



Escuela Politécnica
Superior (EPS)



UNIVERSIDAD DE CÓRDOBA
Departamento de
Informática y
Análisis Numérico

Práctica 4. **Ejemplo de creación de un árbol de decisión**

Juan Carlos Fernández Caballero (jfcaballero@uco.es)

Introducción al Aprendizaje Automático (IAA)

3º de Grado de Ingeniería Informática

Especialidad en Computación

mayo de 2020



GRUPO DE INVESTIGACIÓN AYRNA
APRENDIZAJE Y REDES NEURONALES ARTIFICIALES
uco.es/ayrna

Entrenamiento de los árboles de decisión

Para crear y entrenar un árbol se usa el conjunto de entrenamiento o ***training***, y se realizan **procesos iterativos** (***Greedy*** o **Voraz**, **Divide y Vencerás**, etc) sobre los datos, que van obteniendo ramas y subconjuntos de datos homogéneos respecto a la clase. Estos procesos iterativos se basan en la **ganancia de información**, que depende de la **Entropía**.

- **Entropía:** Medida de incertidumbre, relacionada con:
 - ▶ **1) Pureza:** Cómo de cerca está un conjunto de pertenecer a una misma clase.
 - ▶ **2) Impureza (desorden):** Cómo de cerca está un conjunto de la incertidumbre total.
- La Entropía es una medida:
 - ▶ **1) Directamente proporcional** a la impureza, incertidumbre, irregularidad.
 - ▶ **2) Inversamente proporcional** a la pureza, certidumbre, redundancia.
- **Objetivo:** Minimizar la Entropía (desorden o error)!

Entrenamiento de los árboles de decisión

A partir del conjunto de ***entrenamiento*** y de forma general, los pasos son los siguientes:

1. Elección del **nodo raíz**: Se escogerá aquél que provoque una mejor separación de las clases.

Se hace uso de la **Ganancia de información**, que a su vez usa la **Entropía de la clase** y la **Entropía de la clase condicionada a un valor**.

- Cantidad de información mutua o ganancia de información:

$$I(C, X_i) = H(C) - H(C|X_i)$$

- ### ► Entropía de una variable:

$$H(C) = - \sum_{c=1}^n p(c) \log_2 p(c)$$

Entrenamiento de los árboles de decisión

- ▶ Entropía de una variable condicionada a un valor:

$$H(C|X) = - \sum_{i=1}^n p(x|c) \log_2 p(c|x) =$$

$$- \sum_c \sum_x p(x|c) \log_2 p(c|x)$$

- Se elige como primer nodo aquél que maximiza la expresión:
 $I(C, X_i)$, es decir, **minimiza la Entropía (error o desorden)**.

2. Cuando se dividen las ramas, se crean **nuevos subconjuntos de training** para cada una de ellas.
 3. Se repite el proceso (puntos 1 y 2) para cada rama generada.

Ejemplo de entrenamiento (construcción) de un árbol de decisión

A partir de la siguiente tabla con datos de ***training***, usando las expresiones de **Entropía** y **Cantidad de información**, cree un árbol de decisión.

Day	Outlook	Temperature	Humidity	Wind	Decision
1	Rain	Mild	Normal	Weak	Yes
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
4	Overcast	Cool	Normal	Strong	Yes
5	Rain	Cool	Normal	Strong	No
6	Rain	Cool	Normal	Weak	Yes
7	Rain	Mild	High	Weak	Yes
8	Overcast	Hot	High	Weak	Yes
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

Cálculo de la entropía de la variable *decision (clase)*

1. Elección del nodo raíz.

- ▶ Paso 1: Calcular la Entropía de la clase:

$$H(C) = - \sum_{c=1}^n p(c) \log_2 p(c)$$

$$H(Decision) = -p(Yes) \log_2 p(Yes) - p(No) \log_2 p(No)$$

$$H(Decision) = -\left(\frac{6}{10} \log_2 \frac{6}{10}\right) - \left(\frac{4}{10} \log_2 \frac{4}{10}\right) = 0,971$$

Cálculo de la entropía condicionada para la variable *Outlook*

- Paso 2 : Calcular la Entropía de clase condicionada a un valor:

Day	Outlook	Temperature	Humidity	Wind	Decision
1	Rain	Mild	Normal	Weak	Yes
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
4	Overcast	Cool	Normal	Strong	Yes
5	Rain	Cool	Normal	Strong	No
6	Rain	Cool	Normal	Weak	Yes
7	Rain	Mild	High	Weak	Yes
8	Overcast	Hot	High	Weak	Yes
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

$$H(C|X) = - \sum_{i=1}^n p(x|c) \log_2 p(c|x) = \\ - \sum_c \sum_x p(x, c) \log_2 p(c|x)$$

$$H(\text{Decision}| \text{Outlook}) = - p(\text{Rain}, \text{Yes}) \log_2 p(\text{Yes}|\text{Rain}) - p(\text{Rain}, \text{No}) \log_2 p(\text{No}|\text{Rain}) \\ - p(\text{Sun}, \text{Yes}) \log_2 p(\text{Yes}|\text{Sun}) - p(\text{Sun}, \text{No}) \log_2 p(\text{No}|\text{Sun}) \\ - p(\text{Overcast}, \text{Yes}) \log_2 p(\text{Yes}|\text{Overcast}) - p(\text{Overcast}, \text{No}) \log_2 p(\text{No}|\text{Overcast})$$

$$H(\text{Decision}| \text{Outlook}) = - \left(\frac{3}{10} \log_2 \frac{3}{4} \right) - \left(\frac{1}{10} \log_2 \frac{1}{4} \right) \\ - \left(\frac{1}{10} \log_2 \frac{1}{4} \right) - \left(\frac{3}{10} \log_2 \frac{3}{4} \right) \\ - \left(\frac{2}{10} \log_2 \frac{2}{2} \right) - \left(\frac{0}{10} \log_2 \frac{0}{2} \right) = 0,325 + 0,325 + 0 = 0,65$$

Cálculo de la entropía condicionada para la variable *Temperature*

Day	Outlook	Temperature	Humidity	Wind	Decision
1	Rain	Mild	Normal	Weak	Yes
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
4	Overcast	Cool	Normal	Strong	Yes
5	Rain	Cool	Normal	Strong	No
6	Rain	Cool	Normal	Weak	Yes
7	Rain	Mild	High	Weak	Yes
8	Overcast	Hot	High	Weak	Yes
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

$$\begin{aligned}
 H(\text{Decision} | \text{Temperature}) &= - p(\text{Mild}, \text{Yes}) \log_2 p(\text{Yes} | \text{Mild}) - p(\text{Mild}, \text{No}) \log_2 p(\text{No} | \text{Mild}) \\
 &\quad - p(\text{Cool}, \text{Yes}) \log_2 p(\text{Yes} | \text{Cool}) - p(\text{Cool}, \text{No}) \log_2 p(\text{No} | \text{Cool}) \\
 &\quad - p(\text{Hot}, \text{Yes}) \log_2 p(\text{Yes} | \text{Hot}) - p(\text{Hot}, \text{No}) \log_2 p(\text{No} | \text{Hot})
 \end{aligned}$$

$$\begin{aligned}
 H(\text{Decision} | \text{Temperature}) &= - \left(\frac{2}{10} \log_2 \frac{2}{3} \right) - \left(\frac{1}{10} \log_2 \frac{1}{3} \right) \\
 &\quad - \left(\frac{3}{10} \log_2 \frac{3}{4} \right) - \left(\frac{1}{10} \log_2 \frac{1}{4} \right) \\
 &\quad - \left(\frac{1}{10} \log_2 \frac{1}{3} \right) - \left(\frac{2}{10} \log_2 \frac{2}{3} \right) = 0,275 + 0,325 + 0,275 = 0,875
 \end{aligned}$$

