

# **Sistemas Inteligentes**

**Escuela Técnica Superior de Informática**

**Universitat Politècnica de València**

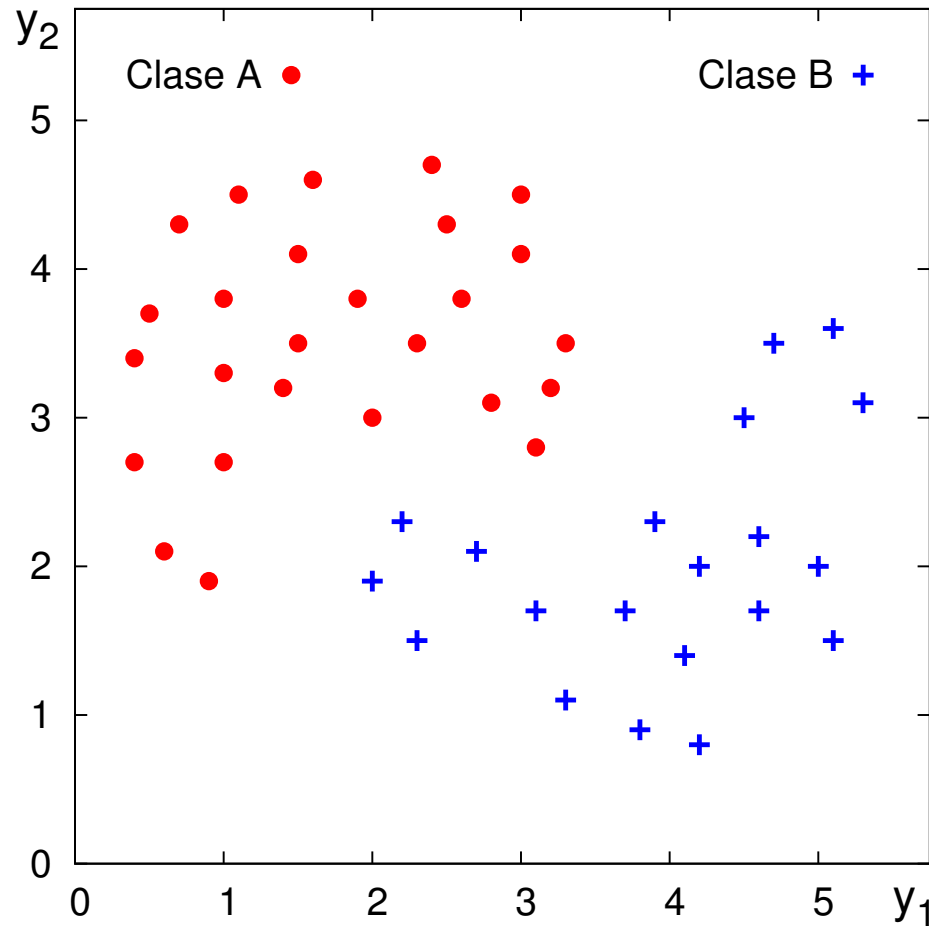
**Tema B2T3**

**Árboles de Clasificación**

# Índice

- 1 *Árboles de Clasificación (ADC)* ▷ 1
- 2 Aprendizaje de ADC ▷ 12
- 3 Bibliografía ▷ 28

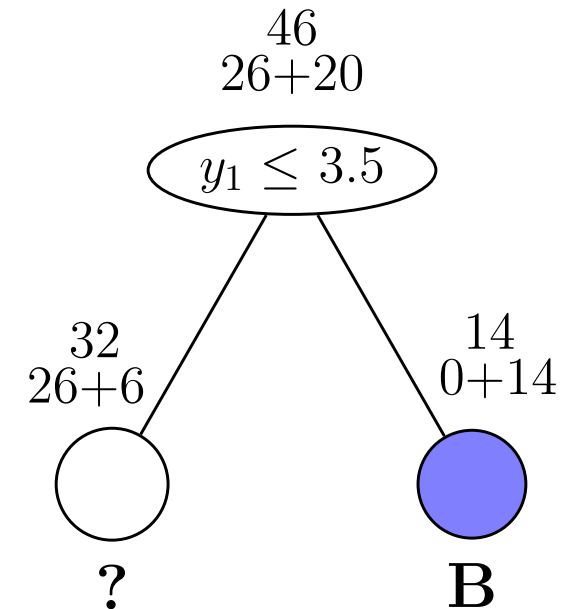
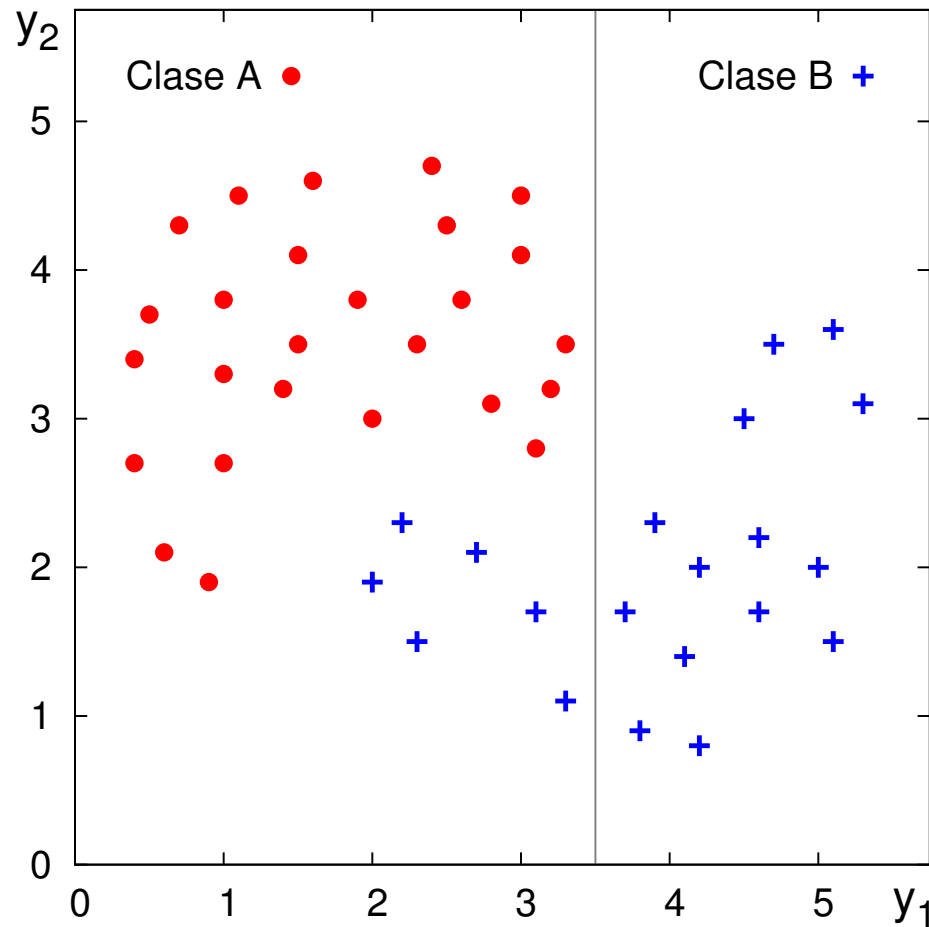
# Ejemplo



Tarea simple ilustrativa:

