

Name : Riddha Majumder

ID : 23341023

Sec : 13

Course : C&E 422

Assignment : 04

Ans to the question no. 1

(a)

Dog_id is not significant for determining whether the dog has a disease, as it is not a feature describing the dog's characteristics. It's just a label

(b)

Here, let's calculate

$$\begin{aligned} E(\text{Disease}) &= -P(\text{yes}) \log_2 P(\text{yes}) - P(\text{No}) \log_2 P(\text{No}) \\ &= -4/8 \log_2 4/8 - 4/8 \log_2 4/8 \Rightarrow 1 \end{aligned}$$

$$\begin{aligned} E(\text{Black}) &= -P(Y|B) \log_2 P(Y|B) - P(N|B) \log_2 P(N|B) \\ &\Rightarrow -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \Rightarrow 0.970 \end{aligned}$$

$$\begin{aligned} E(\text{White}) &= -P(Y|W) \log_2 P(Y|W) - P(N|W) \log_2 P(N|W) \\ &\Rightarrow -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \Rightarrow 0.918 \end{aligned}$$

$$\begin{aligned} \therefore IG(\text{Fur color}) &= \text{Entropy}(D) - P(B) \cdot E(B) - P(W) \cdot E(W) \\ &= 1 - 5/8 \times 0.970 - 3/8 \times 0.918 \Rightarrow 0.0495 \end{aligned}$$

$$\begin{aligned} E(\text{Large}) &= -P(Y|L) \log_2 P(Y|L) - P(N|L) \log_2 P(N|L) \\ &= -1/4 \log_2 1/4 - 3/4 \log_2 3/4 \Rightarrow 0.811 \end{aligned}$$

$$\begin{aligned} E(\text{Small}) &= -P(Y|S) \log_2 P(Y|S) - P(N|S) \log_2 P(N|S) \\ &= -3/4 \log_2 3/4 - 1/4 \log_2 1/4 \Rightarrow 0.811 \end{aligned}$$

$$\begin{aligned} IG(\text{size}) &\Rightarrow E(D) - P(L) E(L) - P(S) E(S) \\ &\Rightarrow 1 - 4/8 \times 0.811 - 4/8 \times 0.811 \Rightarrow 0.189 \end{aligned}$$

∴ The tail length in sorted order will be

1.2, 1.4, 2.2, 2.3, 3.5, 3.8, 4.2, 5.6 [n=8]

Now, we calculate the median $\Rightarrow \frac{\left(\frac{8}{2}\right)^{th} + \left(\frac{8}{2} + 1\right)^{th}}{2}$

$$\Rightarrow \frac{2.3 + 3.5}{2}$$

$$\Rightarrow 2.9$$

As we can see, the value is less than 2.9. So, it will be short, or else long

Dog_ID	Fur color	size	Tail-length	disease
1	Black	Large	Long	No
28	White	Large	Short	Yes
3	Black	Small	Long	Yes
34	Black	Small	Long	Yes
26	Black	Large	Short	No
11	White	Small	Short	No
32	Black	Small	Long	Yes
13	White	Large	Short	No

$$\begin{aligned} \text{Now, } E(\text{long}) &\Rightarrow -P(Y=1) \log_2 P(Y=1) - P(N=1) \log_2 P(N=1) \\ &\Rightarrow -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \Rightarrow 0.811 \end{aligned}$$

$$E(\text{short}) \Rightarrow -\frac{1}{4} \log_2 \left(\frac{1}{4}\right) - \frac{3}{4} \log_2 \left(\frac{3}{4}\right) \Rightarrow 0.811$$

$$\therefore IG(\text{Tail length}) \Rightarrow E(D) - P(L) E(L) - P(S) E(S)$$

$$\Rightarrow 1 - 0.811 \times \frac{4}{8} - 0.811 \times \frac{4}{8}$$

$$\Rightarrow 0.189$$

So, we got 20 IG's. Now, we can say by looking them, two features among three has high information gain and equal. So, I will say, size is the best suited for the root node.

(c)

Given, $n=8$

\therefore Maximum possible value for entropy will be $H(X)=3$ and this only happens when all the possible outcomes have equal probability that is $1/8$.

$$\begin{aligned}\therefore \max H(X) &= \sum_1^8 -1/8 \log_2 1/8 \\ &= 3\end{aligned}$$

(Ans:)