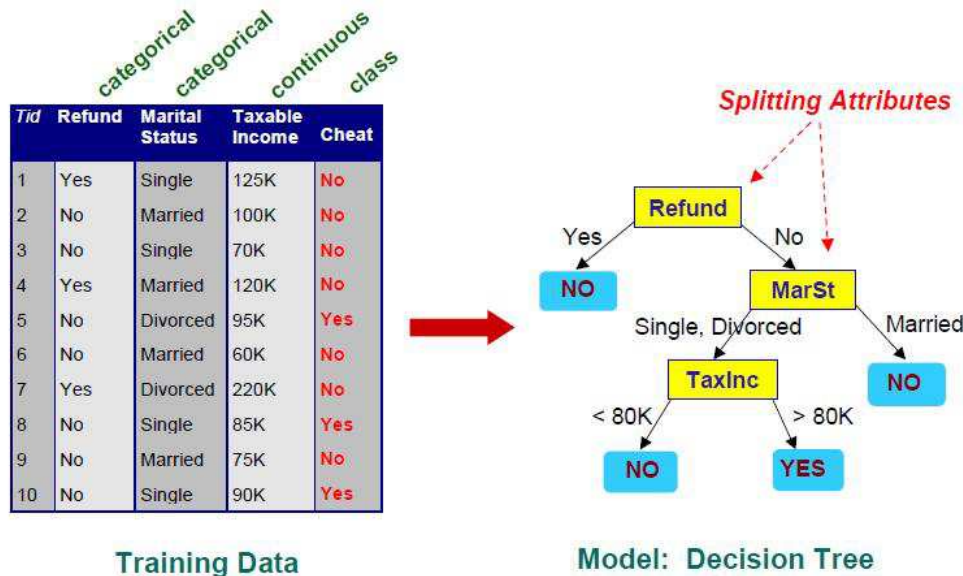


**Ejercicio 1.-** En la siguiente figura se muestra un árbol de decisión para predecir ejemplos acerca de si la persona miente. En el árbol de decisión, la raíz y los nodos internos contienen atributos con condiciones de prueba para separar los registros que tienen características diferentes. A todo nodo terminal se le asigna la clase Sí o No.

Refund = Devolución, Marital Status: Estado Civil, Taxable Income = Renta gravable en miles, Cheat= Miente. Splitting Attribute: Atributos de división



1) ¿Coincide el árbol de decisión con el construido utilizando la reducción de entropía?

$$H(C) = -\sum_c p(c) \log_2 p(c) = -(7/10) \log_2 (7/10) - (3/10) \log_2 (3/10) = 0,8813$$

#### Refund

$$\begin{aligned}
 H(C|R) &= -\sum_c \sum_r p(r,c) \log_2 p(c|r) = -P(NR,NC) \log_2 P(NC|NR) - \\
 &\quad -P(NR,SC) \log_2 P(SC|NR) - P(SC,NR) \log_2 P(NC|SR) - \\
 &\quad -P(SR,SC) \log_2 P(SC|SR) = -(4/10) \log_2 (4/7) - (3/10) \log_2 (3/7) - \\
 &\quad -(3/10) \log_2 (3/3) - (0/10) \log_2 (0/3) = 0,6896
 \end{aligned}$$

#### Marital Status

$$\begin{aligned}
 H(C|MS) &= -\sum_c \sum_r p(ms,c) \log_2 p(c|ms) = -P(S,NC) \log_2 P(NC|S) - \\
 &\quad -P(S,SC) \log_2 P(SC|S) - P(M,NC) \log_2 P(NC|M) - \\
 &\quad -P(M,SC) \log_2 P(SC|M) - P(D,NC) \log_2 P(NC|D) - \\
 &\quad -P(D,SC) \log_2 P(SC|D) = -(2/10) \log_2 (2/4) - (2/10) \log_2 (2/4) - \\
 &\quad -(4/10) \log_2 (4/4) - (0/10) \log_2 (0/4) - (1/10) \log_2 (1/2) - (1/10) \log_2 (1/2) = 0,6
 \end{aligned}$$

#### Taxable Income

$$\begin{aligned}
 H(C|TI) &= -\sum_c \sum_r p(ti,c) \log_2 p(c|ti) = -P(< 80,NC) \log_2 P(NC|<80) - \\
 &\quad -P(< 80,SC) \log_2 P(SC|<80) - P(> 80,NC) \log_2 P(NC|>80) - \\
 &\quad -P(> 80,SC) \log_2 P(SC|>80) = -(3/10) \log_2 (3/3) - (0/10) \log_2 (0/3) - \\
 &\quad -(4/10) \log_2 (4/7) - (3/10) \log_2 (3/7) = 0,6896
 \end{aligned}$$

El nodo raíz es Marital Status dado que es el que menos entropía condicional tiene. De esta forma la nueva base de datos a analizar es por una parte

| Entrenamiento | Refund | Marital Status | Taxable Income | Cheat |
|---------------|--------|----------------|----------------|-------|
| 1             | Yes    | Single         | >80K           | No    |
| 3             | No     | Single         | <80K           | No    |
| 8             | No     | Single         | >80K           | Yes   |
| 10            | No     | Single         | >80K           | Yes   |

Ahora

$$H(C) = -\sum_c p(c) \log_2 p(c) = -(2/4) * \log_2 (2/4) - (2/4) * \log_2 (2/4) = 1$$

### Refund

$$\begin{aligned} H(C | R) &= -\sum_c \sum_r p(r, c) \log_2 p(c | r) = -P(NR, NC) \log_2 P(NC | NR) - \\ &\quad - P(NR, SC) \log_2 P(SC | NR) - P(SR, NC) \log_2 P(NC | SR) \\ &\quad - P(SR, SC) \log_2 P(SC | SR) = -(1/4) \log_2 (1/3) - (2/4) \log_2 (2/3) - \\ &\quad -(1/4) \log_2 (1) - (0/4) \log_2 (0/1) = (3/4) \log_2 3 - (1/2) \log_2 2 \end{aligned}$$

### Taxable Income

$$\begin{aligned} H(C | TI) &= -\sum_c \sum_r p(ti, c) \log_2 p(c | ti) = -P(< 80, NC) \log_2 P(NC | <80) - \\ &\quad - P(< 80, SC) \log_2 P(SC | <80) - P(> 80, NC) \log_2 P(NC | >80) \\ &\quad - P(> 80, SC) \log_2 P(SC | >80) = -(1/4) \log_2 (1/1) - (0/4) \log_2 (0/1) - \\ &\quad -(1/4) \log_2 (1/3) - (2/4) \log_2 (2/3) = (3/4) \log_2 3 - (1/2) \log_2 2 \end{aligned}$$

Luego el siguiente nodo por la rama de la izquierda puede ser cualquiera de los dos. Si Consideramos la rama de la derecha del árbol, tenemos

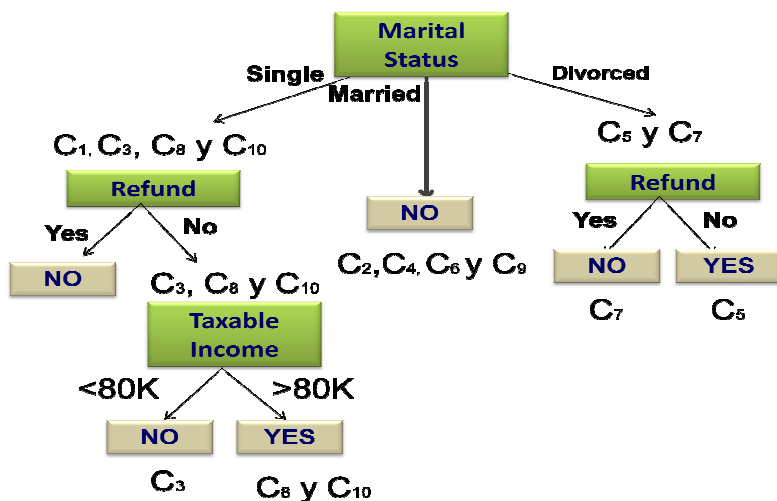
| Entrenamiento | Refund | Marital Status | Taxable Income | Cheat |
|---------------|--------|----------------|----------------|-------|
| 5             | No     | Divorced       | >80K           | Yes   |
| 7             | Yes    | Divorced       | >80K           | No    |

