

\* Decision Tree :-

Instance	q1	q2	q3	classification
1	True	Hot	High	No
2	True	Hot	High	No
3	False	Hot	High	Yes
4	False	Cool	Normal	Yes
5	False	Cool	Normal	Yes
6	True	Cool	High	No
7	True	Hot	High	No
8	True	Hot	Normal	Yes
9	False	Cool	Normal	Yes
10	False	Cool	High	Yes

For Entropy of all data  
Distinct values in

classification	Total
Yes	6
No	4
	<u>10</u>

$$\text{Entropy}(D) = -\frac{6}{10} \log_2 \left( \frac{6}{10} \right) - \frac{4}{10} \log_2 \left( \frac{4}{10} \right)$$

$$= 0.9709 \left[ \because \log_y x = \frac{\log_{10} x}{\log_{10} y} \right]$$

⇒ chain of q1



$$\text{Gain}(D, a_1) = \text{Entropy}(D) - \text{Entropy}(D|a_1)$$

Distinct values in $a_1$	yes	no	Total
True	1	4	5
False	5	0	5
			<u>10</u>

$$\text{Entropy}(a_1) = \frac{5}{10} \times \left[ -\frac{1}{5} \log_2 \left( \frac{1}{5} \right) - \frac{4}{5} \log_2 \left( \frac{4}{5} \right) \right]$$

$$\frac{5}{10} \times \left[ -\frac{5}{5} \log_2 \left( \frac{5}{5} \right) - \frac{0}{5} \log_2 \left( \frac{0}{5} \right) \right]$$

$$= 0.7219 \times \frac{5}{10}$$

$$= 0.3609$$

$$\therefore \text{Gain}(D, a_1) = 0.9709 - 0.3609$$

$$= 0.6099$$

$\Rightarrow$  Gain of  $a_2$  :

Distinct value in $a_2$	yes	no	Total
Hot	2	3	5
Cool	4	1	5
			<u>10</u>



$$\text{Gain}(D, a_2) = 0.9709 - \left\{ \frac{5}{10} \left[ -\frac{2}{5} \log_2 \left( \frac{2}{5} \right) - \frac{3}{5} \log_2 \left( \frac{3}{5} \right) \right] \right.$$

$$\left. + \frac{5}{10} \left[ -\frac{4}{5} \log_2 \left( \frac{4}{5} \right) - \frac{1}{5} \log_2 \left( \frac{1}{5} \right) \right] \right\}$$

$$= 0.1245$$

⇒ Gain of  $a_3$

distinct values in $a_3$	yes	no	Total
high	2	4	6
normal	4	0	4
			<u>10</u>

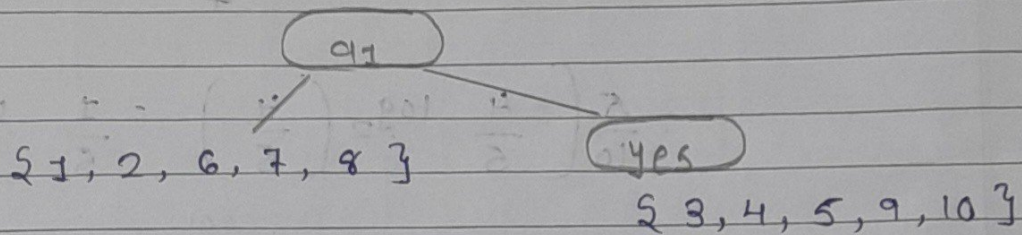
$$\text{Gain}(D, a_3) = 0.9709 - \left\{ \frac{6}{10} \left[ -\frac{2}{6} \log_2 \left( \frac{2}{6} \right) - \frac{4}{6} \log_2 \left( \frac{4}{6} \right) \right] \right.$$

$$\left. + \frac{4}{10} \left[ -\frac{4}{4} \log_2 \left( \frac{4}{4} \right) \right] \right\}$$



$$= 0.4200$$

- $\text{gain}(D, a_1) = 0.6099 \rightarrow \text{Maximum}$
- $\text{gain}(D, a_2) = 0.1245$
- $\text{gain}(D, a_3) = 0.4200$



- New Data

Instance	$a_2$	$a_3$	classification
1	Hot	High	No
2	Hot	High	No
6	Cool	High	No
7	Hot	High	No
8	Hot	Normal	Yes

$$\text{Entropy}(D) = -\frac{1}{5} \log_2 \left( \frac{1}{5} \right) - \frac{4}{5} \log_2 \left( \frac{4}{5} \right)$$

$$= 0.7219$$

Dist. Value in classification	Count
yes	1
no	4
	5



⇒ gain of  $q_2$  :

$$\text{gain}(D, q_2) = \text{Entropy}(D) - \text{Entropy}(q_2)$$

$$= 0.7219 - \left\{ \frac{4}{5} \left[ \frac{1}{4} \log_2 \left( \frac{1}{4} \right) - \frac{3}{4} \log_2 \left( \frac{3}{4} \right) \right] \right.$$

$$\left. + \frac{1}{5} \left[ -\frac{1}{1} \log_2 \left( \frac{1}{1} \right) \right] \right\}$$

$$= 0.0729$$

Disti. Values in $q_2$	yes	no	total
Hot	1	3	4
Cool	1	0	$\frac{1}{5}$

⇒ gain for  $q_3$

Disti. Value	yes	no	total
High	0	4	4
Normal	1	0	$\frac{1}{5}$

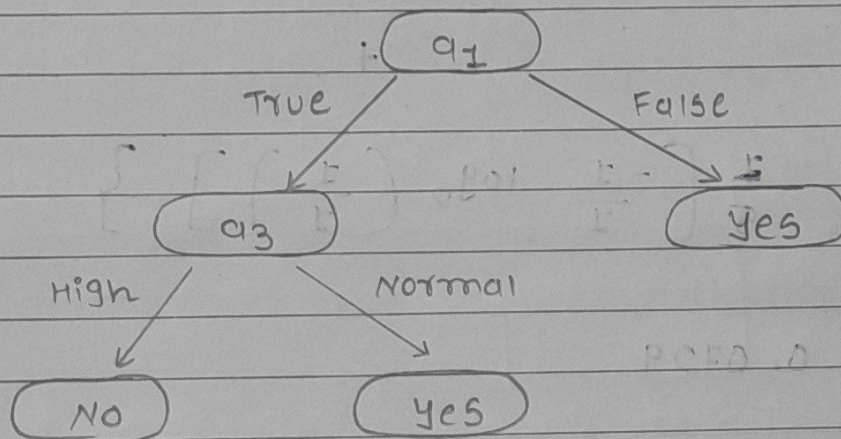


$$Gain(D, a_3) = 0.7219 - \left\{ \frac{4}{5} \left[ -\frac{4}{4} \log_2 \left( \frac{4}{4} \right) \right] \right.$$

+

$$\left. \frac{1}{5} \left[ -\frac{1}{1} \log_2 \left( \frac{1}{1} \right) \right] \right\}$$

$$= 0.7219 - \text{Max}$$



{1, 2, 6, 7}

{8}