

Prem A. Jethwa

1001861810

Assignment - 4 -

* Task 1 :

Part a :

$$\begin{aligned}
 & P(\text{Color} = \text{Green} \mid \text{Vehicle} = \text{Truck}) \\
 &= P(\text{Color} = \text{Red} \mid \text{Vehicle} = \text{Truck}) + P(\text{Color} = \text{Blue} \mid \text{Vehicle} = \text{Truck}) \\
 &= P(\text{Color} = \text{Red} \wedge \text{Vehicle} = \text{Truck}) + P(\text{Color} = \text{Blue} \wedge \text{Vehicle} = \text{Truck}) \\
 &\quad P(\text{Vehicle} = \text{Truck}) \qquad \qquad \qquad P(\text{Vehicle} = \text{Truck}) \\
 &= \frac{0.1554}{0.1554 + 0.1680 + 0.0966} + \frac{0.0966}{0.1554 + 0.1680 + 0.0966} \\
 &= \underline{\underline{0.60}}
 \end{aligned}$$

Part b :

Proove 'Vehicle' & 'color' are independent of each other.

To do so, we need to proove,

$$P(\text{Vehicle} \mid \text{Color}) = P(\text{Vehicle})$$

or

$$P(\text{Color} \mid \text{Vehicle}) = P(\text{Color})$$

For Color = Green and Vehicle = Car :

$$P(\text{Color} = \text{Green} \mid \text{Vehicle} = \text{Truck}) = \underline{\underline{0.60}} - (\text{Part a})$$

$$P(\text{Color} = \text{Green}) = 1 - P(\text{Color} = \text{Green})$$

$$= 1 - (0.1280 + 0.0480 + 0.1680 + 0.0560)$$

$$= \underline{\underline{0.60}}$$

We can see that,

$$P(\text{Color} = \text{Green} \mid \text{Vehicle} = \text{Truck}) = P(\text{Color} = \text{Green})$$

As both the values are same, \therefore variable color & vehicle are independent of each other.

* Task 2:

Part a:

Given:

$A \rightarrow 7$ values of $B_n \rightarrow 8$ values each where $n = 1$ to 10 .

\therefore No. of values to be stored for Joint P.D. (theoretically)

$$= 7 \times 8^{10} = 7.51 \times 10^9 \text{ (approx)}$$

And practically $= (7 \times 8^{10}) - 1$

Part b:

It is given that B_j is conditionally independent of B_i (where $j \neq i$) given A .

$$\therefore P(B_1, B_2, B_3, \dots, B_{10} | A) = \prod_{j=1}^{10} P(B_j | A) - ①$$

\therefore Using product rule we can say that,

$$\begin{aligned} P(A, B_1, B_2, B_3, \dots, B_{10}) &= P(B_1, B_2, B_3, \dots, B_{10} | A) \cdot P(A) \\ &= \left(\prod_{j=1}^{10} P(B_j | A) \right) \cdot P(A) \end{aligned}$$

- from ①

$$\left(\prod_{j=1}^{10} P(B_j|A) \right), P(A)$$

Each B_i variable can take 8 values & A can take 7 values.

\therefore No. of values to store for each $P(B_j|A) = 8 \times 7$

And as there are 10 B_i variables.

\therefore Total no. of values to be stored for $P(B_j|A)$

$$= 8 \times 7 \times 10 — ②$$

$$= 7 \times 7 \times 10 — ④$$

A variable can take 7 values

\therefore No. of values to store for $P(A) = 7 — ③$

$$= 6 — ⑤$$

\therefore Total no. of values to be stored = $(8 \times 7 \times 10) + 7$ — from ② & ③