Cálculo de la entropía condicionada para las variables *Humidity* y *Wind*

Day	Outlook	Temperature	Humidity	Wind	Decision
1	Rain	Mild	Normal	Weak	Yes
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
4	Overcast	Cool	Normal	Strong	Yes
5	Rain	Cool	Normal	Strong	No
6	Rain	Cool	Normal	Weak	Yes
7	Rain	Mild	High	Weak	Yes
8	Overcast	Hot	High	Weak	Yes
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

$$H(\text{Decision}|\text{Humidity}) = - p(\text{Normal}, \text{Yes}) \log_2 p(\text{Yes}|\text{Normal}) - p(\text{Normal}, \text{No}) \log_2 p(\text{No}|\text{Normal}) \\ - p(\text{High}, \text{Yes}) \log_2 p(\text{Yes}|\text{High}) - p(\text{High}, \text{No}) \log_2 p(\text{No}|\text{High})$$

$$H(\text{Decision}|\text{Humidity}) = - \left(\frac{4}{10} \log_2 \frac{4}{5} \right) - \left(\frac{1}{10} \log_2 \frac{1}{5} \right) \\ - \left(\frac{2}{10} \log_2 \frac{2}{5} \right) - \left(\frac{3}{10} \log_2 \frac{3}{5} \right) = 0,361 + 0,485 = 0,846$$

$$H(\text{Decision}|\text{Wind}) = - p(\text{Weak}, \text{Yes}) \log_2 p(\text{Yes}|\text{Weak}) - p(\text{Weak}, \text{No}) \log_2 p(\text{No}|\text{Weak}) \\ - p(\text{Strong}, \text{Yes}) \log_2 p(\text{Yes}|\text{Strong}) - p(\text{Strong}, \text{No}) \log_2 p(\text{No}|\text{Strong})$$

$$H(\text{Decision}|\text{Wind}) = - \left(\frac{5}{10} \log_2 \frac{5}{7} \right) - \left(\frac{2}{10} \log_2 \frac{2}{7} \right) \\ - \left(\frac{1}{10} \log_2 \frac{1}{3} \right) - \left(\frac{2}{10} \log_2 \frac{2}{3} \right) = 0,604 + 0,275 = 0,879$$

Cálculo de la ganancia de información

- ▶ Cantidad de información mutua
o ganancia de información (a maximizar):

$$I(C, X_i) = H(C) - H(C|X_i)$$

$$I(\text{Decision}, \text{Outlook}) = H(\text{Decision}) - H(\text{Decision}|\text{Outlook}) = 0,971 - 0,650 = 0,321$$

$$I(\text{Decision}, \text{Temperature}) = H(\text{Decision}) - H(\text{Decision}|\text{Temperature}) = 0,971 - 0,875 = 0,096$$

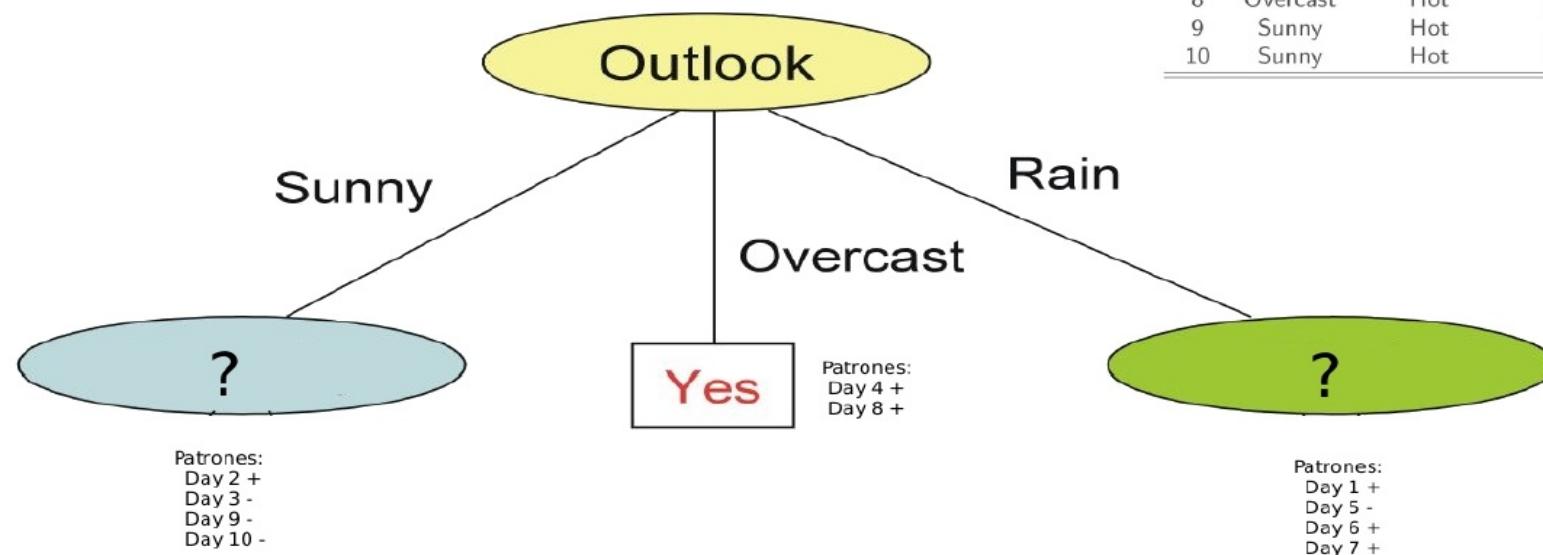
$$I(\text{Decision}, \text{Humidity}) = H(\text{Decision}) - H(\text{Decision}|\text{Humidity}) = 0,971 - 0,846 = 0,125$$

