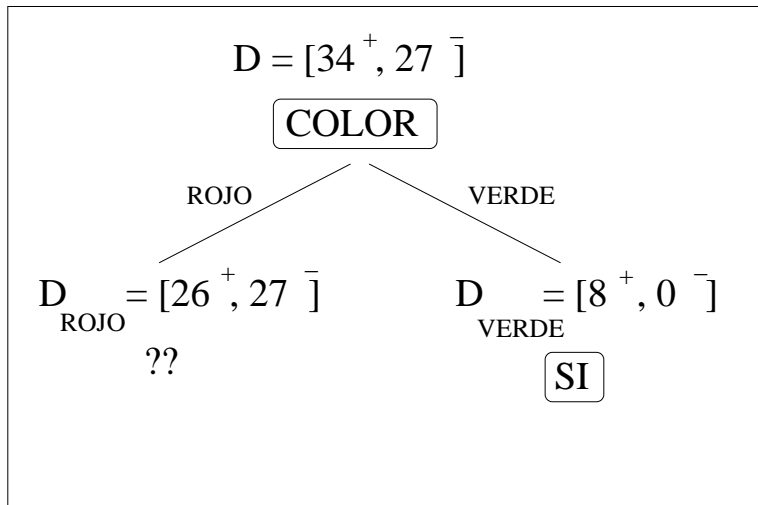
ROJO

VERDE

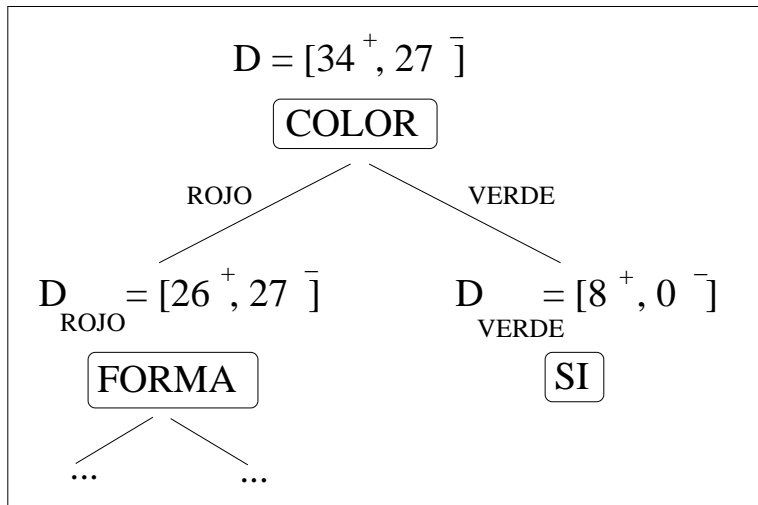
# Árboles de decisión



# Árboles de decisión

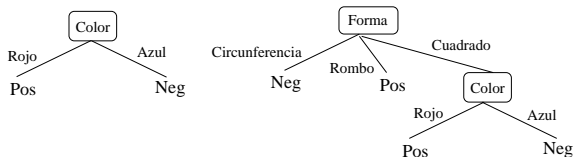


# Árboles de decisión



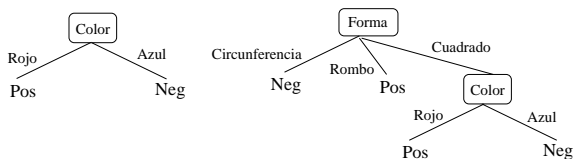
# Clasificadores

Color	Forma	Clasificación
Rojo	Cuadrado	Pos
Rojo	Rombo	Pos
Azul	Circunferencia	Neg
Azul	Cuadrado	Neg



# Clasificadores

Color	Forma	Clasificación
Rojo	Cuadrado	Pos
Rojo	Rombo	Pos
Azul	Circunferencia	Neg
Azul	Cuadrado	Neg



$\langle \text{Azul}, \text{Rombo} \rangle$  ???

# La navaja de Occam

Guillermo de Occam (1288-1349)

## *Lex parsimoniae*

*Entia non sunt multiplicanda praeter necessitatem* (No ha de presumirse la existencia de más cosas que las absolutamente necesarias)



Guillermo de Occam

## La navaja de Occam

En igualdad de condiciones la solución más sencilla es probablemente la correcta



# Algoritmo ID3

## Algoritmo ID3

ID3(Ejemplos, Atributo-objetivo, Atributos)

1. Si todos los Ejemplos son positivos, devolver un nodo etiquetado con +
2. Si todos los Ejemplos son negativos, devolver un nodo etiquetado con -
3. Si Atributos está vacío, devolver un nodo etiquetado con el valor más frecuente de Atributo-objetivo en Ejemplos.
4. En otro caso:
  - 4.1. Sea A el atributo de Atributos que MEJOR clasifica Ejemplos
  - 4.2. Crear Árbol, con un nodo etiquetado con A.
  - 4.3. Para cada posible valor v de A, hacer:
    - \* Añadir un arco a Árbol, etiquetado con v.
    - \* Sea Ejemplos(v) el subconjunto de Ejemplos con valor del atributo A igual a v.
    - \* Si Ejemplos(v) es vacío:
      - Entonces colocar debajo del arco anterior un nodo etiquetado con el valor más frecuente de Atributo-objetivo en Ejemplos.
      - Si no, colocar debajo del arco anterior el subárbol ID3(Ejemplos(v), Atributo-objetivo, Atributos-{A}).
  - 4.4 Devolver Árbol

## ¿Cómo saber qué atributo clasifica mejor?

- Entropía de un conjunto de ejemplos  $D$  (resp. de una clasificación):

$$Ent(D) = -\frac{|P|}{|D|} \cdot \log_2 \frac{|P|}{|D|} - \frac{|N|}{|D|} \cdot \log_2 \frac{|N|}{|D|}$$

donde  $P$  y  $N$  son, respectivamente, los subconjuntos de ejemplos positivos y negativos de  $D$

- Notación:  $Ent([p^+, n^-])$ , donde  $p = |P|$  y  $n = |N|$
- Intuición:
  - Mide la ausencia de “homogeneidad” de la clasificación
  - Teoría de la Información: cantidad media de información (en bits) necesaria para codificar la clasificación de un ejemplo
- Ejemplos:
  - $Ent([9^+, 5^-]) = -\frac{9}{14} \cdot \log_2 \frac{9}{14} - \frac{5}{14} \cdot \log_2 \frac{5}{14} = 0,94$
  - $Ent([k^+, k^-]) = 1$  (ausencia total de homogeneidad)
  - $Ent([p^+, 0^-]) = Ent([0^+, n^-]) = 0$  (homogeneidad total)