$$= \underline{\underline{570}} \text{ (theoretically).}$$

And practically = $(7 \times 7 \times 10) + 6$ — from ④ & ⑤

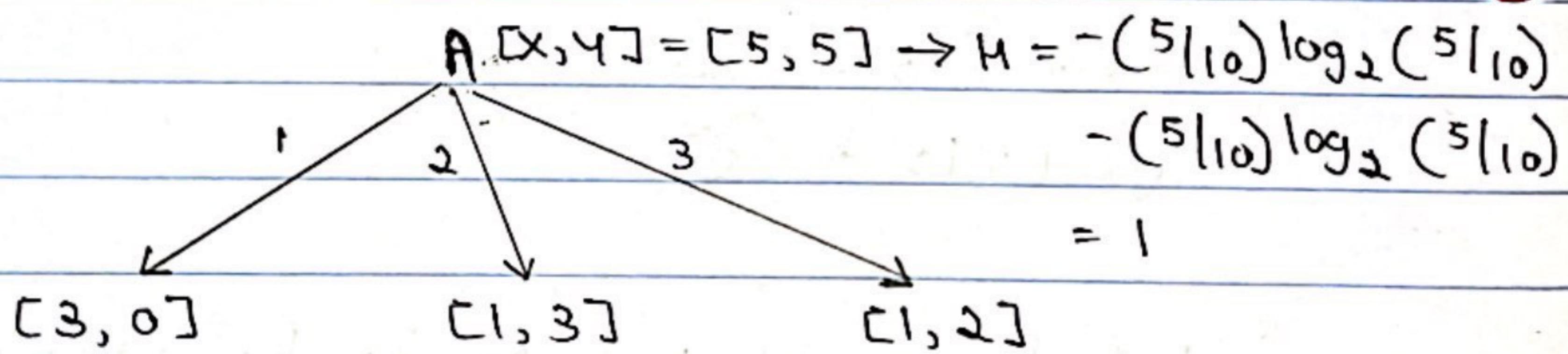
$$= \underline{\underline{496}} \text{ (Practically).}$$

Part c:

Here we know that the variables B_1, B_2, \dots, B_{10} are independent of each other given A. In a way, we can say that the effects $(B_1, B_2, \dots, B_{10})$ are independent of each other given cause (A). Thus we can say that this scenario follows the Naive Bayes Model.

* Task 4:

Attribute A:



$$H_1 =$$

$$-(3/3) \log_2(3/3)$$

$$-(0/3) \log_2(0/3)$$

$$= 0$$

$$H_2 =$$

$$-(1/4) \log_2(1/4)$$

$$-(3/4) \log_2(3/4)$$

$$= 0.81125$$

$$H_3 =$$

$$-(1/3) \log_2(1/3)$$

$$-(2/3) \log_2(2/3)$$

$$= 0.918.$$

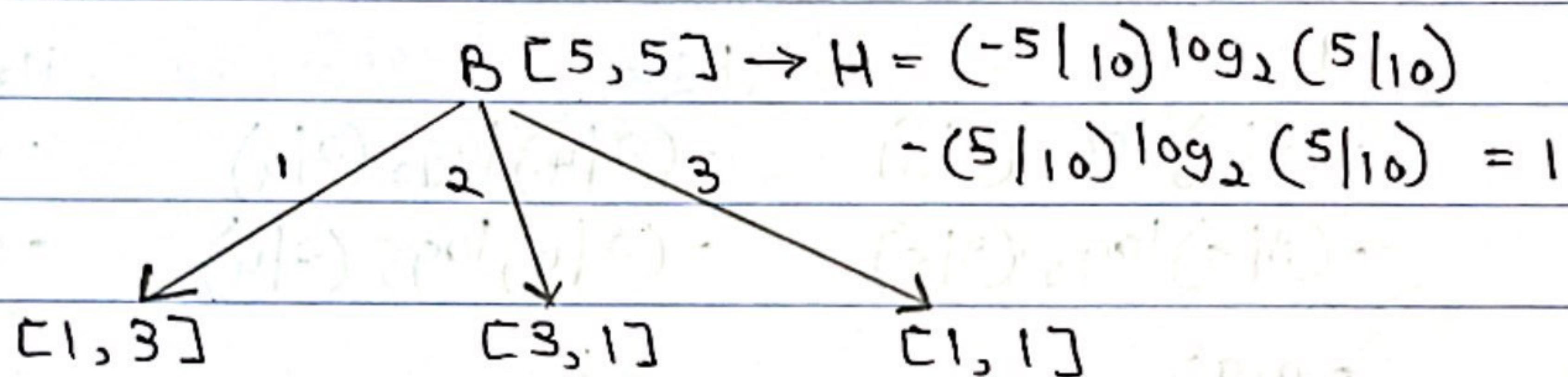
Information gain, for attribute A

$$= H - \sum_{j=1}^L (R_j | R) \cdot H_j$$

$$= 1 - \left[\left(\frac{3}{10} \right) (0) + \left(\frac{4}{10} \right) (0.8112) + \left(\frac{3}{10} \right) (0.918) \right]$$

$$= 0.40012 \approx 0.4$$

Attribute B:



$$H_1 = -(1/4) \log_2 (1/4) = 0.8112$$

$$H_2 = -(3/4) \log_2 (3/4) = 0.8112$$

$$H_3 = -(1/2) \log_2 (1/2) = 1$$

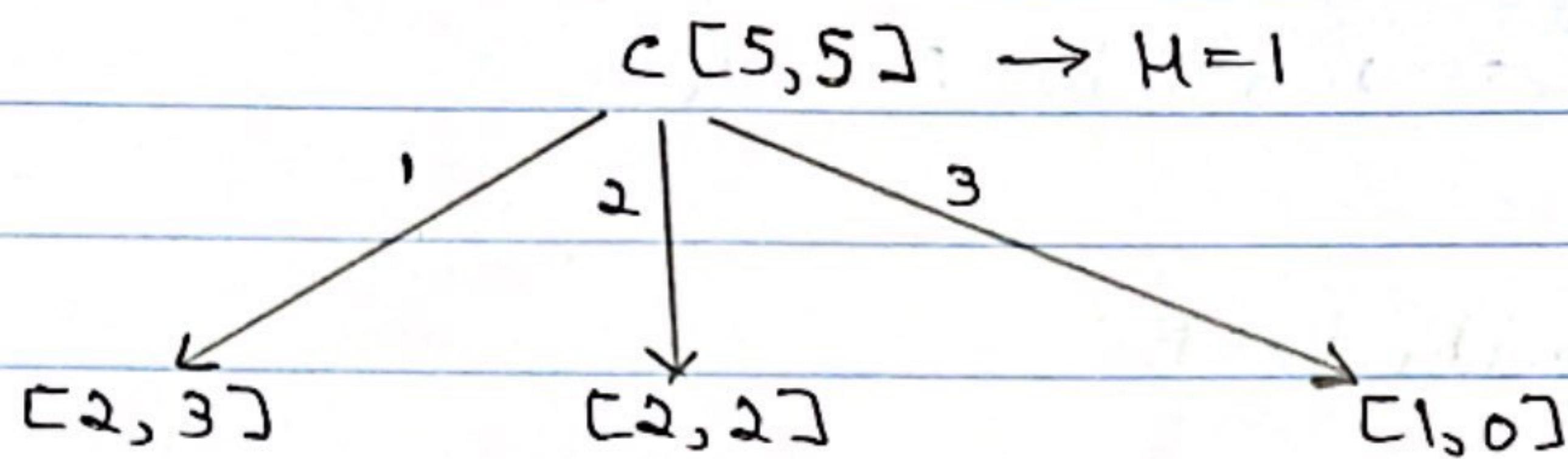
\therefore Information gain for attribute B,

$$= H - \sum_{i=1}^L (K_i/K) \cdot H_i$$

$$= 1 - \left[\left(\frac{4}{10} \right) \cdot (0.8112) + \left(\frac{4}{10} \right) \cdot (0.8112) + \left(\frac{2}{10} \right) (1) \right]$$

$$= 0.15104$$

Attribute c:



$$H_1 =$$

$$\begin{aligned} & -(2|5) \log_2 (2|5) \\ & -(3|5) \log_2 (3|5) \\ & = 0.97 \end{aligned}$$

$$H_2 =$$

$$\begin{aligned} & -(2|4) \log_2 (2|4) \\ & -(2|4) \log_2 (2|4) \\ & = 1 \end{aligned}$$

$$H_3 =$$

$$\begin{aligned} & -(1|1) \log_2 (1|1) \\ & -(0|2) \log_2 (0|2) \\ & = 0 \end{aligned}$$

∴ Information gain for attribute c

$$\begin{aligned} I &= 1 - \left[\left(\frac{5}{10}\right) \cdot (0.97) + \left(\frac{4}{10}\right) \cdot (1) + \left(\frac{1}{10}\right) \cdot (0) \right] \\ &= 0.115 \end{aligned}$$

Attribute	Information Gain
A	0.4
B	0.15
C	0.11

∴ A has maximum information gain, we can select it at root.

