

Can you prove math

Info gain

entropy

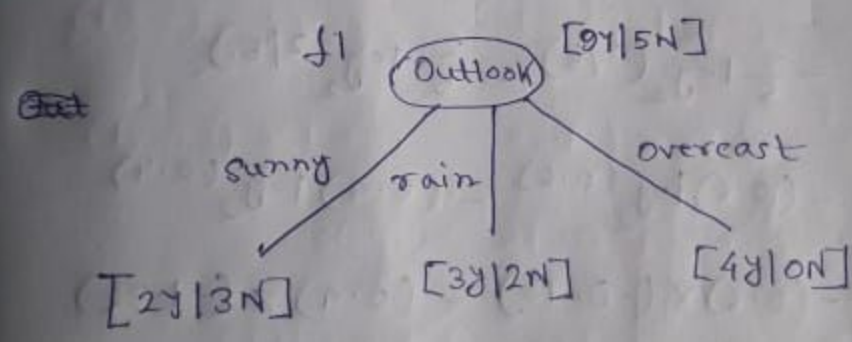
Gini impurity

② entropy vs Gini-impurity [5 diff]

③ take this data [tenis play] build DT from scratch (mathematical representation)

DTR | RFR | RFC

④ Making of Decision Tree (Tennis Play Dataset) and Mathematical Representation



$H(S) \Rightarrow$  root featuring entropy

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

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

$$\sim (0.4) \log_2(0.4) - (0.6) \log_2(0.6)$$

$$\sim (0.4) \times (-1.3) - (0.6) \times (-0.7)$$

Prove it

$$\begin{aligned}
 &= -\frac{2}{14} \log_2(2/14) - \frac{5}{14} \log_2(5/14) \\
 &= -(0.4) \log_2(0.64) - (0.35) \log_2(0.35) \\
 &= -(0.64) \times (-0.64) - (0.35) \times (-1.51) \\
 &= 0.40 + 0.52 \\
 &\approx 0.93
 \end{aligned}$$

$$\begin{aligned}
 28/3N &\Rightarrow -\frac{2}{5} \log_2(2/5) - \frac{3}{5} \log_2(3/5) \\
 &= -(0.4) \log_2(0.4) - (0.6) \log_2(0.6) \\
 &= -(0.4) \times (-1.32) - (0.6) \times (-0.73) \\
 &= 0.52 + 0.43 \\
 &= 0.95
 \end{aligned}$$

$$\begin{aligned}
 32/2N &\Rightarrow -\frac{3}{5} \log_2(3/5) - \frac{2}{5} \log_2(2/5) \\
 &= -(0.6) \log_2(0.6) - (0.4) \log_2(0.4) \\
 &= -(0.6) \times (-0.73) - (0.4) \times (-1.32) \\
 &= 0.43 + 0.52 \\
 &= 0.95
 \end{aligned}$$

$$\begin{aligned}
 48/10N &\Rightarrow -\frac{4}{4} \log_2(4/4) - \frac{0}{4} \log_2(0/4) \\
 &= -1 \log_2 1 - 0 \\
 &= 0
 \end{aligned}$$

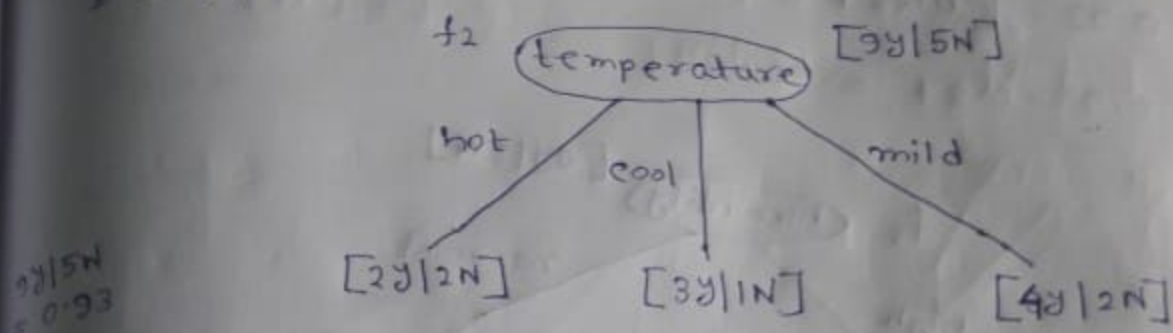
$$\text{Gain}(s, f_1) = H_s - \sum_{v \in \text{val}} \frac{|s_v|}{|s|} \times H(s_v)$$

$$= 0.93 - \left[ \frac{5}{14} \times 0.95 + \frac{5}{14} \times 0.95 + \frac{4}{14} \times 0 \right]$$

$$= 0.93 - [0.33 + 0.33]$$

$$= 0.93 - 0.66$$

$$= 0.27$$



$$2y|2N \Rightarrow -\frac{2}{4} \log_2(2/4) - \frac{2}{4} \log_2(2/4)$$

$$= 1$$

$$3y|1N \Rightarrow -\frac{3}{4} \log_2(3/4) - \frac{1}{4} \log_2(1/4)$$

$$= -(0.75) \times \log_2(0.75) - (0.25) \log_2(0.25)$$

$$= -(0.75) \times (-0.41) - (0.25) \times (-2)$$

$$= 0.30 + 0.5 = 0.80$$

$$4y|2N \Rightarrow -\frac{4^2}{36} \log_2(4/6) - \frac{2^1}{36} \log_2(2/6)$$

$$= -\frac{2}{3} \log_2(2/3) - \frac{1}{3} \log_2(1/3)$$

$$= -(0.66) \log_2(0.66) - (0.33) \log_2(0.33)$$

$$= -(0.66) \times (-0.59) - (0.33) \times (-1.59)$$

$$= 0.38 + 0.52$$

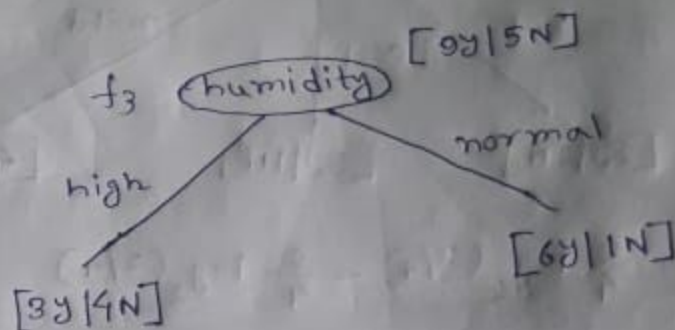
$$= 0.90$$

$$\text{Gain}(S, f_2) = H_S - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} \times H(S_v)$$

$$= 0.93 - \left[ \frac{4^2}{7 \times 1} \times 1 + \frac{4^2}{7 \times 1} \times 0.80 + \frac{4^3}{7 \times 1} \times 0.90 \right]$$



$$\begin{aligned}
 &= 0.93 - \left[ \frac{2}{7} + \frac{2}{7} \times 0.80 + \frac{3}{7} \times 0.90 \right] \\
 &= 0.93 - [0.28 + 0.28 \times 0.80 + (0.42) \times 0.90] \\
 &= 0.93 - [0.28 + 0.22 + 0.37] \\
 &= 0.93 - 0.87 = 0.06
 \end{aligned}$$



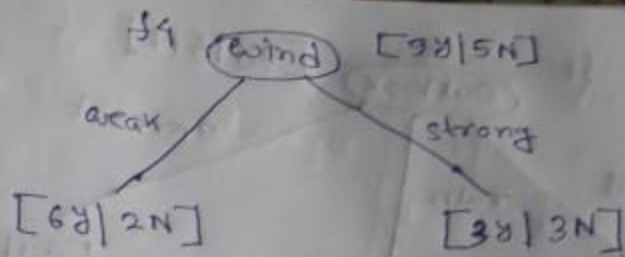
$$9y|5N = 0.93$$

$$\begin{aligned}
 3y|4N &\Rightarrow -\frac{3}{7} \log_2(3/7) - \frac{4}{7} \log_2(4/7) \\
 &= -0.42 \log_2(0.42) - 0.57 \log_2(0.57) \\
 &= -0.42 \times (-1.25) - 0.57 \times (-0.81) \\
 &= 0.52 + 0.46 = 0.98
 \end{aligned}$$

$$\begin{aligned}
 6y|1N &\Rightarrow -\frac{6}{7} \log_2(6/7) - \frac{1}{7} \log_2(1/7) \\
 &= -(0.85) \log_2(0.85) - (0.14) \log_2(0.14) \\
 &= -(0.85) \times (-0.23) - (0.14) \times (-2.83) \\
 &= 0.19 + 0.39 \\
 &= 0.58
 \end{aligned}$$

$$\text{Gain}(S, f_3) = H_S - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} \times H(S_v)$$

