Assignment 10

Task 1

Part (a)    Entropy :

$$H(A) = H\left(\frac{80}{100}, \frac{20}{100}\right)$$

$$= -\frac{80}{100} \log_2\left(\frac{80}{100}\right) - \frac{20}{100} \log_2\left(\frac{20}{100}\right)$$

$$= -0.8 \log(0.8) - .2 \log(.2)$$

$$= -0.8 \times (-0.322) - (0.2)(-2.322)$$

$$= 0.2575 + 0.4643$$

$$\boxed{H(A) = 0.7218}$$

Part. (b)    Info Gain $= H(A) - \frac{35}{100} \times H\left(\frac{20}{35}, \frac{15}{35}\right)$

$$- \frac{65}{100} \times H\left(\frac{5}{65}, \frac{60}{65}\right)$$

$$= 0.7218 - \frac{35}{100} \times \left(-\frac{20}{35} \log_2\left(\frac{20}{35}\right) - \frac{15}{35} \log_2\left(\frac{15}{35}\right)\right)$$

$$- \frac{65}{100} \times \left(-\frac{60}{65} \log\left(\frac{60}{65}\right) - \frac{5}{65} \log\left(\frac{5}{65}\right)\right)$$

$$= 0.7218 - \left(0.35 \times \left(-0.571 \times -0.80735\right)\right.$$

$$- \left(0.4285 \times -1.2239\right)$$

$$- 0.65\left(-0.923 \times (-0.1154) - 0.0769 \times (-3.700)\right))$$

$$= 0.7218 - 0.344 - 0.2543$$

$$= 0.7218 - 0.5990$$

$$\boxed{\text{Info Gain} = 0.1224}$$

Part(C) | The Info Gain would be 0, it's repeated, hence no change would be observed.
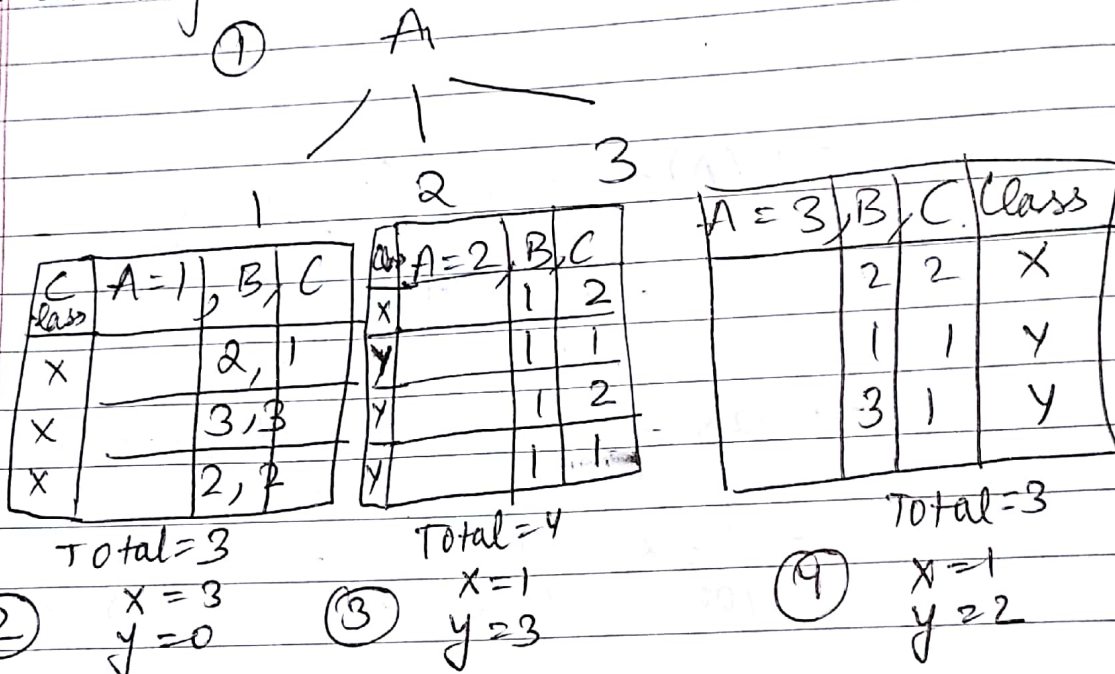
$$A \rightarrow C \rightarrow F$$

Part(D). deaf node : F
output : will wait

$$A \rightarrow B \rightarrow E \rightarrow H$$

Part (E) deaf node : H
output : will not wait

Task 2 | Considering 1st case with [Rootnode = A]

① A

1     2     3

| Class | A=1 | B | C |
|-------|-----|---|---|
| X | | 2, | 1 |
| X | | 3, 3 | |
| X | | 2, 2 | |

Total = 3

| Class | A=2 | B | C |
|-------|-----|---|---|
| X | | 1 | 2 |
| Y | | 1 | 1 |
| Y | | 1 | 2 |
| Y | | 1 | 1 |

Total = 4

| A = 3 | B | C | Class |
|-------|---|---|-------|
| | 2 | 2 | X |
| | 1 | 1 | Y |
| | 3 | 1 | Y |

Total = 3

② X = 3
    Y = 0

③ X = 1
    Y = 3

④ X = 1
    Y = 2

① $H(E) = 1$

② $H(E_1) = 0$

③ $H(E_2) = -\frac{1}{4} \log_2 \left(\frac{1}{4}\right) - \frac{3}{4} \log_2 \left(\frac{3}{4}\right)$

$$= 0.5 + 0.31127$$
$$= 0.81127$$

④ $H(E_3) = -\frac{1}{3} \log_2 \left(\frac{1}{3}\right) - \frac{2}{3} \log_2 \left(\frac{2}{3}\right)$