## Ganancia de información

- Preferimos nodos con menos entropía (árboles pequeños)
- Entropía esperada después de usar un atributo  $A$  en el árbol:

$$\sum_{v \in \text{Valores}(A)} \frac{|D_v|}{|D|} \cdot \text{Ent}(D_v)$$

donde  $D_v$  es el subconjunto de ejemplos de  $D$  con valor del atributo  $A$  igual a  $v$

- Ganancia de información esperada después de usar un atributo  $A$ :

$$\text{Ganancia}(D, A) = \text{Ent}(D) - \sum_{v \in \text{Valores}(A)} \frac{|D_v|}{|D|} \cdot \text{Ent}(D_v)$$

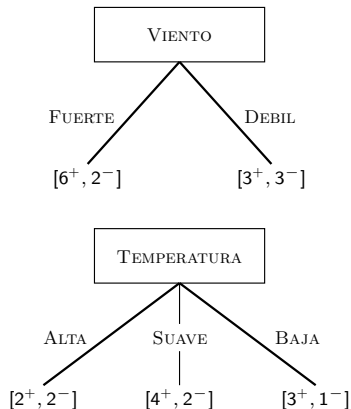
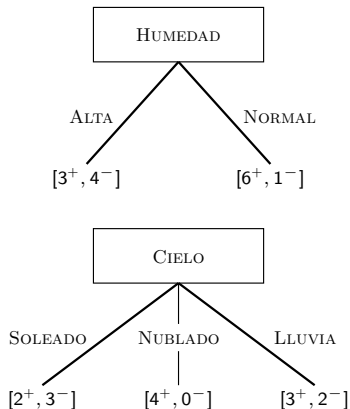
- En el algoritmo ID3, en cada nodo usamos el atributo con mayor ganancia de información (considerando los ejemplos correspondientes al nodo)

# Algoritmo ID3 (ejemplo 1)

- Conjunto de entrenamiento:

EJ.	CIELO	TEMPERATURA	HUMEDAD	VIENTO	JUGAR TENIS
$D_1$	SOLEADO	ALTA	ALTA	DÉBIL	-
$D_2$	SOLEADO	ALTA	ALTA	FUERTE	-
$D_3$	NUBLADO	ALTA	ALTA	DÉBIL	+
$D_4$	LLUVIA	SUAVE	ALTA	DÉBIL	+
$D_5$	LLUVIA	BAJA	NORMAL	DÉBIL	+
$D_6$	LLUVIA	BAJA	NORMAL	FUERTE	-
$D_7$	NUBLADO	BAJA	NORMAL	FUERTE	+
$D_8$	SOLEADO	SUAVE	ALTA	DÉBIL	-
$D_9$	SOLEADO	BAJA	NORMAL	DÉBIL	+
$D_{10}$	LLUVIA	SUAVE	NORMAL	DÉBIL	+
$D_{11}$	SOLEADO	SUAVE	NORMAL	FUERTE	+
$D_{12}$	NUBLADO	SUAVE	ALTA	FUERTE	+
$D_{13}$	NUBLADO	ALTA	NORMAL	DÉBIL	+
$D_{14}$	LLUVIA	SUAVE	ALTA	FUERTE	-

# Algoritmo ID3 (ejemplo 1)

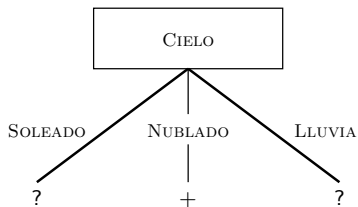


## Algoritmo ID3 (ejemplo 1)

- Entropía inicial:  $Ent([9^+, 5^-]) = 0,94$
- Selección del atributo para el nodo raíz:
  - $Ganancia(D, HUMEDAD) =$   
 $0,94 - \frac{7}{14} \cdot Ent([3^+, 4^-]) - \frac{7}{14} \cdot Ent([6^+, 1^-]) = 0,151$
  - $Ganancia(D, VIENTO) =$   
 $0,94 - \frac{8}{14} \cdot Ent([6^+, 2^-]) - \frac{6}{14} \cdot Ent([3^+, 3^-]) = 0,048$
  - $Ganancia(D, CIELO) =$   
 $0,94 - \frac{5}{14} \cdot Ent([2^+, 3^-]) - \frac{4}{14} \cdot Ent([4^+, 0^-])$   
 $- \frac{5}{14} \cdot Ent([3^+, 2^-]) = 0,246$  (mejor atributo)
  - $Ganancia(D, TEMPERATURA) =$   
 $0,94 - \frac{4}{14} \cdot Ent([2^+, 2^-]) - \frac{6}{14} \cdot Ent([4^+, 2^-])$   
 $- \frac{4}{14} \cdot Ent([3^+, 1^-]) = 0,02$
- El atributo seleccionado es CIELO

# Algoritmo ID3 (ejemplo 1)

- Árbol parcialmente construido:



## Algoritmo ID3 (ejemplo 1)

- Selección del atributo para el nodo CIELO=SOLEADO
- $D_{\text{SOLEADO}} = \{D_1, D_2, D_8, D_9, D_{11}\}$  con entropía  $Ent([2^+, 3^-]) = 0,971$ 
  - $Ganancia(D_{\text{SOLEADO}}, \text{HUMEDAD}) = 0,971 - \frac{3}{5} \cdot 0 - \frac{2}{5} \cdot 0 = 0,971$  (mejor atributo)
  - $Ganancia(D_{\text{SOLEADO}}, \text{TEMPERATURA}) = 0,971 - \frac{2}{5} \cdot 0 - \frac{2}{5} \cdot 1 - \frac{1}{5} \cdot 0 = 0,570$
  - $Ganancia(D_{\text{SOLEADO}}, \text{VIENTO}) = 0,971 - \frac{2}{5} \cdot 1 - \frac{3}{5} \cdot 0,918 = 0,019$
- El atributo seleccionado es HUMEDAD

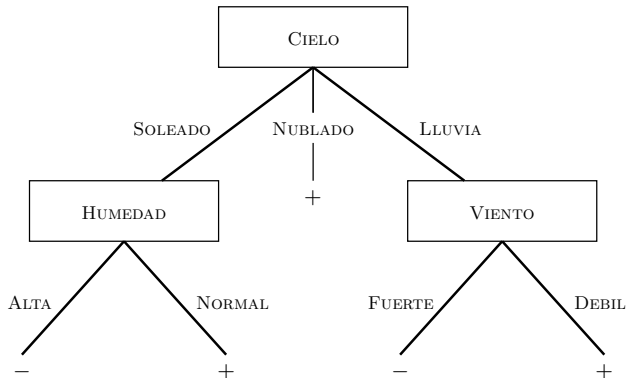


## Algoritmo ID3 (ejemplo 1)

- Selección del atributo para el nodo CIELO=LLUVIA:
- $D_{LLUVIA} = \{D_4, D_5, D_6, D_{10}, D_{14}\}$  con entropía  $Ent([3^+, 2^-]) = 0,971$ 
  - $Ganancia(D_{LLUVIA}, HUMEDAD) = 0,971 - \frac{2}{5} \cdot 1 - \frac{3}{5} \cdot 0,918 = 0,020$
  - $Ganancia(D_{LLUVIA}, TEMPERATURA) = 0,971 - \frac{3}{5} \cdot 0,918 - \frac{2}{5} \cdot 1 = 0,020$
  - $Ganancia(D_{LLUVIA}, VIENTO) = 0,971 - \frac{3}{5} \cdot 0 - \frac{2}{5} \cdot 0 = 0,971$  (mejor atributo)
- El atributo seleccionado es VIENTO

