

$$\log_2 3 = \log_2 \frac{3}{10} = \log_2 3 - \log_2 10 = 0.477 - 1.030 = -0.553$$

2.

(Căştigul de informaţie — determinarea celui mai „bun” atribut; arbori de decizie consistenţi cu un set de date)



Fie attributele binare X_1, X_2, X_3, X_4 (valorile lor pot fi doar T sau F), precum şi două tipuri de etichete, 0 şi 1. Veţi considera cele 8 instanţe din tabelul alăturat.

Eticheta instanţei	X_1	X_2	X_3	X_4
1	T	T	T	F
1	T	T	T	F
1	F	T	T	F
0	T	T	F	F
0	T	T	F	F
0	F	T	F	F
0	F	T	F	F
0	F	T	T	F

a. Vrem să învăţăm un arbore de decizie din acest set de exemple. În vederea selectării celui mai bun candidat pentru nodul rădăcină, calculaţi câştigul de informaţie pentru fiecare atribut X_i , cu $i = 1, \dots, 4$. Ce atribut veţi selecta?

b. Există oare un arbore de decizie care poate clasifica în mod perfect instanţele date? În cazul afirmativ, desenaţi acel arbore de decizie. În caz contrar, daţi o explicaţie simplă.

a) $P_k = \text{câştigul de informaţie}$, nu beluie

$$H(E) = H[X_1, X_2, X_3, X_4] = \frac{1}{8} \log_2 8 = 1 \text{ (câştig de inf e 0)}$$

$$I_G(E, X_1) = H(E) - \left(\frac{1}{2} H[2+2] + \frac{1}{2} H[4+4-] \right) = 1 - \left(\frac{1}{2} + \frac{1}{2} \right) = 0$$

$$I_G(E, X_2) = H(E) - \left(\frac{1}{2} H[1+3] + \frac{1}{2} H[4+4-] \right) = 1 - 1 = 0$$

$$I_G(E, X_3) = H(E) - \left(\frac{1}{2} H[1+3] + \frac{1}{2} H[4+4-] \right) =$$

$$= 1 - \left(\frac{1}{2} \left(\frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4} \right) + \frac{1}{2} \left(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4} \right) \right) =$$

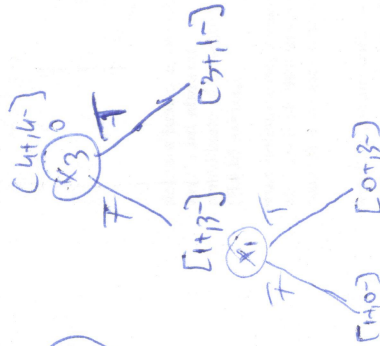
$$= 1 - \left(\frac{1}{2} \left(\frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4} \right) + \frac{1}{2} \left(\frac{1}{2} + \frac{3}{2} \left(\log_2 \frac{1}{2} - \log_2 \frac{3}{4} \right) \right) \right) =$$

$$= 1 - \left(\frac{1}{2} + \frac{3}{4} \left(2 - 1.58 \right) \right) = 1 - \left(\frac{1}{2} + \frac{3}{4} \cdot 0.42 \right) =$$

$$= 1 - (0.5 + 0.315) = 0.185$$

$$I_G(E, X_4) = H(E) - \left(\frac{1}{10} H[4+4-] + \frac{9}{10} H[1+3] \right) = 0$$

Alegem X_3 ca rădăcină



— Alurech de decizie nu poate clasifica perfect
datele de antrenament deoarece datale nu
sunt consistente adică:

E	X_1	X_2	X_3	X_4
0	F	T	F	F
1	F	T	F	F

din instanţele v şi w etichetate

— nodul rădăcină nu poate clasifica perfect
datele de antrenament deoarece datale nu
sunt consistente adică: