# Assignment-I – Entrophy

| Instance | $a_1$ | $a_2$ | $a_3$ | classification |
|----------|-------|-------|--------|----------------|
| 1 | T | Hot | High | NO |
| 2 | T | Hot | High | NO |
| 3 | F | Hot | High | Yes |
| 4. | F | Cool | normal | Yes |
| 5. | F | Cool | normal | Yes |
| 6 | T | Cool | High | NO |
| 7 | T | Hot | High | NO |
| 8 | T | Hot | normal | Yes |
| 9. | F | Cool | normal | Yes |

| 10. | F | cool | High | | Yes. |
|-----|---|------|------|--|------|

For calculating Entrophy:-

Step 1:- overall Entrophy for the Data set

Step 2: find Individual Entrophy for each attribute

Step 3: For that attribute calculate the Information gain formula.

Step 4: Repeat 3 & 4 for all attributes.

Step 5: Sum the Entrophy finded from all each attribute.

Calculation:-

1) overall Entrophy:-

Given

No. of DS = 10

Positive = $6/10$

negative = $4/10$

$= -1.322$

$4/10 = 0.4$

$6/10 = -0.6$

$= -1.322$

$= -0.737$

Entrophy formula:-

$$E(S) = -P_S \ Log_2 \ P_S - N_S \ Log_2 \ N_S$$

$$E(S) = -6/10 \ Log_2 \ 6/10 - 4/10 \ Log_2 \ 4/10$$
$$= 0.5288 + 0.4422$$

$$E(S) = 0.971$$

2) Entrophy & IGF $(a_1)$

values - True, False (2)

True                          False

Total - 5                     Total - 5

$T_s$ - 4/5                   $F_s$ - 5/5

$T_N$ - 1/5                   $F_N$ - 0

                                       are yes

$E(T) = -4/5 \log_2 4/5 - 1/5 \log_2 1/5$     $\boxed{E(F) = 0}$

       $= 0.25752 + 0.4644$

$\boxed{E(T) = 0.72192}$

INFORMATION GAIN FORMULA FOR $a_1$ :-

$$\boxed{G(a_1) = Entropy - \frac{True}{a_1} * E(T) - \frac{False}{a_1} * E(F)}$$

$G(a_1) = 0.971 - 5/10 * 0.72192 - 5/10 * 0$

       $= 0.971 - 0.36096 - 0$

$\boxed{G(a_1) = 0.61}$

3) Entrophy & IGF ($a_2$) :-

values - Hot, cool

Hot

Total - 5    3/5 = 0.6

Hs - 3/5    L=0.73

Hn - 2/5    2/5 = 0.4
            L=-1.32

$E(H) = -3/5 \log_2 3/5 - 2/5 \log_2 2/5$

= 0.4422 & 0.5288

= 0.971

cool     4/5

Total - 5

Cs - 4/5   1/5

Cn - 1/5

$E(C) = -4/5 \log_2 4/5 - 1/5 \log_2 1/5$

= 0.

= 0.72192

IFG Fa $a_2$ :-

$G(a_2) = 0.971 - 5/10 * 0.971 - 5/10 * 0.72192$

= 0.971 - (0.5 * 0.971) - (0.5 * 0.72192)

            - 0.4855        - 0.36096

= 0.971

= 0.971 - 0.84646

G($a_2$) = 0.124

3) a₃ :-

value = 2

**H**

$T = 6$

$\frac{2}{6} = 0.33$

$L = -1.585$

$Hs = \frac{2}{6}$

$\frac{4}{6} = 0.67$

$HN = \frac{4}{6}$

$L = -0.585$

$= -\frac{2}{6} \log_2 \frac{2}{6} - \frac{4}{6} \log_2 \frac{4}{6}$

$= 0.52305 + 0.39195$

$= 0.915$

**N**

$T = 4$

$Ns = \frac{4}{4}$

$NN = 0$

$= 0$

$= 0$

**IGF :-**

$= 0.917 - \frac{6}{10} * 0.915 - 0$

$= 0.917 - 0.549$

$\boxed{G(a_3) = 0.368}$

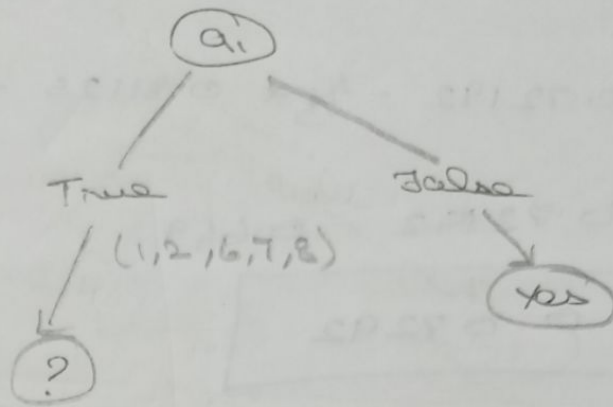$q_1 = 0.61 \;\wedge$

$q_2 = 0.124$

$q_3 = 0.368$

(✗) Fuom the calculated IG we can also to find the Root node

Decision tree

Root node



(*) From False we get cconclusion, But in true we have to again form the dataset belongs to true.

| | $a_1$ | $a_2$ | $a_3$ | Classification |
|---|---|---|---|---|
| 1 | T | Hot | High | NO |
| 2 | T | Hot | High | NO |
| 6 | T | cool | High | NO |
| 7 | T | Hot | High | NO |
| 8 | T | Hot | Normal | Yes |

By using this dataset again we have to calculate next?

Isr by attributes

$1/5 = 0.2$, $L = -2.322$, $4/5 = 0.8$, $L = 0.321$

$$E(S) = -1/5 \log_2 1/5 - 4/5 \log_2 4/5$$

$$= 0.4644 + 0.25752$$

$$E(S) = 0.72192$$

$a_2$ :- Entrophy

values - Hot, cool

Hot

T - 4

$H_S - 3/4$

$H_N - 1/4$

$$E(H) = -3/4 \log_2 3/4 - 1/4 \log_2 1/4$$

$$= 0.31125 + 0.5$$

$$= 0.81125$$

cool

T - 1

$C_S - 0$

$C_N - 1/1$

$$E(C) = 0$$

$3/4 = 0.75$

$L = 0.415$

$1/4 = 0.25$

$L = 2$

IGGF  $a_2$

$$= 0.72192 - 4/5 * 0.81125 - 0$$

$$= 0.72192 - 0.649$$

$$\boxed{G(a_2) = 0.07292}$$

---

$a_3$ :-

Entrophy :-

values - High, normal

High - 4

HS - 0

HN - 4/4

$$\boxed{E(H) = 0}$$

NO

normal - 1

NS - 1/1

NN - 0

$$\boxed{E(N) = 0}$$

yes

IGF :-

$$= 0.72192 - 0 - 0$$

$$\boxed{G(a_3) = 0.72192}$$

---

$$a_2 = 0.07292$$

$$a_3 = 0.72192$$

# Decision Tree

$a_1$

True $(1,2,6,7,8)$     False

$a_2$                  yes

High $(1,2,6,7)$     normal $(8)$

NO            yes