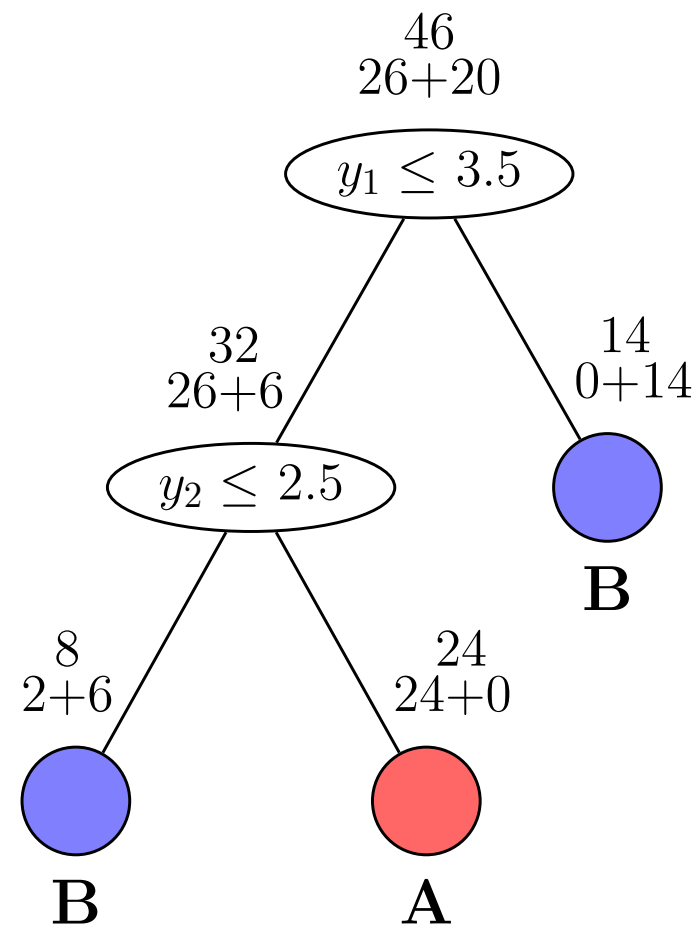
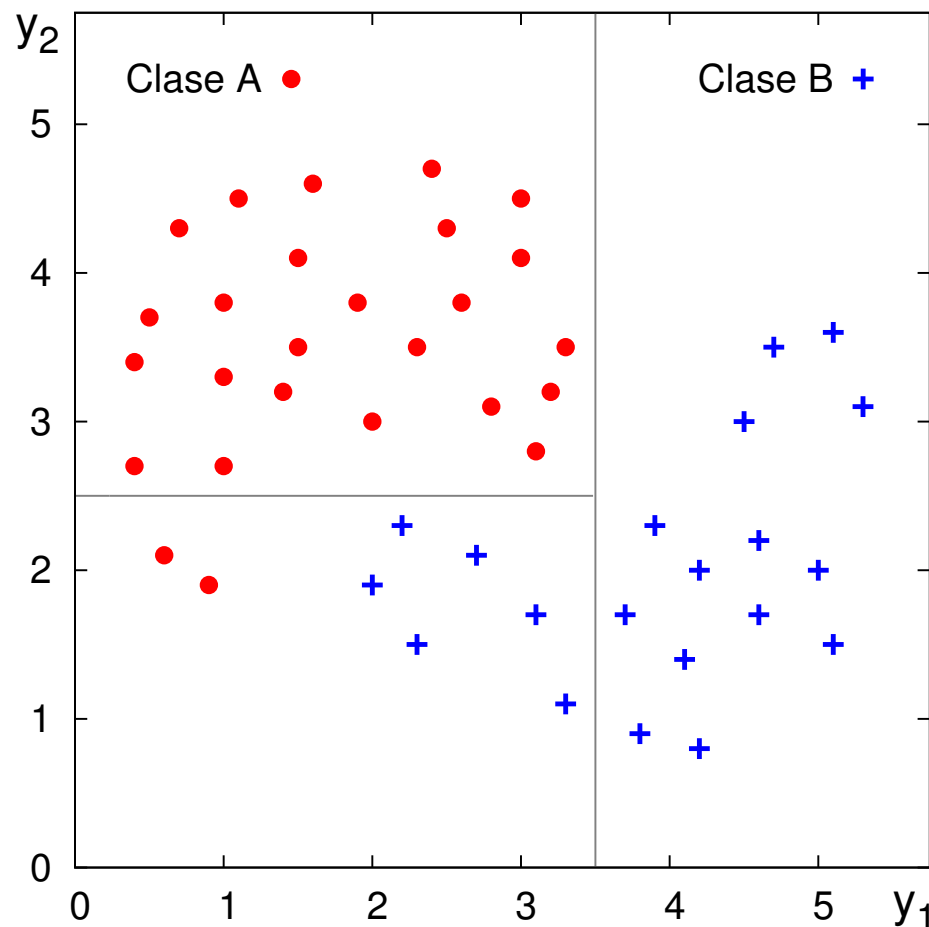
- representación en dos dimensiones ( $E = \mathbb{R}^2$ )
- clasificación en 2 clases *no* separables linealmente
- 46 datos (vectores)
  - 26 de clase A
  - 20 de clase B

## Ejemplo: Primera partición



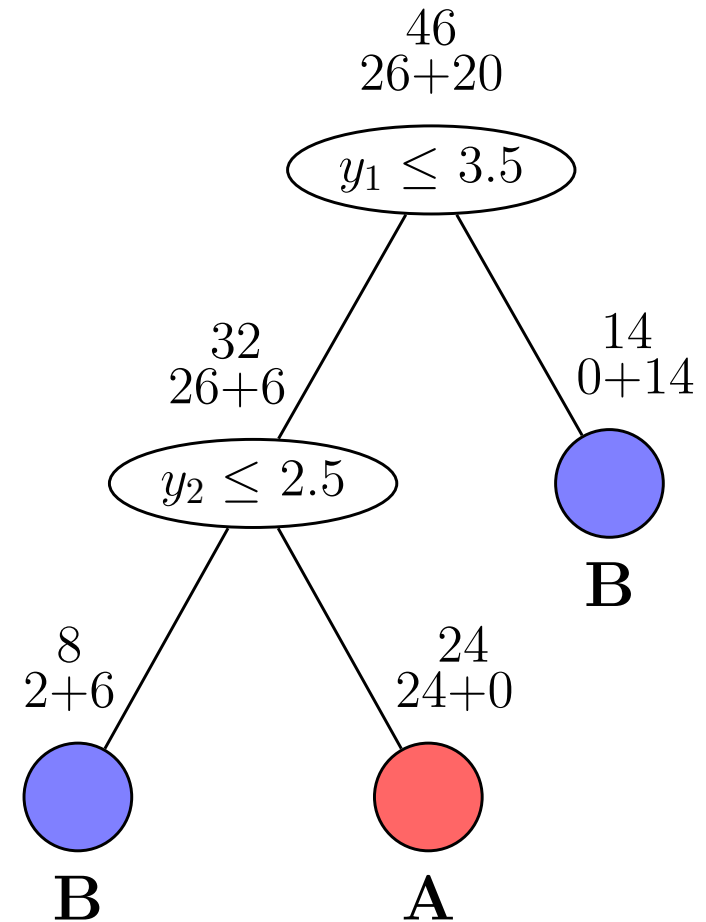
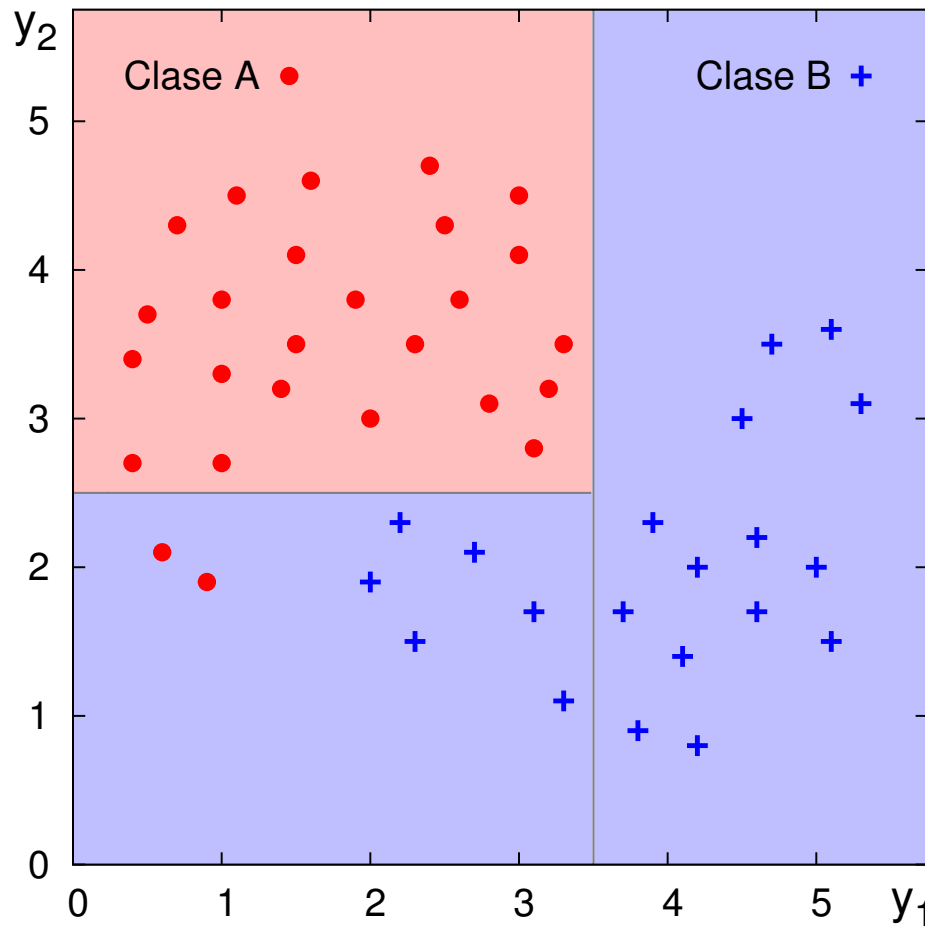
- El nodo raíz evalúa todos los datos con ¿ $y_1 \leq 3.5$ ?
- El nodo derecho es *puro* porque todos los datos de ese nodo son de clase B.