$$I(\text{Decision}, \text{Wind}) = H(\text{Decision}) - H(\text{Decision}|\text{Wind}) = 0,971 - 0,879 = 0,092$$

- ▶ Por tanto, el atributo que maximiza la ganancia de información es ***Outlook***

Day	Outlook	Temperature	Humidity	Wind	Decision
1	Rain	Mild	Normal	Weak	Yes
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
4	Overcast	Cool	Normal	Strong	Yes
5	Rain	Cool	Normal	Strong	No
6	Rain	Cool	Normal	Weak	Yes
7	Rain	Mild	High	Weak	Yes
8	Overcast	Hot	High	Weak	Yes
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

Árbol parcial obtenido



Day	Outlook	Temperature	Humidity	Wind	Decision
1	Rain	Mild	Normal	Weak	Yes
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
4	Overcast	Cool	Normal	Strong	Yes
5	Rain	Cool	Normal	Strong	No
6	Rain	Cool	Normal	Weak	Yes
7	Rain	Mild	High	Weak	Yes
8	Overcast	Hot	High	Weak	Yes
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

Nuevos conjuntos de datos

2. Se crean ***nuevos subconjuntos de training*** para cada rama.

Conjunto de datos para la rama ***Sunny***

Day	Outlook	Temperature	Humidity	Wind	Decision
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

Conjunto de datos para la rama ***Rain***

Day	Outlook	Temperature	Humidity	Wind	Decision
1	Rain	Mild	Normal	Weak	Yes
5	Rain	Cool	Normal	Strong	No
6	Rain	Cool	Normal	Weak	Yes
7	Rain	Mild	High	Weak	Yes

Cálculo de la entropía de la variable *decision (clase)*, rama *Sunny*

1. Elección del **nodo raíz rama Sunny**.

- ▶ Entropía de la clase:

$$H(C) = - \sum_{c=1}^n p(c) \log_2 p(c)$$

$$H(Decision) = -p(\text{Yes}) \log_2 p(\text{Yes}) - p(\text{No}) \log_2 p(\text{No})$$

$$H(Decision) = -\left(\frac{1}{4} \log_2 \frac{1}{4}\right) - \left(\frac{3}{4} \log_2 \frac{3}{4}\right) = 0,811$$

Day	Outlook	Temperature	Humidity	Wind	Decision
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

Cálculo de la entropía condicionada para la variable *Temperature*

- Entropía de la clase condicionada a un valor:

Day	Outlook	Temperature	Humidity	Wind	Decision
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

$$H(C|X) = - \sum_{i=1}^n p(x|c) \log_2 p(c|x) = \\ - \sum_c \sum_x p(x, c) \log_2 p(c|x)$$

$$H(\text{Decision} | \text{Temperature}) = - p(\text{Mild}, \text{Yes}) \log_2 p(\text{Yes} | \text{Mild}) - p(\text{Mild}, \text{No}) \log_2 p(\text{No} | \text{Mild}) \\ - p(\text{Cool}, \text{Yes}) \log_2 p(\text{Yes} | \text{Cool}) - p(\text{Cool}, \text{No}) \log_2 p(\text{No} | \text{Cool}) \\ - p(\text{Hot}, \text{Yes}) \log_2 p(\text{Yes} | \text{Hot}) - p(\text{Hot}, \text{No}) \log_2 p(\text{No} | \text{Hot})$$

