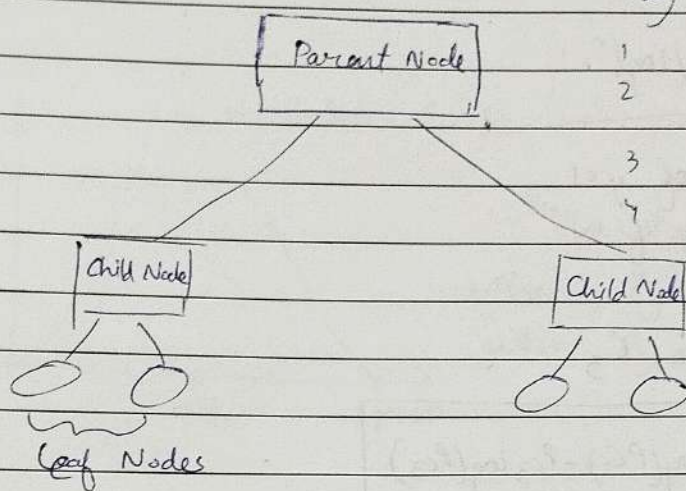


DECISION - TREE

Dataset

Day	Outlook	Temp	Humidity	
1	Sunny	hot	High	Yes
2	Sunny	hot	low	No
3	Overcast	cold	High	No
4	rainfall	cold	low	Yes



① ID3

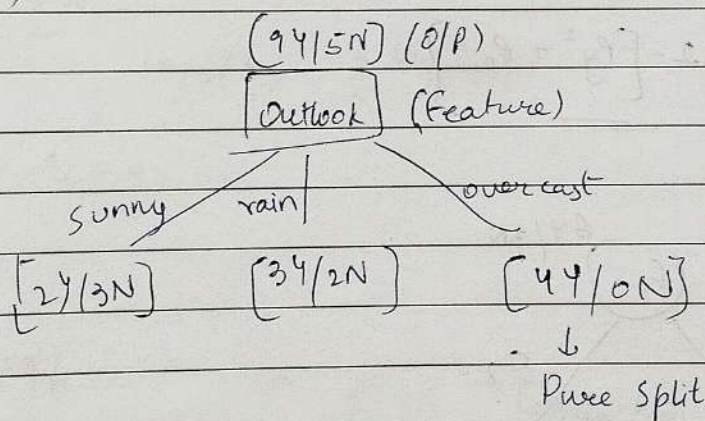
Iterative Decomposition

(entropy)

② CART

Classification & Regression trees

Gini Impurity.



⇒ For checking the purity of the feature :-

① Entropy

② Gini coefficient / Gini impurity

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

Binary Classification

↓
 $-P_y \log(P_y) - P_N \log(P_N)$

P_y → Probability of yes
 P_N → Probability of No

3 Class → C_1, C_2, C_3

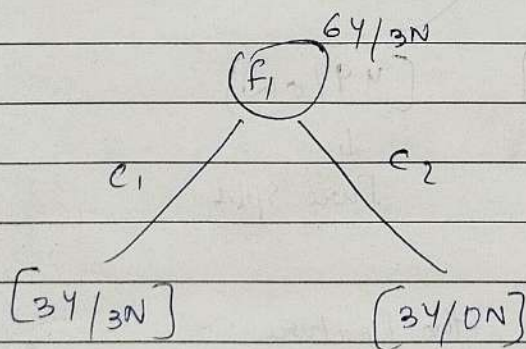
Multi class Classification

$-P_{C_1} \log(P_{C_1}) - P_{C_2} \log(P_{C_2}) - P_{C_3} \log(P_{C_3})$

⇒ Gini coeff → $1 - \sum_{i=1}^n P_i^2$

→ $1 - [P_y^2 + P_N^2]$

Eg:

