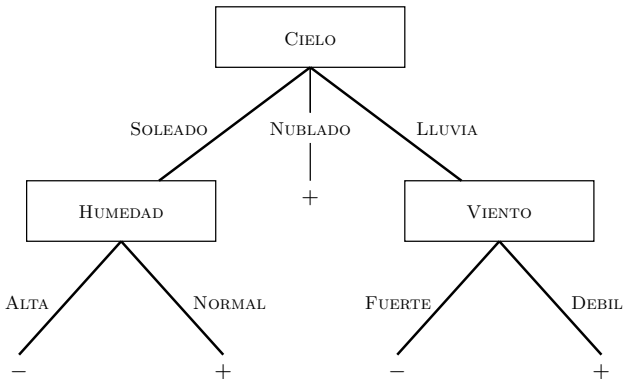


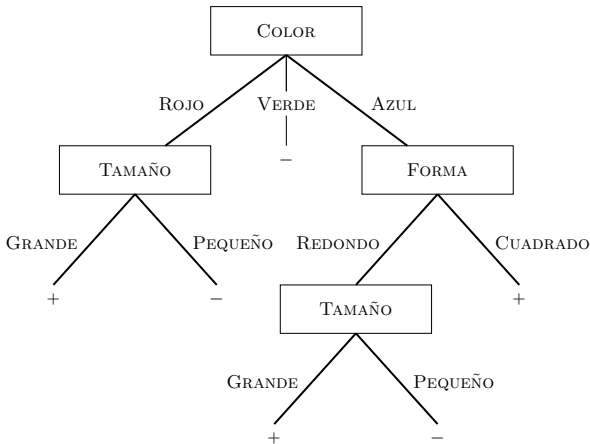
Árboles de decisión

- 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



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:

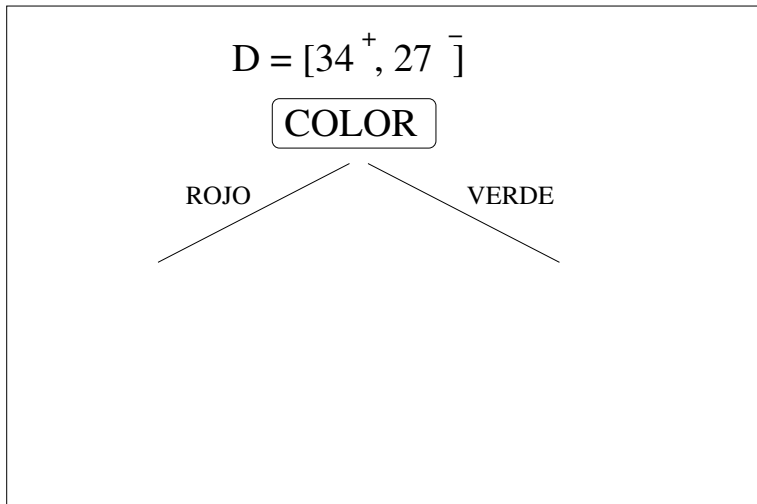
EJ.	CIELO	TEMPERATURA	HUMEDAD	VIENTO	JUGAR TENIS
D_1	SOLEADO	ALTA	ALTA	DÉBIL	-
D_2	SOLEADO	ALTA	ALTA	FUERTE	-
D_3	NUBLADO	ALTA	ALTA	DÉBIL	+
D_4	LLUVIA	SUAVE	ALTA	DÉBIL	+
...					