# Algoritmo ID3 (ejemplo 1)

- Árbol finalmente aprendido:



## Algoritmo ID3 (ejemplo 2)

- Conjunto de entrenamiento:

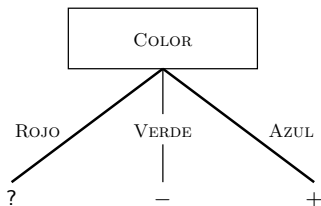
EJ.	COLOR	FORMA	TAMAÑO	CLASE
$O_1$	ROJO	CUADRADO	GRANDE	+
$O_2$	AZUL	CUADRADO	GRANDE	+
$O_3$	ROJO	REDONDO	PEQUEÑO	-
$O_4$	VERDE	CUADRADO	PEQUEÑO	-
$O_5$	ROJO	REDONDO	GRANDE	+
$O_6$	VERDE	CUADRADO	GRANDE	-

## Algoritmo ID3 (ejemplo 2)

- Entropía inicial en el ejemplo de los objetos,  $Ent([3^+, 3^-]) = 1$
- Selección del atributo para el nodo raíz:
  - $Ganancia(D, \text{COLOR}) = 1 - \frac{3}{6} \cdot Ent([2^+, 1^-]) - \frac{1}{6} \cdot Ent([1^+, 0^-]) - \frac{2}{6} \cdot Ent([0^+, 2^-]) = 0,543$
  - $Ganancia(D, \text{FORMA}) = 1 - \frac{4}{6} \cdot Ent([2^+, 2^-]) - \frac{2}{6} \cdot Ent([1^+, 1^-]) = 0$
  - $Ganancia(D, \text{TAMAÑO}) = 1 - \frac{4}{6} \cdot Ent([3^+, 1^-]) - \frac{2}{6} \cdot Ent([0^+, 2^-]) = 0,459$
- El atributo seleccionado es COLOR

## Algoritmo ID3 (ejemplo 2)

- Árbol parcialmente construido:

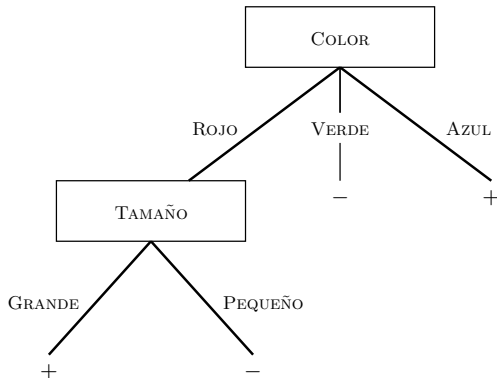


## Algoritmo ID3 (ejemplo 2)

- Selección del atributo para el nodo COLOR=ROJO:
- $D_{\text{ROJO}} = \{O_1, O_3, O_5\}$  con entropía  $Ent([2^+, 1^-]) = 0,914$ 
  - $Ganancia(D_{\text{ROJO}}, \text{FORMA}) =$   
 $0,914 - \frac{1}{3} \cdot Ent([1^+, 0^-]) - \frac{2}{3} \cdot Ent([1^+, 1^-]) = 0,247$
  - $Ganancia(D_{\text{ROJO}}, \text{TAMAÑO}) =$   
 $0,914 - \frac{2}{3} \cdot Ent([2^+, 0^-]) - \frac{1}{3} \cdot Ent([0^+, 1^-]) = 0,914$
- El atributo seleccionado es TAMAÑO

## Algoritmo ID3 (ejemplo 2)

- Árbol finalmente aprendido:



# Algunas cuestiones prácticas a resolver en aprendizaje automático

- Validar la hipótesis aprendida
  - ¿Podemos *cuantificar* la bondad de lo aprendido respecto de la explicación real?
- Sobreajuste
  - ¿Se ajusta *demasiado* lo aprendido al conjunto de entrenamiento?

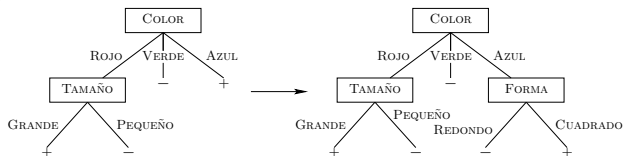


# Medida del rendimiento del aprendizaje

- Conjuntos de entrenamiento y prueba (*test*)
  - Aprender con el conjunto de entrenamiento
  - Medir el rendimiento en el conjunto de prueba:
    - proporción de ejemplos bien clasificados en el conjunto de prueba
- Repetición de este proceso
  - Curva de aprendizaje
  - Estratificación: cada clase correctamente representada en el entrenamiento y en la prueba
- Si no tenemos suficientes ejemplos como para apartar un conjunto de prueba: validación cruzada
  - Dividir en  $k$  partes, y hacer  $k$  aprendizajes, cada uno de ellos tomando como prueba una de las partes y entrenamiento el resto. Finalmente hacer la media de los rendimientos.
  - En la práctica: validación cruzada, con  $k = 10$  y estratificación

## Sobreajuste y ruido

- Una hipótesis  $h \in H$  *sobreajusta* los ejemplos de entrenamiento si existe  $h' \in H$  que se ajusta peor que  $h$  a los ejemplos pero actúa mejor sobre la distribución completa de instancias.
- *Ruido*: ejemplos incorrectamente clasificados. Causa sobreajuste
- Ejemplo: supongamos que por error, se incluye el ejemplo  $\langle \text{AZUL}, \text{REDONDO}, \text{PEQUEÑO} \rangle$  como ejemplo negativo
- El árbol aprendido en este caso sería (sobrejustado a los datos):



# Sobreajuste y ruido

- Otras causas de sobreajuste:
  - Atributos que en los ejemplos presentan una aparente regularidad pero que no son relevantes en realidad
  - Conjuntos de entrenamiento pequeños
- Maneras de evitar el sobreajuste en árboles de decisión:
  - Parar el desarrollo del árbol antes de que se ajuste perfectamente a todos los datos
  - Podar el árbol *a posteriori*
- Poda *a posteriori*, dos aproximaciones:
  - Transformación a reglas, podado de las condiciones de las reglas
  - Realizar podas directamente en el árbol
  - Las podas se producen siempre que reduzcan el error sobre un conjunto de prueba

