

Decision Tree

Play Tennis Dataset

- To Find out the root node we will follow 2 approaches.

① Information Gain \rightarrow Max Value

② Gini Impurity \rightarrow Min Value

Information Gain

$$\text{Entropy} = \sum -P \log_2(P)$$

Label \rightarrow Play Tennis

Output \rightarrow $Y=9$ $N=5$

$$E(L) = -P(Y) \log P(Y) - P(N) \log P(N)$$

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

$$= -\frac{9}{14} \times (-0.64) - \frac{5}{14} \times (-1.49)$$

$$= 0.41 + 0.54$$

$$= \underline{\underline{0.95}}$$

Entropy of Individual class in feature

<u>Outlook</u>	<u>Play Tennis</u>	<u>TV</u>	<u>P(Y)</u>	<u>P(N)</u>
Sunny	$Y=2$ $N=3$	5	$\frac{2}{5}$	$\frac{3}{5}$
Overcast	$Y=4$ $N=0$	4	$\frac{4}{4}$	$\frac{0}{4}$
Rainfall	$Y=3$ $N=2$	5	$\frac{3}{5}$	$\frac{2}{5}$

$$E(\text{Sunny}) = -\left(\frac{2}{5}\right) \log\left(\frac{2}{5}\right) - \left(\frac{3}{5}\right) \log\left(\frac{3}{5}\right)$$

$$= -0.4 \times (-1.32) - 0.6 \times (-0.74)$$

$$= 0.53 + 0.44$$

$$= \underline{\underline{0.97}}$$

$$E(\text{Overcast}) = -\left(\frac{4}{4}\right) \log\left(\frac{4}{4}\right) - \frac{0}{4} \log\left(\frac{0}{4}\right)$$

$$= 0$$

$$\begin{aligned}
 E(\text{each Fall}) &= -\left(\frac{2}{5}\right) \log\left(\frac{2}{5}\right) - \left(\frac{2}{5}\right) \log\left(\frac{2}{5}\right) \\
 &= -0.6 \times (-0.74) - 0.4 \times (-1.32) \\
 &= 0.42 + 0.53 \\
 &= 0.97
 \end{aligned}$$

$$E(\text{class}) = \sum \frac{\text{No of observations}}{\text{Total no. of observations}} \times E_i$$

$$\begin{aligned}
 E(\text{outlook}) &= \frac{5}{14} \times 0.97 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.97 \\
 &= 0.35 + 0.35 \\
 &= 0.70
 \end{aligned}$$

$$\text{Information Gain} = E_{\text{before}} - E_{\text{after}}$$

$$\begin{aligned}
 IG &= 0.95 - 0.70 \\
 &= 0.25
 \end{aligned}$$

Temperature	Play Tennis	TV	P(Y)	P(N)
hot	Y=2 N=2	4	$2/4 = \frac{1}{2}$	$2/4 = \frac{1}{2}$
mild	Y=4 N=2	6	$4/6 = \frac{2}{3}$	$2/6 = \frac{1}{3}$
cool	Y=3 N=1	4	$3/4$	$1/4$

$$\begin{aligned}
 E(\text{mild}) &= -\frac{2}{3} \log\left(\frac{2}{3}\right) - \frac{1}{3} \log\left(\frac{1}{3}\right) \\
 &= -\frac{2}{3} \times (-0.59) - \frac{1}{3} \times (-1.59) \\
 &= 0.39 + 0.53 = 0.92
 \end{aligned}$$

$$\begin{aligned}
 E(\text{hot}) &= -\frac{1}{2} \log\left(\frac{1}{2}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right) \\
 &= -\frac{1}{2} (-1) - \frac{1}{2} (-1) \\
 &= 1
 \end{aligned}$$

$$\begin{aligned}
 E(\text{cool}) &= -\frac{3}{4} \log\left(\frac{3}{4}\right) - \frac{1}{4} \log\left(\frac{1}{4}\right) \\
 &= -\frac{3}{4} (-0.42) - \frac{1}{4} (-2) \\
 &= 0.32 + 0.5 \\
 &= 0.82
 \end{aligned}$$

$$E(\text{Temperature}) = \frac{4}{14} \times 1 + \frac{6}{14} \times 0.92 + \frac{4}{14} \times 0.82$$

$$= 0.28 + 0.39 + 0.23$$

$$= 0.90$$

$$I.G = 0.95 - 0.90$$

$$= \underline{\underline{0.05}}$$

<u>Humidity</u>	<u>Play Tennis</u>	<u>TV</u>	<u>PCY</u>	<u>PCN</u>
high	Y=3 N=4	7	3/7	4/7
normal	Y=6 N=1	7	6/7	1/7

