# Decision tree

① ID3

Iterative Decotomiser.

(Entropy)

② CART

classification and Regression Tree.

(Gini impurity)

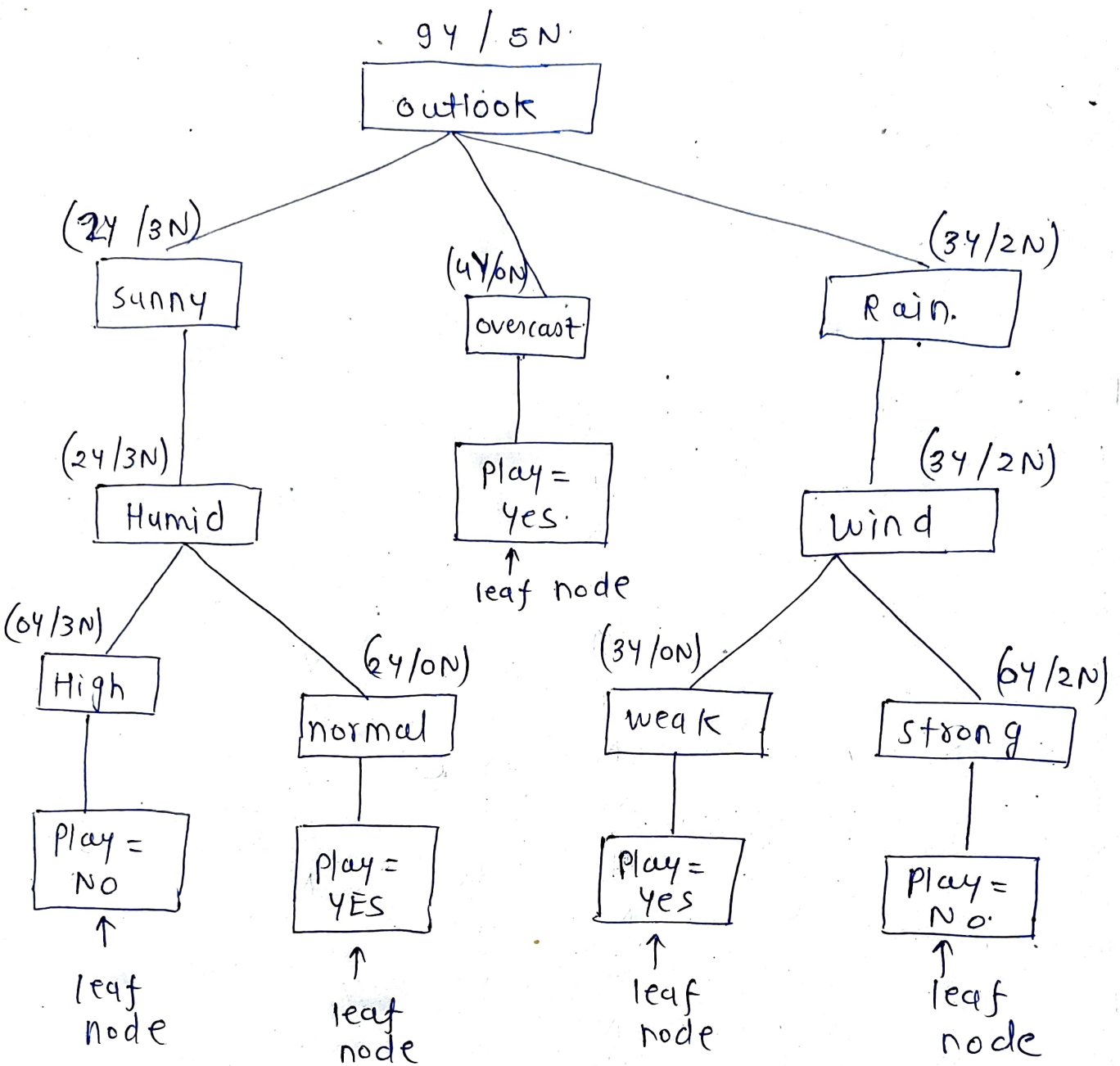| Day | outlook | Temp | Humid | wind | Decision. |
|-----|---------|------|-------|------|-----------|
| 1 | S | H | H | W | N |
| 2 | S | H | H | S | N |
| 3 | O | H | H | W | Y |
| 4 | R | M | H | W | Y. |
| 5 | R | C | N | W | Y. |
| 6 | R | C | N | S | N |
| 7 | O | C | N | S | Y. |
| 8 | S | M | H | W | N |
| 9 | S | C | N | W | Y. |
| 10 | R | M | N | W | Y. |
| 11 | S | M | N | W | Y. |
| 12 | O | M | H | S | Y. |
| 13 | O | H | N | W | Y. |
| 14 | R | M | H | S | N. |