# Podado de árboles

## Algoritmo de poda para reducir el error

1. Dividir el conjunto de ejemplos en Entrenamiento y Prueba
2. Árbol=árbol obtenido por ID3 usando Entrenamiento
3. Continuar=True
4. Mientras Continuar:
  - \* Medida = proporción de ejemplos en Prueba correctamente clasificados por Árbol
  - \* Por cada nodo interior N de Árbol:
    - Podar temporalmente Árbol en el nodo N y sustituirlo por una hoja etiquetada con la clasificación mayoritaria en ese nodo
    - Medir la proporción de ejemplos correctamente clasificados en el conjunto de prueba.
  - \* Sea K el nodo cuya poda produce mejor rendimiento
  - \* Si este rendimiento es mejor que Medida, entonces  
Árbol = resultado de podar permanentemente Árbol en K
  - \* Si no, Continuar=Falso
5. Devolver Árbol

## Atributos con valores continuos

- Reemplazamos los atributos continuos por atributos booleanos que se crean dinámicamente, introduciendo umbrales  $C$ .
  - $A$  continuo
  - $A_{<C}$  booleano, toma el valor SI cuando el valor es menor que  $C$  y NO en otro caso.

## Atributos con valores continuos

Temperatura	50	52	60	68	70	78	84
Clase	+	+	-	-	-	+	+

- Los candidatos a umbral  $C$  son los valores adyacentes con distinta clase
  - $56 = (52 + 60)/2$  y  $74 = (70 + 78)/2$ .
- Seleccionamos el umbral con máxima ganancia
- El nuevo atributo  $A_{<C}$  compite con los restantes.
- El proceso se realiza a cada paso eligiendo el mejor umbral en el conjunto de entrenamiento.

# Otras cuestiones prácticas del algoritmo ID3

- Extensiones del algoritmo:
  - Otras medidas para seleccionar atributos
  - Otras estimaciones de error
  - Atributos sin valores
  - Atributos con coste
- Algoritmos C4.5 y C5.0 (Quinlan)

# Aplicaciones

Journal of Theoretical Biology 357 (2014) 21–25



Contents lists available at ScienceDirect

Journal of Theoretical Biology

journal homepage: [www.elsevier.com/locate/yjtbi](http://www.elsevier.com/locate/yjtbi)



## Decision trees for the analysis of genes involved in Alzheimer's disease pathology



Sonia L. Mestizo Gutiérrez<sup>a</sup>, Marisol Herrera Rivero<sup>b</sup>, Nicandro Cruz Ramírez<sup>c</sup>,  
Elena Hernández<sup>d</sup>, Gonzalo E. Aranda-Abreu<sup>d,\*</sup>

<sup>a</sup> Doctorado en Investigaciones Cerebrales, Universidad Veracruzana, Av. Luis Castelazo Ayala S/N, Xalapa, Veracruz 91190, Mexico

<sup>b</sup> Doctorado en Ciencias Biomédicas, Universidad Veracruzana, Av. Luis Castelazo Ayala S/N, Xalapa, Veracruz, Mexico

<sup>c</sup> Departamento de Inteligencia Artificial, Universidad Veracruzana, Sebastián Camacho 5, Centro, Xalapa, Veracruz 91000, Mexico

<sup>d</sup> Centro de Investigaciones Cerebrales, Cuerpo Académico de Neuroquímica, Universidad Veracruzana, Av. Luis Castelazo Ayala S/N, Xalapa, Veracruz, Mexico

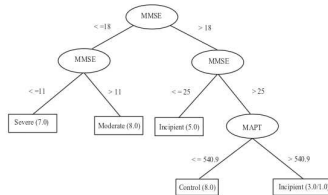
Decision trees for the analysis of genes involved in Alzheimer's disease pathology Sonia L. Mestizo Gutiérrez, Marisol Herrera Rivero, Nicandro Cruz Ramírez, Elena Hernández, Gonzalo E. Aranda-Abreu. *Journal of Theoretical Biology*. Volume 357, 21 September 2014, Pages 21?25.



# Aplicaciones

## 2.3.3. Decision trees

We tested two disease classifiers with different sets of genes, using the J48 algorithm in the Waikato Environment for Knowledge Analysis (WEKA) (Hall et al., 2009; Waikato Environment for Knowledge Analysis, 2013) an implementation of the C4.5 algorithm, with 10-fold cross-validation. The first set of genes tested consisted of those found to significantly differ by ANOVA, and the second set consisted of the AD-related genes we added to the list: *APP*, *APOE*, *BACE1*, *NCSTN*, *PSEN1*, *PSEN2* and *MAPT*. Expression data as well as MMSE and NFT scores were included.



Decision trees for the analysis of genes involved in Alzheimer's disease pathology Sonia L. Mestizo Gutiérrez,  
Marisol Herrera Rivero, Nicandro Cruz Ramírez, Elena Hernández, Gonzalo E. Aranda-Abreu. *Journal of Theoretical Biology*. Volume 357, 21 September, 2014, pages 21-25.

# Aplicaciones

## Diagnosis of gastric cancer using decision tree classification of mass spectral data

Yahui Su,<sup>1\*</sup> Jing Shen,<sup>1\*</sup> Honggang Qian,<sup>1</sup> Huachong Ma,<sup>2</sup> Jiafu Ji,<sup>1</sup> Hong Ma,<sup>1</sup> Longhua Ma,<sup>3</sup> Weihua Zhang,<sup>3</sup> Ling Meng,<sup>1</sup> Zhenfu Li,<sup>1</sup> Jian Wu,<sup>1</sup> Genglin Jin,<sup>1</sup> Jianzhi Zhang<sup>1</sup> and Chengchao Shou<sup>1,4</sup>

<sup>1</sup>Peking University School of Oncology and Beijing Cancer Hospital and Institute, Haidian, Beijing 100036; <sup>2</sup>Beijing Chaoyang Hospital, Chaoyang, Beijing 100020, China; and <sup>3</sup>Ciphergen Biosystems, Fremont, CA 94555, USA

(Received July 15, 2006/Revised September 2, 2006/Accepted September 4, 2006/Online publication October 19, 2006)

Cancer Science, Volume 98, Issue 1, pages 37-43, January 2007

Computers in Biology and Medicine 42 (2012) 195–204



Contents lists available at SciVerse ScienceDirect

Computers in Biology and Medicine

journal homepage: [www.elsevier.com/locate/cbm](http://www.elsevier.com/locate/cbm)



## Using partial decision trees to predict Parkinson's symptoms: A new approach for diagnosis and therapy in patients suffering from Parkinson's disease

