



3. Describir con precisión cómo se eligen, en el algoritmo ID3, los atributos que se van colocando en cada momento en el árbol.

4. Supongamos que tenemos un conjunto de observaciones  $\mathbf{d}$  y que tenemos un conjunto de hipótesis  $H = \{h_1, \dots, h_n\}$  que podrían explicar esas observaciones. Definir con precisión qué entendemos por hipótesis MAP y por hipótesis ML.

Apellidos: .....

Nombre: .....

**Ejercicio 2** (1.5 puntos): (**Aprendizaje inductivo**)

Una empresa suministra a un cliente información sobre novelas, las cuales a veces compra y a veces no. La siguiente tabla muestra las últimas ofertas de la empresa y si el cliente compró o no. Las ofertas dependen de los atributos *Idioma*, *Género*, *Precio* y *Edición*.

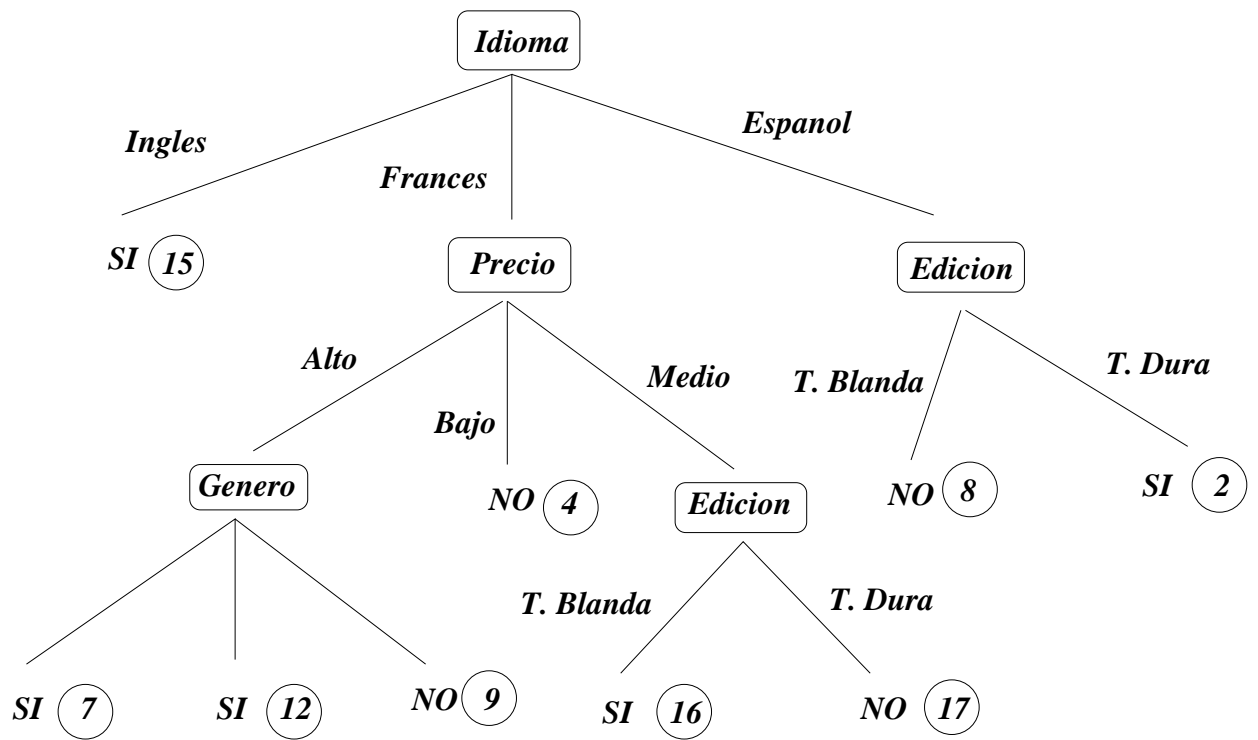
Ej.	IDIOMA	GÉNERO	PRECIO	EDICIÓN	Compra
1	Español	Aventuras	Alto	Tapa dura	Sí
2	Francés	Policíaco	Alto	Tapa blanda	Sí
3	Inglés	Aventuras	Medio	Tapa blanda	No
4	Francés	Histórico	Bajo	Tapa blanda	Sí
5	Francés	Aventuras	Alto	Tapa blanda	No
6	Español	Aventuras	Bajo	Tapa blanda	No
7	Francés	Histórico	Bajo	Tapa blanda	Sí
8	Inglés	Policíaco	Medio	Tapa blanda	No
9	Español	Aventuras	Bajo	Tapa dura	No
10	Francés	Aventuras	Bajo	Tapa blanda	No

- (a) Se pide usar el clasificador *Naive Bayes* con *m-estimación* usando 10 como *tamaño de muestreo equivalente* para clasificar la siguiente instancia

(Inglés, Histórico, Bajo, Tapa dura)

**Nota:** Hacer las cuentas tomando las 4 primeras cifras decimales.

- (b) ¿Cuál es el mínimo valor de  $m$  para que el algoritmo anterior mantuviera la clasificación obtenida en el apartado anterior? Razona tu respuesta.
- (c) Supongamos que hemos dividido un conjunto de 100 ejemplos en dos conjuntos. El primero, con 90 ejemplos lo hemos usado para crear un árbol de decisión y el segundo *Prueba*, con 10 ejemplos, es el que aparece en la tabla anterior. Supongamos que el árbol obtenido tras aplicar el algoritmo ID3 utilizando el conjunto de 90 ejemplos es