## Ejemplo: segunda partición y fronteras de decisión



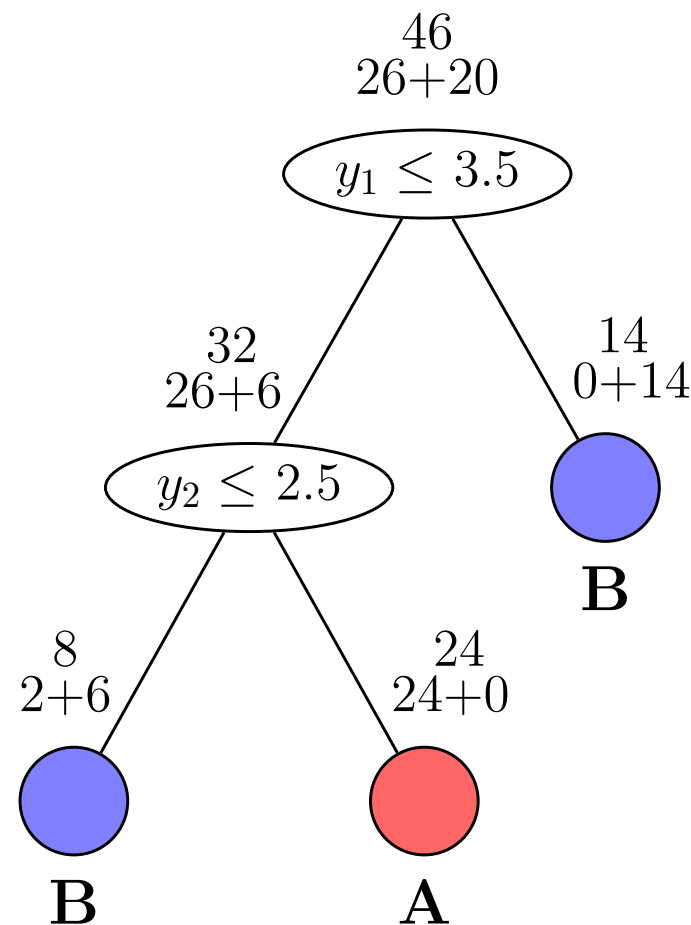
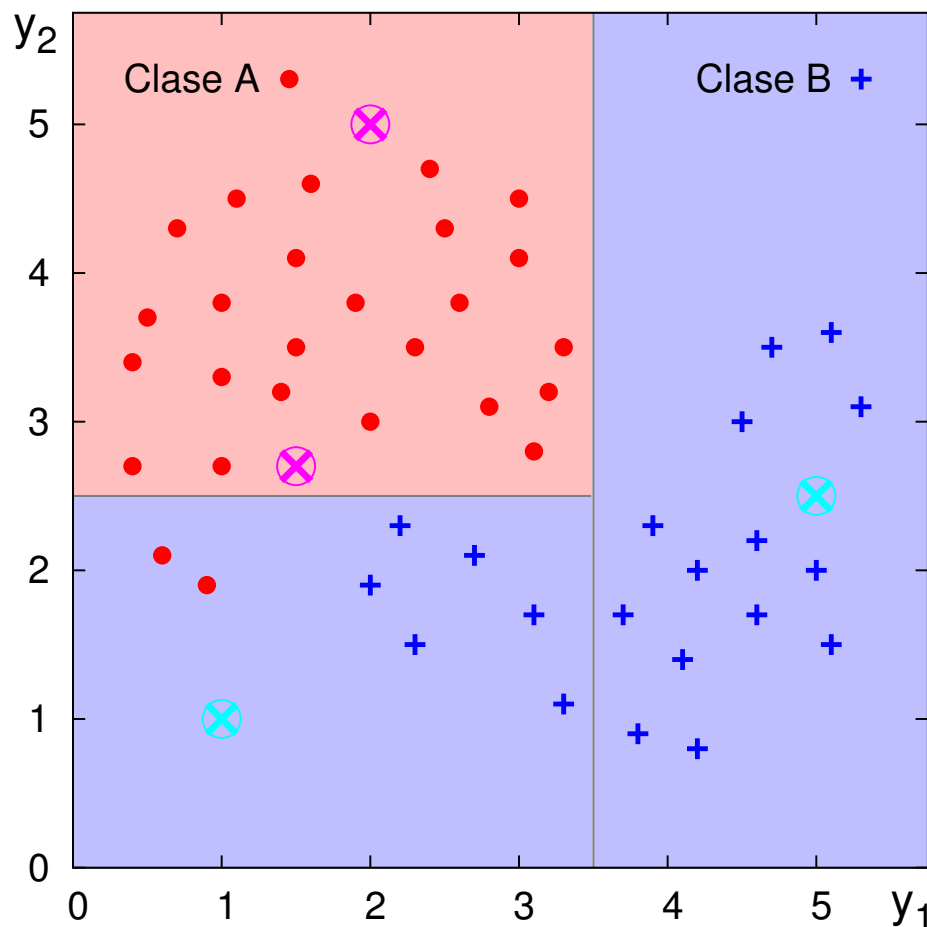
- Los datos del nodo izquierdo se particiona con  $y_2 \leq 2.5$ ?
- El nodo-hijo derecho es *puro* y se etiqueta como clase A
- El izquierdo no se particiona más y se etiqueta como B (clase mayoritaria)

## Ejemplo: regiones de decisión



- Fronteras de decisión paralelas a los ejes
- Regiones de decisión rectangulares.
- La probabilidad de error estimada por resustitución es  $2/46 = 0.0435 \rightarrow 4.35\%$

## Ejemplo: clasificación de nuevos datos



El árbol de decisión obtenido permite clasificar nuevos datos:

$(1.0, 1.0)^t$ :  $y_1 \leq 3.5, y_2 \leq 2.5 \rightarrow$  clase B

$(5.0, 2.5)^t$ :  $y_1 > 3.5 \rightarrow$  clase B

...

$(1.5, 2.7)^t$ :  $y_1 \leq 3.5, y_2 > 2.5 \rightarrow$  clase A

$(2.0, 5.0)^t$ :  $y_1 \leq 3.5, y_2 > 2.5 \rightarrow$  clase A

...

# Árboles de clasificación (ADC)

- Estructura resultante de la *partición recursiva* del *espacio de representación* a partir de una *datos de aprendizaje*.
- Cada nodo interno evalúa con una pregunta un atributo concreto
- Cada nodo hoja contiene una etiqueta de clase (clasificación)
- Los datos de test se clasifican mediante una serie de preguntas sobre los valores de sus atributos, empezado por el nodo raíz y siguiendo el camino determinado por las respuestas a las preguntas de los nodos internos, hasta llegar a un nodo hoja que clasifica la muestra en una clase
- Veremos CART<sup>1</sup> caracterizado por una partición binaria de los nodos basada en criterios estadísticos sólidos.

---

<sup>1</sup>Classification And Regression Trees



# Notación

- Espacio de representación:  $E \equiv \mathbb{R}^D$ ;  $\mathbf{y} = (y_1, y_2, \dots, y_D)^t \in E$
- Muestra de aprendizaje:  $N$  vectores, con su correcta clasificación:  $(\mathbf{y}_1, c_1), \dots, (\mathbf{y}_N, c_N)$ ,  $\mathbf{y}_i \in E$ ,  $c_i \in \mathbb{C} = \{1, 2, \dots, C\}$ ,  $1 \leq i \leq N$
- Un árbol se denota por  $T$  (“*Tree*”), un nodo por  $t$ , sus hijos izquierdo y derecho por  $t_L, t_R$ , respectivamente y el conjunto de nodos-hoja o terminales por  $\tilde{T}$
- Una partición binaria (“*split*”) se denota por  $s$  y el conjunto de particiones admisibles por  $S$

# Estimación de probabilidades asociadas a los nodos de un ADC

Sean:  $N$ ,  $N_c$ ,  $N(t)$ ,  $N_c(t)$ , respectivamente, el número total de datos de la muestra de aprendizaje, el número de estos datos de la clase  $c$ , el número de los que están representados en el nodo  $t$ , y el número de estos últimos que son de la clase  $c$ .