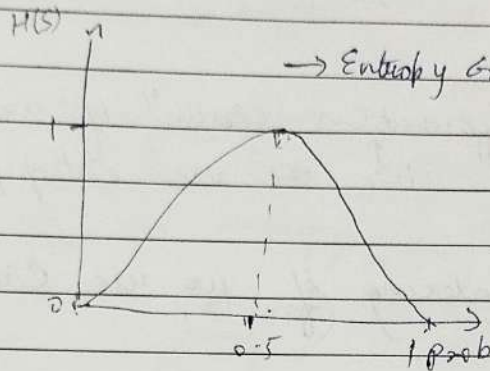


Entropy = $H(S) = -\sum_{i=1}^n P_i \times \log_2(P_i)$

$= -0.5 \log_2(0.5) - 0.5 \log_2(0.5)$
 $\nearrow P_y$ $\nearrow P_N$

$= -1 \log_2(0.5) = -1 \times -1 = 1$

Entropy for other node = $-\log_2(1) - 0\log_2(0)$
 $= 0$



→ Entropy Graph w.r.t probability.

When; $P_Y = 0.5, P_N = 0.5, H(S) = 1$
 When; $P_Y = 1, P_N = 0, H(S) = 0$
 When; $P_Y = 0, P_N = 1, H(S) = 0$

Note :-

$H(S) = 1 \rightarrow$ very impure split
 $H(S) = 0 \rightarrow$ pure split

$2Y/3N$
 $H(S) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$

$H(S) = 0.92$

② Gini coeff.

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

① $3Y/3N \rightarrow$ Entropy-high $\rightarrow H(S)=1 \Rightarrow$ Very impure split

② $3Y/0N \rightarrow$ Entropy-0 \rightarrow Gini index=0

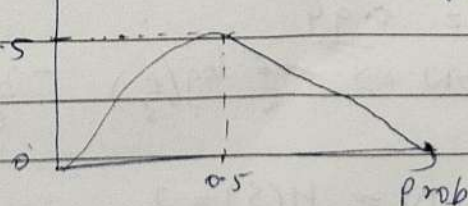
③ $2Y/3N$

$3Y/3N = 1 - (0.5^2 + 0.5^2)$
 $= 0.5$

Gini impurity

(Gini) 0.5

Entropy = $[0, 1]$
 Gini coeff = $[0, 0.5]$



$\frac{0.44}{9 \times \frac{40}{36}}$

$\frac{44}{60N} = 1 - \left(\frac{1}{9} + \frac{4}{9}\right)$

$\frac{32}{48}$
 $\frac{8}{12}$
 $= 1 - \frac{5}{9} = \frac{4}{9} = 0.44$

$$\rightarrow [84/2N] = 1 - \left[\left(\frac{4}{5}\right)^2 + \left(\frac{1}{5}\right)^2 \right]$$

$$= 1 - \left[\frac{16+1}{25} \right] = \frac{25-17}{25} = \frac{8}{25} = 0.32$$

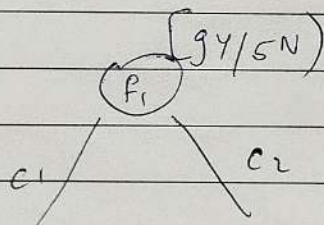
$$\begin{array}{r} 0.32 \\ 25 \overline{) 8.00} \\ \underline{75} \\ 50 \\ \underline{50} \\ 00 \end{array}$$

→ we have to calculate "Information Gain." If we consider ID3 algorithm to calculate IG, we use entropy.

→ we have to use Gini indexing if we use CART algorithm to find "Information Gain."

ID3 $\left\{ \begin{array}{l} \text{Entropy} \\ \text{I.G.} \end{array} \right\}$ Gini Impurity $\left\{ \begin{array}{l} \text{I.G.} \end{array} \right\}$ CART \rightarrow Nowadays this is used.

$$\Rightarrow \boxed{\text{Gain}(S, f_i) = H(S) - \sum_{S \in S_i} \frac{|S_i|}{|S|} (H(S_i))}$$



$$[64/2N]$$

$$[34/3N]$$

$H(S) \Rightarrow$ root Feature Entropy.

$$- H(S) = - \sum_{i=1}^2 p_i (\log(p_i))$$

$$= - \frac{9}{14} \log\left(\frac{9}{14}\right) - \frac{5}{14} \log\left(\frac{5}{14}\right)$$

$$= 0.94$$

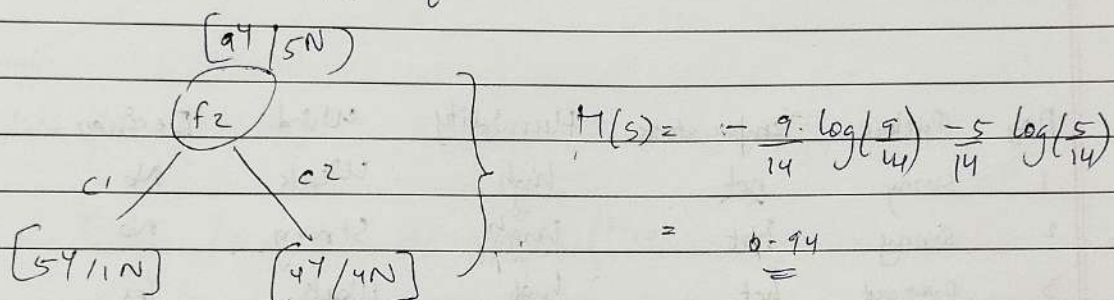
$$- 64/2N \Rightarrow - \frac{6}{8} \log\left(\frac{6}{8}\right) - \frac{2}{8} \log\left(\frac{2}{8}\right) = 0.81$$

$$- 34/3N \Rightarrow H(S) = 1$$

$$\text{Gain}(S, F_1) = 0.94 - \left(\frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right)$$

$F_1 = 0.32$ } Information Gain w.r. to feature 1.

S_{SV} represents total responses of child node $(64/2N) \rightarrow 2$
 $S_S \rightarrow$ total responses
 $H_{SV} \rightarrow$ entropy of child Node.



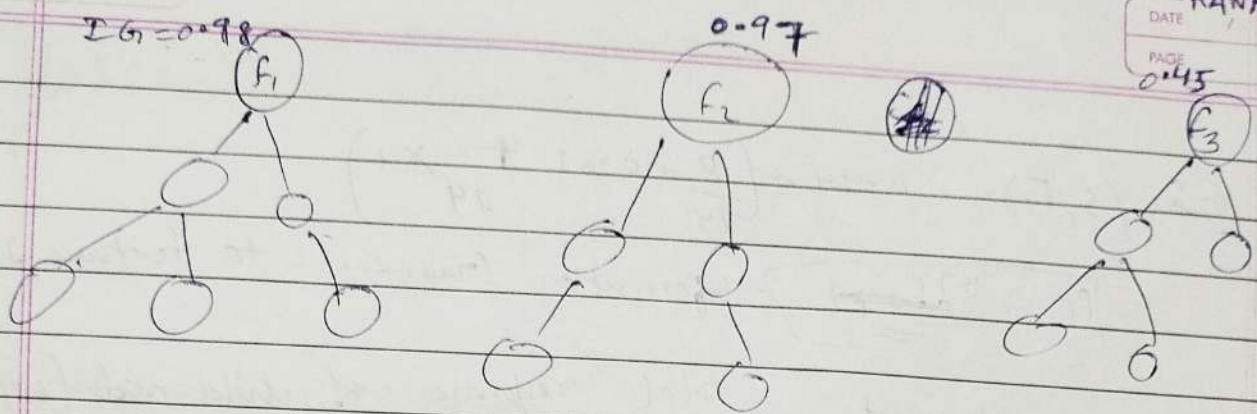
$$\text{Gain} = 0.94 - \left(\frac{5}{14} \times 0.65 - \left(\frac{8}{14} \right) \times 1 \right)$$

$$\text{Gain}(F_2) = 1 - 0.8$$

I.G of f_2 is more,

$$\text{IG}(f_2) > \text{IG}(f_1)$$

\downarrow
 $\text{IG}(f_2)$
 This is greater & it is providing more information.