Ahora 
$$H(C) = -\sum_c p(c) \log_2 p(c) = -(1/2) * \log_2 (1/2) - (1/2) * \log_2 (1/2) = 1$$

$$\begin{aligned} H(C | R) &= -\sum_c \sum_r p(r, c) \log_2 p(c | r) = -P(NR, NC) \log_2 P(NC | NR) - \\ &\quad - P(NR, SC) \log_2 P(SC | NR) - P(SR, NC) \log_2 P(NC | SR) \\ &\quad - P(SR, SC) \log_2 P(SC | SR) = -(0/2) \log_2 (0/1) - (1/2) \log_2 (1/1) - \\ &\quad -(1/2) \log_2 (1/1) - (0/2) \log_2 (0/1) = 0 \end{aligned}$$

$$\begin{aligned} H(C | TI) &= -\sum_c \sum_r p(ti, c) \log_2 p(c | ti) = -P(< 80, NC) \log_2 P(NC | <80) - \\ &\quad - P(< 80, SC) \log_2 P(SC | <80) - P(> 80, NC) \log_2 P(NC | >80) \\ &\quad - P(> 80, SC) \log_2 P(SC | >80) = -(0/2) \log_2 (0/0) - (0/2) \log_2 (0/0) - \\ &\quad -(1/2) \log_2 (1/2) - (1/2) \log_2 (1/2) = 1 \end{aligned}$$

luego en este caso el siguiente nodo es Refund, y el árbol puede ser



**Ejercicio 2.-** Con el clasificador de la figura, clasificar los patrones del conjunto de test con las siguientes características

P1.- Refund = No, Marital Status: Married, Taxable income = 80K, Cheat=No

P2.- Refund = Si, Marital Status: Married, Taxable income = 180K, Cheat=Si

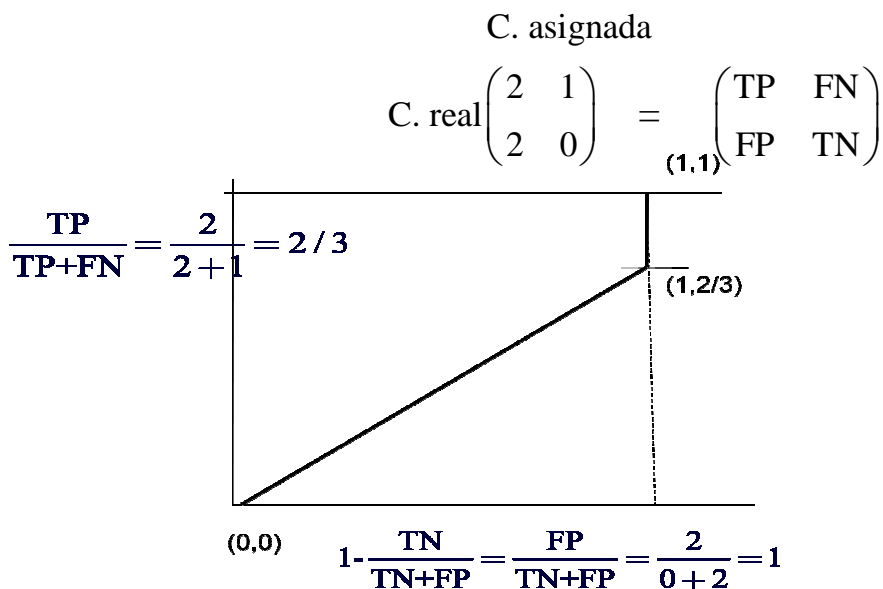
P3.- Refund = No, Marital Status: Divorced, Taxable income = 90K, Cheat=No

P4.- Refund = No, Marital Status: Married, Taxable income = 50K, Cheat=Si

P5.- Refund = Si, Marital Status: Single, Taxable income = 60K, Cheat=No

3) Construir la Matriz de Confusión y a partir de ella construir el área bajo la curva ROC, AUC. En función del resultado cómo calificaría al clasificador.

**Solución.-** Los patrones P1, P2, P4 y P5 el clasificador los clasifica en la clase No, a P3 en la clase Si. La matriz de confusión es de la forma.



Luego el  $AUC = 2/6 = 1/3$ , por lo que el clasificador es peor que un clasificador aleatorio.