Árboles de decisión

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

COLOR

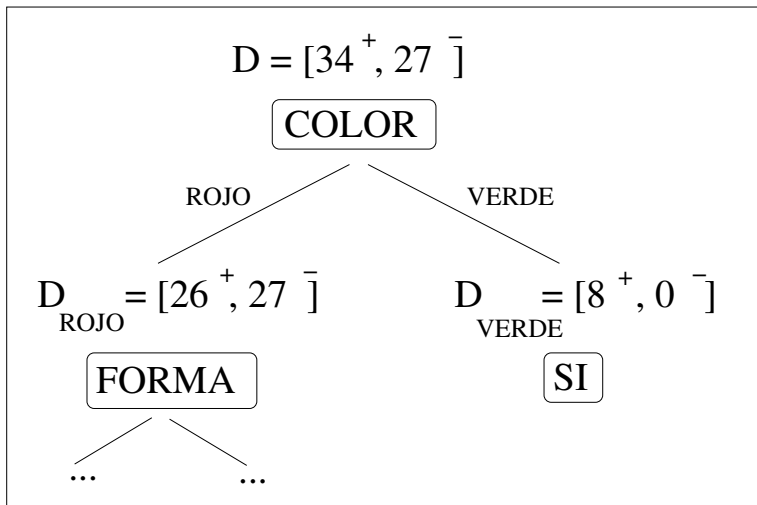
Árboles de decisión





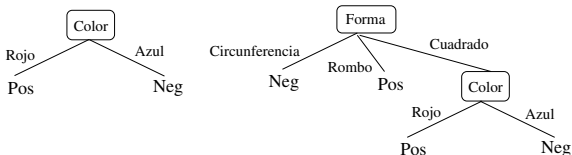


Árboles de decisión



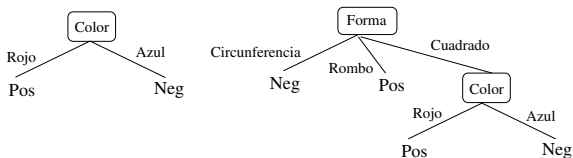
Clasificadores

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 prae-
ter necessitatem* (No ha de presu-
mirse 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)

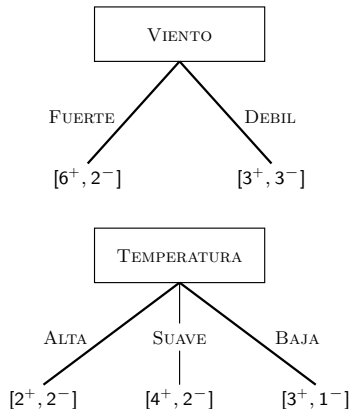
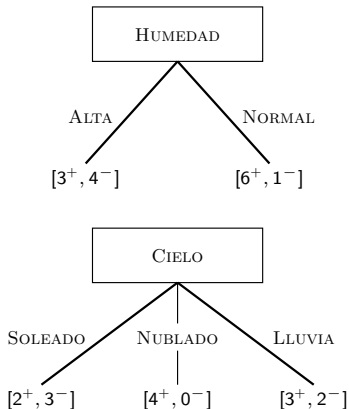
1. Si todos los Ejemplos son positivos, devolver un nodo etiquetado con +
2. Si todos los Ejemplos son negativos, devolver un nodo etiquetado con -
3. Si Atributos está vacío, devolver un nodo etiquetado con el valor más frecuente de Atributo-objetivo en Ejemplos.
4. En otro caso:
 - 4.1. Sea A el atributo 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

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)

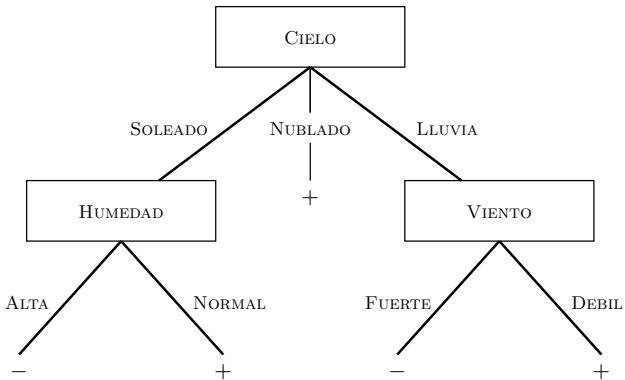
- Selección del atributo para el nodo CIELO=SOLEADO
- $D_{\text{SOLEADO}} = \{D_1, D_2, D_8, D_9, D_{11}\}$ con entropía
 $Ent([2^+, 3^-]) = 0,971$
 - $Ganancia(D_{\text{SOLEADO}}, \text{HUMEDAD}) =$
 $0,971 - \frac{3}{5} \cdot 0 - \frac{2}{5} \cdot 0 = 0,971$ (mejor atributo)
 - $Ganancia(D_{\text{SOLEADO}},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