outlook
{
S → Sunny
O → overcast
R → Rain
}

Humidity
{
H → High
N → Normal
}

Temp
{
H → Hot
M → mild
C → cold
}

wind
{
W → weak
S → strong
}

Play   Y - Yes
       N - No

9Y / 5N

outlook

(2Y /3N)

sunny

(4Y/ON)

overcast

(3Y/2N)

Rain.

(2Y/3N)

Humid

play =
yes.

↑
leaf node

(3Y/2N)

wind

(0Y/3N)

High

(2Y/0N)

normal

(3Y/0N)

weak

(6Y/2N)

strong

Play =
No

↑
leaf
node

play =
YES

↑
leaf
node

Play=
yes

↑
leaf
node

play =
No.

↑
leaf
node

* I have provided detailed mathematical solution of above problem using entropy at last of these notes.

* Entropy: measure of randomness (Disorder) or measure of purity or impurity.

$$E = -\sum_{i=1}^{n} P_i \times \log_2 (P_i)$$

for 2 classes (Y/N)

$$E_2 = -P_Y \log_2(P_Y) - P_N \cdot \log_2(P_N).$$

for 3 classes $(C_1, C_2, C_3)$

$$E_3 = -P_{C_1} \log_2(P_{C_1}) - P_{C_2} \log_2(P_{C_2}) - P_{C_3} \log_2(P_{C_3})$$
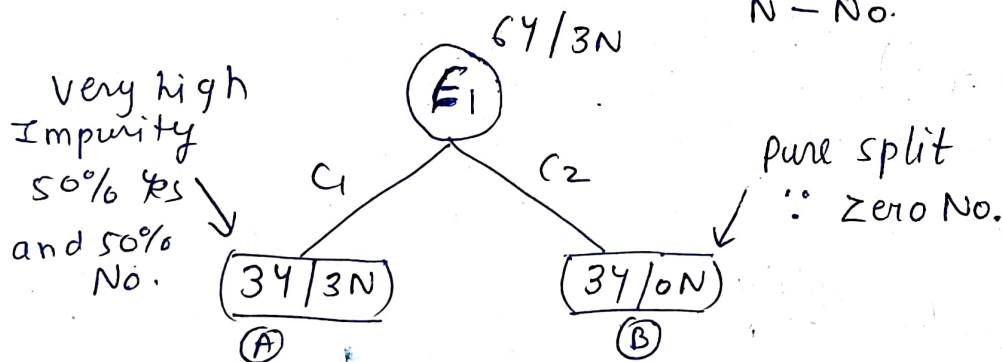
## Gini Impurity

$$G = 1 - \sum_{i=1}^{n} P_i^2$$

$n \equiv$ no. of classes

for 2 classes (Y, N).

$$G_2 = 1 - (P_Y^2 + P_N^2)$$

## Example

Y — Yes
N — No.



very high Impurity 50% Yes and 50% No.

$6Y/3N$ at $E_1$, branches $C_1$ to $(3Y/3N)$ (A) and $C_2$ to $(3Y/0N)$ (B)

pure split
$\therefore$ Zero No.

| f1 | Decision |
|----|----------|
| $C_1$ | N |
| $C_2$ | Y |
| $C_1$ | Y |
| $C_1$ | N |
| $C_1$ | Y |
| $C_2$ | Y |
| $C_1$ | Y |
| $C_1$ | N |
| $C_2$ | Y |

checking Impurity of f1.

$$H(s) = E_A = -\sum_{i=1}^{n} P_i \times \log(P_i)$$

$$= -P_Y \log_2(P_Y) - P_N \log_2(P_N)$$

$$= -\frac{3}{6} \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right)$$

$$= -2 \times \frac{3}{6} \log_2\left(\frac{3}{6}\right) = -\log_2\frac{1}{2} = 1$$