Probabilidad a priori de la clase  $c$ :  $\hat{P}(c) = \frac{N_c}{N}$

Probabilidad a posteriori de clase en el nodo  $t$ :  $\hat{P}(c | t) = \frac{N_c(t)}{N(t)}$

Probabilidad de un nodo *terminal*,  $t \in \tilde{T}$ :  $\hat{P}(t) = \frac{N(t)}{N}$

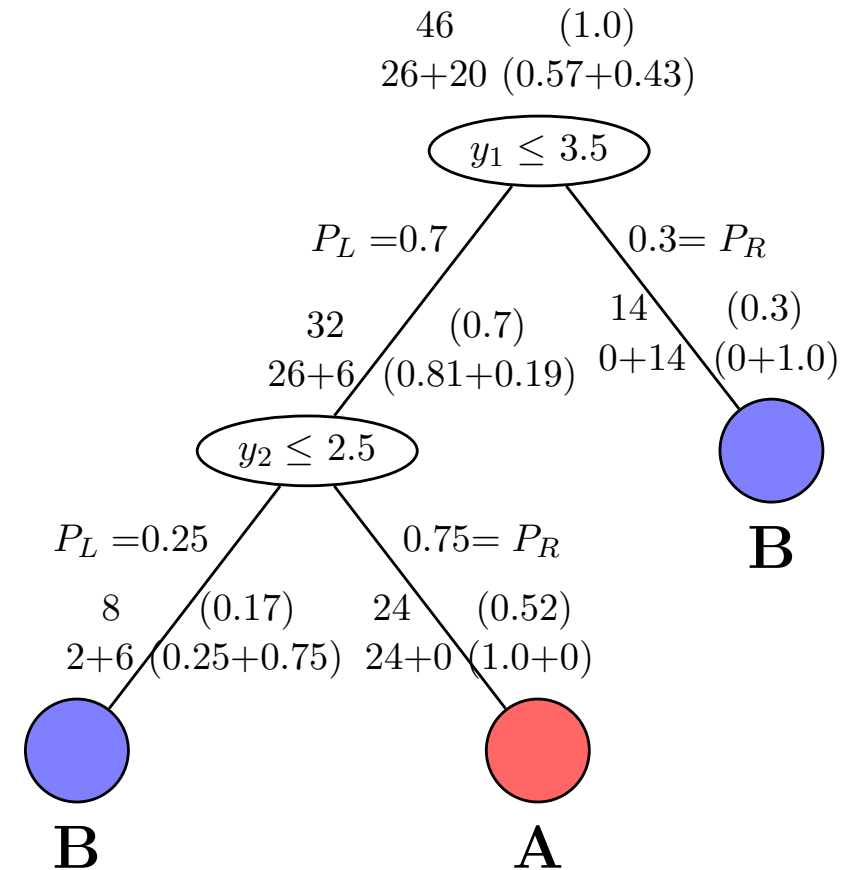
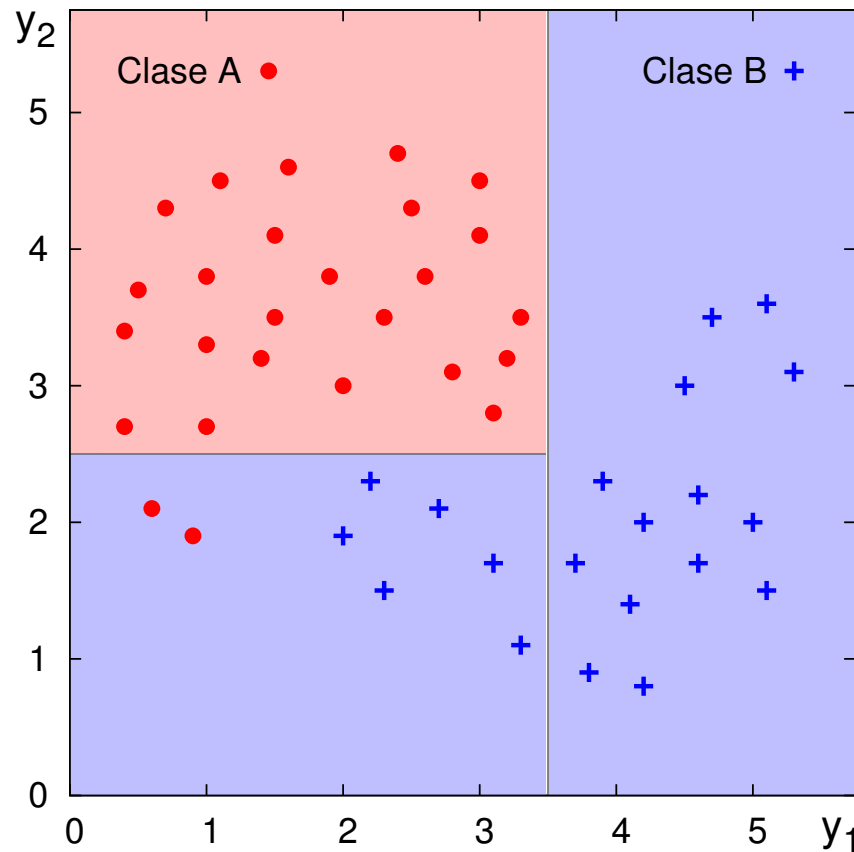
Probabilidad de decisión por el hijo izquierdo de  $t$ :  $\hat{P}_L(t) = \frac{N(t_L)}{N(t)}$

Probabilidad de decisión por el hijo derecho de  $t$ :  $\hat{P}_R(t) = \frac{N(t_R)}{N(t)}$

Todo nodo no terminal  $t$  verifica las siguientes propiedades:

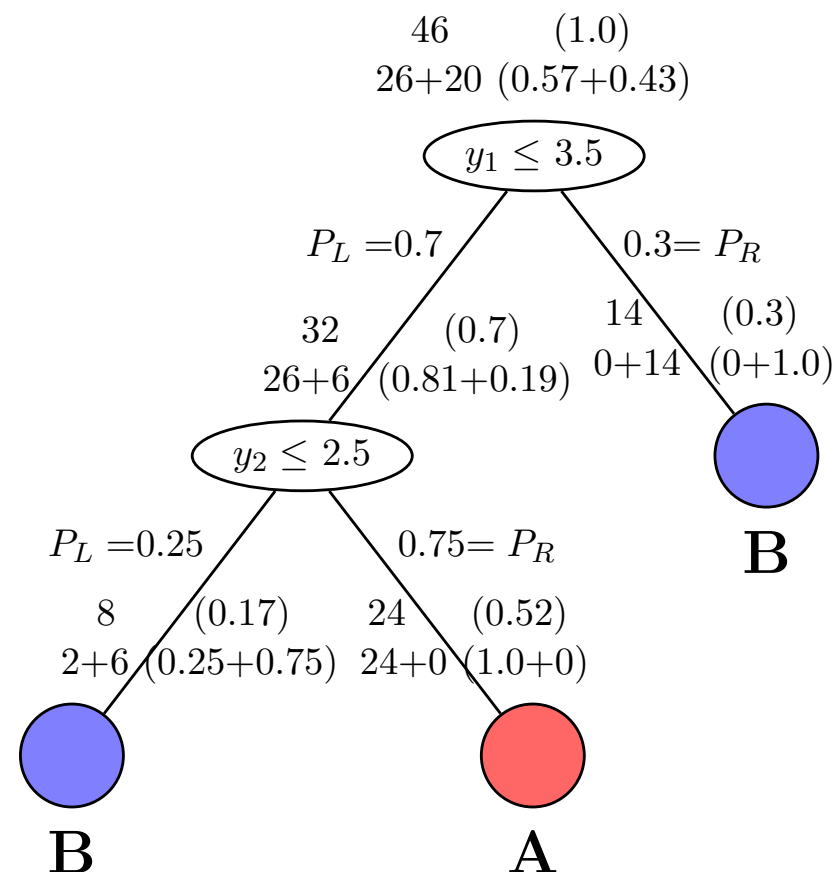
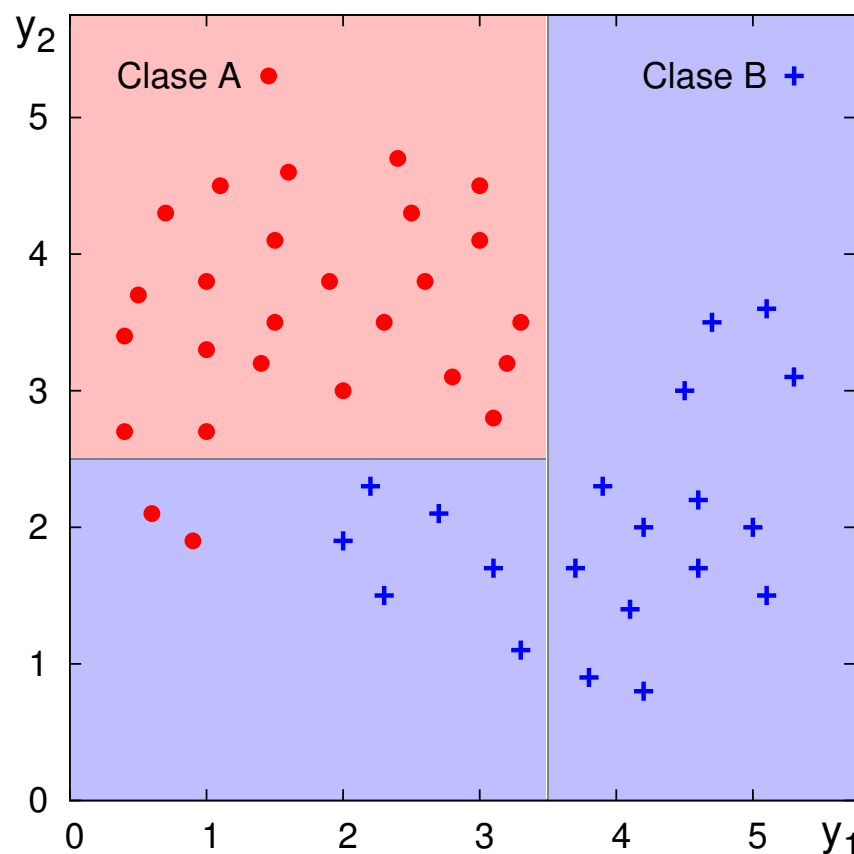
$$\hat{P}_L(t) + \hat{P}_R(t) = 1$$