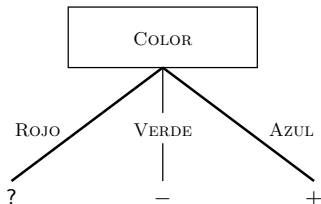
EJ.	COLOR	FORMA	TAMAÑO	CLASE
O_1	ROJO	CUADRADO	GRANDE	+
O_2	AZUL	CUADRADO	GRANDE	+
O_3	ROJO	REDONDO	PEQUEÑO	-
O_4	VERDE	CUADRADO	PEQUEÑO	-
O_5	ROJO	REDONDO	GRANDE	+
O_6	VERDE	CUADRADO	GRANDE	-

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:

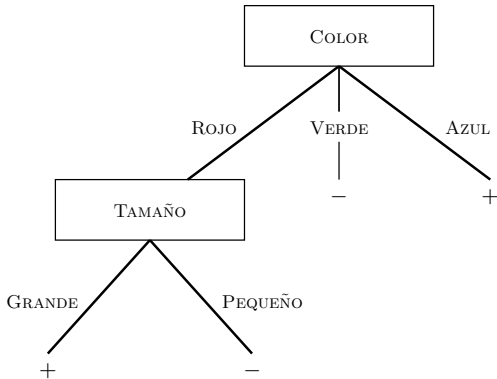


Algoritmo ID3 (ejemplo 2)

- 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 a los datos concretos que se han usado en el aprendizaje?

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
 - Estratificación: cada clase correctamente representada en el entrenamiento y en la prueba
- A veces es necesario usar un tercer conjunto para *validar* el ajuste de parámetros del modelo
- Si no tenemos suficientes ejemplos como para apartar un conjunto de prueba: validación cruzada
 - Dividir en k partes, y hace 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

- El sobreajuste es el gran enemigo a batir en aprendizaje automático
- Causas de sobreajuste:
 - Ruido (errores en los datos)
 - 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:
 - Evitar modelos excesivamente complejos (penalizar la complejidad)
 - Medir el rendimiento sobre conjuntos de validación independientes, para comprobar la capacidad de generalización de lo aprendido

Podado de árboles

Algoritmo de poda para reducir el error

1. Dividir el conjunto de ejemplos en Entrenamiento, Validación (y Prueba)
2. Árbol=árbol obtenido por ID3 usando Entrenamiento
3. Continuar=True
4. Mientras Continuar:
 - * Medida = proporción de ejemplos en conjunto de Validación 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 Validación.
 - * 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 (y su rendimiento sobre Prueba)

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 $A_{<C}$ compete con los restantes.
- El proceso se realiza a cada paso eligiendo el mejor umbral en el conjunto de entrenamiento.

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



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



Sonia L. Mestizo Gutiérrez^a, Marisol Herrera Rivero^b, Nicandro Cruz Ramírez^c,
Elena Hernández^d, Gonzalo E. Aranda-Abreu^{d,*}

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

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

^eDepartamento de Inteligencia Artificial, Universidad Veracruzana, Sebastián Camacho 5, Centro, Xalapa, Veracruz 91000, Mexico

^d 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 Gutiérrezz, 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,^{1*} Jing Shen,^{1*} Honggang Qian,¹ Huachong Ma,² Jiafu Ji,¹ Hong Ma,¹ Longhua Ma,³ Weihua Zhang,³ Ling Meng,¹ Zhenfu Li,¹ Jian Wu,¹ Genglin Jin,¹ Jianzhi Zhang¹ and Chengchao Shou^{1,4}

¹Peking University School of Oncology and Beijing Cancer Hospital and Institute, Haidian, Beijing 100036; ²Beijing Chaoyang Hospital, Chaoyang, Beijing 100020, China; and ³Ciphergen 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



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^a, Alexandros T. Tzallas^a, Dina Baga^a, Dimitra Chaloglou^b, Dimitrios I. Fotiadis^{a,*}, Sofia Tsouli^c, Maria Diakou^c, Spyros Konitsiotis^c