## Entropy impurity

$$E_B = - P_Y \log_2(P_Y) - P_N \log_2(P_N)$$

$$= - \frac{3}{3} \log_2(1) - 0 \times \log_2(0).$$

$$= - \log_2(1).$$

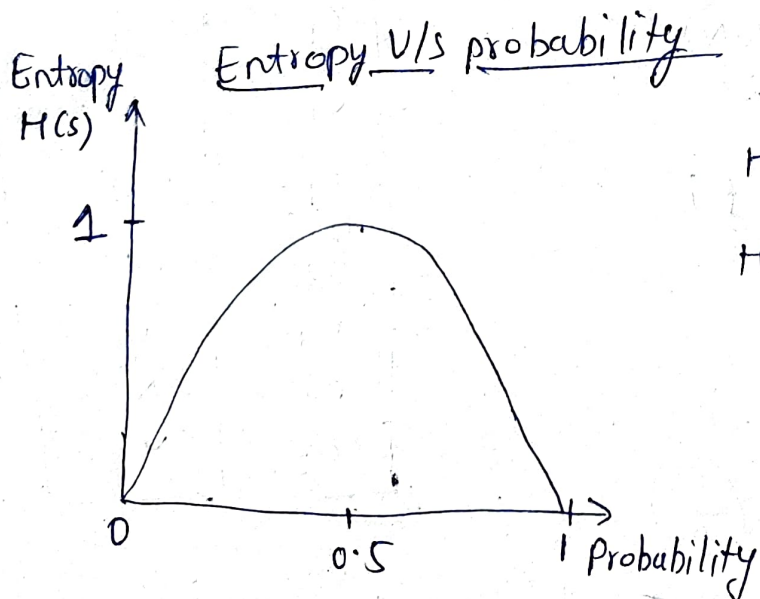$$\boxed{E_B = -0.}$$

## Gini impurity

$$G_A = 1 - \sum_{i=1}^{n} P_i^2 = 1 - (P_Y^2 + P_N^2)$$

$$= 1 - \left[\left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2\right]$$

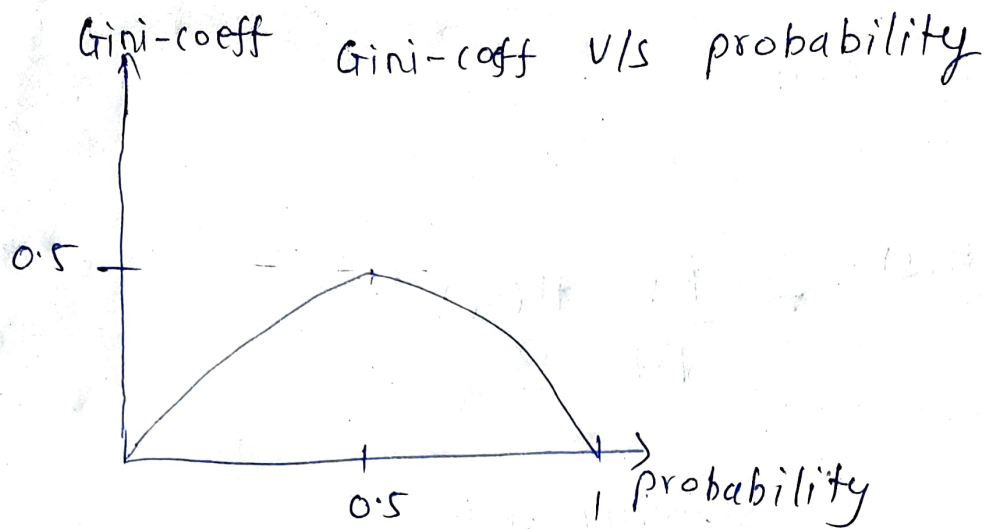$$= 1 - \left(\frac{1}{4} + \frac{1}{4}\right) = 1 - \frac{1}{2} = \frac{1}{2}.$$

$$G_B = 1 - \left[P_Y^2 + P_N^2\right]$$

$$= 1 - (1^2 + 0^2)$$

$$G_B = 0$$

## Entropy V/s probability



$H(s) = 1$ very impure split

$H(s) = 0$ pure split

Gini-coeff

Gini-coeff v/s probability



Entropy $\Rightarrow$ 0 to 1

Gini-coeff $\Rightarrow$ 0 to 0.5

44/8N

$$G = 1 - \left[\left(\frac{4}{12}\right)^2 + \left(\frac{8}{12}\right)^2\right] = 1 - \left(\frac{1}{9} + \frac{4}{9}\right)$$