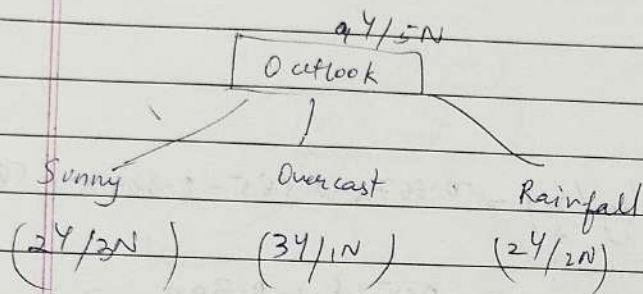
IG = 0.98 is good

Day	Outlook	Temperature	Humidity	Wind	Decision
1	sunny	hot	high	Weak	No
2	sunny	hot	high	Strong	No
3	overcast	hot	high	Weak	Yes
4	rainfall	mild	high	Weak	Yes
5	rainfall	cool	Normal	Weak	Yes
6	rainfall	cool	Normal	Strong	No
7	overcast	cool	Normal	Strong	Yes
8	sunny	mild	High	weak	No
9	sunny	cool	Normal	weak	Yes
10	rainfall	mild	Normal	Weak	Yes
11	sunny	mild	Normal	strong	Yes
12	overcast	mild	High	strong	Yes
13	overcast	hot	Normal	weak	Yes
14	rainfall	mild	High	Strong	No

4/25
4/25

Assignment

- 1) Diff b/w Gini Impurity & Entropy.
 - 2) Take this tennis play dataset & build Decision Tree from so
- (2) Information Gain



$$H_{\text{train}}(S, F_i) = H_S - \sum_{S \in S} \frac{|S|}{N} (H_{S|S})$$

$$H_S = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = -0.64 \times -0.6374 - 0.35 \times -1.40$$

$$= 0.407 + 0.5194$$

$$= 0.9264$$

$$H_{S|S_{\text{Sunny}}} = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) = -0.4 \times -1.322 - 0.6 \times -0.737$$

$$= 0.5288 + 0.4422$$

$$= 0.971$$

$$(H_{S|S})_{\text{Overcast}} = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) = -0.75 \times -0.415 - 0.25 \times -2$$

$$= 0.31125 + 0.5 = 0.81125$$

$$(H_{S|S})_{\text{Rainfall}} = -0.5 \log_2(0.5) - 0.5 \log_2(0.5) =$$

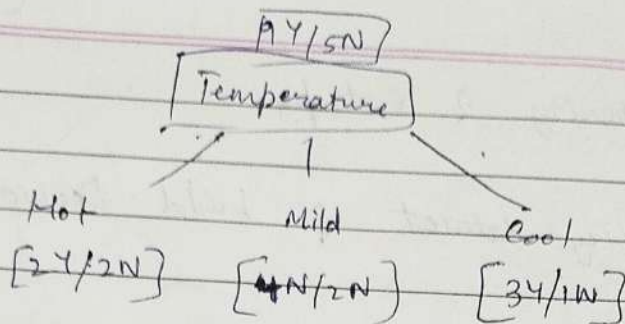
$$= 2 \times -0.5 \times -1 = 2 \times 0.5 = 1$$

$$\text{Gain} = 0.9264 - \left[\frac{5}{14} \times 0.971 + \frac{4}{14} \times 0.81125 + \frac{4}{14} \times 1 \right]$$

$$= 0.9264 - [0.346 + 0.231 + 0.285] = 0.9264 - 0.8627$$

$$= 0.0637$$

⇒



$$H_S = 0.9264$$

$$H_{Sv} = 2$$

$$H_{Sv}^{mid} = -\frac{4}{6} \log\left(\frac{4}{6}\right) - \frac{2}{6} \log\left(\frac{2}{6}\right) = 0.667 \times -0.585 - 0.334 \times -1.585$$

$$= 0.528 + 0.390$$

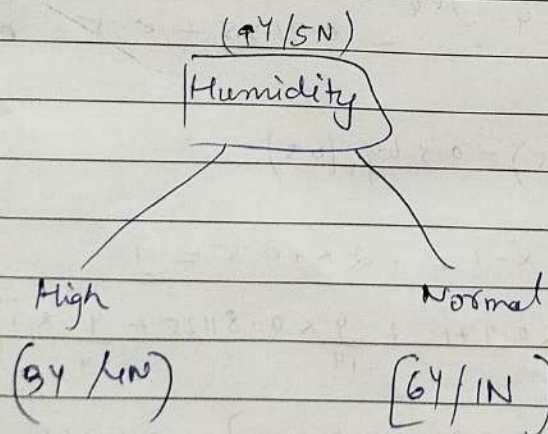
$$H_{Sv}^{cool} = -\frac{3}{4} \log\left(\frac{3}{4}\right) - \frac{1}{4} \log\left(\frac{1}{4}\right) = 0.81125$$

