# Decision Tree

A Decision Tree is a decision based hierarchical model that uses a tree-like model of decisions and their possible consequences.

## Decision Tree Classifier

$$\underset{\text{i/p}}{P_1 (Num/Cat)} \longrightarrow \underset{\text{o/p}}{categorical}$$

Different forms of DTC (Decision tree clasifier)
① $ID_3$ . ② CART

| $ID_3$ | CART |
|---|---|
| (Iterative Dichotomiser 3) | (Classification And Regression Tree) |
| → Entropy | → Gini Impurity. |

In case of Decision Tree



Root Node

Child Node

→ Leaf Node

| Day | f1 outlook | f2 Temperature | f3 humidity | f4 wind | Decision |
|---|---|---|---|---|---|
| 1 | Sunny | hot | high | weak | NO |
| 2 | Sunny | hot | high | strong | NO |
| 3 | overcast | hot | high | weak | Yes |
| 4 | rainfall | mild | ~~normal~~ 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 | rain fall | mild | high . | strong | No |

Here

$Size(O/P) = 14$ (9 yes, 5 No)

Let root node is Outlook



[outlook] [9Y/5N] O/P
~~feature~~

sunny
[2Y/3N]

rainfall
[3Y/2N]

overcast
[4Y/0N]

Leaf Node
because
It is a pure split

## Pure split

If the split contain only one output feature its called Pure split.

→ To find Impurity of a split we use two feature method.

① Entropy    ② ~~Information Goin~~ Gini Inpurity

→ How to choose root node or child node?.

Ans. we have to find the Impurity of all feature.

① Entropy

                            ② Information Goin
                            ② Gini Impurity/coeff

$$-\sum_{i=1}^{n} P_i \times \log(P_i)$$
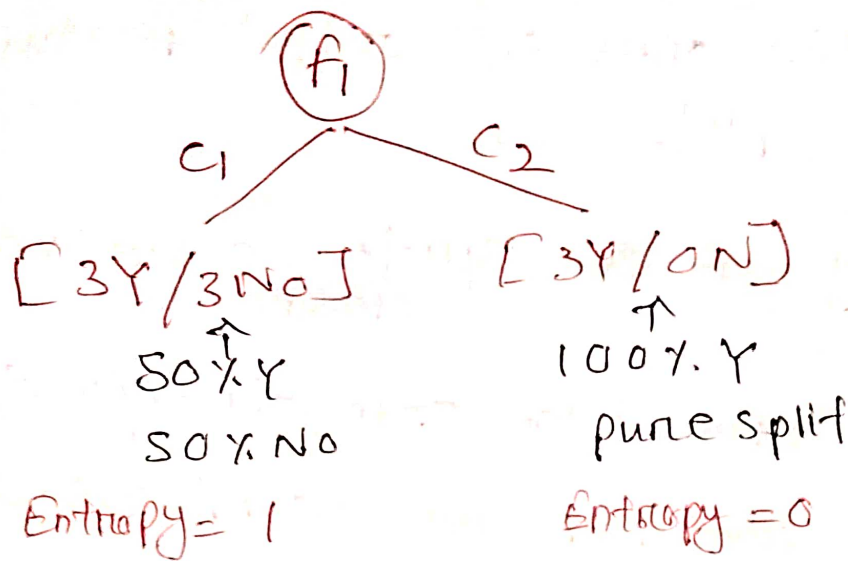$$1-\sum_{i=1}^{n} P_i^2$$

For two class $(y/N)$

Entropy $= -P_y \log(P_y) - P_N \log(P_N)$

for 3 class $(C_1, C_2, C_3)$

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

$\underline{\text{Gini Impurity}}$ (for 2 class $(y/N)$)

$\Rightarrow \quad 1 - [P_y^2 + P_N^2]$

| $f_1$ | $o/p$ |
|-------|-------|
| $C_1$ | Y |
| $C_2$ | Y |
| $C_1$ | Y |
| $C_2$ | Y |
| $C_1$ | Y |
| $C_1$ | No |
| $C_2$ | Y |
| $C_1$ | No |
| $C_1$ | No |

