# Discision Tree

## Description:

→ A Flowchart for making decisions. It starts with question at the root and branches into answers at each node, leading to a final descision at the leaves. Each deeision is based on data, helping the algorithm learn patterns and make predictions or classifications, making it a powerful too like sorting and predicting outcomes.

There are **two types** of decision Tree:

① Decision Tree Classifier [Classification]
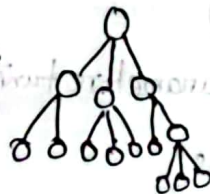
② Decision Tree Regressor [Regression]

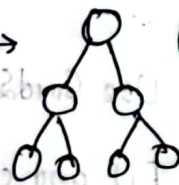## Decision Tree Classifier: (Two types)

① ID3 [Iterative Dichotomiser 3]

② CART [Classification and Regression Tree]  **Scikit learn use this**

ID3 technique →

CART Technique →  (only binary)
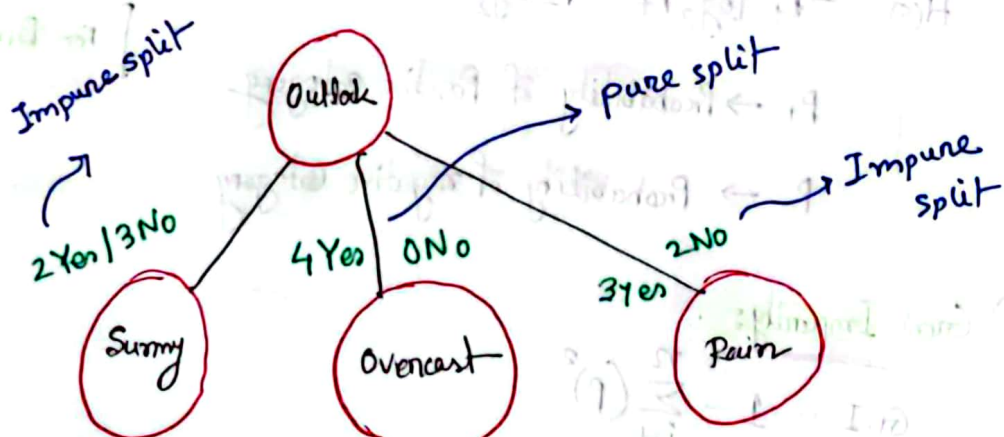
Let's take an example of a dataset of predicting Play tennis or not.

| Day | Outlook | Temp | Humidity | Wind | Play Tennis |
|-----|---------|------|----------|------|-------------|
| 1 | Sunny | Hot | High | Weak | No |
| 2 | Sunny | Hot | High | Strong | No |
| 3 | Overcast | Hot | High | Weak | Yes |
| 4 | Rain | Mild | High | Weak | Yes |
| 5 | Rain | Cool | Normal | Weak | Yes |
| 6 | Rain | Cool | Normal | Strong | No |
| 7 | Overcast | Cool | Normal | Strong | Yes |
| . | . | . | . | . | . |
| . | . | . | . | . | . |
| . | . | . | . | . | . |
| 14 | Rain | Mild | High | Strong | No |

Take one independant feature say → "Outlook" and compare with target feature.



① Purity check:

pure split → Only yes/ only no

impure split → Some yes and some no combination (Need further splitting)

==To check purities== we use ==two techniques:==

1) Entropy   2) Gini Impurity.

Using these techniques to find pure or impure split and decide to decide further split

② ==What Feature you need to select ? to start the split:==

For this, we a use "==information Gain==" technique

This helps to undestistand which independent feature has to be chosen to start with.

==Technique explanation of purity check:==

① ==Entropy:==

$$H(s) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

[For Binary classification]

$P_+ \rightarrow$ Probability of Positive category
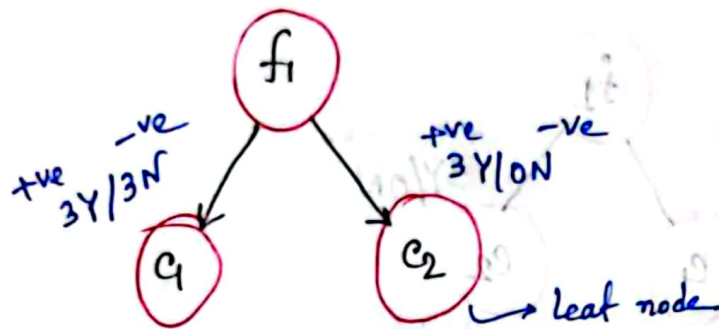
$P_- \rightarrow$ Probability of negative category

② ==Gini Impurity:==

$$G.I = 1 - \sum_{j=1}^{n} (p)^2$$
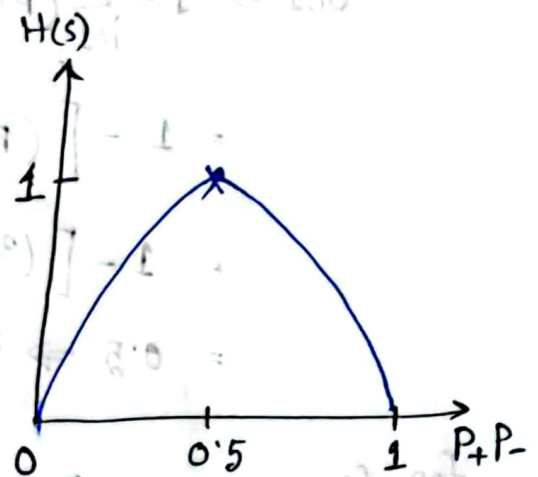
Let's take an example for **Entropy calculation:**