# Ejercicio de cálculo de probabilidades



Nodos:	$\hat{P}(t_i)$	$\hat{P}(A   t_i)$	$\hat{P}(B   t_i)$	$\hat{P}_L(t_i)$	$\hat{P}_R(t_i)$
$t_1$ (raiz)					
$t_2$ (interno)					
$t_3$ (hoja B)					
$t_4$ (hoja B)					
$t_5$ (hoja A)					

# Solución al cálculo de probabilidades



Nodos:	$\hat{P}(t_i)$	$\hat{P}(A   t_i)$	$\hat{P}(B   t_i)$	$\hat{P}_L(t_i)$	$\hat{P}_R(t_i)$
$t_1$ (raíz)	46/46	26/46	20/46	32/46	14/46
$t_2$ (interno)	32/46	26/32	6/32	8/32	24/32
$t_3$ (hoja B)	14/46	0/14	14/14	—	—
$t_4$ (hoja B)	8/46	2/8	6/8	—	—
$t_5$ (hoja A)	24/46	24/24	0/24	—	—

# Índice

- 1 Árboles de Clasificación (ADC) ▷ 1
- 2 *Aprendizaje de ADC* ▷ 12
- 3 Bibliografía ▷ 28

# Aprendizaje de un ADC a partir de una muestra de aprendizaje

Parámetros de un algoritmo de aprendizaje de un ADC:

1. Criterio de partición y evaluación de partición:
  - Condiciones o “preguntas” (“splits”):  $\text{¿}y \in B\text{?, } B \subseteq E$
  - Evaluación y optimización de la calidad de una partición
2. Criterio de *pureza* de un nodo para detener proceso de partición
3. Criterio para asignar una etiqueta de clase a un nodo *terminal*

# Conjunto de preguntas admisibles para formar particiones

- Cada partición involucra a una única componente  $j$  de  $E$ ,  $1 \leq j \leq D$
- Las preguntas  $y \in B$  son de la forma:  $y_j \leq r$ ,  $r \in \mathbb{R}$
- Una partición o “split” es un par  $s = (j, r)$
- Los datos son finitos, por tanto sólo hay un número finito de particiones
- Para un nodo  $t$  con  $N(t)$  elementos:
  - Hay que explorar cada una de las componentes  $j$ ,  $1 \leq j \leq D$ , de  $E$
  - Para cada  $j$ , hay que explorar (al menos)  $N(t)$  posibles valores de  $r$
- Para cada nodo  $t$ , hay que explorar al menos  $\mathcal{O}(D \cdot N(t))$  “splits”

## Evaluación de la calidad de una partición

- Para evaluar las particiones posibles se usa el concepto de “*impureza*”
- La impureza de un nodo  $t$ ,  $\mathcal{I}(t)$ , se mide en función de las probabilidades estimadas de las clases en  $t$ , para lo cual existen varias aproximaciones.
- Una de las más interesantes se basa en el concepto de *entropía* (pág. 16):

$$\mathcal{I}(t) = - \sum_{c=1}^C \hat{P}(c | t) \log_2 \hat{P}(c | t) = - \sum_{c=1}^C \frac{N_c(t)}{N(t)} \log_2 \frac{N_c(t)}{N(t)} \quad (1)$$

Otras definiciones de  $\mathcal{I}(t)$ : *índice Gini* y *probabilidad de error* (ver [1,2,3])

- La calidad de una partición del nodo  $t$  mediante un “*split*”  $s = (j, r)$ , se mide mediante el *decremento de impureza*:

$$\Delta \mathcal{I}(j, r, t) \stackrel{\text{def}}{=} \mathcal{I}(t) - \hat{P}_L(t) \mathcal{I}(t_L) - \hat{P}_R(t) \mathcal{I}(t_R) \quad (2)$$

- La mejor partición es aquella que produce mayor decremento de impureza:

$$(j^*, r^*) = \underset{\substack{1 \leq j \leq D \\ -\infty < r < +\infty}}{\operatorname{argmax}} \Delta \mathcal{I}(j, r, t) \quad (3)$$



# Entropía

- Mide la cantidad de información asociada a una decisión  $k$ -aria:

$$H = - \sum_{i=1}^k P_i \log_2 P_i \quad (0 \log 0 \stackrel{\text{def}}{=} 0)$$

- La unidad es el *bit*: información asociada a tomar una decisión *binaria* en la que las dos alternativas son equiprobables.
- El valor mínimo es 0 y corresponde a una decisión en la que solo hay una alternativa posible.
- El valor máximo es  $+\infty$  que se da para decisiones  $k$ -arias equiprobables con  $k \rightarrow \infty$ :

- Ejemplos:

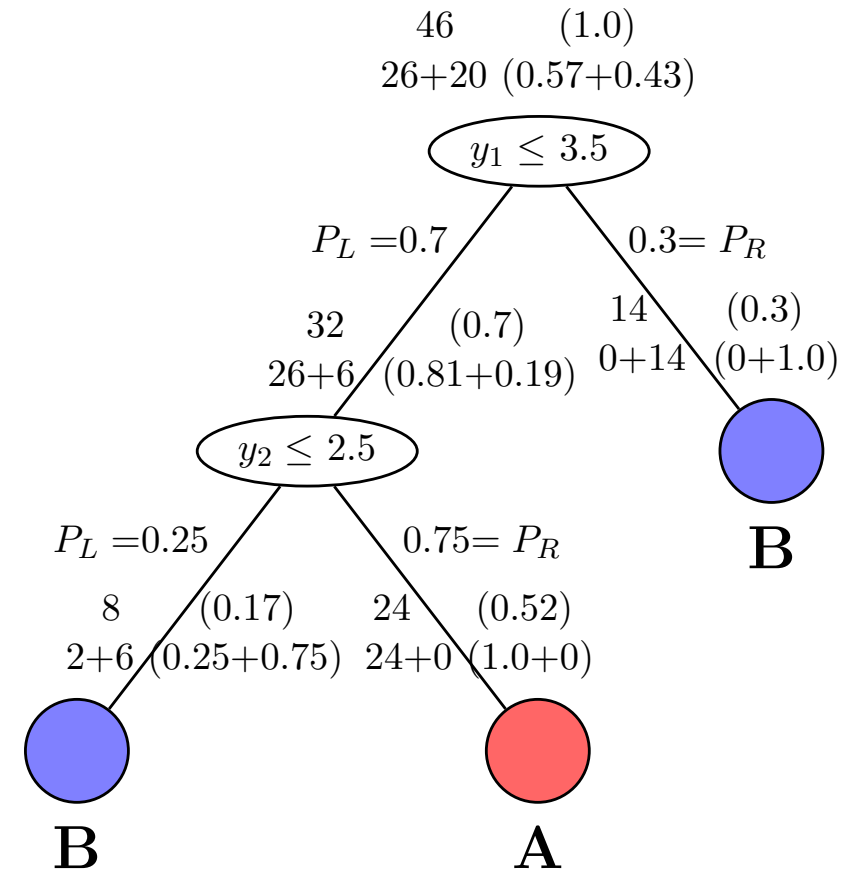
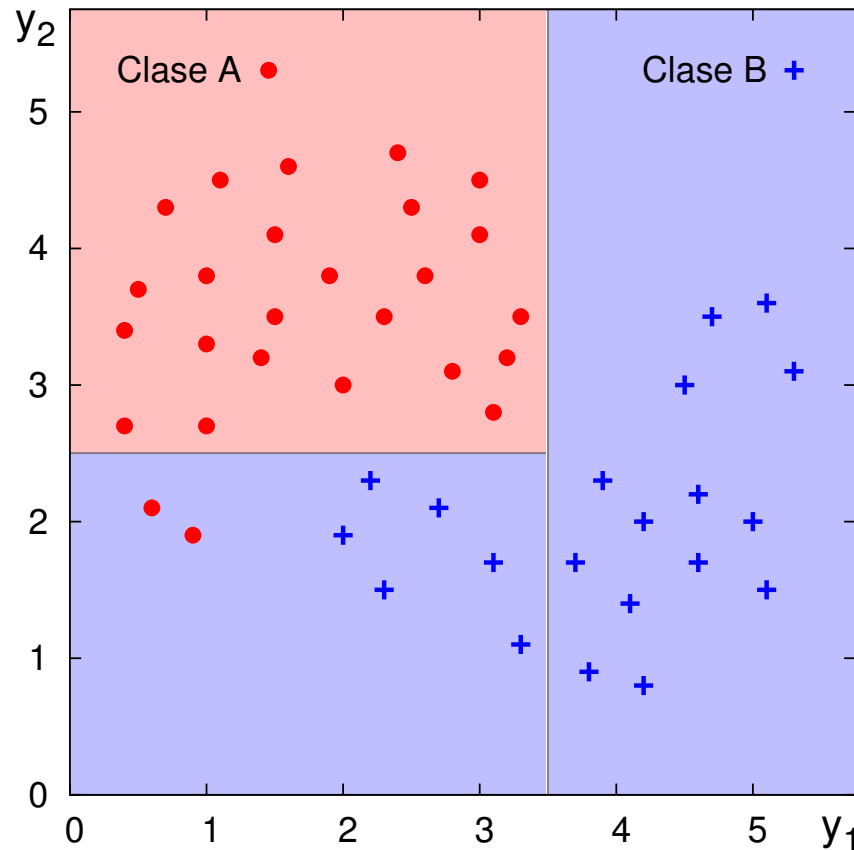
$$\text{Si } P_1 = P_2 = 1/2, \quad H = - (0.5(0 - 1) + 0.5(0 - 1)) = 1$$

