

Notes

Machine learning (Decision Tree)

Decision Tree classifier and Decision Tree Regressor

(DTC)

(DTR)

f_1
↑

f_2
↑

f_3
↑

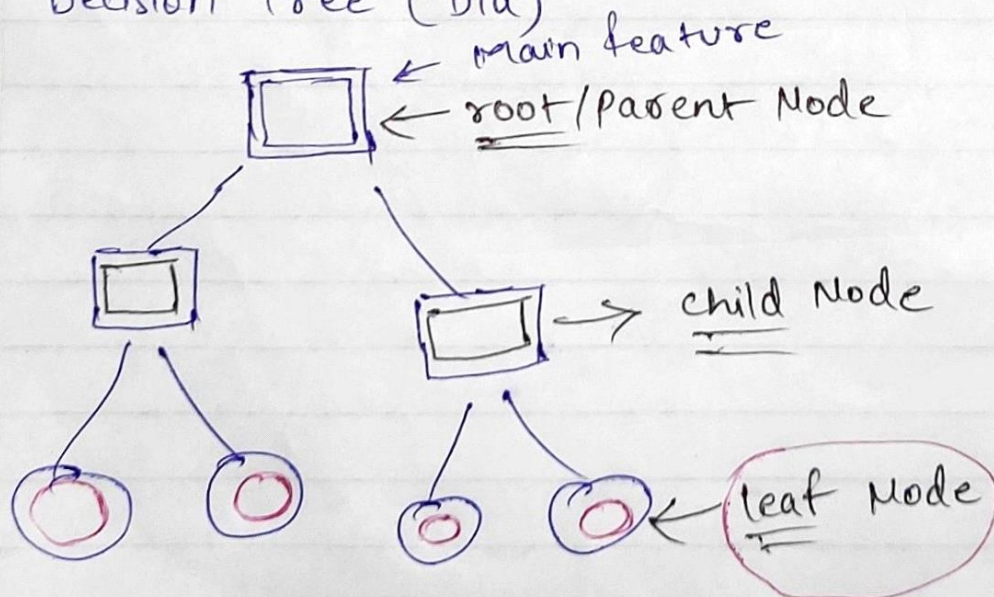
f_4
↑

f_5
↑

O/P.
↑

Day	outlook	temperature	Humidity	wind	yes/No 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	strong	yes
6	overcast	cool	normal	strong	NO
7	sunny	cool	normal	weak	yes
8	rainfall	mild	high	weak	yes

Decision Tree (Dia)



Decision Tree classifier (DTC):

↑
(Num/cat) feature o/p (cate)

Decision Tree Regressor (DTR):

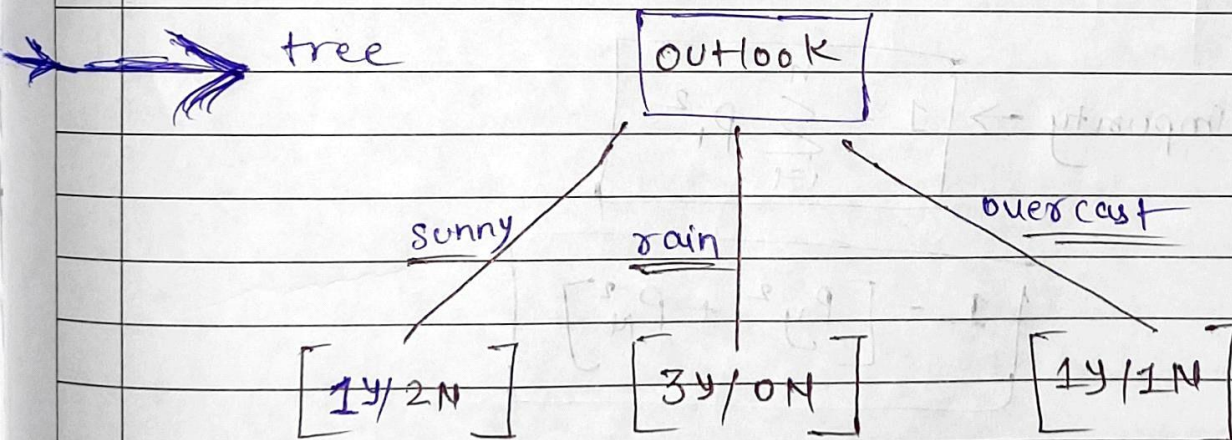
↑
features
(Num/cat) o/p (Num)

