

ÁRBOLES DE DECISIÓN

Dado el siguiente conjunto de datos:

Edad	Antecedentes Familiares	Volumen Próstata (mm ³)	Valor PSA	Cáncer de Próstata
<45	N	< 60	< 4	N
<45	N	< 60	< 4	N
[45,70]	N	> 60	[4,10]	P
>70	N	> 60	> 10	P
>70	S	< 60	[4,10]	N
>70	S	< 60	>10	P
[45,70]	S	> 60	< 4	N
<45	N	< 60	[4,10]	N
<45	S	< 60	< 4	N
>70	S	< 60	< 4	N
<45	S	> 60	[4,10]	P
[45,70]	N	> 60	>10	P
[45,70]	S	< 60	[4,10]	P
>70	N	> 60	> 10	P

Tabla 1

Considerando los datos de la **tabla 1**, se pretende construir un clasificador basado en árboles de decisión.

Se definen las siguientes medidas:

Entropía o Cantidad de Información:

Si un experimento puede tener m resultados distintos: v_1, v_2, \dots, v_m que se pueden producir con probabilidades $P(v_1), P(v_2), \dots, P(v_m)$, la cantidad de información I que se obtiene al conocer el resultado real del experimento es:

$$I[P(v_1), P(v_2), \dots, P(v_m)] = \sum_{i=1}^m -P(v_i) \log_2 P(v_i)$$

Entropía Residual:

La información residual respecto a un atributo A es la cantidad de entropía que se debe reducir para clasificar una instancia después de emplear el atributo A .

$$I_{RES}(A) = \sum_{v \in V(A)} P(v) \times I(E_v)$$

Ganancia de Información:

La ganancia de información es la reducción esperada en la entropía tras la partición de los ejemplos de acuerdo a un atributo A .

$$G(A) = I - I_{RES}(A)$$

Se deben calcular los siguientes elementos:

- Entropía inicial del problema considerando las clases de salida:

$$I = -P(\text{Cáncer de Próstata}=P) \log_2(P(\text{Cáncer de Próstata}=P)) - P(\text{Cáncer de Próstata}=N) \log_2(P(\text{Cáncer de Próstata}=N)) = -(7/14) \log_2(7/14) - (7/14) \log_2(7/14) = 1$$

- Información residual y ganancia para cada atributo

Para el atributo Edad

Se determina en primer lugar las siguientes cantidades que se utilizarán para realizar los cálculos de entropías.

$$\text{Valores del atributo: } V(\text{Edad}) = \{< 45, [45,70], > 70\}$$

$$\text{Distribución de los ejemplos por clase: } E_{TOT} = [7N \quad 7P]$$

Distribución de los ejemplos para cada rango del atributo:

$$E_{<45} = [4N \quad 1P]$$

$$E_{[45,70]} = [1N \quad 3P]$$

$$E_{>70} = [2N \quad 3P]$$

$$I(E_{<45}) = -P(\text{Cáncer de Próstata}=N | \text{Edad}<45) \log_2(P(\text{Cáncer de Próstata}=N | \text{Edad}<45)) - P(\text{Cáncer de Próstata}=P | \text{Edad}<45) \log_2(P(\text{Cáncer de Próstata}=P | \text{Edad}<45)) = -(4/5) \log_2(4/5) - (1/5) \log_2(1/5) = 0.72$$

$$I(E_{[45,70]}) = -P(\text{Cáncer de Próstata}=N | \text{Edad}=[45,70]) \log_2(P(\text{Cáncer de Próstata}=N | \text{Edad}=[45,70])) - P(\text{Cáncer de Próstata}=P | \text{Edad}=[45,70]) \log_2(P(\text{Cáncer de Próstata}=P | \text{Edad}=[45,70])) = -(1/4) \log_2(1/4) - (3/4) \log_2(3/4) = 0.811$$

$$I(E_{>70}) = -P(\text{Cáncer de Próstata}=N | \text{Edad}>70) \log_2(P(\text{Cáncer de Próstata}=N | \text{Edad}>70)) - P(\text{Cáncer de Próstata}=P | \text{Edad}>70) \log_2(P(\text{Cáncer de Próstata}=P | \text{Edad}>70)) = -(2/5) \log_2(2/5) - (3/5) \log_2(3/5) = 0.97$$

$$G(\text{Edad}) = I - I_{RES}(\text{Edad}) = \sum_{v \in V(\text{Edad})} P(v) \times I(E_v) = 1 - (5/14) I(E_{<45}) - (4/14) I(E_{[45,70]}) - (5/14) I(E_{>70}) = 1 - (5/14) \cdot 0.72 - (4/14) \cdot 0.811 - (5/14) \cdot 0.97 = 0.165$$

Para el atributo Antecedentes Familiares

$$V(AF) = \{N, S\}$$

$$E_{TOT} = [7N \quad 7P]$$

$$E_N = [3N \quad 4P]$$

$$E_P = [4N \quad 3P]$$

$$I(E_N) = -P(\text{Cáncer de Próstata}=N | AF=N) \log_2(P(\text{Cáncer de Próstata}=N | AF=N)) - P(\text{Cáncer de Próstata}=P | AF=N) \log_2(P(\text{Cáncer de Próstata}=P | AF=N)) = -(3/7) \log_2(3/7) - (4/7) \log_2(4/7) = 0.985$$

$$I(E_P) = -P(\text{Cáncer de Próstata}=N \mid AF=P) \log_2(P(\text{Cáncer de Próstata}=N \mid AF=P)) - \\ P(\text{Cáncer de Próstata}=P \mid AF=P) \log_2(P(\text{Cáncer de Próstata}=P \mid AF=P)) = -(4/7) \log_2(4/7) - (3/7) \log_2(3/7) = 0.985$$

$$G(\text{Antecedent esFamiliar es}) = I - I_{RES}(\text{Antecedent esFamiliar es}) = \sum_{v \in V(AF)} P(v) \times I(E_v) \\ = 1 - (7/14) I(E_N) - (7/14) I(E_P) = 1 - (7/14) \cdot 0.985 - (7/14) \cdot 0.985 = 0.0147$$

Para el atributo Volumen Próstata

$$V(VP) = \{< 60, > 60\}$$

$$E_{TOT} = [7N \quad 7P]$$

$$E_{<60} = [6N \quad 2P]$$

$$E_{>60} = [1N \quad 5P]$$

$$I(E_{<60}) = -P(\text{Cáncer de Próstata}=N \mid VP<60) \log_2(P(\text{Cáncer de Próstata}=N \mid VP<60)) - \\ P(\text{Cáncer de Próstata}=P \mid VP<60) \log_2(P(\text{Cáncer de Próstata}=P \mid VP<60)) = -(6/8) \log_2(6/8) - (2/8) \log_2(2/8) = 0.811$$

$$I(E_{>60}) = -P(\text{Cáncer de Próstata}=N \mid VP>60) \log_2(P(\text{Cáncer de Próstata}=N \mid VP>60)) - \\ -P(\text{Cáncer de Próstata}=P \mid VP>60) \log_2(P(\text{Cáncer de Próstata}=P \mid VP>60)) = -(1/6) \log_2(1/6) - (5/6) \log_2(5/6) = 0.650$$

$$G(VP) = I - I_{RES}(VP) = \sum_{v \in V(VP)} P(v) \times I(E_v) = 1 - (8/14) I(E_{<60}) - (6/14) I(E_{>60}) = 1 - \\ (8/14) \cdot 0.811 - (6/14) \cdot 0.650 = 0.258$$

Para el atributo PSA

$$V(PSA) = \{< 4, [4,10], > 10\}$$

$$E_{TOT} = [7N \quad 7P]$$

$$E_{<4} = [5N \quad 0P]$$

$$E_{[4,10]} = [2N \quad 3P]$$

$$E_{>10} = [0N \quad 4P]$$

$$I(E_{<4}) = -P(\text{Cáncer de Próstata}=N \mid PSA<4) \log_2(P(\text{Cáncer de Próstata}=N \mid PSA <4)) - \\ P(\text{Cáncer de Próstata}=P \mid PSA <4) \log_2(P(\text{Cáncer de Próstata}=P \mid PSA <4)) = - \\ (5/5) \log_2(5/5) - (0/5) \log_2(0/5) = 0$$

$$I(E_{[4,10]}) = -P(\text{Cáncer de Próstata}=N \mid PSA=[4,10]) \log_2(P(\text{Cáncer de Próstata}=N \mid PSA = [4,10])) - \\ P(\text{Cáncer de Próstata}=P \mid PSA=[4,10]) \log_2(P(\text{Cáncer de Próstata}=P \mid PSA=[4,10])) = -(2/5) \log_2(2/5) - (3/5) \log_2(3/5) = 0.97$$

$$I(E_{>10}) = -P(\text{Cáncer de Próstata}=N \mid PSA >10) \log_2(P(\text{Cáncer de Próstata}=N \mid PSA >10)) - \\ -P(\text{Cáncer de Próstata}=P \mid PSA >10) \log_2(P(\text{Cáncer de Próstata}=P \mid PSA >10)) = -(0/4) \log_2(0/4) - (4/4) \log_2(4/4) = 0$$

$$G(PSA) = I - I_{RES}(PSA) = \sum_{v \in V(PSA)} P(v) \times I(E_v) = 1 - (5/14) I(E_{<4}) - (5/14) I(E_{[4,10]}) - (4/14) I(E_{>10}) = 1 - (5/14) \cdot 0 - (5/14) \cdot 0.97 - (4/14) \cdot 0 = 0.346$$

Considerando los valores de ganancia de información para cada atributo el algoritmo ID3 elegirá PSA como primera variable de separación.

El proceso de seleccionar un nuevo atributo y particionar los ejemplos de entrenamiento es repetido ahora para cada nodo descendiente no-terminal, usando esta vez sólo los ejemplos de entrenamiento asociados con este nodo. Los atributos que han sido incorporados más alto en el árbol son excluidos y por consiguiente, cualquier atributo puede aparecer a lo sumo una vez en cualquier paso del árbol.

El proceso continúa para cada nuevo nodo hoja, hasta que alguna de las siguientes 2 condiciones es alcanzada:

1. Todo atributo ya ha sido incluido en este paso a través del árbol.
2. Todos los ejemplos de entrenamiento asociados con este nodo hoja tienen el mismo valor de atributo objetivo (entropía = 0).

El árbol resultante es el siguiente:

