

Question NO 1:-

(2)

As the training set given to us is.

C_1	C_2	C_3	out
1	1	1	1
1	0	0	1
1	1	0	0
0	0	1	0

As we know entropy is given by

$$E(S) = \sum_{i=1}^K -p_i \log_2 p_i$$

$$\begin{aligned} \text{so } E(C_1) &= \left[\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} \right] \frac{3}{4} + \frac{1}{4} [-1 \log 1] \\ &= (0.918) \frac{3}{4} + 0 = 0.688 \end{aligned}$$

$$\begin{aligned} E(C_2) &= \frac{2}{4} \left[-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right] + \frac{2}{4} \left[-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right] \\ &= \frac{2}{4} [1] + \frac{2}{4} [1] = \frac{2}{4} \\ &= 1 \end{aligned}$$

$$E(c_2) = \frac{2}{4} \left[-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right] + \frac{2}{4} \left[-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right]$$

$$= \frac{2}{4} [1] + \frac{2}{4} [1]$$

$$= 1$$

For total Entropy

$$E(\text{out}) = \left[\frac{2}{4} \log \frac{2}{4} \right] + \left[-\frac{2}{4} \log \frac{2}{4} \right]$$

$$= 1$$

$$\begin{aligned} IG(\text{out}, c_1) &= E(\text{out}) - E(c_1) \\ &= 1 - 0.688 = 0.312 \end{aligned}$$

$$\begin{aligned} IG(\text{out}, c_2) &= E(\text{out}) - E(c_2) \\ &= 1 - 1 = 0 \end{aligned}$$

$$\begin{aligned} IG(\text{out}, c_3) &= E(\text{out}) - E(c_3) \\ &= 1 - 1 = 0 \end{aligned}$$

As gain of C_1 is the best & has the highest information gain so we will select this.

$$E(C_2) = \left[\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right] \frac{2}{3} +$$

$$[-1 \log 1] \frac{1}{3}.$$

$$= 0.666...$$

$$E(C_3) = \frac{2}{3} \left[-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right] + \frac{1}{3} [-1 \log 1]$$

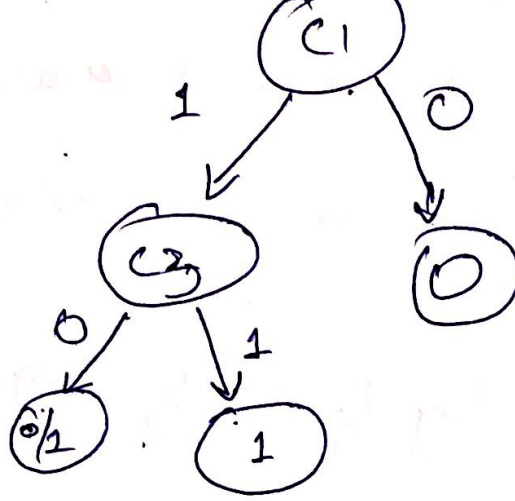
$$= 0.666...$$

So our root node would be having the max info gain which is C_1 & the next attribute we could choose any C_2 or C_3 as they have equal entropy. As if we have $C_1 = 0$ then.

$$E(C_2) = -1 \log 1 = 0$$

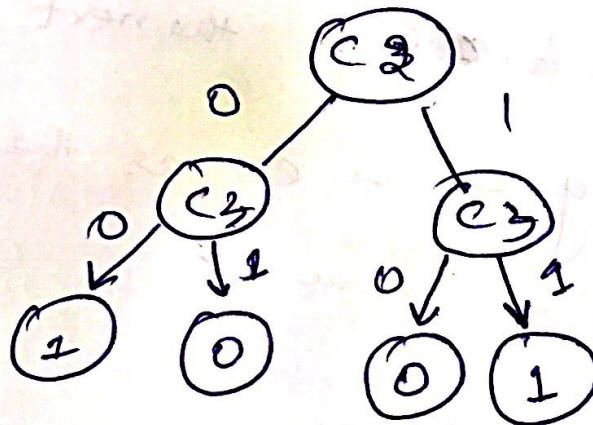
$$E(C_3) = -1 \log 1 = 0.$$

So A-T-G-C



It will give us error in a case when, $C1$ is 1 & $C2$ is zero then is $1/4$ probability that will occur in training data. so this is the training error that could occur.

6) Following is the decision tree of depth 2.



The above tree has training error zero as it satisfies all the given conditions in our data. & has depth 2 as well.