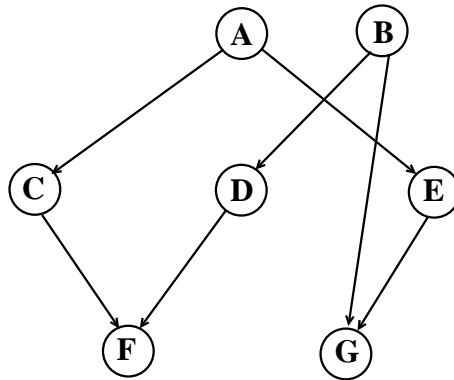
Junto a la clasificación de cada hoja aparece el número de elementos del conjunto  $D$  que verifica la condición, esto es, hay 4 ejemplos con el Idioma *Francés* y Precio *Bajo*, todos ellos con la clasificación *No* y hay 2 ejemplos con Idioma *Español* y Edición en *Tapa Dura*, ambos con la clasificación *Sí*. Se pide usar el ALGORITMO DE PODA PARA REDUCIR EL ERROR sobre el árbol usando el conjunto *Prueba*. **Especificar claramente cuál es el árbol obtenido.**

Apellidos: .....

Nombre: .....

**Ejercicio 3** (2.0 puntos): **(Incertidumbre)**

Consideremos la siguiente red bayesiana que expresa las dependencias existentes entre las variables aleatorias booleanas  $A, B, C, D, E, F$  y  $G$ :



con las siguientes tablas de distribución:

$P(a)$
0.3

$P(b)$
0.6

$A$	$P(c A)$
$a$	0.7
$\neg a$	0.1

$B$	$P(d B)$
$a$	0.1
$\neg a$	0.9

$A$	$P(e A)$
$a$	0.3
$\neg a$	0.8

$C$	$D$	$P(f C, D)$
$c$	$d$	0.9
$c$	$\neg d$	0.7
$\neg c$	$d$	0.5
$\neg c$	$\neg d$	0.2

$B$	$E$	$P(g B, E)$
$b$	$e$	0.2
$b$	$\neg e$	0.8
$\neg b$	$e$	0.2
$\neg b$	$\neg e$	0.9

Se pide:

1. Supongamos que la red se ha ido construyendo siguiendo el algoritmo de construcción de redes bayesianas visto en clase, considerando las variables en orden alfabético. ¿Qué relaciones de dependencia condicional se han ido suponiendo?
2. Según las dependencias *que se deducen* de la estructura de la red, decir si son ciertas o no las siguientes afirmaciones (justificando la respuesta):
  - Sabiendo el valor que toma  $A$ , el grado de creencia en que ocurra  $C$  no se ve actualizado si además sabemos el valor que toma  $G$ .
  - $F$  y  $G$  son condicionalmente dado  $A$
  - $P(F|A, B) = P(F|A, B, G)$
3. Aplicar el algoritmo de eliminación para calcular la probabilidad de que  $A$  sea falso, dado que se ha observado que  $F$  es falso y  $G$  es verdadero.

Apellidos: .....

Nombre: .....

**Ejercicio 4 (1.5 puntos): Procesamiento de lenguaje natural**

En este ejercicio, consideraremos el poema XXIX de *Proverbios y Cantares*, de D. Antonio Machado, publicado en 1912 dentro de su libro *Campos de Castilla*. El poema consta de 10 versos. Para nuestro propósito, consideraremos el poema dividido en cinco documentos, cada uno de ellos formado por dos versos:

 $D_1 \equiv$  Caminante, son tus huellas / el camino y nada más; $D_2 \equiv$  caminante, no hay camino, / se hace camino al andar. $D_3 \equiv$  Al andar se hace camino, / y al volver la vista atrás $D_4 \equiv$  se ve la senda que nunca / se ha de volver a pisar. $D_5 \equiv$  Caminante, no hay camino, / sino estelas en la mar.

Sean  $\vec{W}_1, \dots, \vec{W}_5$  las representaciones vectoriales de los documentos  $D_1, \dots, D_5$ , respecto del conjunto de términos  $T = \{ \text{camino, senda, andar, caminante} \}$ . *Nota* : Caminante y caminante son la misma palabra.

Consideremos el siguiente conjunto de entrenamiento a partir de los documentos  $D_1, \dots, D_4$

$$\mathcal{D} = \{ \langle D_1, \oplus \rangle, \langle D_2, \oplus \rangle, \langle D_3, \ominus \rangle, \langle D_4, \ominus \rangle \}$$

- (a) ¿Cuál es la clasificación del documento  $D_5$  obtenida a partir de  $\mathcal{D}$  mediante el algoritmo  $k$ -NN, con  $k = 3$  usando como medida de similitud de dos documentos el coseno del ángulo que forman sus representaciones vectoriales?
- (b) Aplicar el algoritmo de  $k$ -medias con  $k = 2$  sobre  $\vec{W}_1, \dots, \vec{W}_5$  a partir de los centros iniciales  $c_1 = (0, 0, 0, 0)$  y  $c_2 = (1, 1, 1, 1)$  usando la distancia euclídea.