**Ejercicio 3.-** Un médico con experiencia está construyendo un sistema de razonamiento basado en casos para automatizar una tarea de diagnóstico. Los casos se corresponden con personas individuales, donde sus datos se componen de una serie de características que describen los

posibles síntomas y la parte de solución representa el diagnóstico (clasificación de la enfermedad). La base de casos contiene 10 casos que se pueden ver en la siguiente tabla.

| Entrenamiento | Fiebre | Vómitos | Diarrea | Clasificación              |
|---------------|--------|---------|---------|----------------------------|
| C1            | alta   | si      | si      | Into. por Salmonella, (IS) |
| C2            | media  | no      | si      | Into. por Salmonella, (IS) |
| C3            | no     | si      | si      | Infección intestinal, (II) |
| C4            | media  | si      | no      | Infección intestinal, (II) |
| C5            | no     | no      | si      | Infección intestinal, (II) |
| C6            | alta   | si      | si      | Into. por Salmonella, (IS) |
| C7            | media  | no      | no      | Infección intestinal, (II) |
| C8            | alta   | si      | si      | Into. por Salmonella, (IS) |
| C9            | alta   | no      | si      | Into. por Salmonella, (IS) |
| C10           | alta   | no      | si      | Into. por Salmonella, (IS) |

Construya el árbol de decisión y clasifique un patrón con valores de las características  $q = (\text{alta}; \text{si}; \text{si})$ .

**Solución.-** La entropía asociada a la clase es

$$H(C) = -\sum_c^n p(c) \log_2 p(c) = -(6/10) * \log_2 (6/10) - (4/10) * \log_2 (4/10) = 0,97$$

Las entropías condicionadas a diferentes características de los pacientes son:

**Fiebre**

$$\begin{aligned} H(C | \text{Fiebre}) = & -\sum_c \sum_x p(x, c) \log_2 p(c | x) = -P(\text{Alta}, \text{IS}) \log_2 P(\text{IS} | \text{Alta}) - \\ & -P(\text{Alta}, \text{II}) \log_2 P(\text{II} | \text{Alta}) - P(\text{Media}, \text{IS}) \log_2 P(\text{IS} | \text{Media}) - \\ & -P(\text{Media}, \text{II}) \log_2 P(\text{II} | \text{Media}) - P(\text{No}, \text{IS}) \log_2 P(\text{IS} | \text{No}) - \\ & -P(\text{No}, \text{II}) \log_2 P(\text{II} | \text{No}) = -(5/10) \log_2 (5/5) - (0/10) \log_2 (0/5) - \\ & -(1/10) \log_2 (1/3) - (2/10) \log_2 (2/3) - (2/10) \log_2 (2/2) - (0/10) \log_2 (0/2) = \mathbf{0,27} \end{aligned}$$

**Vómitos**

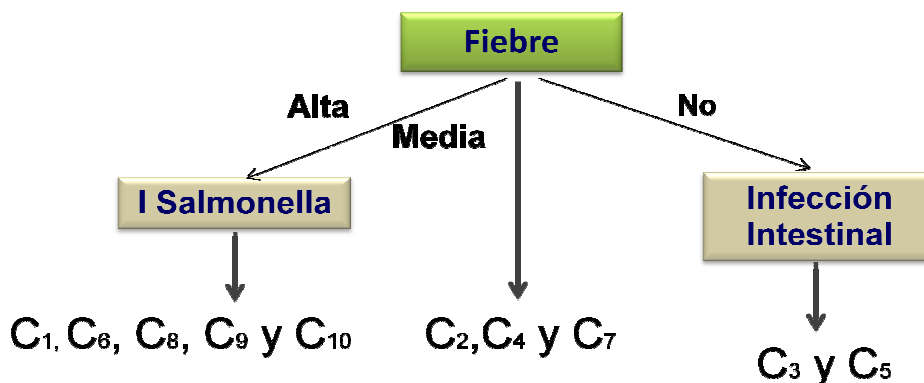
$$\begin{aligned} H(C | \text{Vómitos}) = & -P(\text{SI}, \text{IS}) \log_2 P(\text{IS} | \text{SI}) - P(\text{SI}, \text{II}) \log_2 P(\text{II} | \text{SI}) \\ & -P(\text{NO}, \text{IS}) \log_2 P(\text{IS} | \text{NO}) - P(\text{NO}, \text{II}) \log_2 P(\text{II} | \text{NO}) = \\ & -(3/10) \log_2 (3/5) - (2/10) \log_2 (2/5) - (3/10) \log_2 (3/5) \\ & -(2/10) \log_2 (2/5) = 0,97 \end{aligned}$$