Themis P. Exarchos<sup>a</sup>, Alexandros T. Tzallas<sup>a</sup>, Dina Baga<sup>a</sup>, Dimitra Chaloglou<sup>b</sup>, Dimitrios I. Fotiadis<sup>a,\*</sup>, Sofia Tsouli<sup>c</sup>, Maria Diakou<sup>c</sup>, Spyros Konitsiotis<sup>c</sup>

<sup>a</sup> Unit of Medical Technology and Intelligent Information System, Department of Materials Science and Engineering, University of Ioannina, GR 45110, Ioannina, Greece

<sup>b</sup> ANCO S.A., Athens GR 11742, Greece

<sup>c</sup> Dept. of Neurology, Medical School, University of Ioannina, GR 45110, Ioannina, Greece

# Sección 3

## Sección 3

### Aprendizaje de reglas

# Cambiando la salida: reglas proposicionales

- Reglas de clasificación:
  - R1: Si  $\text{CIELO}=\text{SOLEADO} \wedge \text{HUMEDAD}=\text{ALTA}$   
Entonces  $\text{JUGAR} \text{TENIS} = -$
  - R2: Si  $\text{ASTIGMATISMO} = + \wedge \text{LAGRIMA} = \text{NORMAL}$   
Entonces  $\text{LENTE} = \text{RIGIDA}$
- Ventaja de usar el formalismo de reglas
  - Claridad
  - Modularidad
  - Expresividad: pueden representar cualquier conjunto de instancias
  - Métodos generalizables a primer orden de manera natural
  - Formalismo usado en sistemas basados en el conocimiento
- Reglas y árboles de decisión
  - Fácil traducción de árbol a reglas, pero no a la inversa

## Aprendizaje de reglas

- Objetivo: aprender un conjunto de reglas consistente con los ejemplos
  - Una regla *cubre* un ejemplo si el ejemplo satisface las condiciones
  - Lo cubre *correctamente* si además el valor del atributo en la conclusión de la regla coincide con el valor que el ejemplo toma en ese atributo
- Una medida del ajuste de una regla  $R$  a un conjunto de ejemplos  $D$ :
  - *Frecuencia relativa*:  $p/t$  (donde  $t$  = ejemplos cubiertos por  $R$  en  $D$ ,  $p$  = ejemplos correctamente cubiertos). Notación:  $FR(R, D)$
- Algoritmos de aprendizaje de reglas:
  - ID3 + traducción a reglas
  - Cobertura
  - Algoritmos genéticos

# Un conjunto de entrenamiento

EJ.	EDAD	DIGNOSTICO	ASTIGMATISMO	LAGRIMA	LENTE
$E_1$	JOVEN	MIOPE	-	REDUCIDA	NINGUNA
$E_2$	JOVEN	MIOPE	-	NORMAL	BLANDA
$E_3$	JOVEN	MIOPE	+	REDUCIDA	NINGUNA
$E_4$	JOVEN	MIOPE	+	NORMAL	RÍGIDA
$E_5$	JOVEN	HIPERMÉTROPE	-	REDUCIDA	NINGUNA
$E_6$	JOVEN	HIPERMÉTROPE	-	NORMAL	BLANDA
$E_7$	JOVEN	HIPERMÉTROPE	+	REDUCIDA	NINGUNA
$E_8$	JOVEN	HIPERMÉTROPE	+	NORMAL	RÍGIDA
$E_9$	PREPRESBICIA	MIOPE	-	REDUCIDA	NINGUNA
$E_{10}$	PREPRESBICIA	MIOPE	-	NORMAL	BLANDA
$E_{11}$	PREPRESBICIA	MIOPE	+	REDUCIDA	NINGUNA
$E_{12}$	PREPRESBICIA	MIOPE	+	NORMAL	RÍGIDA
$E_{13}$	PREPRESBICIA	HIPERMÉTROPE	-	REDUCIDA	NINGUNA
$E_{14}$	PREPRESBICIA	HIPERMÉTROPE	-	NORMAL	BLANDA
$E_{15}$	PREPRESBICIA	HIPERMÉTROPE	+	REDUCIDA	NINGUNA
$E_{16}$	PREPRESBICIA	HIPERMÉTROPE	+	NORMAL	NINGUNA

## Un conjunto de entrenamiento

EJ.	EDAD	DIGNOSTICO	ASTIGMATISMO	LAGRIMA	LENTE
$E_{17}$	PRESBICIA	MIOPE	-	REDUCIDA	NINGUNA
$E_{18}$	PRESBICIA	MIOPE	-	NORMAL	NINGUNA
$E_{19}$	PRESBICIA	MIOPE	+	REDUCIDA	NINGUNA
$E_{20}$	PRESBICIA	MIOPE	+	NORMAL	RÍGIDA
$E_{21}$	PRESBICIA	HIPERMÉTROPE	-	REDUCIDA	NINGUNA
$E_{22}$	PRESBICIA	HIPERMÉTROPE	-	NORMAL	BLANDA
$E_{23}$	PRESBICIA	HIPERMÉTROPE	+	REDUCIDA	NINGUNA
$E_{24}$	PRESBICIA	HIPERMÉTROPE	+	NORMAL	NINGUNA

- R2: Si  $\text{ASTIGMATISMO} = + \wedge \text{LAGRIMA} = \text{NORMAL}$   
Entonces  $\text{LENTE} = \text{RÍGIDA}$
- R2 cubre  $E_4, E_8, E_{12}, E_{16}, E_{20}$  y  $E_{24}$ , de los cuales cubre correctamente  $E_4, E_8, E_{12}$  y  $E_{20}$

## Aprendiendo reglas que cubren ejemplos

- Aprender una regla para clasificar  $\text{LENTE}=\text{RIGIDA}$

- Si ?

Entonces  $\text{LENTE}=\text{RIGIDA}$

- Alternativas para ?, y frecuencia relativa de la regla resultante

EDAD=JOVEN	2/8
EDAD=PREPRESBICIA	1/8
EDAD=PRESBICIA	1/8
DIAGNOSTICO=MIOPÍA	3/12
DIAGNOSTICO=HIPERMETROPÍA	1/12
ASTIGMATISMO=−	0/12
ASTIGMATISMO=+	4/12 *
LÁGRIMA=REDUCIDA	0/12
LÁGRIMA=NORMAL	4/12 *

- Regla parcialmente aprendida

- Si  $\text{ASTIGMATISMO}=+$

Entonces  $\text{LENTE}=\text{RIGIDA}$



## Aprendiendo reglas que cubren ejemplos

- Continuamos para excluir ejemplos cubiertos incorrectamente
  - Si  $\text{ASTIGMATISMO} = + \wedge ?$   
Entonces  $\text{LENTE} = \text{RIGIDA}$
  - Alternativas para  $?$ , y frecuencia relativa de la regla resultante