$$H(\text{Decision} | \text{Temperature}) = - \left(\frac{0}{4} \log_2 \frac{0}{1} \right) - \left(\frac{1}{4} \log_2 \frac{1}{1} \right) \\ - \left(\frac{1}{4} \log_2 \frac{1}{1} \right) - \left(\frac{0}{4} \log_2 \frac{0}{1} \right) \\ - \left(\frac{0}{4} \log_2 \frac{0}{2} \right) - \left(\frac{2}{4} \log_2 \frac{2}{2} \right) = 0 + 0 + 0 = 0$$

Cálculo de la entropía condicionada para las variables *Humidity* y *Wind*

Day	Outlook	Temperature	Humidity	Wind	Decision
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

$$H(\text{Decision}|\text{Humidity}) = - p(\text{Normal}, \text{Yes}) \log_2 p(\text{Yes}|\text{Normal}) - p(\text{Normal}, \text{No}) \log_2 p(\text{No}|\text{Normal}) \\ - p(\text{High}, \text{Yes}) \log_2 p(\text{Yes}|\text{High}) - p(\text{High}, \text{No}) \log_2 p(\text{No}|\text{High})$$

$$H(\text{Decision}|\text{Humidity}) = - \left(\frac{1}{4} \log_2 \frac{1}{1} \right) - \left(\frac{0}{4} \log_2 \frac{0}{1} \right) \\ - \left(\frac{0}{4} \log_2 \frac{0}{3} \right) - \left(\frac{3}{4} \log_2 \frac{3}{3} \right) = 0 + 0 = 0$$

$$H(\text{Decision}|\text{Wind}) = - p(\text{Weak}, \text{Yes}) \log_2 p(\text{Yes}|\text{Weak}) - p(\text{Weak}, \text{No}) \log_2 p(\text{No}|\text{Weak}) \\ - p(\text{Strong}, \text{Yes}) \log_2 p(\text{Yes}|\text{Strong}) - p(\text{Strong}, \text{No}) \log_2 p(\text{No}|\text{Strong})$$

$$H(\text{Decision}|\text{Wind}) = - \left(\frac{1}{4} \log_2 \frac{1}{3} \right) - \left(\frac{2}{4} \log_2 \frac{2}{3} \right) \\ - \left(\frac{0}{4} \log_2 \frac{0}{1} \right) - \left(\frac{1}{4} \log_2 \frac{1}{1} \right) = 0,689 + 0 = 0,689$$

Cálculo de la ganancia de información

- ▶ Cantidad de información mutua o ganancia de información:

$$I(C, X_i) = H(C) - H(C|X_i)$$

Day	Outlook	Temperature	Humidity	Wind	Decision
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