$$G_{rain} = 0.9264 - \left[\frac{4 \times 1}{14} + \frac{6}{14} \times 0.918 + \frac{4 \times 0.81125}{14} \right]$$

$$= 0.9264 - [0.393 + 0.231 + 0.285]$$

$$= 0.9264 - 0.9097 = 0.0167$$

⇒



$$H_S = 0.9264$$

$$H_{Sv}^{(High)} = -\frac{3}{7} \log\left(\frac{3}{7}\right) - \frac{4}{7} \log\left(\frac{4}{7}\right)$$

$$= -0.4285 \times -1.2224 - 0.5714$$

$$= 0.4613 + 0.5244 = 0.9857$$

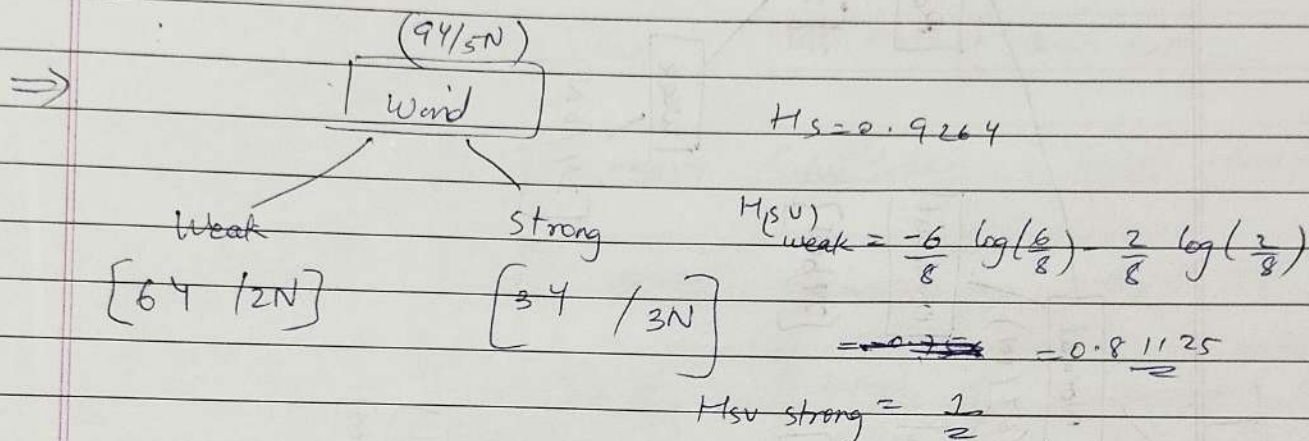
$$H_{Sv}^{(Normal)} = -\frac{6}{7} \log\left(\frac{6}{7}\right) - \frac{1}{7} \log\left(\frac{1}{7}\right)$$

$$= -0.8571 \times -0.2224 - 2.8074 \times 0.1429$$

$$= 0.4010 + 0.1905$$

$$= 0.5915$$

$$\begin{aligned} \text{Gain} &= 0.9264 - \left[0.5 \times 0.9857 + 0.5 \times 0.57159 \right] \\ &= 0.9264 - (0.49285 + 0.2957) \\ &= 0.9264 - 0.78855 = \underline{0.13785} \end{aligned}$$