$$\begin{aligned}
 &= 0.93 - \left[ \frac{5}{14} \times 0.98 + \frac{2}{14} \times 0.58 \right] \\
 &= 0.93 - 0.78 = 0.15
 \end{aligned}$$



$$98|5N = 0.93$$

$$\begin{aligned} 68|2N &\Rightarrow -\frac{3}{4} \log_2 \left(\frac{3}{4}\right) - \frac{1}{4} \log_2 \left(\frac{1}{4}\right) \\ &= -\frac{3}{4} \log_2 (3/4) - \frac{1}{4} \log_2 (1/4) \\ &= -\frac{3}{4} \times 0.75 \times \log_2 (0.75) - 0.25 \times \log_2 (0.25) \\ &= -0.75 \times (-0.41) - 0.25 \times (-2) \\ &= 0.30 + 0.50 \\ &= 0.80 \end{aligned}$$

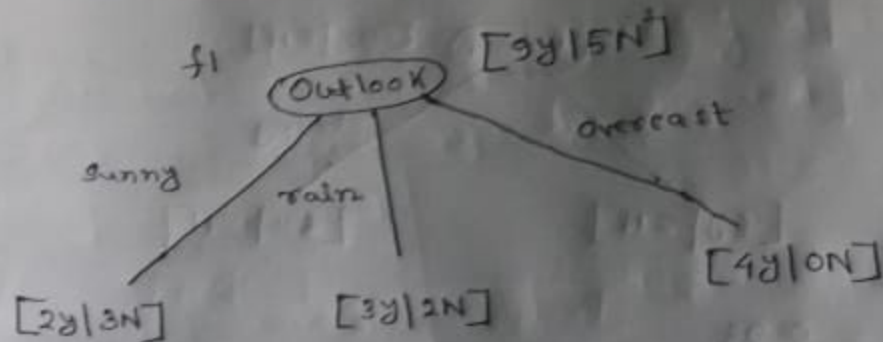
$$38|3N \Rightarrow 1$$

$$\begin{aligned} \text{Gain}(s, f_4) &= H_s - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} \times H(S_v) \\ &= 0.93 - \left[ \frac{2}{4} \times 0.80 + \frac{2}{4} \times 1 \right] \\ &= 0.93 - [0.45 + 0.42] \\ &= 0.93 - 0.87 \\ &= 0.06 \end{aligned}$$

$$\text{Gain}(s, f_1) = 0.27, \quad \text{Gain}(s, f_2) = 0.06$$

$$\text{Gain}(s, f_3) = 0.15, \quad \text{Gain}(s, f_4) = 0.06$$

$$\therefore \text{Root Node} = \text{Gain}(s, f_1) = \text{Outlook}$$



Outlook	temperature	humidity	wind	Decision
sunny	hot	high	weak	N
sunny	hot	high	strong	N
sunny	mild	high	weak	Y
sunny	cool	normal	weak	Y
sunny	mild	normal	strong	Y



$$2y|3N \Rightarrow 0.95$$

$$0y|2N \Rightarrow 0$$

$$1y|1N \Rightarrow 1$$

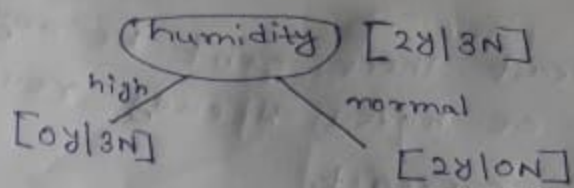
$$1y|0N \Rightarrow 0$$

$$\text{Gain}(\text{temp}) = 0.95 - \left( \frac{2}{5} \times 0 + \frac{2}{5} \times 1 + \frac{1}{5} \times 0 \right)$$

$$= 0.95 - 0.40$$

$$= 0.55$$

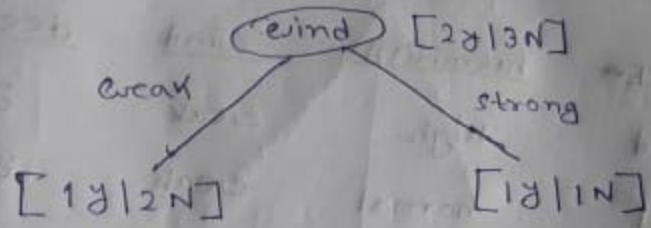




$$08/3N \Rightarrow 0$$

$$28/10N \Rightarrow 0$$

$$\text{Gain}(S, \text{humidity}) = 0.95 - (0+0) = 0.95$$



$$18/2N \Rightarrow -\frac{1}{3} \log_2(1/3) - \frac{2}{3} \log_2(2/3)$$

$$= -(0.33) \log_2(0.33) - (0.66) \log_2(0.66)$$

$$= -(0.33) \times (-1.59) - (0.66) \times (-0.59)$$

$$= 0.52 + 0.38$$

$$\approx 0.91$$

$$18/1N \Rightarrow 1$$

$$\text{Gain}(S, \text{wind}) = 0.95 - \left( \frac{3}{5} \times 0.91 + \frac{1}{5} \times 1 \right)$$