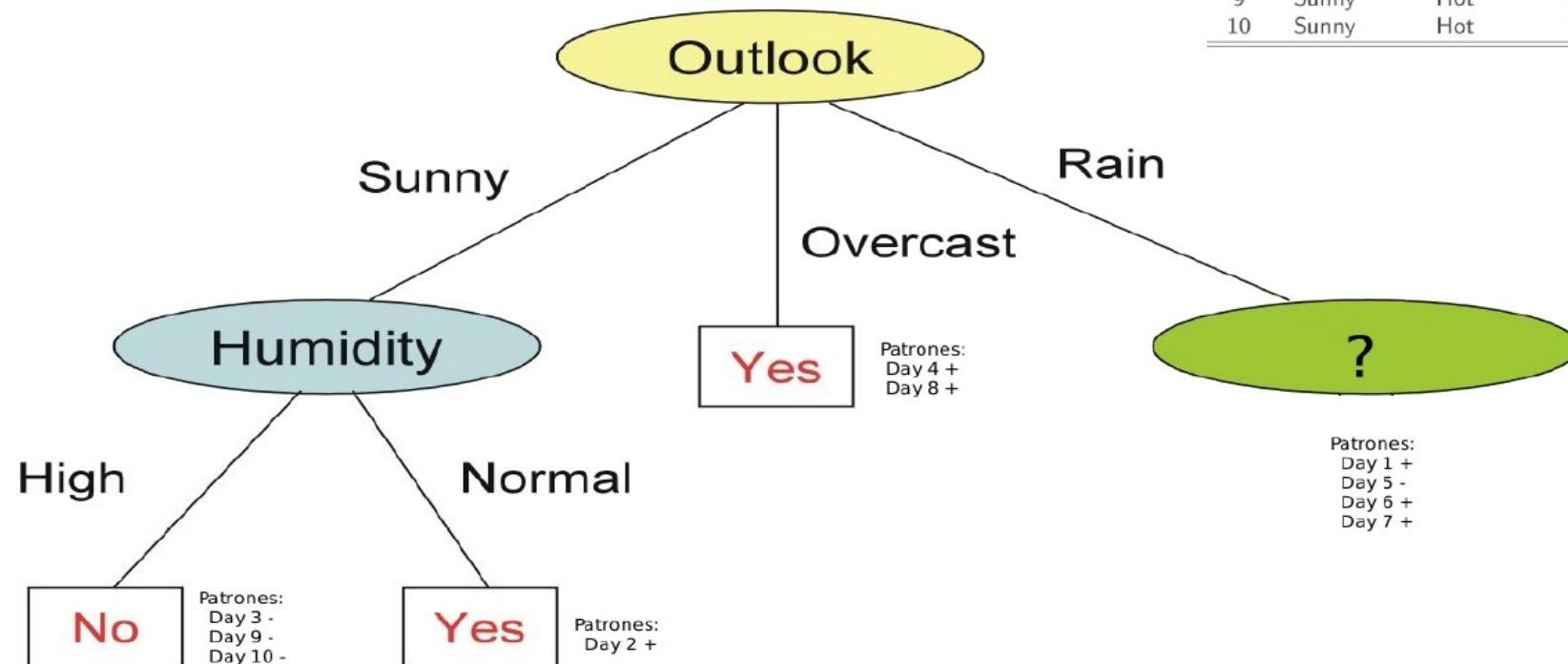
$$I(Decision, Temperature) = H(Decision) - H(Decision|Temperature) = 0,811 - 0 = 0,811$$

$$I(Decision, Humidity) = H(Decision) - H(Decision|Humidity) = 0,811 - 0 = 0,811$$

$$I(Decision, Wind) = H(Decision) - H(Decision|Wind) = 0,811 - 0,689 = 0,122$$

- ▶ En este caso los atributos que maximizan la ganancia de información son **Temperature** y **Humidity**.
- ▶ Se selecciona **Humidity** porque con sólo dos ramas clasifica a todos los patrones.

Nuevo árbol parcial obtenido



Day	Outlook	Temperature	Humidity	Wind	Decision
2	Sunny	Cool	Normal	Weak	Yes
3	Sunny	Mild	High	Weak	No
9	Sunny	Hot	High	Strong	No
10	Sunny	Hot	High	Weak	No

Cálculo de la entropía de la *variable decision (clase)*, rama Rain

Day	Outlook	Temperature	Humidity	Wind	Decision
1	Rain	Mild	Normal	Weak	Yes
5	Rain	Cool	Normal	Strong	No
6	Rain	Cool	Normal	Weak	Yes
7	Rain	Mild	High	Weak	Yes

1. Elección del **nodo raíz rama Rain**.

- Entropía de la clase:

$$H(C) = - \sum_{c=1}^n p(c) \log_2 p(c)$$

$$H(\text{Decision}) = -p(\text{Yes}) \log_2 p(\text{Yes}) - p(\text{No}) \log_2 p(\text{No})$$

$$H(\text{Decision}) = -\left(\frac{3}{4} \log_2 \frac{3}{4}\right) - \left(\frac{1}{4} \log_2 \frac{1}{4}\right) = 0,811$$

Cálculo de la entropía condicionada para la variable *Temperature*

- ▶ Entropía de la clase condicionada a un valor:

$$H(C|X) = - \sum_{i=1}^n p(x|c) \log_2 p(c|x) = \\ - \sum_c \sum_x p(x, c) \log_2 p(c|x)$$