$$E(\text{high}) = -\frac{3}{7} \log\left(\frac{3}{7}\right) - \frac{4}{7} \log\left(\frac{4}{7}\right)$$

$$= -\frac{3}{7}(-1.22) - \frac{4}{7}(-0.81)$$

$$= 0.52 + 0.46$$

$$= 0.98$$

$$E(\text{normal}) = -\frac{6}{7} \log\left(\frac{6}{7}\right) - \frac{1}{7} \log\left(\frac{1}{7}\right)$$

$$= -\frac{6}{7}(-0.22) - \frac{1}{7}(-2.80)$$

$$= 0.19 + 0.4$$

$$= 0.59$$

$$E(\text{Humidity}) = \frac{7}{14} \times 0.98 + \frac{7}{14} \times 0.59$$

$$= 0.49 + 0.295$$

$$= 0.79$$

$$IG = 0.95 - 0.79$$

$$= \underline{\underline{0.16}}$$

Wind	Tennis play	TV	PCV	PCW
Weak	$Y = 6 \quad N = 2$	8	$6/8 = \frac{3}{4}$	$2/8 = \frac{1}{4}$
Strong	$Y = 3 \quad N = 3$	6	$3/6 = \frac{1}{2}$	$3/6 = \frac{1}{2}$

$$\begin{aligned}
 E(\text{Weak}) &= -\frac{3}{4} \log\left(\frac{3}{4}\right) - \frac{1}{4} \log\left(\frac{1}{4}\right) \\
 &= -\frac{3}{4} \log(-0.42) - \frac{1}{4} \log(-2) \\
 &= 0.32 + 0.5 \\
 &= 0.82
 \end{aligned}$$

$$\begin{aligned}
 E(\text{Strong}) &= -\frac{1}{2} \log\left(\frac{1}{2}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right) \\
 &= \frac{1}{2} + \frac{1}{2} = 1
 \end{aligned}$$

$$\begin{aligned}
 E(\text{Wind}) &= \frac{8}{14} \times 0.82 + \frac{6}{14} \times 1 \\
 &= 0.47 + 0.43 \\
 &= 0.90
 \end{aligned}$$

$$\begin{aligned}
 IG &= 0.95 - 0.90 \\
 &= 0.05 \\
 &=
 \end{aligned}$$

- As outlook feature has maximum Information Gain, we will consider it as Root Node.

Gain Impurity

$$G = 1 - \sum_{i=1}^c (P_i)^2$$

$$\begin{aligned} G(U) &= 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 \\ &= 1 - \frac{81}{196} - \frac{25}{196} \\ &= \frac{90}{196} = \underline{\underline{0.46}} \end{aligned}$$

Outlook

$$G(\text{sunny}) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 1 - \frac{9}{25} - \frac{4}{25} = \frac{12}{25} = 0.48$$

$$G(\text{overcast}) = 1 - \left(\frac{4}{4}\right)^2 - \left(\frac{0}{4}\right)^2 = 0$$

$$G(\text{rainfall}) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0.48$$

$$G(\text{class}) = \sum \frac{\text{No. of observations}}{\text{Total no. of observations}} \times G_c$$

$$\begin{aligned} G(\text{outlook}) &= \frac{5}{14} \times 0.48 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.48 \\ &= 0.17 + 0 + 0.17 \\ &= \underline{\underline{0.34}} \end{aligned}$$

$$\begin{aligned} IG &= 0.48 - 0.34 \\ &= \underline{\underline{0.14}} \end{aligned}$$

Temperature

$$\begin{aligned} G(\text{hot}) &= 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 \\ &= 1 - \frac{1}{4} - \frac{1}{4} = \frac{1}{2} = 0.5 \end{aligned}$$

$$\begin{aligned} G(\text{mild}) &= 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 \\ &= 1 - \frac{4}{9} - \frac{1}{9} = \frac{4}{9} = 0.44 \end{aligned}$$

$$\begin{aligned} G(\text{cool}) &= 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \\ &= 1 - \frac{9}{16} - \frac{1}{16} = \frac{6}{16} = 0.375 \end{aligned}$$

$$\begin{aligned} G(\text{temperature}) &= \frac{4}{14} \times 0.5 + \frac{6}{14} \times 0.44 + \frac{4}{14} \times 0.375 \\ &= 0.14 + 0.19 + 0.11 \\ &= \underline{\underline{0.43}} \end{aligned}$$

$$IG = 0.48 - 0.43 = 0.05$$

Humidity

$$\begin{aligned} G(\text{High}) &= 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 \\ &= 1 - \frac{9}{49} - \frac{16}{49} \\ &= \frac{24}{49} = 0.48 \end{aligned}$$

$$\begin{aligned} G(\text{normal}) &= 1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2 \\ &= 1 - \frac{36}{49} - \frac{1}{49} \\ &= \frac{12}{49} = 0.24 \end{aligned}$$

$$\begin{aligned} G(\text{Humidity}) &= \frac{7}{14} \times 0.48 + \frac{7}{14} \times 0.24 \\ &= 0.24 + 0.12 = 0.36 \end{aligned}$$

$$\begin{aligned} IG &= 0.48 - 0.36 \\ &= 0.12 \end{aligned}$$

Wind

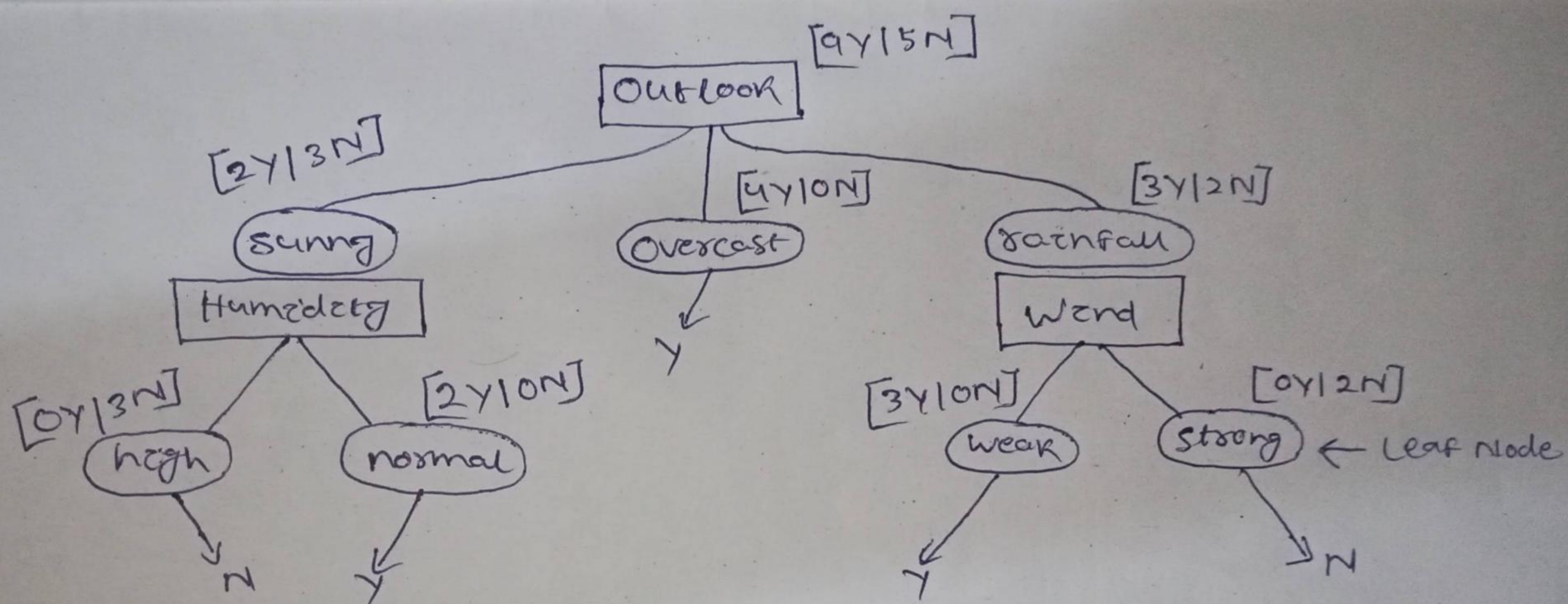
$$\begin{aligned} G(\text{weak}) &= 1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 \\ &= 1 - \frac{36}{64} - \frac{4}{64} \\ &= \frac{64 - 36 - 4}{64} = \frac{24}{64} = 0.375 \end{aligned}$$

$$\begin{aligned} G(\text{strong}) &= 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 \\ &= 1 - \frac{1}{4} - \frac{1}{4} = \frac{2}{4} = 0.5 \end{aligned}$$

$$\begin{aligned} G(\text{Wind}) &= \frac{8}{14} \times 0.375 + \frac{6}{14} \times 0.5 \\ &= 0.21 + 0.21 \\ &= 0.42 \end{aligned}$$

$$\begin{aligned} IG &= 0.48 - 0.42 \\ &= 0.06 \end{aligned}$$

- With respect to the calculation of G , the Impurity and max information gain 'outlook' has least value. So, we will consider it as root node.



- we won't split more, as we already got conclusion. Otherwise our model will be overfitted.