$$= 1 - \frac{5}{9}$$

$$= 4/9$$

$$\boxed{G = 0.44}$$

$(84, 2N)$

$$G = 1 - \left[\left(\frac{8}{10}\right)^2 + \left(\frac{2}{10}\right)^2\right] = 1 - \left(\frac{16}{25} + \frac{1}{25}\right)$$
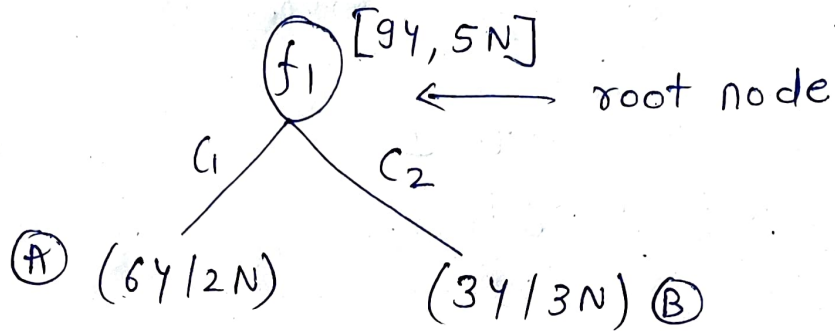
$$= 1 - \frac{17}{25}$$

$$= 8/25$$

$$\boxed{G = 0.32}$$

# Information Gain

1) Entropy

$$Gain(s,f) = H(s) - \sum \frac{|s_v|}{|s|} H(s_v).$$

[9Y, 5N] ← root node

(f1)

$C_1$ / \ $C_2$

(A) (6Y/2N)     (3Y/3N) (B)

$H(s) \equiv$ root feature entropy.

$$H(s) = -P_Y \log_2(P_Y) - P_N (\log_2(P_N))$$

$$= -\frac{9}{14} \log_2(9/14) - \frac{5}{14} \log_2(5/14)$$

$$= 0.41 - (-0.53)$$

$$\boxed{H(s) = 0.94}$$
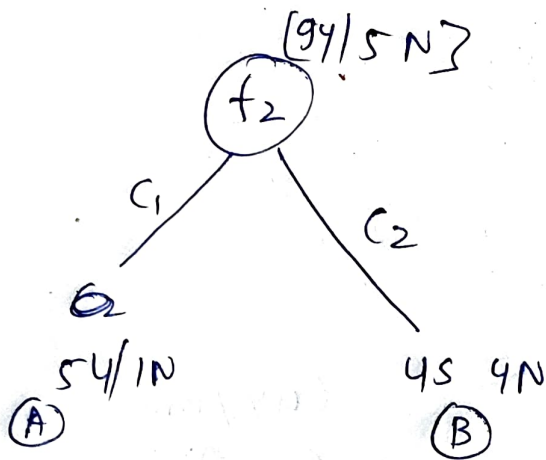
$$E_A = -\frac{6}{8} \log(6/8) - \frac{2}{8} \log(2/8)$$

$$\underline{E_A = 0.81}$$

$$E_B = -\frac{3}{6} \log(3/6) - \frac{3}{6} \log(3/6)$$

$$\boxed{E_B = 1}$$

Total Y/N in A          Total Y/N in B.

$$Gain(s, f_1) = 0.94 - \left[ \frac{8}{14} \times 0.81 + \frac{6}{14} \times 1.0 \right]$$

$$= 0.048$$

$[9+/5N]$

$f_2$

$H(s) = 0.94$

$C_1$     $C_2$

$6_2$

$5+/1N$     $4S \ 4N$

$(A)$     $(B)$

$E_A = -\frac{5}{86} \log_2 (5/86) - \frac{1}{86} \log_2 (1/86)$

$= 0.002 \ 0.65$

$E_B = -4/8 \log_2 (4/8) - 4/8 \log_2 (4/8) = 1$
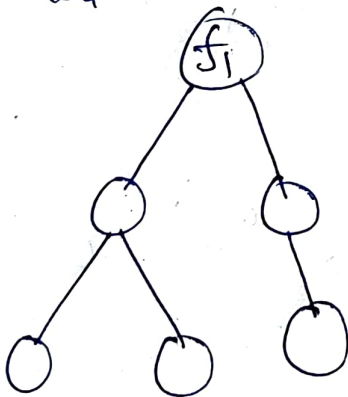
$Gain(s, f_2) = 0.94 - \left( \frac{6}{14} \times 0.65 + \frac{8}{14} \times 1 \right)$

$= 0.09$

∵ $Gain(s, f_2) > Gain(s, f_1)$
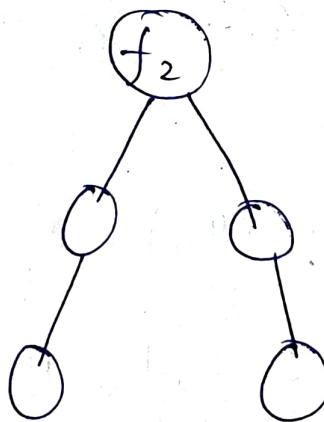
∴ $f_2$ will be better as root feature.