**Diarrea**

$$\begin{aligned} H(C | \text{Diarrea}) = & -P(\text{SI}, \text{IS}) \log_2 P(\text{IS} | \text{SI}) - P(\text{SI}, \text{II}) \log_2 P(\text{II} | \text{SI}) \\ & -P(\text{NO}, \text{IS}) \log_2 P(\text{IS} | \text{NO}) - P(\text{NO}, \text{II}) \log_2 P(\text{II} | \text{NO}) = \\ & -(6/10) \log_2 (6/8) - (2/10) \log_2 (2/8) - (0/10) \log_2 (0/2) \\ & -(2/10) \log_2 (2/2) = 0,66 \end{aligned}$$

| Entrenamiento | Fiebre | Vómitos | Diarrea | Clasificación              |
|---------------|--------|---------|---------|----------------------------|
| C2            | media  | no      | si      | Into. por Salmonella, (IS) |

|    |       |    |    |                            |
|----|-------|----|----|----------------------------|
| C4 | media | si | no | Infección intestinal, (II) |
| C7 | media | no | no | Infección intestinal, (II) |



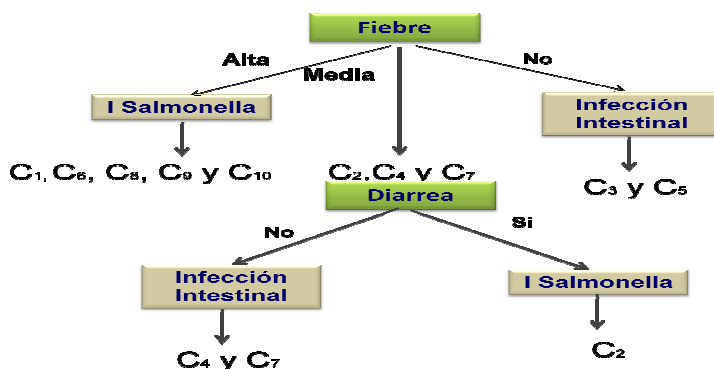
Ahora volvemos a rehacer los cálculos, pero con la nueva estructura de datos

$$H(C) = -\sum_c^n p(c) \log_2 p(c) = -(1/3) * \log_2 (1/3) - (2/3) * \log_2 (2/3) = 0,91$$

$$\begin{aligned}
 H(C|Vómitos) &= -P(SI,II) \log_2 P(II|SI) \\
 &\quad -P(NO,IS) \log_2 P(IS|NO) - P(NO,II) \log_2 P(II|NO) = \\
 &\quad -(1/3) \log_2 (1/2) - (1/3) \log_2 (1/2) - (1/3) \log_2 (1) = 0,66
 \end{aligned}$$

$$\begin{aligned}
 H(C|Diarrea) &= -P(SI,IS) \log_2 P(IS|SI) - P(NO,II) \log_2 P(II|NO) = \\
 &\quad -(1/3) \log_2 (1) - (2/3) \log_2 (1) = 0
 \end{aligned}$$

De nuevo se elige el que tiene menos entropía ya que hará máxima la cantidad de información mutua, en este caso Diarrea, por lo que el árbol de decisión será (Ver figura) y el patrón q = (alta; si; si) se clasificará en I Salmonella dado que su primera característica al ser Alta lo sitúa en la citada clase



**Ejercicio 4.-** Un médico con experiencia está construyendo un sistema de razonamiento basado en casos para automatizar una tarea de diagnóstico. Los casos se corresponden con personas individuales, donde sus datos se componen de una serie de características que describen los posibles síntomas y la parte de solución representa el diagnóstico

(clasificación de la enfermedad). La base de casos contiene 10 casos que se pueden ver en la siguiente tabla.

