

## \*) Mathematical Intuition of Decision Tree \*

→ Decision Tree classifier :- Consider the following dataset :-

	( $f_1$ ) outlook	( $f_2$ ) temperature	( $f_3$ ) humidity	( $f_4$ ) wind	( $o/p$ ) 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.	overcast	cool	high	weak	No
9.	sunny	mild	normal	weak	Yes
10.	sunny	cool	normal	weak	Yes
11.	rainfall	mild	normal	strong	Yes
12.	sunny	mild	high	strong	Yes
13.	overcast	mild	normal	weak	Yes
14.	overcast	hot	normal	weak	Yes
15.	rainfall	mild	high	strong	No

$f_1, f_2, f_3, f_4$  are independent features.

$o/p \rightarrow$  is our dependent feature.

Our independent features can be both numerical or categorical.

For Decision Tree classifier, our dependent feature (or  $o/p$  feature) will be categorical.

Let's talk about 2 Decision Tree classifiers:-

①  $ID_3$

② CART  $\rightarrow$  Classification & Regression Trees.

$ID_3$  ~~Iterative Decision Tree~~

$ID_3$

$\rightarrow$  Iterative Decision Tree

$\rightarrow$  we talk about entropy

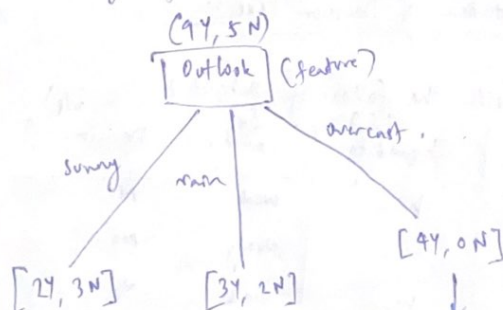
CART

$\rightarrow$  Classification & Regression Trees

$\rightarrow$  we talk about Gini impurity

Let's consider outlook feature as our root node.

Now in our target feature, we have 9 yes & 5 no.



this is a pure split.  
We have only yes values  
& no no values.

this is a leaf node.

Our objective is to get the purity of the feature.  
 To get the purity of our feature, we have 2 methods :-  
 ① Entropy  
 ② Gini impurity / Gini co-efficient

Entropy :-

$$\text{Entropy} = - \sum_{i=1}^n p_i \cdot \log_2(p_i)$$

Gini-coeff :-

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

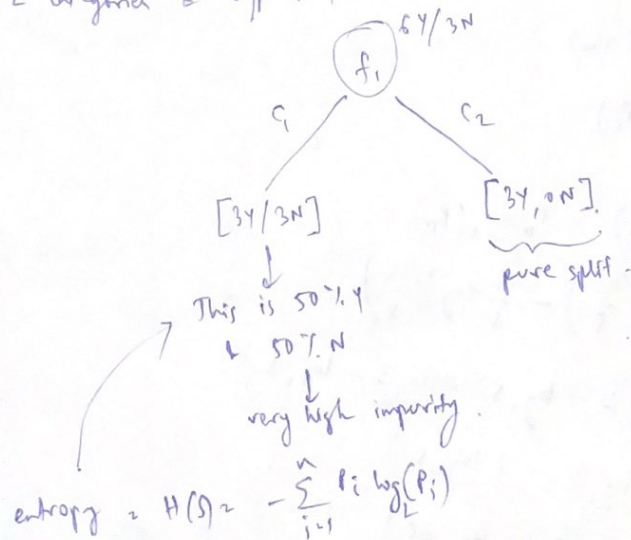
In our dataset, we have 2 classes  $\rightarrow$  Yes & No.  $\rightarrow$  binary classification.

$$\therefore \text{entropy} = - p_Y \log_2(p_Y) - p_N \log_2(p_N)$$

$$= - p_Y \log_2(p_Y) - p_N \log_2(p_N)$$

$$\therefore \text{Gini impurity} = 1 - [p_Y^2 + p_N^2]$$

Say we have a binary classification & based on one feature  $f_1$ , which has 2 categories & o/p  $\rightarrow$  Y or N.



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

$$= - p_Y \log_2(p_Y) - p_N \log_2(p_N)$$

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

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

$$= 1$$

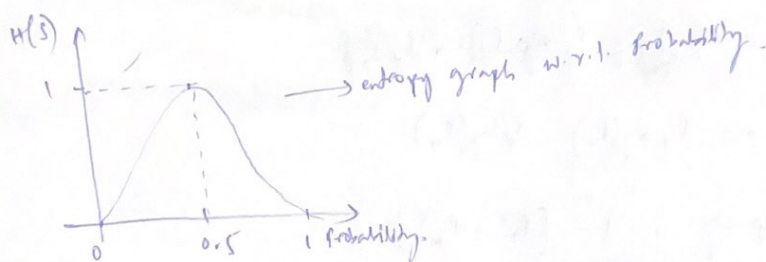
So, for very impure split, we are getting entropy  $H(S) = 1$



Now let's see the entropy of the pure split.

$$\begin{aligned}
 H(s) &= - \sum_{i=1}^n p_i \log_2(p_i) \\
 &= -p_Y \log_2(p_Y) - p_N \log_2(p_N) \\
 &= -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - 0 \\
 &= 0.
 \end{aligned}$$

$\therefore$  for pure split  $H(s) = 0$



Highest value of entropy = 1.  $\rightarrow$  very impure split  
 $H(s) = 0 \rightarrow$  pure split.

Let's consider the split (2Y/3N)

$$\begin{aligned}
 H(s) &= - \sum_{i=1}^n p_i \log_2(p_i) \\
 &= -p_Y \log_2(p_Y) - p_N \log_2(p_N) \\
 &= -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) \\
 &= -(\log_2 2 - \log_2 5 + \log_2 3 - \log_2 5) \\
 &= -(1 - \log_2 3 - 2 \log_2 5) \\
 &= -(1 - \log_2 3 - \log_2 25) \\
 &= 0.97.
 \end{aligned}$$

# ⊗ Gini Coefficient or Gini Impurity ⊗

$$\boxed{\text{Gini coeff} = 1 - \sum_{i=1}^n (p_i)^2}$$

Now let's consider the three splits:-

①  $3Y/3N$

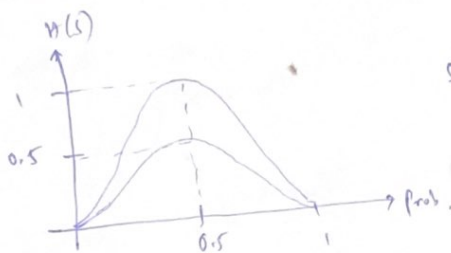
②  $3Y/0N$

③  $2Y/3N$

①  $3Y/3N$

$$\text{Gini coeff} \approx 1 - \left[ \left( \frac{3}{6} \right)^2 + \left( \frac{3}{6} \right)^2 \right]$$

$$\approx 1 - [0.25 + 0.25] \approx 0.5$$



So, range of  $h(s) \rightarrow [0, 1]$   
range of Gini coeff  $\rightarrow [0, 0.5]$

②  $3Y/0N$

$$\text{Gini coeff} \approx 1 - \left[ \left( \frac{3}{3} \right)^2 + \left( \frac{0}{3} \right)^2 \right]$$

$$\approx 0.$$

$$\frac{100}{32}$$

③ consider  $\rightarrow [4Y, 8N]$

$$\text{Gini coeff} \approx 1 - \left[ \left( \frac{4}{12} \right)^2 + \left( \frac{8}{12} \right)^2 \right]$$

$$\approx 1 - \left[ \frac{80}{144} \right] \approx 0.44$$

④ consider  $\rightarrow [8Y/2N]$

$$\text{Gini coeff} \approx 1 - \left[ \left( \frac{8}{10} \right)^2 + \left( \frac{2}{10} \right)^2 \right]$$

$$\approx 1 - \left[ \frac{64}{100} + \frac{4}{100} \right] \approx 1 - \frac{68}{100}$$

$$\approx \frac{100-68}{100} \approx \frac{32}{100} \approx 0.32$$

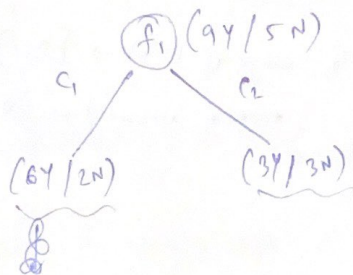
Now, say we have 3 features  $\rightarrow$  feature 1, feature 2, feature 3.

Now, to see Purity, if we find entropy & then the information gain  $\rightarrow$  this is ID3 approach.

If we find Gini impurity & then the information gain  $\rightarrow$  CART approach.  
(this is much faster & suitable for larger datasets)

$$\text{Gain}(S, f_i) = H(S) - \sum \frac{|S_v|}{|S|} H(S_v)$$

Consider this scenario:-



$H(S)$  = root feature entropy.

$$\begin{aligned} &= -p_1 \log_2(p_1) - p_2 \log_2(p_2) \\ &= -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) \\ &= -\log_2\left(\frac{9}{14}\right) - \log_2\left(\frac{5}{14}\right) \\ &= -\log_2 9 - \log_2 5 \\ &\approx 0.94. \end{aligned}$$

for 64/2N

$$H(S) = -\frac{6}{8} \log_2\left(\frac{6}{8}\right) - \frac{2}{8} \log_2\left(\frac{2}{8}\right)$$

$$= -\log_2 6 - \log_2 2$$

$$= -\log_2 6 - 1$$

$$\approx 0.81$$

for 34/3N

$$H(S) \approx 1$$

$$\text{Gain}(S, f_1) = 0.99 - \left[ \frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$= 0.049$$

This is the ID3 approach.

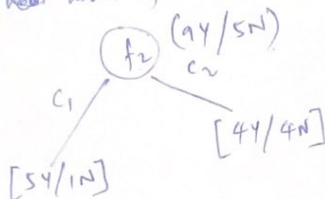
Similarly in the information gain formula, if we replace the entropy with Gini impurity  $\rightarrow$  it will be the CART approach.

But, this is w.r.t. only one feature  $f_1$ .

~~Let's talk about ~~feature~~ feature~~

We have to choose the feature that has less impurity & more information.

Let's talk ~~about~~ about feature 2  $\rightarrow f_2$ .



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

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

$$= -\log_2 9 - \log_2 5 \approx 0.99$$

$$= -3.17 - 2.32$$

$$= -5.49$$

$$\text{gain}_{f_2} = 0.99 - \left[ \frac{6}{14} \times 0.65 + \frac{8}{14} \times 1 \right]$$

$$= 0.09$$

$$\text{Now, } \text{gain}(f_1) = 0.049 \text{ \& } \text{gain}(f_2) = 0.09$$

Clearly, info. gain of  $f_2 >$  info. gain of  $f_1$ .

So, we select  $f_2$  as our root node. because it is greater & is providing more info.