^a Unit of Medical Technology and Intelligent Information System, Department of Materials Science and Engineering, University of Ioannina, GR 45110, Ioannina, Greece

^b ANCO S.A., Athens GR 11742, Greece

^c Dept. of Neurology, Medical School, University of Ioannina, GR 451 10, Ioannina, Greece

Sección 3

Sección 3

Aprendizaje de reglas

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** ASTIGMATISMO=+ \wedge LAGRIMA=NORMAL
Entonces LENTE=RIGIDA
- 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 LENTE=RIGIDA
 - Si ?
 - Entonces LENTE=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 ASTIGMATISMO= +
 - Entonces LENTE=RIGIDA

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*)
 - *Early stopping*: no generar todas las condiciones
 - Podado de las reglas *a posteriori*: eliminar progresivamente condiciones hasta que no se produzca *mejora*

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

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

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

- 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 $k = 1$.

Clasificador Naive Bayes: un ejemplo

- Por tanto, las dos probabilidades (sin normalizar) a posteriori son:
 - $P(+)P(\text{soleado}|+)P(\text{suave}|+)P(\text{alta}|+)P(\text{fuerte}|+) = 0,0070$
 - $P(-)P(\text{soleado}|-)P(\text{suave}|-)P(\text{alta}|-)P(\text{fuerte}|-) = 0,0411$
- Así que el clasificador devuelve la clasificación con mayor probabilidad a posteriori, en este caso la respuesta es $-$ (no es un día bueno para jugar al tenis)

Detalles técnicos sobre las estimaciones: log-probabilidades

- Tal y como estamos calculando las estimaciones, existe el riesgo de que algunas de ellas sean excesivamente bajas
- Si realmente alguna de las probabilidades es baja y tenemos pocos ejemplos en el conjunto de entrenamiento, lo más seguro es que la estimación de esa probabilidad sea 0
- Esto plantea dos problemas:
 - La inexactitud de la propia estimación
 - Afecta enormemente a la clasificación que se calcule, ya que se multiplican las probabilidades estimadas y por tanto si una de ellas es 0, anula a las demás
- Una primera mejora técnica, intentado evitar productos muy bajos: usar logaritmos de las probabilidades.
 - Los productos se transforman en sumas

$$c_{NB} = \underset{c_j \in \mathcal{C}}{\operatorname{argmax}} [\log(P(c_j)) + \sum_i \log(P(a_i|c_j))]$$

Sección 6

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
- Tipos:
 - Clustering de partición estricta
 - Clustering jerárquico
 - Clustering basado en densidad

Aplicaciones

Published online 2017 Feb 24. doi: [10.1371/journal.pone.0171429](https://doi.org/10.1371/journal.pone.0171429)

Clustering cancer gene expression data by projective clustering ensemble

Xianxue Yu, Guoxian Yu, and Jun Wang*

Guy N Brock, Editor

[Author information](#) ► [Article notes](#) ► [Copyright and License information](#) ►

This article has been cited by other articles in PMC.

Go to:

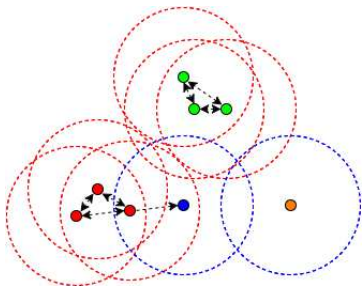
Gene expression data analysis has paramount implications for gene treatments, cancer diagnosis and other domains. Clustering is an important and promising tool to analyze gene expression data. Gene expression

Clustering cancer gene expression data by projective clustering ensemble Xianxue Yu, Guoxian Yu and Jun Wang

PLoS One. 2017; 12(2): doi: 10.1371/journal.pone.0171429

DBSCAN

Density Based Spatial Clustering of Applications with Noise



Ejemplo de salida de DBSCAN con $minPts = 2$. Se han obtenido dos clusters, uno con tres puntos y otro con cuatro. Los puntos rojos y verdes son puntos núcleo, el punto azul es frontera y el punto naranja es ruido.

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)