$IG = 0.98$          $IG = 0.97$          $IG = 0.45$

$f_1$                    $f_2$                    $f_3$

∴ Root node ⟹ $f_1$

⇒ solution of problem on first page.
   Total Y and N ⇒ D. Y = 9     N = 5

① outlook feature.

$$(Y=9, N=5)$$
Outlook

(2Y/3N)                       (4Y/0N)
Sunny       Rain          overcast·
(2Y/3N)      (3Y/2N).       (4Y/0N)

$H(S) = -\Sigma P_i \log_2 (P_i)$

$\quad = - P_Y \log_2 (P_Y) - P_N \log_2 (P_N)$

$\quad = - \dfrac{9}{814} \log_2 (9/14) - \dfrac{5}{14} \log_2 (5/14)$

$H(S) = 0.94$

$E_{sunny} = - P_Y \log_2 (P_Y) - P_N \log_2 (P_N)$

$\quad = - \dfrac{2}{5} \log_2 (2/5) - \dfrac{3}{5} \log_2 (3/5)$

$\quad = 0.971$

$E_{Rain} = - \dfrac{3}{5} \log_2 (3/5) - \dfrac{2}{5} \log_2 (2/5)$

$\quad = 0.971$

$E_{overcast} = - \dfrac{4}{4} \log_2 (4/4) - \dfrac{0}{4} \log_2 (0/4) = 0$

$Gain (s, outlook) = H(S) - \Sigma \dfrac{|S_v|}{|S|} \times H(S_v).$

$\quad = 0.94 - \left( \dfrac{5}{14} \times 0.971 + \dfrac{5}{14} \times 0.971 + \dfrac{4}{14} \times 0 \right)$

$Gain (s, outlook) = 0.246$

②

94, 5N

$$\text{Temp}$$

$$\overset{\text{Hot}}{(2Y/2N)} \qquad \overset{\text{mild}}{(4Y/2N)} \qquad \overset{\text{cold}}{(3Y, 1N)}$$

$H(s) = 0.94.$    for $9Y$ and $5N$.

$$E_{Hot} = -P_Y \log_2(P_Y) - P_N \log_2(P_N)$$

$$= -\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right)$$

$$= 1$$

$$E_{mild} = -\frac{4}{6} \log\left(\frac{4}{6}\right) - \frac{2}{6} \log\left(\frac{2}{6}\right)$$

$$= 0.918$$

$$E_{cold} = -\frac{3}{4} \log\left(\frac{3}{4}\right) - \frac{1}{4} \log\left(\frac{1}{4}\right)$$

$$= 0.811$$

$$\text{Gain}(s, \text{Temp}) = H_{st} - \sum \frac{|S_v|}{|S|} \times H(S_v)$$

$$= 0.94 - \left(\frac{4}{14} \times 1 + \frac{6}{14} \times 0.918 + \frac{4}{14} \times 0.811\right)$$

$$= 0.029$$

③       $\underline{\text{Humid}}$ 9Y/5N

High                        Normal.

(3Y/4N)                  (6Y/21N)

$H(s) = 0.9y$

$$E_{High} = -P_y \log_2(P_y) - P_N \log_2(P_m)$$

$$= -\frac{3}{7} \log_2(3/7) - \frac{4}{7} \log_2(4/7)$$

$$= 0.985$$

$$E_{Normal} = -\frac{6}{7} \log_2(6/7) - \frac{1}{7} \log_2(1/7)$$

$$= 0.592$$

$$Gain(S, Humid) = 0.9y - \left(\frac{7}{14} \times 0.985 + \frac{7}{14} \times 0.592\right)$$

$$= 0.151$$

④      $\underline{\text{wind}}$ (9Y/5N)

