

3.

(CCCCCalcularea unor entropii;
aplicarea algoritmului ID3)

Centrul [de Medicină] pentru Controlul și Prevenția Maladiilor a fost sesizat în legătură cu o creștere surprinzătoare a aparițiilor de vampiri. Acest centru a strâns anumite date preliminare privind atât caracteristicile vampirilor [deja] cunoscuți cât și ale non-vampirilor, și dorește să construiască un arbore de decizie ca să-i ajute pe cetățeni să identifice noi vampiri. Datele culese sunt prezentate în tabelul următor:

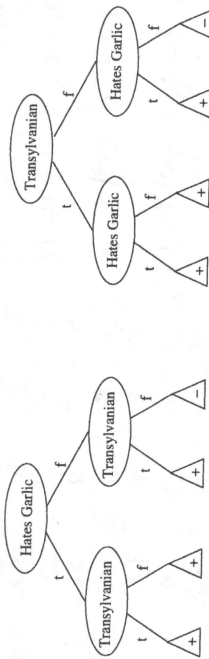
V	G	T	Nr. apariții
(Vampire)	(Hates Garlic)	(Transylvanian)	
+	t	t	9
+	t	f	4
+	f	t	3
+	f	f	0
-	t	t	1
-	t	f	3
-	f	t	1
-	f	f	6

Fiecare linie indică ce caracteristici au fost „observate”, și de câte ori a fost observată fiecare combinație de caracteristici. De exemplu, combinația (+, t, t) a fost observată de 9 ori, pe când combinația (-, f, f) n-a fost observată niciodată. (Simbolii 't' și 'f' au fost folosiți în locul lui True și False, pentru a evita confuzia cu atributul T.)

a. Calculați entropia condițională medie $H(V|G)$. (Toate calculele intermediare trebuie făcute cu o precizie de cel puțin 4 zecimale, pentru a ne asigura că răspunsul final are primele 3 zecimale corecte.)

b. Calculați entropia condițională medie $H(V|T)$. (Din nou, toate calculele intermediare trebuie făcute cu o precizie de cel puțin 4 zecimale.)

c. Care dintre arborii de mai jos reprezintă rezultatul învățării realizate de algoritmul ID3 pe aceste date?



d. Adevărat sau Fals: Arborele produs de către ID3 va clasifica o persoană căreia nu-i place usturoiul (engl., hates garlic) dar nu este transilvănean ca fiind vampir.

Indicație: Este posibil să aveți nevoie de următoarele valori pentru entropia $H(p)$ unei variabile aleatoare Bernoulli de parametru p : $H(11/27) = 0.9751$, $H(4/17) = 0.7871$, $H(3/10) = 0.8812$, $H(1/7) = 0.5916$, $H(4/13) = 0.8904$.

a.
$$H(V|G) = \frac{1}{8} H(2+2-3) + \frac{1}{8} H(2+2-3) = 1$$

$$[2+2-3] \uparrow \text{gărit}$$

a.
$$H(V|G) = \frac{17}{27} \cdot H(13+4-7) + \frac{10}{27} \cdot H(3+7-7) = \frac{17}{27} \cdot \left(\frac{13}{17} \cdot \log_2 \frac{17}{13} + \frac{4}{17} \cdot \log_2 \frac{17}{4} \right) + \frac{10}{27} \cdot \left(\frac{3}{10} \cdot \log_2 \frac{10}{3} + \frac{7}{10} \cdot \log_2 \frac{10}{7} \right) = \frac{13 \log_2 \frac{17}{13} + 4 \log_2 \frac{17}{4} + 3 \log_2 \frac{10}{3} + 7 \log_2 \frac{10}{7}}{27} = 0.822$$

b.
$$H(V|T) = \frac{14}{27} H(12+2-7) + \frac{13}{27} H(4+9-7) = \frac{14}{27} \cdot \left(\frac{12}{14} \cdot \log_2 \frac{14}{12} + \frac{2}{14} \cdot \log_2 \frac{14}{2} \right) + \frac{13}{27} \cdot \left(\frac{4}{13} \cdot \log_2 \frac{13}{4} + \frac{9}{13} \cdot \log_2 \frac{13}{9} \right) = \frac{12 \log_2 \frac{14}{12} + 2 \log_2 \frac{14}{2} + 4 \log_2 \frac{13}{4} + 9 \log_2 \frac{13}{9}}{27} = 0.7356$$

c. Conform subpunctelor a și b, în rădăcina va fi ales atributul T (Transylvanian) deoarece entropia sa este mai mică decât a atributului G , ceea ce înseamnă că va avea câștig de informație mai mare. Adărar arborele ID3 este cel de-al doilea.

d. Adevărat. Se mărește pe arbore ramura f pentru atributul Transylvanian și ramura t pentru atributul Hates Garlic.