EDAD=JOVEN	2/4
EDAD=PREPRESBICIA	1/4
EDAD=PRESBICIA	1/4
DIAGNOSTICO=MIOPÍA	3/6
DIAGNOSTICO=HIPERMETROPÍA	1/6
LÁGRIMA=REDUCIDA	0/6
LÁGRIMA=NORMAL	4/6 *

- Regla parcialmente aprendida
  - Si  $\text{ASTIGMATISMO} = + \wedge \text{LAGRIMA} = \text{NORMAL}$   
Entonces  $\text{LENTE} = \text{RIGIDA}$

## Aprendiendo reglas que cubren ejemplos

- Continuamos para excluir ejemplos cubiertos incorrectamente
  - Si  $\text{ASTIGMATISMO} = + \wedge \text{LAGRIMA} = \text{NORMAL} \wedge ?$   
Entonces  $\text{LENTE} = \text{RIGIDA}$
  - Alternativas para  $?$ , y frecuencia relativa de la regla resultante
 

EDAD=JOVEN	2/2 *
EDAD=PREPRESBICIA	1/2
EDAD=PRESBICIA	1/2
DIAGNOSTICO=MIOPIA	3/3 *
DIAGNOSTICO=HIPERMETROPIA	1/3
- Regla finalmente aprendida (no cubre incorrectamente ningún ejemplo)
  - Si  $\text{ASTIGMATISMO} = + \wedge \text{LAGRIMA} = \text{NORMAL} \wedge$   
 $\text{DIAGNOSTICO} = \text{MIOPIA}$   
 Entonces  $\text{LENTE} = \text{RIGIDA}$

## Aprendiendo reglas que cubren ejemplos

- Queda un ejemplo con  $\text{LENTE}=\text{RIGIDA}$  no cubierto por R2
  - Comenzamos otra vez con “Si ? Entonces  $\text{LENTE}=\text{RIGIDA}$ ”, pero ahora con  $D' = D \setminus \{E_4, E_{12}, E_{20}\}$
- Reglas finalmente aprendidas para  $\text{LENTE}=\text{RIGIDA}$ :
  - R1: Si  $\text{ASTIGMATISMO} = + \wedge \text{LAGRIMA} = \text{NORMAL} \wedge \text{DIAGNOSTICO} = \text{MIOPIA}$   
Entonces  $\text{LENTE}=\text{RIGIDA}$
  - R2: Si  $\text{EDAD} = \text{JOVEN} \wedge \text{ASTIGMATISMO} = + \wedge \text{LAGRIMA} = \text{NORMAL}$   
Entonces  $\text{LENTE}=\text{RIGIDA}$
  - Cubren correctamente los 4 ejemplos de  $\text{LENTE}=\text{RIGIDA}$  (y se solapan)
- Ahora se podría continuar para aprender reglas que clasifiquen:
  - $\text{LENTE}=\text{BLANDA}$
  - $\text{LENTE}=\text{NINGUNA}$

# Algoritmo de aprendizaje de reglas por cobertura

## Algoritmo de aprendizaje por cobertura

**Aprendizaje-por-Cobertura( $D, \text{Atributo}, v$ )**

1. Hacer Reglas-aprendidas igual a vacío
2. Hacer  $E$  igual a  $D$
3. Mientras  $E$  contenga ejemplos cuyo valor de Atributo es  $v$ , hacer:
  - 3.1 Crear una regla  $R$  sin condiciones y conclusión  $\text{Atributo}=v$
  - 3.2 Mientras que haya en  $E$  ejemplos cubiertos por  $R$  incorrectamente y queden atributos que usar, hacer:
    - 3.2.1 Elegir la **MEJOR** condición  $A=w$  para añadir a  $R$ , donde  $A$  es un atributo que no aparece en  $R$  y  $w$  es un valor de los posibles que puede tomar  $A$
    - 3.2.2 Actualizar  $R$  añadiendo la condición  $A=w$  a  $R$
  - 3.3 Incluir  $R$  en Reglas-aprendidas
  - 3.4 Actualizar  $E$  quitando los ejemplos cubiertos por  $R$
4. Devolver Reglas-Aprendidas

- Algoritmo para aprender un conjunto de reglas (a partir de  $D$ )
  - Reglas para predecir situaciones en las que un Atributo dado toma un valor  $v$

## Control en el algoritmo de cobertura

- Bucle externo:
  - Añade reglas (la hipótesis se *generaliza*)
  - Cada regla añadida cubre algunos ejemplos correctamente
  - Elimina en cada vuelta los ejemplos cubiertos por la regla añadida
  - Y se añaden reglas mientras queden ejemplos sin cubrir
- Bucle interno:
  - Añade condiciones a la regla (la regla se *especializa*)
  - Cada nueva condición excluye ejemplos cubiertos incorrectamente
  - Y esto se hace mientras haya ejemplos incorrectamente cubiertos
- Cobertura frente a ID3
  - Aprende una regla cada vez, ID3 lo hace simultáneamente
  - ID3: elecciones de atributos
  - Cobertura: elecciones de parejas atributo-valor

## Algoritmo de cobertura (propiedades)

- Diferentes criterios para elegir la mejor condición en cada vuelta del bucle interno:
  - Se añade la condición que produzca la regla con *mayor frecuencia relativa* (como en el ejemplo)
  - Se añade la que produzca *mayor ganancia de información*:
$$p \cdot (\log_2 \frac{p'}{t'} - \log_2 \frac{p}{t})$$
donde  $p'/t'$  es la frecuencia relativa *después* de añadir la condición y  $p/t$  es la frecuencia relativa *antes* de añadir la condición
- Las reglas aprendidas por el algoritmo de cobertura se ajustan *perfectamente* al conjunto de entrenamiento (*peligro de sobreajuste*)
  - Podado de las reglas *a posteriori*
  - Eliminar progresivamente condiciones hasta que no se produzca *mejora*
  - Criterio probabilístico para decidir la mejora

# Sección 4

## Sección 4

### Aprendizaje basado en instancias: kNN

## Clasificación mediante vecino más cercano

- Una técnica alternativa a construir el modelo probabilístico es calcular la clasificación directamente a partir de los ejemplos (*aprendizaje basado en instancias*)
- Idea: obtener la clasificación de un nuevo ejemplo a partir de las categorías de los ejemplos más “ceranos”.
  - Debemos manejar, por tanto, una noción de “distancia” entre ejemplos.
  - En la mayoría de los casos, los ejemplos serán elementos de  $R^n$  y la distancia, la euclídea.
  - Pero se podría usar otra noción de distancia
- Ejemplo de aplicación: clasificación de documentos



# El algoritmo $k$ -NN

- El algoritmo  $k$ -NN (de “ $k$  nearest neighbors”):
  - Dado un conjunto de entrenamiento (vectores numéricos con una categoría asignada) y un ejemplo nuevo
  - Devolver la categoría mayoritaria en los  $k$  ejemplos del conjunto de entrenamiento más cercanos al ejemplo que se quiere clasificar

## Distancias para $k$ -NN

- Posibles distancias usadas para definir la “cercanía”:
  - Euclídea:  $d_e(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$
  - Manhattan:  $d_m(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n |x_i - y_i|$
  - Hamming: número de componentes en las que se difiere.
- La euclídea se usa cuando cada dimensión mide propiedades similares y la Manhattan en caso contrario; la distancia Hamming se puede usar aún cuando los vectores no sean numéricos.
- Normalización: cuando no todas las dimensiones son del mismo orden de magnitud, se normalizan las componentes (restando la media y dividiendo por la desviación típica)

## Algunas observaciones sobre $k$ -NN

- Elección de  $k$ :
  - Usualmente, basándonos en algún conocimiento específico sobre el problema de clasificación
  - También como resultado de pruebas en conjuntos más pequeños (conjuntos de validación)
  - Si la clasificación es binaria, preferiblemente impar, para intentar evitar empates ( $k=5$ , por ejemplo)

## Variantes del algoritmo $k$ -NN

### Algoritmo $k$ -NN con pesos

- En esta variante se consideran los  $k$  vecinos más cercanos  $\{a_1, \dots, a_k\}$  al objeto  $x$  que queremos clasificar.
- A cada uno de los  $k$  vecinos más cercanos se le asigna un peso  $w_i$ ; se les asigna el peso

$$w_i = \frac{1}{\text{dist}(a_i, x)}$$

- Sumamos los pesos de cada una de las posibles clasificaciones.
- El valor asignado a  $x$  será el que obtenga un mayor peso.
- Así un ejemplo  $a_i$  cuenta más cuanto más cercano esté a  $x$ .

## Variantes del algoritmo $k$ -NN

### NCC (Nearest Centroid Classifier)

- Dado un conjunto entrenamiento formado por puntos de  $R^n$  junto con su clasificación y un nuevo punto  $x$ , NCC asigna al nuevo punto la clasificación de la clase cuyo centroide este mas cercano al punto.
- En otras palabras, para cada valor de clasificación calculamos el centroide de los puntos asociados y luego aplicamos  $k$ -NN sobre el conjunto de centroides con  $k = 1$ .

## Variantes del algoritmo $k$ -NN

### $k$ -NN con rechazo

- En esta variante, además del conjunto de entrenamiento  $D$ , el valor  $k$  y el objeto a clasificar  $x$ , necesitamos un umbral  $\mu$ .
- El algoritmo toma los  $k$  puntos del conjunto de entrenamiento mas cercanos a  $x$  y devuelve el valor de clasificacion mayoritario en esos  $k$  puntos **solo si el numero de puntos con esa clasificacion supera el umbral**.
- Por ejemplo, si consideramos  $k=12$  y de los 12 puntos mas cercanos hay 8 con clasificacion  $A$ , 3 con clasificacion  $B$  y 1 con clasificacion  $C$ , y el umbral es  $\mu = 7$ , entonces el valor de clasificacion devuelto es  $A$ . En cambio, si el umbral es  $\mu = 9$ , el resultado debe ser **NO CLASIFICADO**.

## Sección 5

### Sección 5 *Clustering*

# Clustering

*... in cluster analysis a group of objects is split up into a number of more or less homogeneous subgroups on the basis of an often subjectively chosen measure of similarity (i.e., chosen subjectively based on its ability to create “interesting” clusters), such that the similarity between objects within a subgroup is larger than the similarity between objects belonging to different subgroups.*  
(Backer & Jain, 1981)



# Clustering

- Se trata de dividir un conjunto de datos de entrada en subconjuntos (*clusters*), de tal manera que los elementos de cada subconjunto compartan cierto patrón o características a priori desconocidas
- Aprendizaje *no supervisado*: no tenemos información sobre qué cluster corresponde a cada dato.
- Aplicaciones de clustering:
  - Minería de datos
  - Procesamiento de imágenes digitales
  - Bioinformática
- Dos tipos:
  - Clustering de partición estricta
  - Clustering jerárquico

# Clustering

## Clustering de partición estricta

Dado un conjunto de ejemplos  $D = \{\vec{x}_1, \dots, \vec{x}_j, \dots, \vec{x}_N\}$  con  $\vec{x}_j = (x_{j1}, \dots, x_{jd}) \in \mathbb{R}^d$ , el clustering de partición estricta *Hard partitional clustering* busca una partición de  $D$  en  $K$  clusters,  $\mathcal{P} = \{C_1, \dots, C_K\}$  con  $K \leq N$  tal que

- $C_i \neq \emptyset$  para  $i \in \{1, \dots, K\}$
- $\bigcup_{i=1}^K C_i = D$
- $C_i \cap C_j = \emptyset$  para todo  $i, j \in \{1, \dots, K\}$  con  $i \neq j$ .

# Clustering

## Clustering jerárquico

Dado un conjunto de ejemplos  $D = \{\vec{x}_1, \dots, \vec{x}_j, \dots, \vec{x}_N\}$  con  $\vec{x}_j = (x_{j1}, \dots, x_{jd}) \in \mathbb{R}^n$ , el clustering jerárquico *Hierarchical clustering* busca construir un conjunto de particiones anidadas de  $D$  con estructura de árbol,  $\mathcal{H} = \{P_1, \dots, P_Q\}$  con  $Q \leq N$  (donde cada  $P_j$  es una partición) tal que si  $C_i \in P_m$  y  $C_j \in P_l$  con  $m > l$  entonces  $C_i \subset C_j$  o  $C_i \cap C_j = \emptyset$ , para todo  $i, j, m, l \in \{1, \dots, Q\}$ ,  $i \neq j$ . Cada  $P_i$  es una partición de  $D$ .

# Clustering

- Como hemos visto, el **clustering de partición estricta** divide al conjunto de datos en clusters que no tienen ninguna estructura interna.
- En cambio, el **clusteing jerárquico** da estructura interna a los clusters. De hecho, de manera recursiva, cada cluster está dividido en clusters internos, en una estructura anidada que va desde un cluster general conteniendo a todos los individuos, hasta clusters que contienen un único elemento.

## El concepto de distancia

- La idea básica del *clustering* consiste en agrupar las instancias según su *proximidad*, esto es, dos instancias pertenecerán al mismo cluster si están *próximas* y pertenecerán a clusters distintos si están *lejanas*.
- Distancias:
  - Euclídea:  $d_e(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$
  - Manhattan:  $d_m(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n |x_i - y_i|$
  - Hamming: número de componentes en las que se difiere.
- Problemas:
- Incompatibilidad de las unidades de medida: Metros, gramos, litros, ... (Normalización); Tipos de variables: Booleanas, nominales, numéricas, ...; Datos incompletos ...