weak

(6Y/2N)                    Strong

                             3Y/3N

$H(s) = 0.9y$

$$E_{weak} = -\frac{6}{8} \log_2(6/8) - \frac{2}{8} \log_2(2/8)$$

$$= 0.811$$

$$E_{strong} = -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right)$$
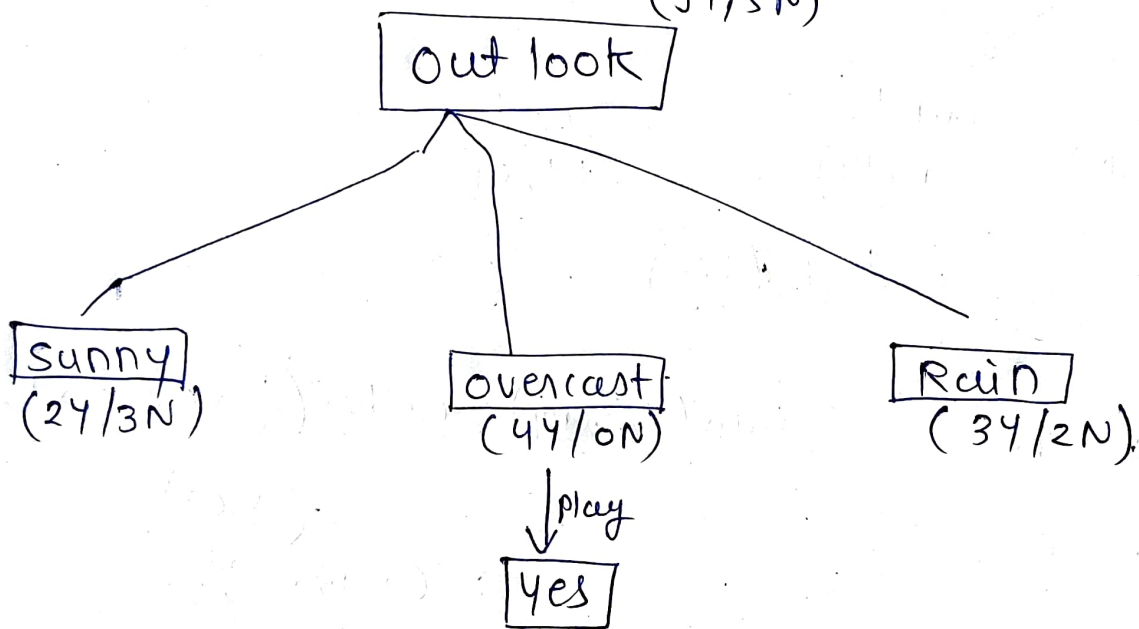
$$= 1$$

$$Gain \,(s,wind) = 0.94 - \left(\frac{8}{14} \times 0.811 + \frac{6}{14} \times 1\right).$$

$$= 0.048$$

$$\therefore \quad Gain \,(S, outlook) \text{ is Highest ie } 0.246$$

$\therefore$ our root node is outlook

(9Y/5N)

| Out look |



| sunny | | Overcast | | Rain |
(2Y/3N) | | (4Y/0N) | | (3Y/2N)

↓ play

| yes |

* Now Doing calculations for Decision Node (sunny)
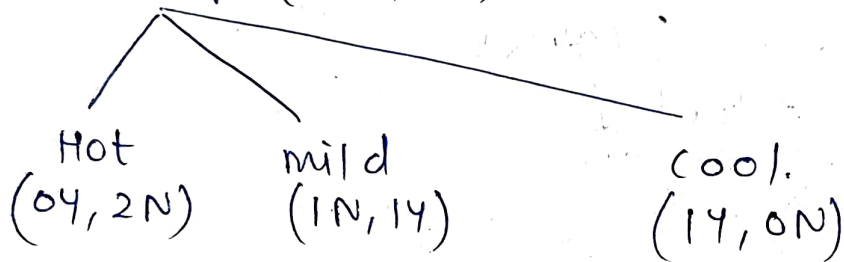* for 2Y and 3N.

$$H(s) = -P_y \log_2 (P_y) - P_N \log_2 (P_N)$$

$$= -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right)$$

$$H(s) = 0.97$$

| outlook | Temp | Humid | wind | Decision |
|---|---|---|---|---|
| S | H | H | W | N |
| S | H | H | S | N |
| S | M | H | W | N |
| S | C | N | W | Y |
| S | M | N | S | Y |

① Temp ( 2Y, 3N)