$$= 0.95 - (0.54 + 0.4)$$

$$= 0.95 - 0.94$$

$$= 0.01$$

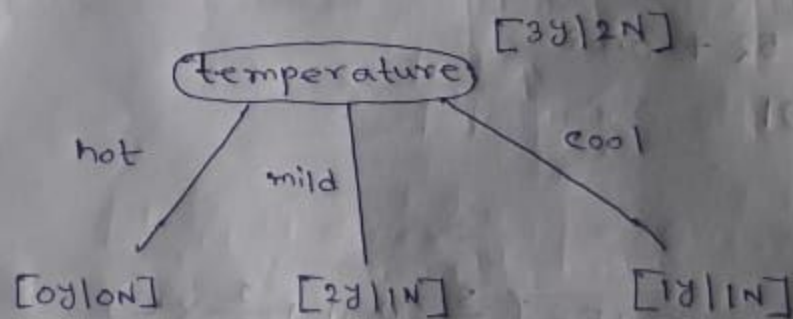
$$\text{Gain}(S, \text{humidity}) > \text{Gain}(S, \text{temperature}) > \text{Gain}(S, \text{wind})$$

∴ Gain(S, humidity) is highest i.e. (0.95)  
 ∴ Humidity is decision node for sunny.

\* for overcast there can't be any decision node because we have already reached the leaf node for overcast.

\* Now Doing calculations for Decision Node (Rain)  
(3Y and 2N)

<u>Outlook</u>	<u>temp</u>	<u>humidity</u>	<u>wind</u>	<u>decision</u>
rainfall	mild	high	weak	yes
rainfall	cool	normal	weak	yes
rainfall	cool	normal	strong	no
rainfall	mild	normal	weak	yes
rainfall	mild	high	strong	no



$$0Y|0N \Rightarrow 0$$

$$2Y|1N \Rightarrow -\frac{2}{3} \log_2 \left( \frac{2}{3} \right) - \frac{1}{3} \log_2 \left( \frac{1}{3} \right)$$

$$\approx 0.91$$

$$1Y|1N \Rightarrow 1$$

$$\text{Gain}(S, \text{temp}) = 0.95 - \left[ 0 + \frac{3}{5} \times 0.91 + \frac{2}{5} \times 1 \right]$$

$$= 0.95 - [0 + 0.54 + 0.40]$$

$$= 0.95 - 0.94 = 0.01$$



node  
leaf

(Rain)