$$\text{Si } P_1 = 1, P_2 = 0, \quad H = - 1 \cdot 0 + 0 = 0$$

$$\text{Si } P_i = 1/k, 1 \leq i \leq k, \quad H = \log_2 k; \quad H \rightarrow \infty \text{ si } k \rightarrow \infty$$

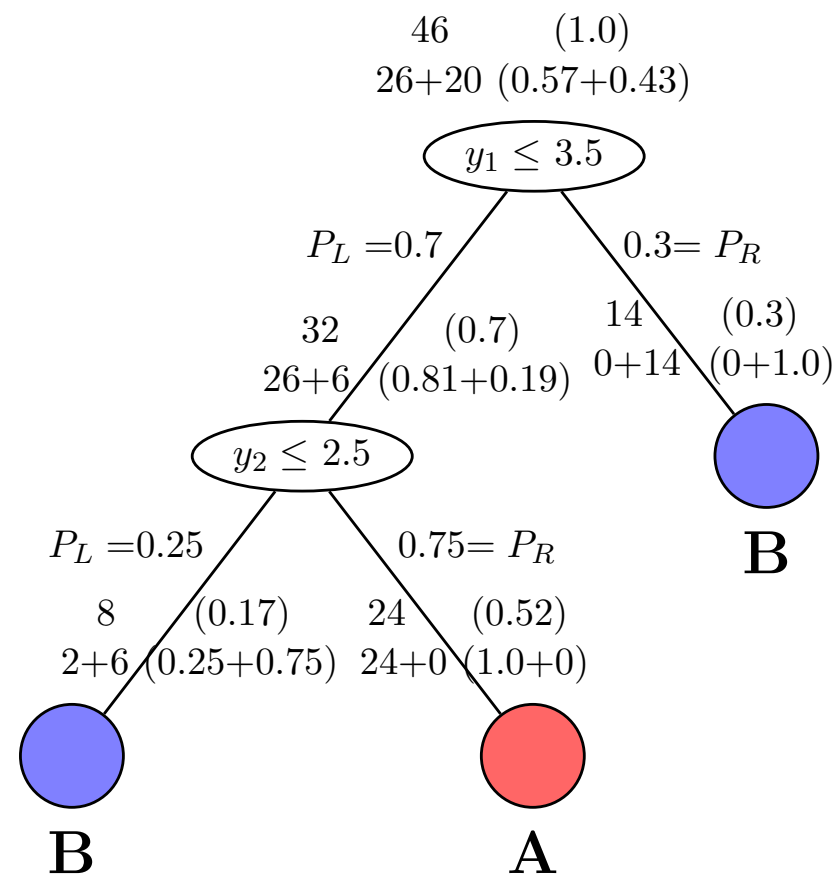
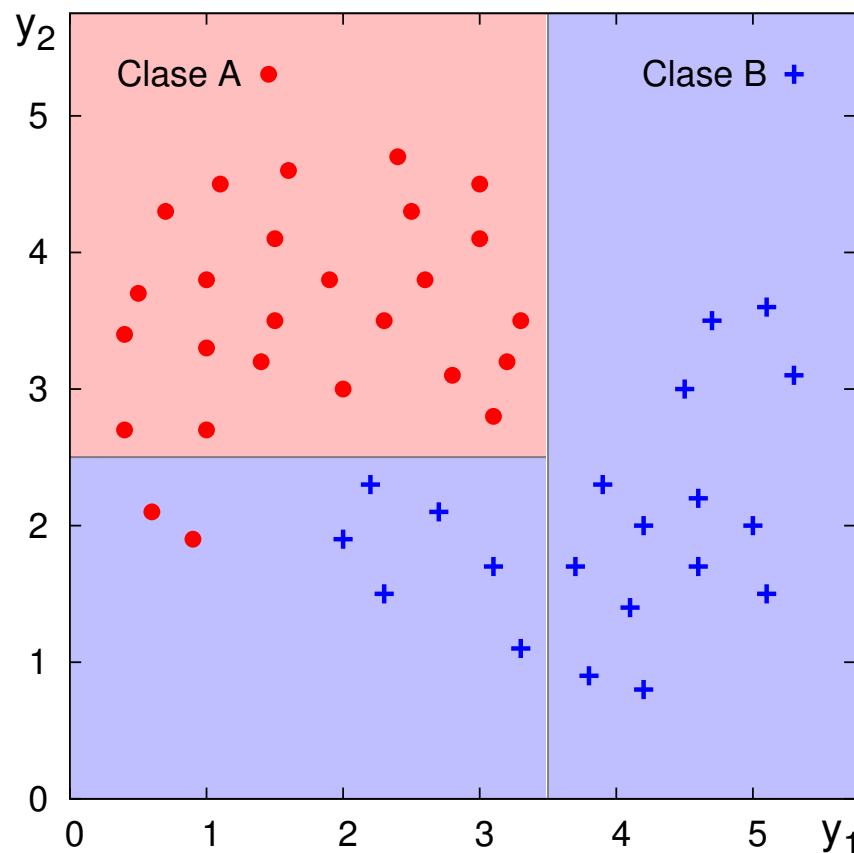
*Ejercicio:* según Eq.(1), calcular  $\mathcal{I}(t) \forall t$  en el árbol de la página 5.

# Ejercicio de cálculo de impureza



Nodos:	$\hat{P}(t_i)$	$\hat{P}(A   t_i)$	$\hat{P}(B   t_i)$	$\hat{P}_L(t_i)$	$\hat{P}_R(t_i)$	$\mathcal{I}(t_i)$
$t_1$ (raiz)	46/46	26/46	20/46	32/46	14/46	
$t_2$ (interno)	32/46	26/32	6/32	8/32	24/32	
$t_3$ (hoja B)	14/46	0/14	14/14	—	—	
$t_4$ (hoja B)	8/46	2/8	6/8	—	—	
$t_5$ (hoja A)	24/46	24/24	0/24	—	—	

# Solución al cálculo de impureza



Nodos:	$\hat{P}(t_i)$	$\hat{P}(A   t_i)$	$\hat{P}(B   t_i)$	$\hat{P}_L(t_i)$	$\hat{P}_R(t_i)$	$\mathcal{I}(t_i)$
$t_1$ (raiz)	46/46	26/46	20/46	32/46	14/46	<b>0.988</b>
$t_2$ (interno)	32/46	26/32	6/32	8/32	24/32	<b>0.695</b>
$t_3$ (hoja B)	14/46	0/14	14/14	—	—	<b>0.000</b>
$t_4$ (hoja B)	8/46	2/8	6/8	—	—	<b>0.811</b>
$t_5$ (hoja A)	24/46	24/24	0/24	—	—	<b>0.000</b>

## Ejercicio (para hacer en clase)

Con respecto al ejemplo de la página 5:

- Calcular el decremento de impureza que se produce al dividir cada uno de los 2 nodos *no* terminales según Eq (2)
- Las dos particiones que se muestran en este ejemplo son solo ejemplos ilustrativos basados en pura intuición geométrica (es decir, *no* son el resultado de la optimización de ninguna expresión de impureza).