Day	Outlook	Temperature	Humidity	Wind	Decision
1	Rain	Mild	Normal	Weak	Yes
5	Rain	Cool	Normal	Strong	No
6	Rain	Cool	Normal	Weak	Yes
7	Rain	Mild	High	Weak	Yes

$$H(\text{Decision} | \text{Temperature}) = - p(\text{Mild}, \text{Yes}) \log_2 p(\text{Yes} | \text{Mild}) - p(\text{Mild}, \text{No}) \log_2 p(\text{No} | \text{Mild}) \\ - p(\text{Cool}, \text{Yes}) \log_2 p(\text{Yes} | \text{Cool}) - p(\text{Cool}, \text{No}) \log_2 p(\text{No} | \text{Cool})$$

$$H(\text{Decision} | \text{Temperature}) = - \left(\frac{2}{4} \log_2 \frac{2}{2} \right) - \left(\frac{0}{4} \log_2 \frac{0}{2} \right) \\ - \left(\frac{1}{4} \log_2 \frac{1}{2} \right) - \left(\frac{1}{4} \log_2 \frac{1}{2} \right) = 0 + 0,5 = 0,5$$

Cálculo de la entropía condicionada para la variables *Humidity* y *Wind*

Day	Outlook	Temperature	Humidity	Wind	Decision
1	Rain	Mild	Normal	Weak	Yes
5	Rain	Cool	Normal	Strong	No
6	Rain	Cool	Normal	Weak	Yes
7	Rain	Mild	High	Weak	Yes

$$H(\text{Decision}|\text{Humidity}) = - p(\text{Normal}, \text{Yes}) \log_2 p(\text{Yes}|\text{Normal}) - p(\text{Normal}, \text{No}) \log_2 p(\text{No}|\text{Normal}) \\ - p(\text{High}, \text{Yes}) \log_2 p(\text{Yes}|\text{High}) - p(\text{High}, \text{No}) \log_2 p(\text{No}|\text{High})$$

$$H(\text{Decision}|\text{Humidity}) = - \left(\frac{2}{4} \log_2 \frac{2}{3} \right) - \left(\frac{1}{4} \log_2 \frac{1}{3} \right) \\ - \left(\frac{1}{4} \log_2 \frac{1}{1} \right) - \left(\frac{0}{4} \log_2 \frac{0}{1} \right) = 0,689 + 0 = 0,689$$

$$H(\text{Decision}|\text{Wind}) = - p(\text{Weak}, \text{Yes}) \log_2 p(\text{Yes}|\text{Weak}) - p(\text{Weak}, \text{No}) \log_2 p(\text{No}|\text{Weak}) \\ - p(\text{Strong}, \text{Yes}) \log_2 p(\text{Yes}|\text{Strong}) - p(\text{Strong}, \text{No}) \log_2 p(\text{No}|\text{Strong})$$

$$H(\text{Decision}|\text{Wind}) = - \left(\frac{3}{4} \log_2 \frac{3}{3} \right) - \left(\frac{0}{4} \log_2 \frac{0}{3} \right) \\ - \left(\frac{0}{4} \log_2 \frac{0}{1} \right) - \left(\frac{1}{4} \log_2 \frac{1}{1} \right) = 0 + 0 = 0$$

Cálculo de la ganancia de información

- ▶ Cantidad de información mutua o ganancia de información:

$$I(C, X_i) = H(C) - H(C|X_i)$$

Day	Outlook	Temperature	Humidity	Wind	Decision
1	Rain	Mild	Normal	Weak	Yes
5	Rain	Cool	Normal	Strong	No
6	Rain	Cool	Normal	Weak	Yes
7	Rain	Mild	High	Weak	Yes

$$I(\text{Decision}, \text{Temperature}) = H(\text{Decision}) - H(\text{Decision}|\text{Temperature}) = 0,811 - 0,5 = 0,311$$

$$I(\text{Decision}, \text{Humidity}) = H(\text{Decision}) - H(\text{Decision}|\text{Humidity}) = 0,811 - 0,689 = 0,122$$

$$I(\text{Decision}, \text{Wind}) = H(\text{Decision}) - H(\text{Decision}|\text{Wind}) = 0,811 - 0 = 0,811$$

- ▶ Por tanto, el atributo que maximiza la ganancia de información es **Wind**

Árbol definitivo obtenido

Árbol construido a partir del conjunto de entrenamiento o *training*.