$f_1$

$C_1$        $C_2$

$[3Y/3No]$        $[3Y/oN]$

$50\% Y$        $100\% Y$

$50\% No$        pure split

Entropy = 1        Entropy = 0

entropy for $C_1 = -\sum_{i=1}^{n} P_i \log(P_i)$

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

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

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

$= -\frac{1}{2} \left[ \log(1) - \log_2 2 + \log(1) - \log_2 2 \right]$

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

entropy for $C_2$
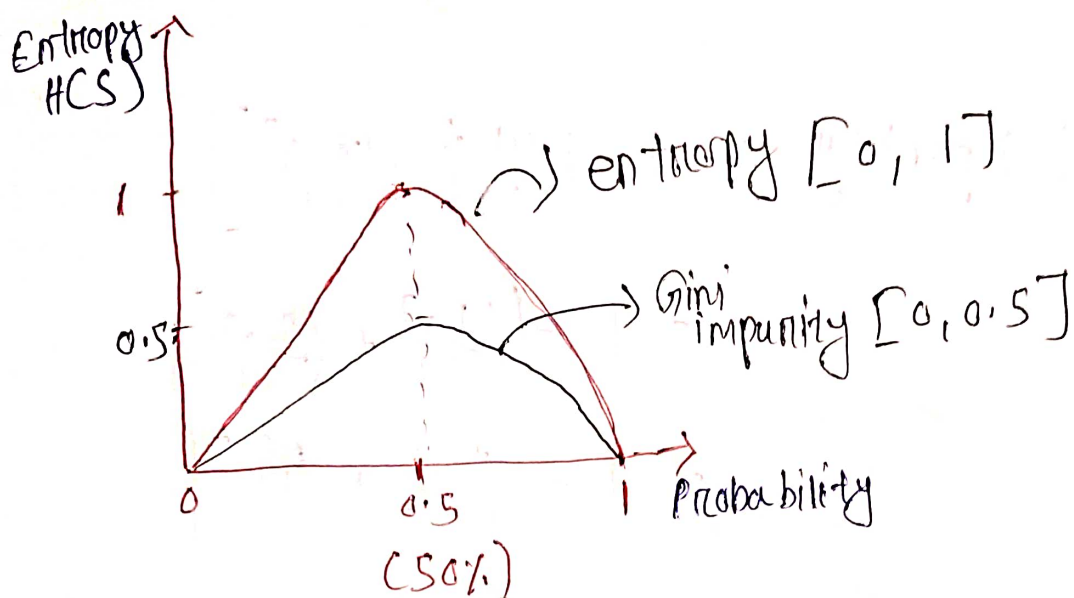
$H(S) = -P_y \log(P_y) - P_{No} \log(P_{No})$

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

$= -1 \log(1)$

$= 0$

Entropy HCS)

→ entropy [0, 1]

→ Gini impunity [0, 0.5]

Probability (50%)

HCS) = 1, very impure split

HCS) = 0, pure split.

Let o/p   2 yes / 3 No.

$$HCS) = -\frac{2}{5} \log_2 (2/5) - \frac{3}{5} \log_2 (3/5)$$

$$= 0.97$$

**By using Gini-Impurity (for Dexter)**

$$\text{Gini Impurity} = 1 - \sum_{i=1}^{n} (p_i^2) \quad (\text{for } c_1)$$

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

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

$$= 1 - \frac{1}{2} = 0.5$$

for $c_2$

$$\text{gini impurity} = 1 - \left[ \left(\frac{3}{3}\right)^2 + \left(\frac{0}{3}\right)^2 \right]$$

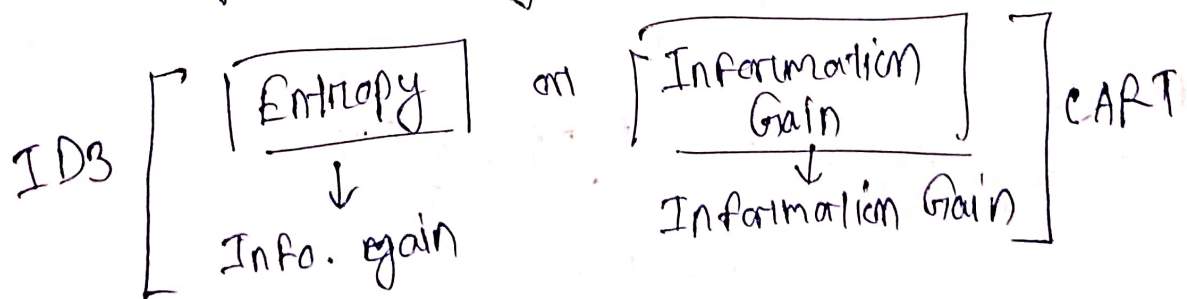$$= 1 - 1 = 0$$

**Q.** For [4Y / 8N] gini impurity ?.

Ans

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

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

$= 1 - \left[ \frac{1}{9} + \frac{4}{9} \right]$

$= 1 - \frac{5}{9} = \frac{4}{9} = 0.44$

**Q.** [8Y / 2N]

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

$= 1 - \left[ \left(\frac{4}{5}\right)^2 + \left(\frac{1}{5}\right)^2 \right]$
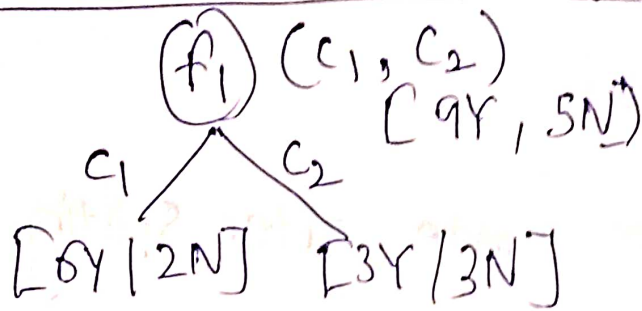
$= 1 - \frac{17}{25}$

$= \frac{8}{25} = 0.32$

$\rightarrow$ feature1 , feature-2 , feature-3
which feature to choose first ?.