Según la Eq. (3), analizar la segunda partición (la que en el ejemplo se resuelve con  $s = (2, 2.5)$ ; es decir,  $y_2 \leq 2.5$ ), y determinar si alguna de las siguientes particiones es mejor para ese nodo:  $(y_1 \leq 1.95)$ ,  $(y_2 \leq 1.8)$

Nodos:	$\hat{P}(t_i)$	$\hat{P}(A   t_i)$	$\hat{P}(B   t_i)$	$\hat{P}_L(t_i)$	$\hat{P}_R(t_i)$	$\mathcal{I}(t_i)$
$t_1$ (raiz)	46/46	26/46	20/46	32/46	14/46	<b>0.988</b>
$t_2$ (interno)	32/46	26/32	6/32	8/32	24/32	<b>0.695</b>
$t_3$ (hoja B)	14/46	0/14	14/14	—	—	<b>0.000</b>
$t_4$ (hoja B)	8/46	2/8	6/8	—	—	<b>0.811</b>
$t_5$ (hoja A)	24/46	24/24	0/24	—	—	<b>0.000</b>

## Solución al ejercicio

- Decrementos de impureza para  $t_1$  y  $t_2$ :

Nodos:	$\hat{P}(t_i)$	$\hat{P}(A   t_i)$	$\hat{P}(B   t_i)$	$\hat{P}_L(t_i)$	$\hat{P}_R(t_i)$	$\mathcal{I}(t_i)$	$\Delta\mathcal{I}(t_i)$
$t_1$ (raiz)	46/46	26/46	20/46	32/46	14/46	<b>0.504</b>	
$t_2$ (interno)	32/46	26/32	6/32	8/32	24/32	<b>0.492</b>	
$t_3$ (hoja B)	14/46	0/14	14/14	—	—	—	
$t_4$ (hoja B)	8/46	2/8	6/8	—	—	—	
$t_5$ (hoja A)	24/46	24/24	0/24	—	—	—	

- Splits* alternativos en  $t_2$ :

$$(y_1 \leq 1.95) : \mathcal{I}(t_L) = 0, \mathcal{I}(t_R) = -(11/17) \log(11/17) - (6/17) \log(6/17) = 0.937$$

$$(y_2 \leq 1.80) : \mathcal{I}(t_L) = 0, \mathcal{I}(t_R) = -(26/29) \log(26/29) - (3/29) \log(3/29) = 0.480$$

$$\Delta\mathcal{I}(1, 1.95, t_2) = 0.695 - 0 - (17/32) \cdot 0.937 = 0.197 < 0.492$$

$$\Delta\mathcal{I}(2, 1.80, t_2) = 0.695 - 0 - (29/32) \cdot 0.480 = 0.260 < 0.492$$

Por tanto, ninguno de estos *splits* habría sido mejor que  $(y_2 \leq 2.5)$ .

## Criterios de suficiente “pureza” en nodos terminales

- Criterio de parada de particionamiento, cuando máximo decremento de impureza posible es demasiado pequeño:

$$\max_{\substack{1 \leq j \leq D \\ -\infty < r < +\infty}} \Delta \mathcal{I}(j, r, t) < \epsilon \quad (4)$$

donde  $\epsilon$  es una constante pequeña a determinar empíricamente.

- Otro criterio es exigir que los nodos terminales sean totalmente puros
- Este último criterio tiene problemas de generalización y árboles grandes

## Asignación de etiquetas de clase a nodos terminales

- Criterio simple y eficaz: asignar a cada nodo terminal la clase mayoritaria en sus elementos:

$$c^*(t) = \operatorname{argmax}_{1 \leq c \leq C} \hat{P}(c | t), \quad \forall t \in \tilde{T} \quad (5)$$

# Algoritmo ADC

*Árbol* CreaNodo(*clase*  $c$ , *componente*  $j$ , *umbral*  $u$ , *nodos*  $t_L, t_R$ ) // crea un árbol ( $\uparrow$ nodo)  
*Árbol* ADC(*muestra*  $\mathcal{Y} \equiv (\mathbf{y}_1, c_1), \dots, (\mathbf{y}_n, c_n)$ ) { // aprende un Árbol de Clasificación  
      $(j^*, r^*, \delta) = \text{MejorPartición}(\mathcal{Y})$  // según Eq. (1,2,3);  $\delta$  es el decremento de impureza  
     **if**  $(\delta < \epsilon)$  { // si  $\delta$  es demasiado pequeño – según Eq. (4)  
          $c = \text{ClaseMayoritaria}(\mathcal{Y})$  // según Eq. (5)  
         **return** CreaNodo( $c, -, -, \text{NULL}, \text{NULL}$ ) // crea nodo terminal y le asigna la clase  $c$   
     } **else** {  
          $\mathcal{Y}_L = \mathcal{Y}_R = \emptyset$  //  $\emptyset$  es el conjunto vacío  
          $\forall (\mathbf{y}, c) \in \mathcal{Y}$  { // realiza la partición en función de  $j^*, r^*$   
             **if**  $(y_{j^*} \leq r^*)$   $\mathcal{Y}_L = \mathcal{Y}_L \cup \{(\mathbf{y}, c)\}$   
             **else** /\*  $(y_{j^*} > r^*)$  \*/  $\mathcal{Y}_R = \mathcal{Y}_R \cup \{(\mathbf{y}, c)\}$   
         }  
          $t_L = \text{ADC}(\mathcal{Y}_L)$  // crea recursivamente el subárbol izquierdo  
          $t_R = \text{ADC}(\mathcal{Y}_R)$  // crea recursivamente el subárbol derecho  
         **return** CreaNodo( $0, j^*, r^*, t_L, t_R$ )  
     }  
 }

# Estimación por resustitución del error de clasificación

- Según la *teoría de la decisión estadística*, la probabilidad de error de un nodo (terminal)  $t$ , estimada por *resustitución*, es:

$$\hat{P}_e(t) = 1 - \max_{1 \leq c \leq C} \hat{P}(c | t)$$

Y para un árbol  $T$ :

$$\hat{P}_e(T) = \sum_{t \in \tilde{T}} \hat{P}(t) \hat{P}_e(t)$$

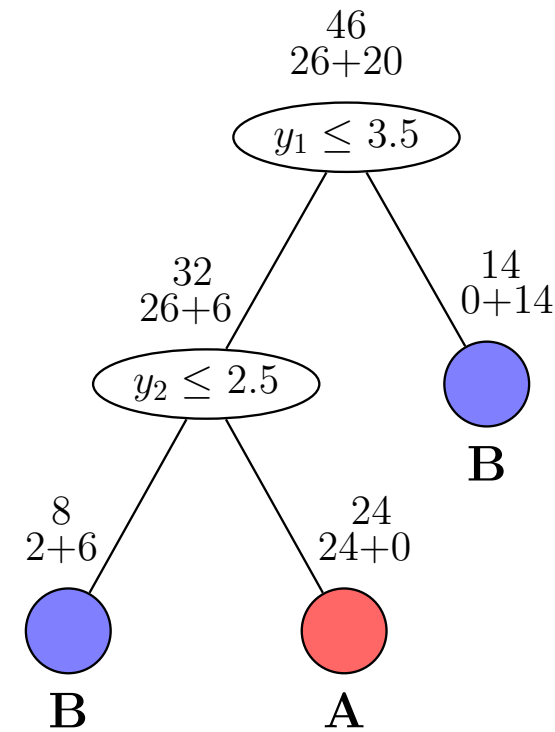
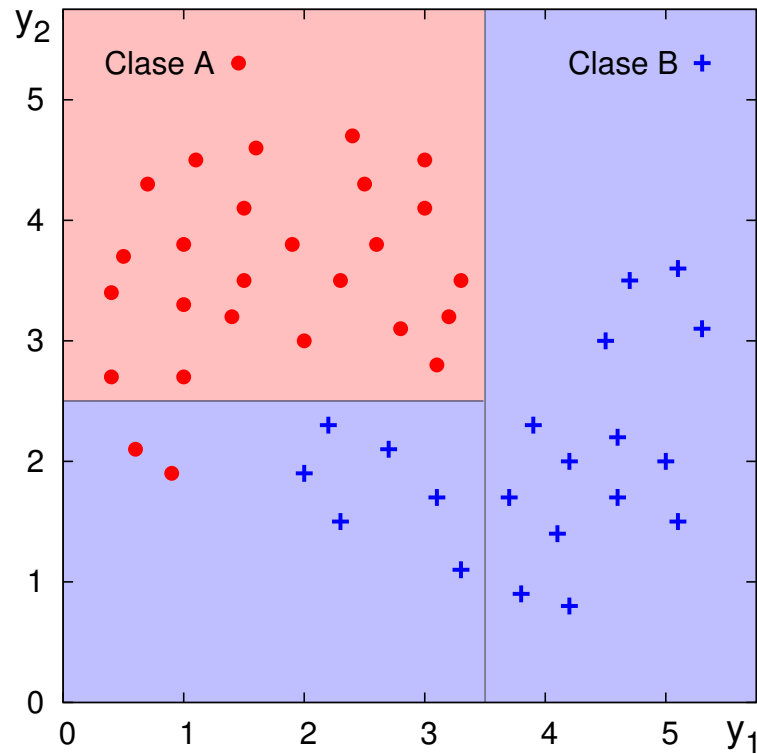
- Si se hace crecer el árbol hasta que los nodos terminales sean totalmente puros, el error estimado será nulo, ya que en este caso:

$$\hat{P}_e(t) = 0 \quad \forall t \in \tilde{T}.$$

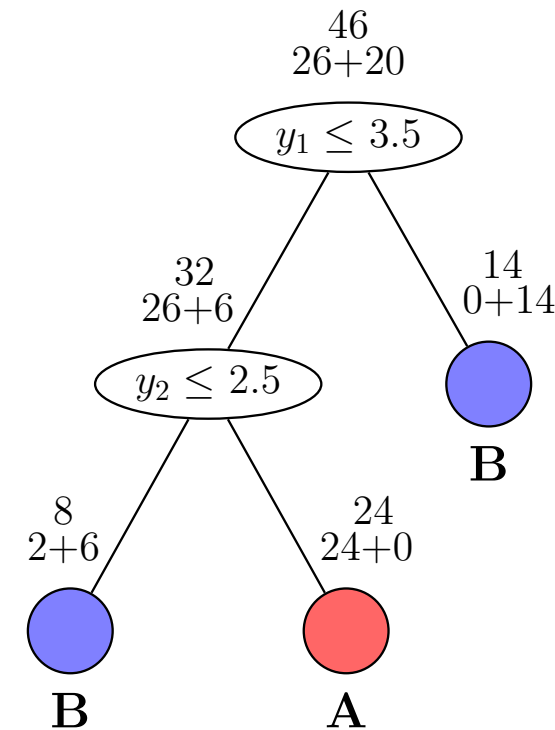
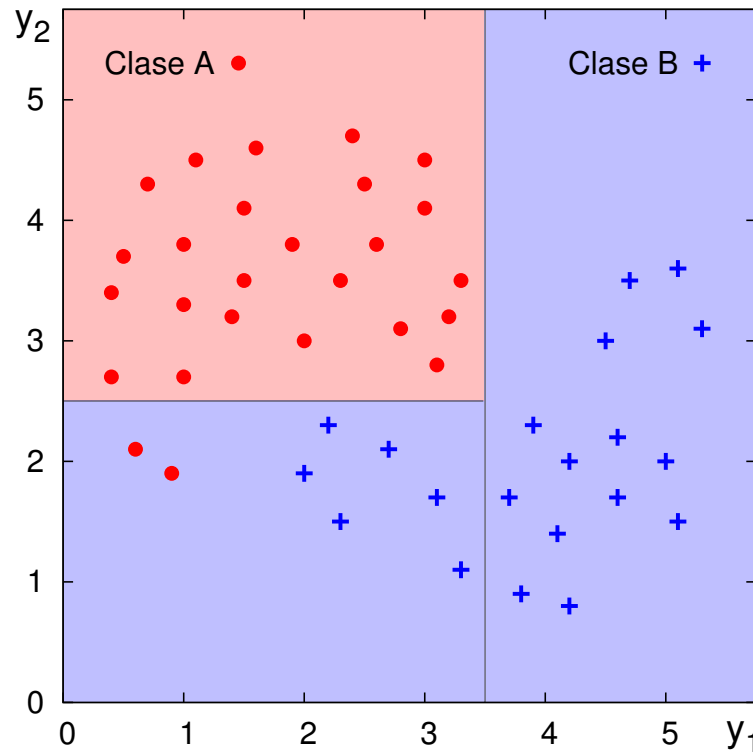
- Esto conlleva un *sobreaprendizaje* que generalmente no es deseable  $\Rightarrow$  esencialmente el árbol se convierte en mero almacén de la muestra de aprendizaje, sin capacidad de generalización ante nuevos datos



# Estimación por resustitución del error de clasificación



# Estimación por resustitución del error de clasificación



En el ejemplo de la página 5, (3 nodos terminales, 2 totalmente puros):

$$\begin{aligned}
 \hat{P}_e(T) &= \hat{P}(t_3)\hat{P}_e(t_3) + \hat{P}(t_4)\hat{P}_e(t_4) + \hat{P}(t_5)\hat{P}_e(t_5) \\
 &= \frac{14}{46} \cdot 0 + \frac{8}{46} \cdot \frac{2}{8} + \frac{24}{46} \cdot 0 = \frac{2}{46} \approx 0.0435 \rightarrow 4.35\%
 \end{aligned}$$

## Algoritmo ADC: Ejercicio en clase

En un problema de clasificación en 3 clases, para objetos representados mediante vectores de características bidimensionales tenemos la siguiente muestra:

	Entrenamiento					Test			
$n$	1	2	3	4	5	6	7	8	9
$x_{n1}$	1	1	2	2	3	3	3	0	0
$x_{n2}$	0	1	0	1	0	1	-1	1	0
$c_n$	A	B	B	B	C	B	C	B	A

1. Realiza una traza de ejecución del algoritmo ADC sobre el conjunto de entrenamiento siendo  $\epsilon = 0$ . Muestra el árbol resultante, los valores de impureza de los nodos generados y los decrementos de impureza de las particiones obtenidas.
2. Representa gráficamente las particiones obtenidas y las regiones de decisión tras la ejecución del algoritmo ADC, así como las muestras de entrenamiento.
3. Calcula la probabilidad de error en test con intervalos de confianza al 95 % y discute el resultado.
4. Muestra el árbol resultante si  $\epsilon = 0.5$  y calcula la probabilidad de error en test en este caso.

## Ejercicio de seguimiento: Algoritmo ADC

En un problema de clasificación en 3 clases, para objetos representados mediante vectores de características tridimensionales tenemos la siguiente muestra:

	Entrenamiento									Test					
$n$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$x_{n1}$	0	1	0	2	0	1	1	2	2	1	0	0	1	1	2
$x_{n2}$	0	1	1	1	0	1	0	1	1	0	1	0	0	1	0
$x_{n3}$	0	0	1	0	2	3	2	2	3	0	2	1	3	1	2
$c_n$	A	A	A	A	B	B	C	C	C	A	A	B	B	C	C

1. A partir del conjunto de entrenamiento, enumera el conjunto de particiones posibles para el nodo raíz en las que ambos nodos hijos tienen al menos un elemento.
2. Calcula la mejor partición inicial (nodo raíz) de acuerdo al decremento de impureza.
3. Calcula la probabilidad de error en test tras aplicar el algoritmo ADC con  $\epsilon = 0.35$
4. Considera sólo las muestras de entrenamiento de las clases B y C. ¿Es posible definir una única partición con error de resustitución igual a cero? ¿Por qué?
5. En el caso del apartado anterior, ¿crees que el algoritmo Perceptron conseguiría un error de resustitución igual a cero? ¿Por qué?

# Índice

- 1 Árboles de Clasificación (ADC) ▷ 1
- 2 Aprendizaje de ADC ▷ 12
- 3 *Bibliografía* ▷ 28

# Bibliografía

- [1] R.O. Duda, D.G. Stork, P.E. Hart. Pattern Classification. Wiley, 2001.
- [2] A. R. Webb, K. D. Copsey. Statistical Pattern Recognition. Wiley, tercera ed., 2011.
- [3] Classification and Regression Trees by L. Breiman, J.H. Friedman, R.A. Olshen y C.J. Stone. Chapman & Hall, 1984.

El material de este tema se basa principalmente en [1] y [2].