Hot
(0Y, 2N)

mild
(1N, 1Y)

cool
(1Y, 0N)

$E_{Hot} = 0 - \frac{2}{2} \log_2 (\frac{2}{2}) = 0$

$E_{mild} = 1$

$E_{cool} = 0$

Gain $(S, Temp) = 0.97 - (\frac{2}{5} \times 0 + \frac{2}{5} \times 1 + \frac{1}{5} \times 0.)$

$= 0.97 - \frac{2}{5}$

$= 0.57$

② Humid (2Y, 3N)



High
(0Y, 3N)

Normal
(2Y, 0N)

$E_{High} = 0$

$E_{Normal} = 0$

Gain $(S, Humid) = 0.97 - (0 + 0)$

$= 0.97$

③ wind (2Y, 3N)

weak         strong
(1Y, 2N)      (1Y, 1N)

$E_{weak} = -\frac{1}{3} \log_2 \left(\frac{1}{3}\right) - \frac{2}{3} \log_2 \left(\frac{2}{3}\right) = 0.918$

$E_{strong} = 1$

$Gain(s, wind) = 0.97 - \left(\frac{3}{5} \times 0.918 + \frac{2}{5} \times 1\right)$

$= 0.0192$

∴ Gain(s, Humid) is Highest ie (0.97)

∴ Humid is    our    Decision

Node for Sunny.

② for overcast there won't be any
Decision node because we have
already reached leaf node for
overcast.
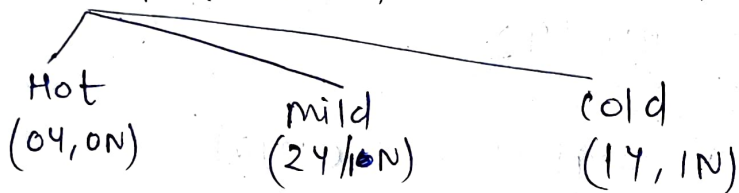
③ Now Doing calculations for Decision
Node (Rain)
(3Y and 2N)

| outlook | Temp | Humid | wind | Decision |
|---------|------|-------|------|----------|
| R | M | H | W | Y |
| R | C | N | W | Y |
| R | C | N | S | N |
| R | M | N | W | Y |
| R | M | H | S | N |

*for 3Y and 2N

$$H(s) = -\frac{3}{5} \log_2 (3/5) - \frac{2}{5} \log_2 (2/5) = 0.97.$$

1) Temp ( 3Y, 2N)

Hot
(0Y, 0N)

mild
(2Y/1N)

Cold
(1Y, 1N)

$E_{Hot} = 0.$

$$E_{mild} = -\frac{2}{3} \log_2 (2/3) - \frac{1}{3} \log_2 (1/3) = \underline{0.918}$$

$E_{cold} = 0.$

$$Gain (s, Temp) = 0.97 - \left(0 + \frac{3}{5} \times 0.918 + 0\right)$$
$$= \underline{0.4192}$$

② Humid ( 3Y, 2N)

High
(1Y, 1N)

Normal
(2Y, 1N)

$E_{High} = 0$

$E_{Normal} = 0.918$

$$Gain (s, Humid) = 0.97 - \left(0 + \frac{3}{5} \times 0.918\right) = \underline{0.4192}$$

③

$$\text{wind } (3Y, 2N)$$

weak — $(3Y, 0N)$
strong — $(0Y, 2N)$

$\bar{E}_{weak} = 0$

$E_{strong} = 0$

$\text{Gain}(s, wind) = 0.97 - (0 + 0)$

$$= 0.97$$

∴ Gain(s, wind) is Highest ∴ wind is our decision node for Rain.

| Final Solution | ⇒

$(9Y / 5N)$

| Outlook | ← Root node

$(2Y/3N)$

| SUNNY | ← Decision node

$(Y4/0N)$

| overcast |

Decision node →

$(3Y/2N)$

| Rain. |

$(2Y/3N)$

| Humid |

$(2Y/3N)$

| Play = yes. |
↑ Leaf node

$(3Y/2N)$

| Wind |

$(0Y/3N)$

| High |

$(2Y, 0N)$

| Normal |

$(3Y/0N)$

| weak |

$(0Y/2N)$

| strong |

| Play = No |
↑ Leaf node

| Play = YES |
↑ Leaf node

| Play = Yes |
↑ Leaf node

| Play = No. |
↖ Leaf node