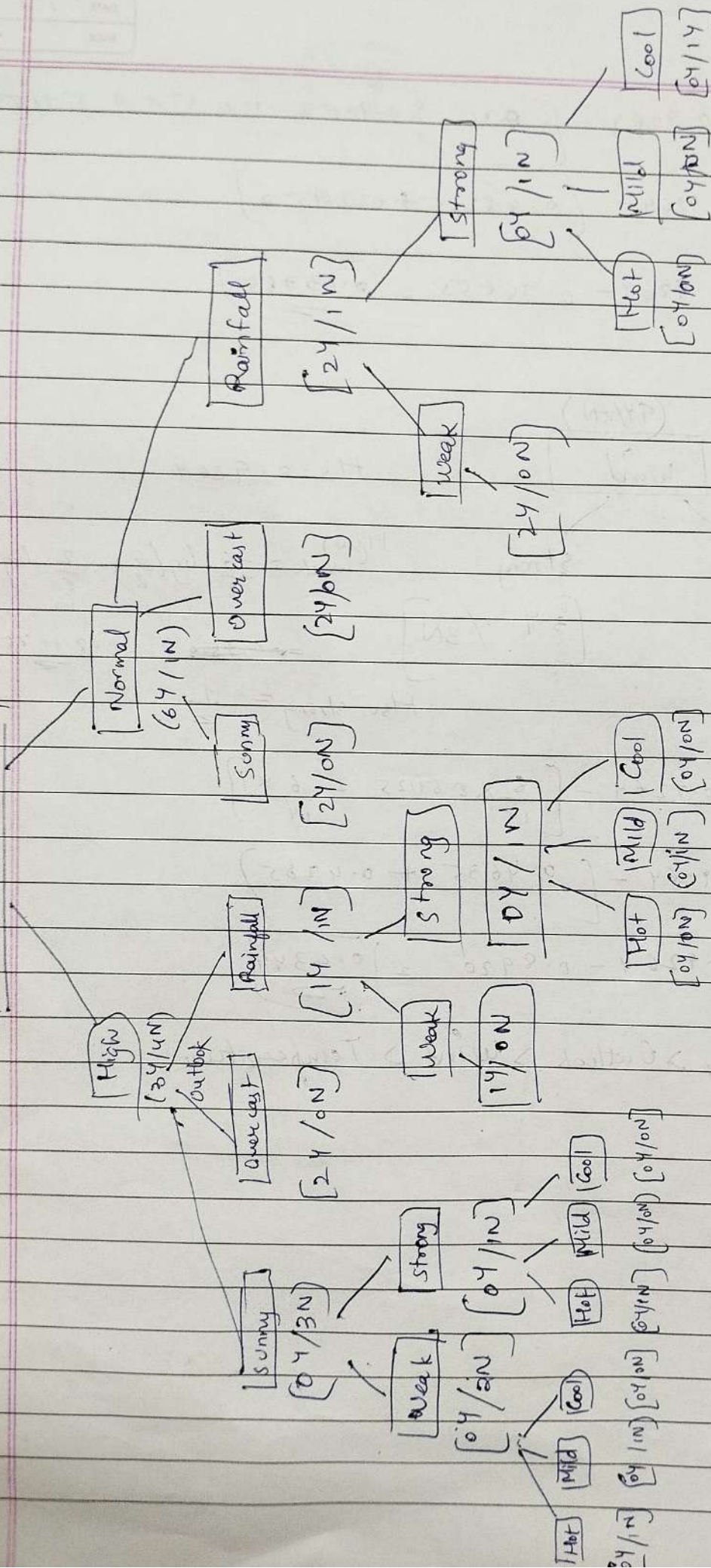


$$\begin{aligned} \text{Gain} &= 0.9264 - \left[\frac{8}{14} \times 0.81125 + \frac{6}{14} \times 1 \right] \\ &= 0.9264 - (0.4635 + 0.4285) \\ &= 0.9264 - 0.8920 = \underline{0.0344} \end{aligned}$$

I-G

Humidity > Outlook > Wind > Temperature

94/5N
HUMIDITY



Entropy

1) Range $\rightarrow [0-1]$

0 \rightarrow Perfectly pure, all data pts. belong to one class

1 \rightarrow Max. impurity \rightarrow equal distrib. -ution of classes.

Gini Index

1) Range $- 0 - 0.5$

0 \rightarrow Perfectly pure (Data pts belong to one)

0.5 \rightarrow Max Impurity (equal distribution)

2) More Robust than Gini Index.

less Robust than Entropy.

3) Entropy is logarithmic measure Gini Index is linear measure.

4) Used in ID3 & C4.5 algorithms. Gini Index is used in CART

5) Entropy measures the amount of uncertainty or randomness in a set. Probability of misclassifying a randomly chosen element in a set