* Task 3:

From the given network.

$$P(\neg \text{not}(\text{Baseball Game on TV}) \mid \text{not}(\text{George Feeds cat}))$$

$$= P\left(\frac{\neg b_{\text{GTV}}}{\neg G_{\text{FC}}}\right)$$

$$= \frac{P(\neg b_{\text{GTV}} \wedge \neg G_{\text{FC}})}{P(\neg G_{\text{FC}})}$$

Inference by Enumeration;

$$\alpha < P(\neg b_{\text{GTV}} \wedge \neg G_{\text{FC}})$$

$$P(\neg b_{\text{GTV}} \wedge \neg G_{\text{FC}})$$

$$= P(\neg b_{\text{GTV}} \wedge \text{G}_{\text{WTV}} \wedge \text{O}_{\text{CF}} \wedge \neg G_{\text{FC}}) +$$

$$P(\neg b_{\text{GTV}} \wedge \text{G}_{\text{WTV}} \wedge \neg \text{O}_{\text{CF}} \wedge \neg G_{\text{FC}}) +$$

$$P(\neg b_{\text{GTV}} \wedge \neg \text{G}_{\text{WTV}} \wedge \text{O}_{\text{CF}} \wedge \neg G_{\text{FC}}) +$$

$$P(\neg b_{\text{GTV}} \wedge \text{G}_{\text{WTV}} \wedge \neg \text{O}_{\text{CF}} \wedge \neg G_{\text{FC}})$$

$$= 0.695890411 \times 0.118110236 \times 0.169863014 \times$$

$$0.958333333$$

$$+ 0.695890411 \times 0.118110236 \times 0.830136986 \times$$

$$0.293577982$$

$$\begin{aligned}
 & + 0.695890411 \times 0.881889764 \times 0.169863014 \times \\
 & \hspace{10em} 0.68421054 \\
 & + 0.695890411 \times 0.881889764 \times 0.830136986 \times \\
 & \hspace{10em} 0.41237113 \\
 & = \underline{0.125744303}
 \end{aligned}$$

Similarly, $P(BGTV \wedge GFC)$

$$\begin{aligned}
 & = 0.304109589 \times 0.927927928 \times 0.169863014 \times \\
 & \hspace{10em} 0.958333333 \\
 & + 0.304109589 \times 0.927927928 \times 0.830136906 \times \\
 & \hspace{10em} 0.293577982 \\
 & + 0.304109589 \times 0.72072072 \times 0.169863014 \times \\
 & \hspace{10em} 0.684210526 \\
 & + 0.304109589 \times 0.72072072 \times 0.830136906 \times \\
 & \hspace{10em} 0.041237113 \\
 & = \underline{0.118007272}
 \end{aligned}$$

tuple,

$\alpha < 0.125744303$	$0.118007272 >$
$= < 0.51587073$	$0.48412927 >$

* Task 5 -

Entropy Before Split

$$H = -\frac{5}{10} \log_2 \left(\frac{5}{10}\right) - \frac{5}{10} \log_2 \left(\frac{5}{10}\right)$$

$$= 1$$

Using A threshold A15,

$$x = 5, y = 3$$

$$H_{A>15} = -\frac{5}{8} \log_2 \left(\frac{5}{8}\right) - \frac{3}{8} \log_2 \left(\frac{3}{8}\right)$$

$$= 0.9544$$

$$H_{A<15} \rightarrow x = 0, y = 2$$

$$= -\frac{0}{2} \log \left(\frac{0}{2}\right) - \frac{2}{2} \log \left(\frac{2}{2}\right)$$

$$= 0$$

$$(\text{Infogain A}) = 1 - 0.9544 + 4 \times \frac{8}{10} + 0 = 0.7635$$

$$= 1 - 0.7635$$

$$= 0.2365$$

Using B threshold 15,

$$B > 15 \rightarrow x = 3, y = 5$$

$$H_{B>15} = -\frac{3}{8} \log_2 \left(\frac{3}{8}\right) - \frac{5}{8} \log_2 \left(\frac{5}{8}\right)$$
$$= 0.9544$$

$$B < 15 \rightarrow x = 2, y = 0$$

$$H_{B<15} = -\frac{2}{2} \log_2 \left(\frac{2}{2}\right) - \frac{0}{2} \log_2 \left(\frac{0}{2}\right)$$
$$= 0$$

$$(Information Gain)_B = 1 - 0 + 0.9544 \times \frac{8}{10}$$

$$= 1 - 0.7635$$

$$= 0.2365$$

Using C threshold 15,

$$C \geq 15 \rightarrow x = 5, y = 4$$

$$H_{C \geq 15} = -\frac{5}{9} \log_2 \left(\frac{5}{9}\right) - \frac{4}{9} \log_2 \left(\frac{4}{9}\right)$$
$$= 0.9910$$

$$C_{\geq 15} \rightarrow X=0, Y=1$$

$$H_{C \geq 15} = -\frac{0}{1} \log_2(0) - \frac{1}{1} \log_2(1)$$
$$\underline{= 0}$$