① ID₃

iterative Decotomiser

② CART

↳ classification and
Regression tree



Purity \rightarrow ① entropy

② Gini-coeff, gini impurity

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

Binary classification $\xrightarrow{2 \text{ class} \rightarrow Y/N}$

$Y \rightarrow \underline{\text{Yes}} \quad N \rightarrow \underline{\text{No}}$

$$\Rightarrow \boxed{-P_Y \log_2(P_Y) - P_N \log_2(P_N)} \Rightarrow 2 \text{ class}$$

$$\text{Gini coeff} \rightarrow 1 - \sum_{i=1}^n P_i^2$$

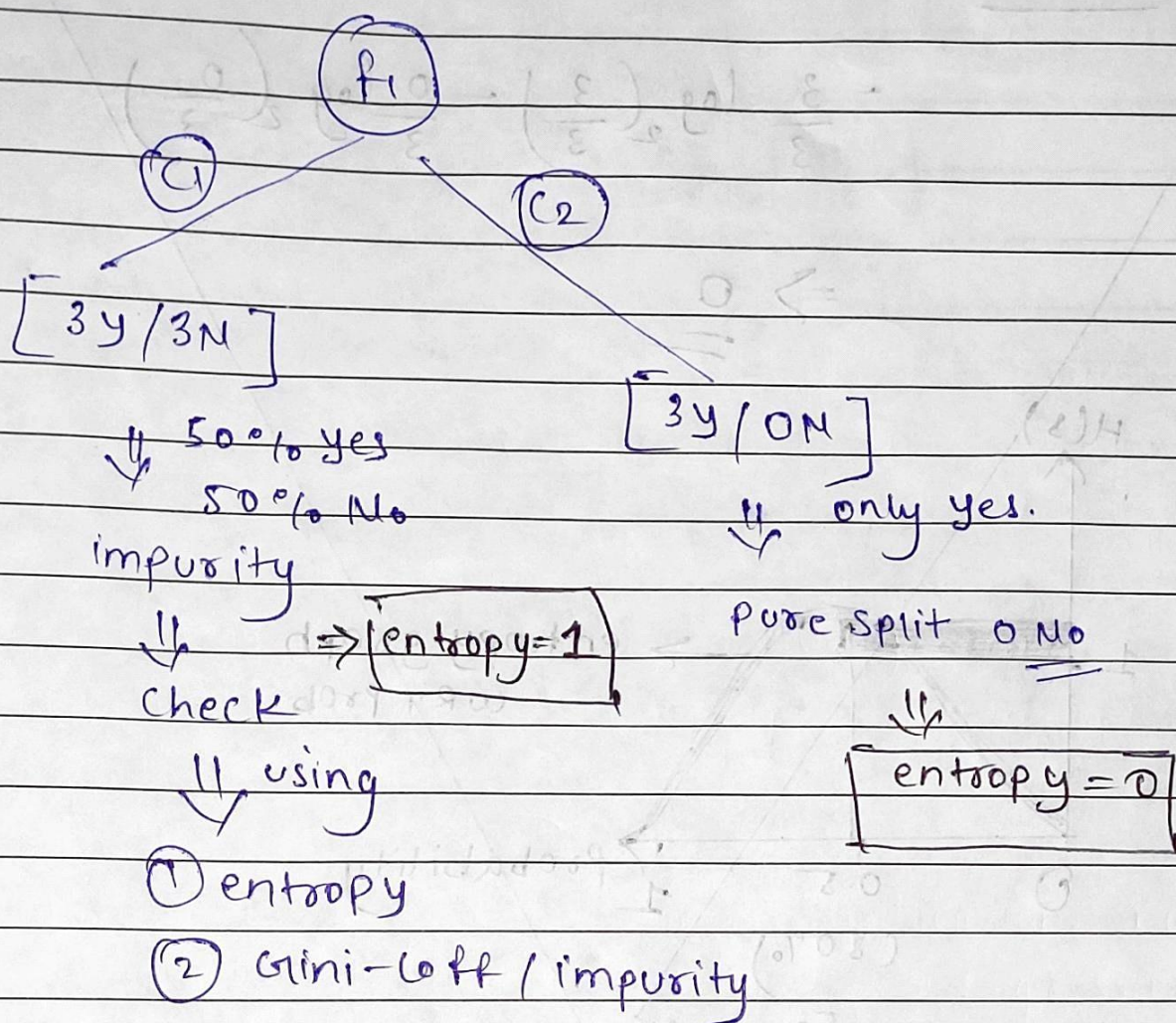
Multi-classification
 \uparrow

3-class Entropy $\rightarrow C_1, C_2, C_3$

$$= P_{C_1} \log_2(P_{C_1}) - P_{C_2} \log_2(P_{C_2}) - P_{C_3} \log_2(P_{C_3})$$

$$\text{Gini-impurity} \rightarrow \boxed{1 - \sum_{i=1}^n P_i^2}$$

$$\boxed{1 - [P_Y^2 + P_N^2]}$$



class 1 \leftarrow n

$$\text{entropy} = \sum_{i=1}^n p_i \log_2(p_i)$$

$$\begin{cases} n = \text{No} \\ y = \text{yes} \end{cases}$$

$$= -p_y \log_2(p_y) - p_n \log_2(p_n)$$

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

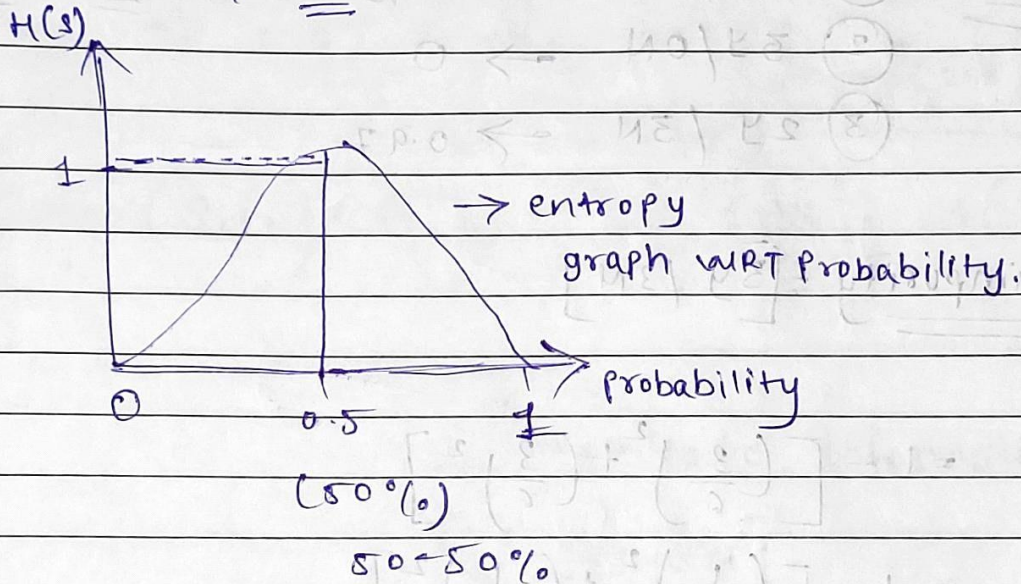
$$= -\frac{1}{2} [\log_2(1) - \log_2(2)] - \frac{1}{2} [\log_2(1) - \log_2(2)]$$

$$= -\frac{1}{2} [0 - 1] - \frac{1}{2} [0 - 1]$$

$$= \frac{1}{2} + \frac{1}{2} = 1$$

class 2 :

$$-\frac{3}{3} \log_2 \left(\frac{3}{3} \right) - \frac{0}{3} \log_2 \left(\frac{0}{3} \right)$$



$H(s) = 1 \rightarrow$ very impure split

$H(s) = 0 \rightarrow$ Pure split

2 yes / 3 no

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

$$H(s) = 0.97$$

Gini-impurity or Gini-coeff :

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

Example : (1) $3Y / 3N \rightarrow \text{Entropy} \Rightarrow H(s) = 1$ (very impure)

(2) $3Y / 0N \rightarrow 0$

(3) $2Y / 3N \rightarrow 0.97$

(1) Gini impurity : $[3Y / 3N]$

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

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

$$\Rightarrow 1 - \left[\frac{1}{4} + \frac{1}{4} \right]$$

$$\Rightarrow 1 - \left[\frac{2}{4} \right] = 1 - \frac{1}{2} = \frac{1}{2} = 0.5$$

(2) $[4Y / 8N] = \text{Gini} \Rightarrow 1 - \left[\left(\frac{4}{12} \right)^2 + \left(\frac{8}{12} \right)^2 \right]$

$$\text{Gini} \Rightarrow 0.44$$

$$③ [8Y/2N]$$

$$\Rightarrow 1 - \left[\left(\frac{8}{10} \right)^2 + \left(\frac{2}{10} \right)^2 \right]$$

$$\underline{\underline{Gini}} \Rightarrow \underline{\underline{0.32}}$$

$$\boxed{\text{Gain}(S, f_r) = H(S) - \sum_{s \in S} \frac{|S_s|}{|S|} H(S_s)}$$

1st split and 2nd split / or no. of split

↓
To calculate information gain

$$f_r [9Y/5N]$$

①

②

step 1+ $[6Y/2N]$

$$[3Y/3N]$$

$$① [9Y/5N] :$$

$H(S) \Rightarrow$ root feature entropy

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

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

$$= -(0.64) \log(0.64) - (0.35) \log(0.35)$$

$$\approx \underline{\underline{0.94}}$$

Teacher's Signature.....

(2) $6Y / 2N$

$$= -\frac{6}{8} \log\left(\frac{6}{8}\right) - \frac{2}{8} \log\left(\frac{2}{8}\right)$$

$$= 0.81$$

(3) $3Y / 8N \Rightarrow H(S) = 1$

$|S_u| \Rightarrow$ total no. of sample after splitting.

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

$$= 0.94 [0.462 + 0.42]$$

$$\approx 0.049$$

\Rightarrow entropy.

\Rightarrow total no. of sample after splitting

$$\text{Gain}(S, h) = H(S) - \sum_{|S|} \frac{|S_u|}{|S|} H(S_u)$$

\Downarrow

ID₃ Method

$$\text{Gain}(f_1) \Rightarrow \underline{\underline{0.049}}$$

step 2: $(f_2) [94/5N]$

$$[54/1N]$$

$$[44/4N]$$

$$\text{Gain} = 0.94 = \left(\frac{6}{14}\right) \times 0.65 - \left(\frac{8}{14}\right) \times 1$$

$$\text{Gain}(f_2) = 0.09 \Rightarrow \text{I.G.}(f_2)$$

$$\boxed{* \text{I.G.}(f_2) > \text{I.G.}(f_1)}$$

this is greater and it is providing more info.