$f_1$

+ve -ve
3Y/3N

+ve -ve
3Y/0N

$c_1$

$c_2$ → leaf node

$$H(c_1) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

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

$$= 1 \Rightarrow \text{Impure Split}$$

$$H(c_2) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

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

$$= 0 \Rightarrow \text{Pure split}$$

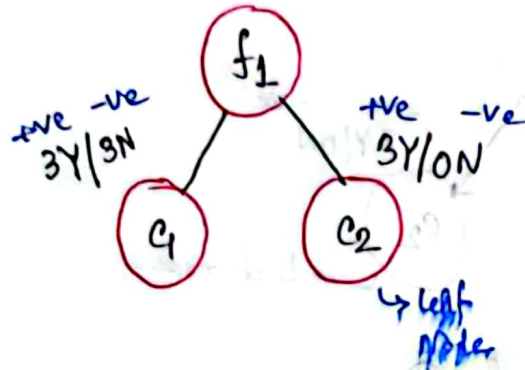That was for Binary classification (Yes/no)

$c_1$  $c_2$  $c_3$
For "Multi class Classification" (Yes/No/Maybe)

Entropy → $H(s) = -P_{c_1} \log_2 P_{c_1} - P_{c_2} \log_2 P_{c_2} - P_{c_3} \log_2 P_{c_3}$

H(s)

1

0        0.5        1    $P_+ P_-$

P → Probability
f → feature
c → Category

Now let's take an example for ==Gini Impurity Calculation:==



$f_1$

+ve  -ve
3Y/3N

+ve  -ve
3Y/0N

$c_1$  $c_2$

→ leaf node

for, c1,

$$G.I \Rightarrow 1 - \sum_{i=1}^{n} (p)^2$$

$$= 1 - \left[ (P_+)^2 + (P_-)^2 \right]$$

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

$$= 0.5 \Rightarrow \text{impure split}$$

for, $c_2$,

$$G.I = 1 - \sum_{i=1}^{n} (p)^2$$
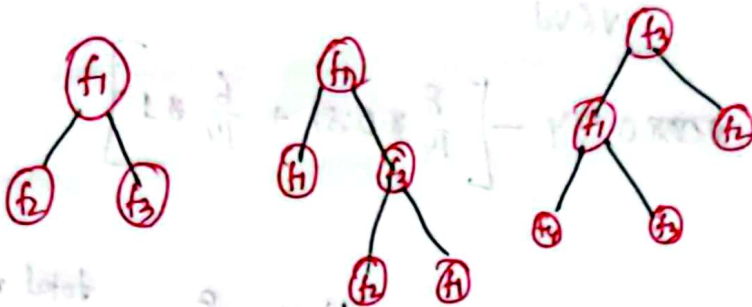
$$= 1 - \left[ (P_+)^2 + (P_-)^2 \right]$$

$$= 1 - \left[ (3/3)^2 + (0/3)^2 \right]$$

$$= 1 - 1$$

$$= 0 \Rightarrow \text{Pure split}$$



$H(s)$

1

0.5

→ Entropy

→ Gini

0      0.5    1    → P-P

## Explanation of "Information Gain" :



which feature should be selected at first, then after splitting which new feature should be selected, that is decided by information gain.
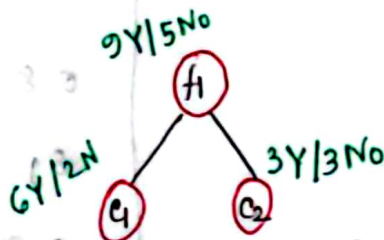
## Information Gain:

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

$H_s \rightarrow$ Entropy of root node

$S \rightarrow$ Sample

$v \rightarrow$ value



$$H(s) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$= -\frac{9}{14} \log \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$= 0.94$$

using the same calculation formula, $H(S_v) = H(c_1), H(c_2)$

$$\Rightarrow H(c_1) = 0.81$$
$$H(c_2) = 1$$

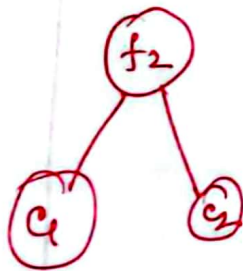$$\text{Gain}\left(S, f_1\right) = H(s) - \sum_{v \in Val} \frac{|sv|}{|s|} H(sv)$$

$$= 0.94 - \left[\frac{8}{14} * 0.81 + \frac{6}{14} * 1\right]$$

$\therefore$ Gain $(s, f_1) = 0.049$

Here, $\frac{8}{14} \rightarrow$ $\dfrac{\text{total num of Y an N in } c_1}{\text{total No of Y and N in } f_1}$

Then suppose we take another feature and measure the Information gain

Here, $\frac{6}{14} =$ $\dfrac{\text{total No of Y and N in } c_2}{\text{total no of Y and N in } f_1}$



Gain$(s, f_2) \rightarrow$ Gain found $\rightarrow 0.051$

$0.8 = H(c_1)$

$1 = H(c_2)$

$0.94 = H(s) = H(f_1)$

==The greater the Gain value,==

==that feature will be used by us.==

So, we will use "==feature f2.=="

==Entropy vs Gini Impurity:== (Which to use When)

Whenever dataset ==small== (1000, 2000 records) → ==Use Entropy==

Whenever dataset ==large== (1M ; 100000, more) → ==Use Gini Impurity==