QUESTION NO 2 :-

(3)

As we know we will use InfoEntropy Algorithm.

$$E(S) = -p \log_2 p.$$

$$\begin{aligned} E(S) &= -\frac{8}{16} \log_2 \frac{8}{16} - \frac{8}{16} \log_2 \frac{8}{16} \\ &= 1 \end{aligned}$$

$$\begin{aligned} E(\text{col}) &= \frac{8}{16} [-\frac{2}{8} \log_2 \frac{2}{8} - \frac{6}{8} \log_2 \frac{6}{8}] \\ &\quad + \frac{8}{16} [-\frac{4}{8} \log_2 \frac{4}{8} - \frac{3}{8} \log_2 \frac{3}{8}] \\ &= \frac{8}{16} (0.9) + \frac{8}{16} (0.8) \\ &= 0.88. \end{aligned}$$

$$\begin{aligned} E(\text{size}) &= \frac{8}{16} [-\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8}] + \\ &\quad \frac{8}{16} [-\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}]. \\ &= 0.88. \end{aligned}$$

$$E(\text{Act.}) = \frac{8}{16} \left[-\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8} \right] +$$

$$\left[-\frac{3}{8} \log \frac{3}{8} - \frac{5}{8} \log \frac{5}{8} \right] \frac{8}{16} = 0.8$$

$$E(\text{Age}) = \frac{8}{16} \left[-\frac{3}{8} \log \frac{3}{8} - \frac{5}{8} \log \frac{5}{8} \right] + \frac{8}{16} \left[-\frac{2}{8} \log \frac{2}{8} - \frac{6}{8} \log \frac{6}{8} \right]$$

$$= 0.8$$

$$IG(S, \text{Col}) = 1 - 0.8$$

$$= 0.12$$

$$IG(S, \text{Size}) = 1 - 0.8$$

$$= 0.12$$

$$IG(S, \text{Act}) = 1 - 0.8$$

$$= 0.12$$

$$IG(S, \text{Age}) = 1 - 0.8$$

$$= 0.12$$

As Information gain in all cases are equal so we can take any of it.

we are choosing size as our attribute -
when size is small.

$$E(\text{col}) = \frac{4}{8} [-1 \log 1] + \frac{4}{8} [-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4}]$$
$$= 0.405 \quad \text{--- (i)}$$

$$E(\text{act}) = \frac{4}{8} [-\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4}] + \frac{4}{8} [-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4}]$$
$$= 0.965 \quad \text{--- (ii)}$$

$$E(\text{Age}) = \frac{4}{8} [-\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4}] +$$
$$\frac{4}{8} [-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4}]$$
$$= 0.905 \quad \text{--- (iii)}$$

As now we can see that eq i i.e.
 $E(\text{col})$ has the smallest value so we will
select this.

Now we will do the process for when we have size large.

$$E(\text{Age}) = \frac{4}{8} (-1 \log 1) + \frac{4}{8} \left[-\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4} \right] \\ = 0.5$$

$$E(\text{col}) = \frac{4}{8} \left[-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \right] + \\ \frac{4}{8} \left[-\frac{1}{4} \log \frac{1}{4} - \frac{3}{4} \log \frac{3}{4} \right] \\ = 0.81$$

$$E(\text{Act}) = \frac{4}{8} (-1 \log 1) + \frac{4}{8} \left[-\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4} \right] \\ = 0.5$$

As we have 2 options here i.e. Age & Act as both have same value so we can select any but let us choose age as attribute.

(5)

col = purple

Age = Adult.

$$E(\text{Act}) = -1 \log 1 = 0.$$

$$E(\text{size}) = \frac{2}{4} \left[-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right] + \frac{2}{4} \left[\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right].$$

$$= 1.$$

size = large
col = yellow.

$$E(\text{Act}) = \frac{2}{4} (-1 \log 1) + \frac{2}{4} \left[-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right].$$

$$= 0.5.$$

$$E(\text{Age}) = \frac{2}{4} (-1 \log 1) + \frac{2}{4} \left[\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right].$$

$$= 0.5.$$

color = purple

Age = child.

$$E(\text{Act}) = -1 \log 1 = 0.$$

$$E(\text{size}) = -1 \log 1 = 0.$$

size = small
Col = yellow.

$$E(\text{Act}) = 0.$$

$$E(\text{Age}) = 0.$$

So from the above we can easily draw.

3 - level Decision tree for baloon classification