$$(\text{Information Gain}) C_{15} = 1 - 0 + 0.9910 \times \frac{9}{10} = 0.8919$$

$$= 1 - 0.8919$$
$$\underline{= 0.1081}$$

$$\text{Wing A threshold} = 20$$

$$A \geq 20 \rightarrow X=5, Y=2$$

$$H_{A \geq 20} = -\frac{5}{7} \log_2\left(\frac{5}{7}\right) - \frac{2}{7} \log_2\left(\frac{2}{7}\right)$$
$$\underline{= 0.8631}$$

$$H_{A \geq 20} \rightarrow X=0, Y=3$$

$$H_{A \geq 20} = \frac{0}{3} \log_2\left(\frac{0}{3}\right) - \frac{3}{3} \log_2\left(\frac{3}{3}\right)$$
$$\underline{= 0}$$

Using B Threshold = 20,

$$B \geq 20 \rightarrow x = 3, y = 5$$

$$H_{B \geq 20} = -\frac{3}{8} \log_2 \left(\frac{3}{8}\right) - \frac{5}{8} \log_2 \left(\frac{5}{8}\right)$$
$$= 0.9544$$

$$B < 20 \rightarrow x = 2, y = 0$$

$$H_{B < 20} = -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right)$$
$$= 0$$
$$(Info Gain B)_{20} = 1 - 0.9544$$
$$= 0.2365$$

Using C Threshold = 20,

$$C \geq 20 \rightarrow x = 5, y = 1$$

$$H_{C \geq 20} = -\frac{5}{6} \log_2 \left(\frac{5}{6}\right) - \frac{1}{6} \log_2 \left(\frac{1}{6}\right)$$
$$= 0.6500$$

$$C < 20 \rightarrow x = 0, y = 4$$

$$H_{C < 20} = -\frac{0}{4} \log_2 \left(\frac{0}{4}\right) - \frac{4}{4} \log_2 \left(\frac{4}{4}\right)$$
$$= 0$$

$$(\text{Info Gain})_{C_{20}} = 0 + 0.6500 \times \left(\frac{6}{10} \right)$$

$$= 1 - 0.3900 = \underline{\underline{0.6100}}$$

Using A threshold = 25,

$$A \geq 25 \rightarrow X = 3, Y = 0$$

$$H_{A \geq 25} = -\frac{3}{3} \log_2 \left(\frac{3}{3} \right) - \frac{0}{3} \log_2 \left(\frac{0}{3} \right)$$

$$= \underline{\underline{0}}$$

$$H_{A \geq 25} = X = 2, Y = 5$$

$$= -\frac{2}{7} \log_2 \left(\frac{2}{7} \right) - \frac{5}{7} \log_2 \left(\frac{5}{7} \right)$$

$$= \underline{\underline{0.8631}}$$

$$(\text{Info Gain})_{A \geq 25} = 0.8631 \times \frac{7}{10}$$

$$= 1 - [0.8631] = \underline{\underline{0.3959}}$$

Using B threshold = 25,

$$B \geq 25 \rightarrow X = 0, Y = 5$$

$$H_{B \geq 25} = -\frac{0}{5} \log_2 \left(\frac{0}{5} \right) - \frac{5}{5} \log_2 \left(\frac{5}{5} \right)$$

$$= \underline{\underline{0}}$$

$$B \geq 25 \rightarrow x=5, y=0$$

$$H_B \geq 25 = 0$$

$$(\text{Info gain } B)_{\geq 25} = 1 - 0 = \underline{\underline{1}}$$

Using c threshold = 25,

$$C \geq 25 \rightarrow x=3, y=0$$

$$H_C \geq 25 = -\frac{3}{3} \log_2 \left(\frac{3}{3}\right) - \log_2 \left(\frac{0}{3}\right)$$

$$= \underline{\underline{0}}$$

$$C \leq 25 \rightarrow x=2, y=5$$

$$H_C \leq 25 = -\frac{2}{7} \log_2 \left(\frac{2}{7}\right) - \frac{5}{7} \log_2 \left(\frac{5}{7}\right)$$
$$\approx \underline{\underline{0.8631}}$$

$$(\text{Info gain})_{C \geq 25} = 1 - (0 + 0.8631 \times 7/10)$$
$$= \underline{\underline{0.3959}}$$

Highest information is obtained when using B with threshold 25.