Ans. For checking impurity

$$ID3 \left[ \boxed{Entropy} \quad \text{or} \quad \boxed{\underset{\downarrow}{\begin{array}{c} Information \\ Gain \end{array}} \atop Information\ Gain} \right] CART$$

$ID3 \rightarrow$ Entropy $\downarrow$ Info. gain

# Information Gain

$$\boxed{\text{Gain}(S, f_1) = H(S) - \sum \frac{|S_v|}{|S|} H(S_v)}$$

$f_1$ $(c_1, c_2)$ $[9Y, 5N]$

$c_1$ $[6Y|2N]$    $c_2$ $[3Y/3N]$

$H(S) \Rightarrow$ root feature entropy

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

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

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

$$\approx 0.94$$

for $c_1$

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

$$= 0.81$$

for $c_2$

$$H(S) = -\frac{3}{6} \log\left(\frac{3}{6}\right) - \frac{3}{6} \{\log\left(\frac{3}{6}\right) = 1$$

$$\text{Gain}(S, f_1) = 0.94 - \left(\frac{\overbrace{8}^{H(S)(C_1)}}{14} \times 0.81 + \frac{\overbrace{6}^{H(S)C_2}}{14} \times 1\right)$$

$$= 0.94 - (0.462 + 0.42)$$

$$= 0.049$$

$|Sv| =$ total no. of sample after spliting
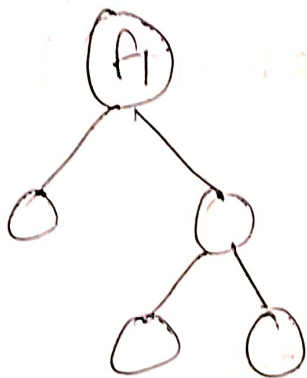
$|S| =$ before spliting total no. of sample.

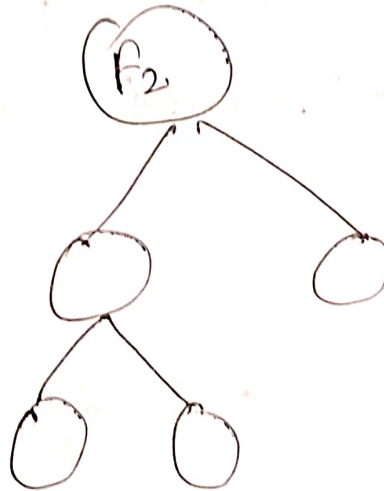Similarly we can calculate Info. Gain for all feature. Then we need to comparle.
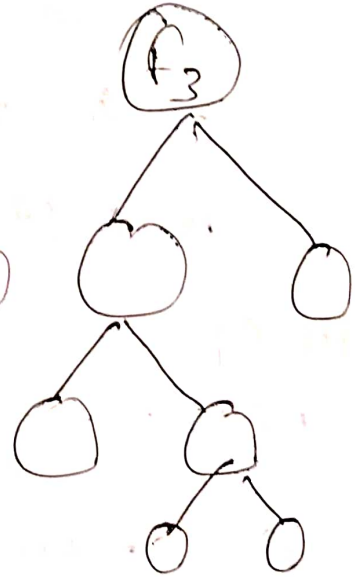
Let    I Gain = 0.78          Gain = 0.62          Gain = 0.20



conclusion

-) $f_1$ is best.

we will select the feature which has highest information gain then we split the node based on that feature.

## Note

Decision tree is robust to _outliers_ and
no scaling of feature is _required_. because
Decision tree works of condition based approach.

Q. when to stop spliting?

In real-world datasets have large no.
of features which results in large no.
of splits. which in turn gives a large tree.
Such trees take time to build and can
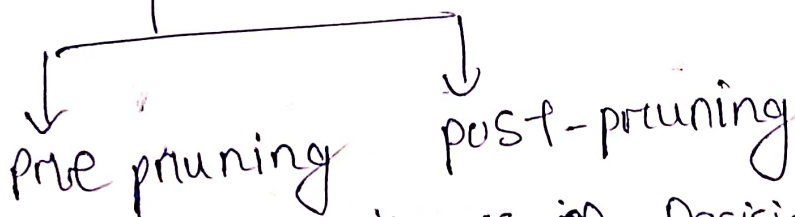lead to overfitting. (traing accuracy
herry high and test accuracy verry low)

There are many ways to tackle this problem
(hyperparameter tuning)

-) Max_depth (no. of depth split should be)

-) min_samples_split ( split the node if
   it contains samples more then min_sample_split)

-) max_feature , min_sample_leaf

-) Pruning (cutting the tree)

                    |
          _____|_____
         |                     |
         ↓                     ↓
    pre pruning          post-pruning

-) more we will discuss in Decision tree regression