00

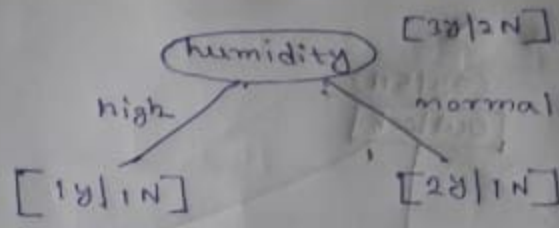
5

0

11/5

0000

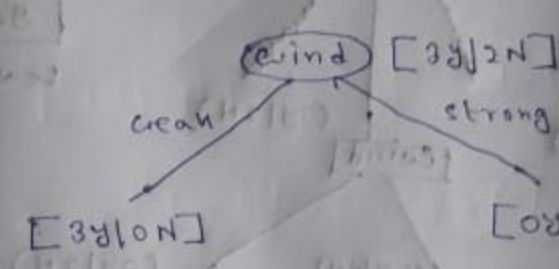
$\times 1$



$$18|1N \Rightarrow 1$$

$$28|1N \Rightarrow 0.01$$

$$\text{Gain}(S_{\text{humidity}}) = 0.01$$



$$38|0N = 0$$

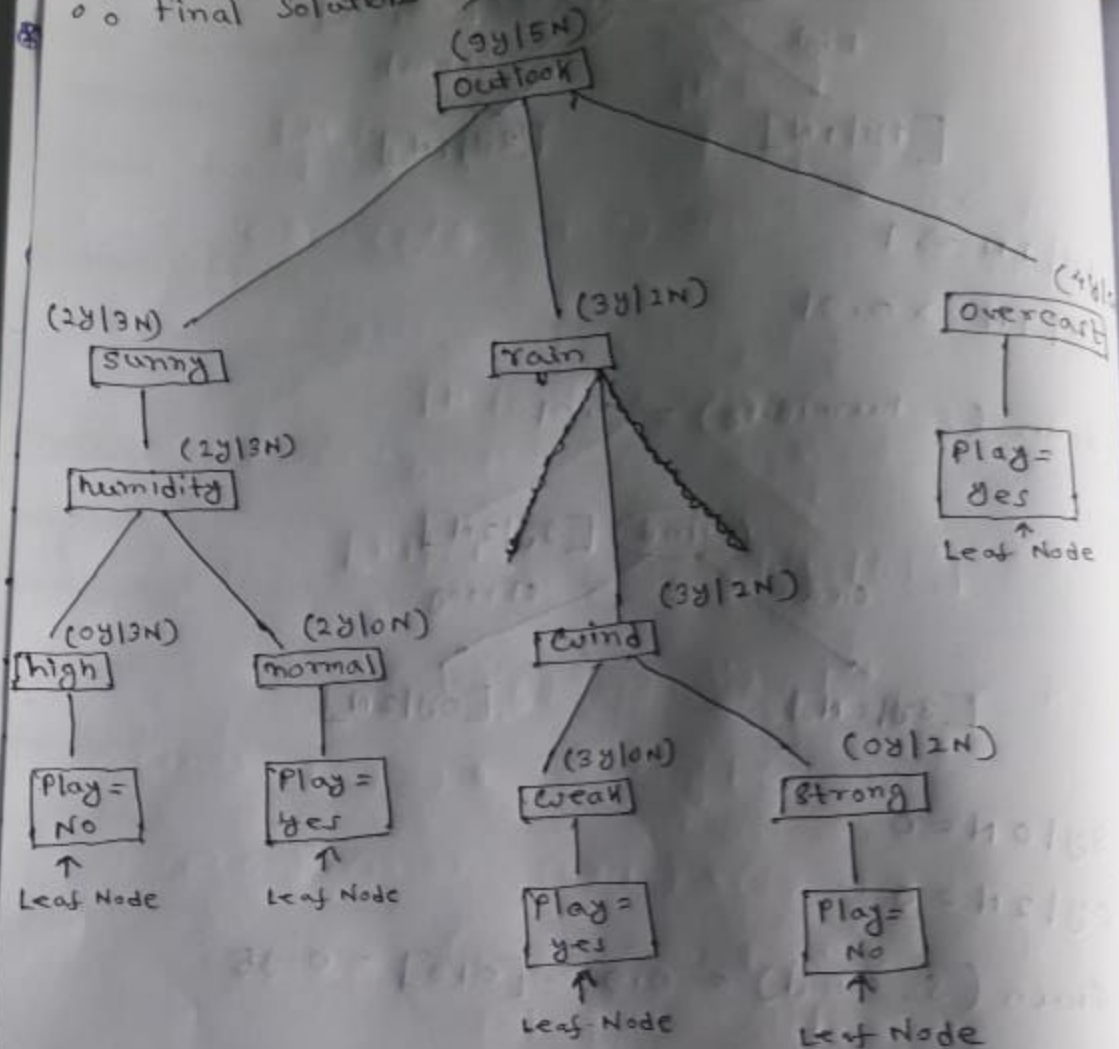
$$08|2N = 0$$

$$\text{Gain}(S_{\text{wind}}) = 0.95 - [0+0] = 0.95$$

∴  $\text{Gain}(S_{\text{wind}})$  is highest

∴ Wind is our decision node for Rain.

Final Solution  $\Rightarrow$



### ⊗ Entropy Vs Gini Impurity

① If we use ID3 approach for calculating information gain we will use calculate entropy.

If we use CART approach for calculating information gain then we will calculate gini impurity

② Gini impurity is faster than Entropy because in Entropy we use log functions to calculate

entropy and in Gini impurity we don't use log function to calculate Gini impurity

② Entropy should be used when less number of features are present in the data set and Gini Impurity should be used when more number of ~~the~~ features are present in the data set.

③ The range of Entropy lies in between 0 to 1 and the range of Gini impurity lies in between 0 to 0.5.

④ formula of ~~Gini~~ Entropy is

$$\sum_{i=1}^n P_i \times \log(P_i)$$

formula of Gini Impurity is  $1 - \sum_{i=1}^n P_i^2$