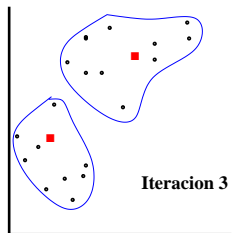
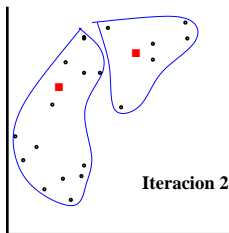
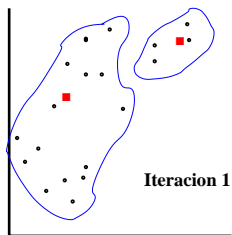
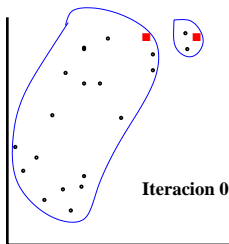
# Un algoritmo de partición estricta: $k$ -medias

- Entrada: un número  $k$  de clusters, un conjunto de datos  $\{x_i\}_{i=1}^N$  y una función de distancia
- Salida: un conjunto de  $k$  centros  $m_1, \dots, m_k$

## $k$ -medias( $k$ ,datos,distancia)

1. Inicializar  $m_i$  ( $i=1, \dots, k$ ) (aleatoriamente o con algún criterio heurístico)
2. REPETIR (hasta que los  $m_i$  no cambien):
  - 2.1 PARA  $j=1, \dots, N$ , HACER:  
Calcular el cluster correspondiente a  $x_j$ , escogiendo, de entre todos los  $m_i$ , el  $m_h$  tal que  $\text{distancia}(x_j, m_h)$  sea mínima
  - 2.2 PARA  $i=1, \dots, k$  HACER:  
Asignar a  $m_i$  la media aritmética de los datos asignados al cluster  $i$ -ésimo
3. Devolver  $m_1, \dots, m_k$

# Idea gráfica intuitiva en el algoritmo de $k$ -medias



## Ejemplo en el algoritmo $k$ -medias

- Datos sobre pesos de la población: 51, 43, 62, 64, 45, 42, 46, 45, 45, 62, 47, 52, 64, 51, 65, 48, 49, 46, 64, 51, 52, 62, 49, 48, 62, 43, 40, 48, 64, 51, 63, 43, 65, 66, 65, 46, 39, 62, 64, 52, 63, 64, 48, 64, 48, 51, 48, 64, 42, 48, 41
- El algoritmo, aplicado con  $k = 2$  y distancia euclídea, encuentra dos centros  $m_1 = 63,63$  y  $m_2 = 46,81$  en tres iteraciones
- 19 datos pertenecen al primer cluster y 32 al segundo cluster



## Diversas cuestiones sobre el algoritmo $k$ -medias

- Inicialización: aleatoria o con alguna técnica heurística (por ejemplo, partir los datos aleatoriamente en  $k$  clusters y empezar con los centros de esos clusters)
- En la práctica, los centros con los que se inicie el algoritmo tienen un gran impacto en la calidad de los resultados que se obtengan

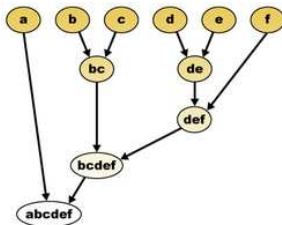
# Clustering jerárquico

Dependiendo del orden en que se cree la jerarquía de clusters, los algoritmos de clustering jerárquico se dividen en dos tipos:

- **Clustering jerárquico aglomerativo:** Partimos de clusters conteniendo un único ejemplo y vamos agrupando los clusters, obteniendo agrupamientos cada vez de mayor tamaño hasta obtener un cluster final con todos los individuos.
- **Clustering jerárquico divisor:** Empezamos con un cluster con todos los individuos y vamos realizando particiones de los clusters obtenidos hasta obtener clusters conteniendo un único ejemplo.

## Clustering jerárquico

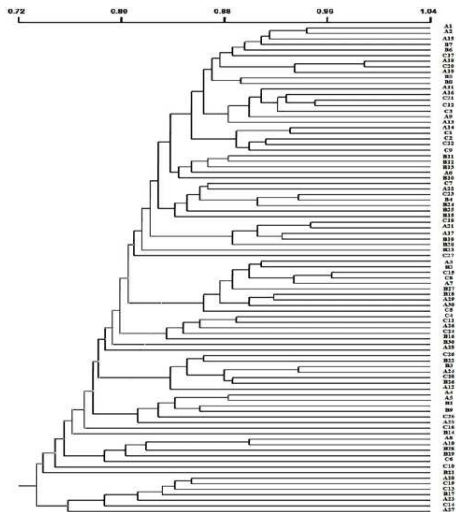
Con independencia de que se use clustering aglomerativo o divisor, los resultados suelen representarse con una estructura de árbol llamada *dendrodrama*



La raíz representa el conjunto completo y cada una de las hojas representa una instancia.

# Clustering jerárquico

## Ejemplo de dendrograma



# Clustering jerárquico

## Complejidad

Consideremos un problema de clustering jerárquico con  $N$  puntos

- Los métodos de *Clustering jerárquico divisor* tienen que considerar  $2^{N-1} - 1$  posibles maneras de dividir el conjunto inicial en dos subconjuntos.
- Los métodos de *Clustering jerárquico aglomerativo* deben considerar la distancia entre pares de puntos de orden  $O(N^2)$ .
- Por tanto *Clustering jerárquico aglomerativo* se usan mucho más que los de tipo *divisor*.

## Clustering jerárquico aglomerativo

El esquema general de *Clustering jerárquico aglomerativo* es el siguiente:

1. Inicializamos el algoritmo con  $N$  clusters individuales. Calculamos la matriz de proximidad (basada en alguna definición de distancia) para los  $N$  clusters.
2. En la matriz de proximidad, buscamos la menor distancia entre clusters. Según definamos la *distancia entre clusters*, tendremos diferentes algoritmos. Combinamos en un único cluster aquellos que estén a distancia mínima.
3. Actualizamos la matriz de proximidad considerando los nuevos clusters.
4. Repetimos los pasos 2 y 3 hasta que quede un único cluster.



## Bibliografía

- Mitchell, T.M. *Machine Learning* (McGraw-Hill, 1997)
  - Caps. 3,6,8 y 10
- Russell, S. y Norvig, P. *Artificial Intelligence (A Modern Approach)* (3rd edition) (Prentice Hall, 2010)
  - Seccs. 18.1, 18.2, 18.3, 20.1 y 20.2
- Witten, I.H. y Frank, E. *Data mining* (Third edition) (Morgan Kaufmann Publishers, 2011)
  - Cap. 3, 4, 5 y 6.
- Alpaydin, E. *Introduction to Machine Learning* (third edition) (The MIT Press, 2014)
- Xu, R y Wunsch II, D.C. *Clustering* (IEEE Press, 2009)