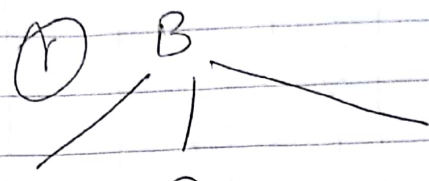$$= 0.5278 - 0.3956$$
$$= 0.9234$$

⑧ ① Information Gain (A) $= H(E) - \frac{4}{10}(H(E_2)) - \frac{3}{10}H(E_3)$

$= 1 - 0.4 \times 0.811 - 0.3 \times 0.9234$

$= 0.39858$

Considering $2^{nd}$ Case with Rootnode = B



① $H(E) = 1$

② $H(E_1) = -\frac{1}{4}\log_2\left(\frac{1}{4}\right) - \frac{3}{4}\log_2\left(\frac{3}{4}\right) = 0.81127$

③ $H(E_2) = -\frac{3}{4}\log_2\left(\frac{3}{4}\right) - \frac{1}{4}\log_2\left(\frac{1}{4}\right) = 0.81127$

④ $H(E_3) = 1$

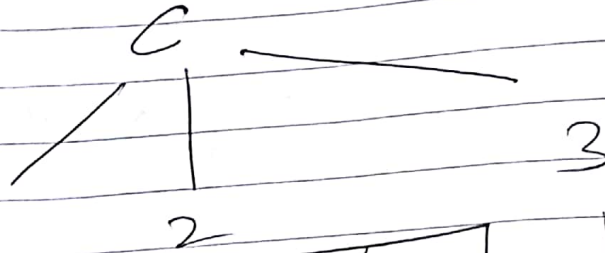Information Gain (B) $= 1 - \frac{4}{10}H(E_1) - \frac{4}{10}H(E_2) - \frac{2}{10}(E_3)$

$= 1 - 0.4 \times 0.811 - 0.4 \times 0.811$

$\qquad\qquad - 0.2 \times 1$

$= 0.1512$

Considering 3rd case with Rootnode = C

C

1                  2                  3

A B

$C_H = $

| | B | C | Class |
|---|---|---|---|
| 1 | 1 | 2 | X |
| 1 | 2 | 1 | Y |
| 2 | 3 | 1 | Y |
| 3 | 2 | 3 | Y |
| 2 | 2 | 1 | Y |

Total = 5

X = 1

y = 4

| C=2 | A | B | Class |
|---|---|---|---|
| | 2 | 1 | X |
| | 3 | 2 | X |
| | 1 | 2 | X |
| | 2 | 2 | Y |

Total = 4

X = 3

y = 1

| C=3 | A | B | Class |
|---|---|---|---|
| | 1 | 0 | X |

Total = 1

X = 1

$H(C_E) = 1$

$H(E_1) = \dfrac{-1}{5} \log_2\left(\dfrac{1}{5}\right) - \dfrac{4}{5} \log_2\left(\dfrac{4}{5}\right)$

$= -0.2 \log_2(0.2) - 0.8 \log_2(0.8)$

$= 0.4643 + 0.2575$

$= 0.7128 \quad = 0.7218$

$H(E_2) = -\dfrac{3}{4} \log_2\left(\dfrac{3}{4}\right) - \dfrac{1}{4} \log_2\left(\dfrac{1}{4}\right)$

$= 0.811127$

$H(E_3) = 0$

Information Gain(C) $= H(E) - \frac{5}{10} H(E_1) - \frac{4}{10} H(E_2) - \frac{1}{10} H(E_3)$

$= 1 - \left(\frac{1}{2} \times 0.7218\right) - \left(\frac{2}{5} \times 0.81127\right)$

$= 1 - 0.3609 - 0.3245$

$= 0.3146$

Since, Info Gain (A) $>$ Info Gain(C) $>$ Info Gain(B)

$0.39858 > 0.3146 > 0.1512$

Hence, A receives the highest Info Gain

Task 3 (a) Highest entropy is incase of even distribution.

$\phantom{i.e \quad 1000 \to} A \quad B \quad C \quad D$

i.e $\quad 1000 \to 250, 250, 250, 250$

Hence Entropy would be $= H(E_A) + H(E_B) + H(E_C) + H(E_D)$

$= 4 \times \left(- \frac{250}{1000} \times \log_2\left(\frac{250}{1000}\right)\right)$

$= -4 \times 0.25 \times \log_2 0.25$

$= 2$

lowest entropy would mean, putting all the examples in a single class. i.e $\quad 1000 \to \overset{A}{1000}$

Hence, Entropy would be $= H(E_A)$

$= -\frac{1000}{1000} \times \log_2\left(\frac{1000}{1000}\right)$

$= 0$

Date    /    /

**Task 3**

(b) Highest possible Entropy is 2

Lowest possible Entropy is 0

**Task 4**

We can improve the classifier depending on true cases in the data set. i.e we would need a more vibrant and diverse database set.

Completely depends on the Data Set.
We cannot guarantee 60%. It would depend on the training Data.

**Task 5**

The total no. of distinct decision trees with n boolean attributes is equal to the distinct truth table with $2^n$ rows will be $2^{2^n}$

Hence, with 5 boolean variables $= 2^{2^5}$

$$= 2^{32}$$

$$= 4294967296$$