| Entrenamiento   | Fiebre | Vómitos | Clasificación              |
|-----------------|--------|---------|----------------------------|
| C <sub>1</sub>  | alta   | si      | Into. por Salmonella, (IS) |
| C <sub>2</sub>  | media  | no      | Into. por Salmonella, (IS) |
| C <sub>3</sub>  | no     | si      | Otros (O)                  |
| C <sub>4</sub>  | media  | si      | Infección intestinal, (II) |
| C <sub>5</sub>  | no     | no      | Infección intestinal, (II) |
| C <sub>6</sub>  | media  | no      | Infección intestinal, (II) |
| C <sub>7</sub>  | alta   | si      | Into. por Salmonella, (IS) |
| C <sub>8</sub>  | alta   | no      | Otros (O)                  |
| C <sub>9</sub>  | no     | si      | Into. por Salmonella, (IS) |
| C <sub>10</sub> | alta   | no      | Into. por Salmonella, (IS) |

Construya el árbol de decisión y clasifique un patrón con valores de las características q = (alta; no).

#### Solución.-

La entropía asociada a la clase es

$$H(C) = -\sum_c^n p(c) \log_2 p(c) = -(5/10) * \log_2 (5/10) - (2/10) * \log_2 (2/10) - (3/10) * \log_2 (3/10) = 2,701$$

#### Fiebre

$$\begin{aligned} H(C | F) &= -\sum_c \sum_r p(f, c) \log_2 p(c | f) = -P(alta, IS) \log_2 P(IS|alta) - \\ &\quad - P(media, IS) \log_2 P(IS|media) - P(no, IS) \log_2 P(IS|no) \\ &\quad - P(alta, O) \log_2 P(O|alta) - P(media, O) \log_2 P(O|media) \\ &\quad - P(no, O) \log_2 P(O|no) - P(alta, II) \log_2 P(II|alta) \\ &\quad - P(media, II) \log_2 P(II|media) - P(no, II) \log_2 P(II|no) = \\ &\quad -(3/10) \log_2 (3/4) - (1/10) \log_2 (1/3) - (1/10) \log_2 (1/3) - (1/10) \log_2 (1/4) \\ &\quad - 0 - (1/10) \log_2 (1/3) - 0 - (2/10) \log_2 (2/3) - (1/10) \log_2 (1/3) = \mathbf{1,34} \end{aligned}$$

#### Vómitos

$$\begin{aligned} H(C | V) &= -\sum_c \sum_r p(v, c) \log_2 p(c | v) = -P(Si, IS) \log_2 P(IS|Si) - \\ &\quad - P(Si, O) \log_2 P(O|Si) - P(Si, II) \log_2 P(II|Si) \\ &\quad - P(No, IS) \log_2 P(IS|No) - P(No, O) \log_2 P(O|No) \\ &\quad - P(No, II) \log_2 P(II|No) = -(3/10) \log_2 (3/5) - (1/10) \log_2 (1/5) - \\ &\quad - (1/10) \log_2 (1/5) - (2/10) \log_2 (2/5) - (1/10) \log_2 (1/5) - (2/10) \log_2 (2/5) = \mathbf{1,445} \end{aligned}$$

El nodo raíz es Fiebre dado que es el que menos entropía condicional tiene. De esta forma la nueva base de datos a analizar es para **Fiebre alta**

| Entrenamiento | Fiebre      | Vómitos   | Clasificación                     |
|---------------|-------------|-----------|-----------------------------------|
| $C_1$         | <b>alta</b> | <b>si</b> | <b>Into. por Salmonella, (IS)</b> |
| $C_7$         | <b>alta</b> | <b>si</b> | <b>Into. por Salmonella, (IS)</b> |
| $C_8$         | alta        | no        | Otros (O)                         |
| $C_{10}$      | alta        | no        | Into. por Salmonella, (IS)        |

Para fiebre media es

| Entrenamiento | Fiebre       | Vómitos   | Clasificación                     |
|---------------|--------------|-----------|-----------------------------------|
| $C_2$         | media        | no        | Into. por Salmonella, (IS)        |
| $C_4$         | <b>media</b> | <b>si</b> | <b>Infección intestinal, (II)</b> |
| $C_6$         | media        | no        | Infección intestinal, (II)        |

y para fiebre NO es

| Entrenamiento | Fiebre    | Vómitos   | Clasificación                     |
|---------------|-----------|-----------|-----------------------------------|
| $C_3$         | no        | si        | Otros (O)                         |
| $C_5$         | <b>no</b> | <b>no</b> | <b>Infección intestinal, (II)</b> |
| $C_9$         | no        | si        | Into. por Salmonella, (IS)        |