

# Decision Tree

The Decision Tree algorithm is a popular machine learning algorithm used for both classification and regression tasks. It is a tree-like model that makes decisions based on feature values, where each internal node represents a feature, each branch represents a decision rule, and each leaf node represents the outcome or prediction.

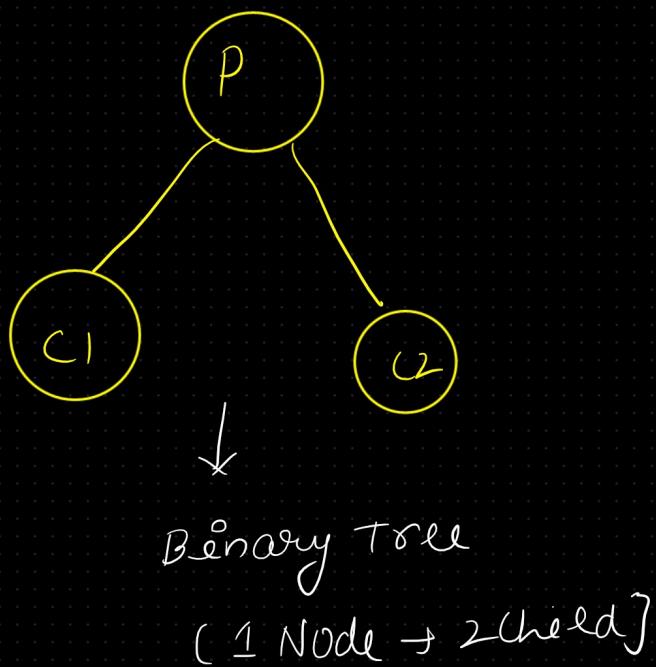
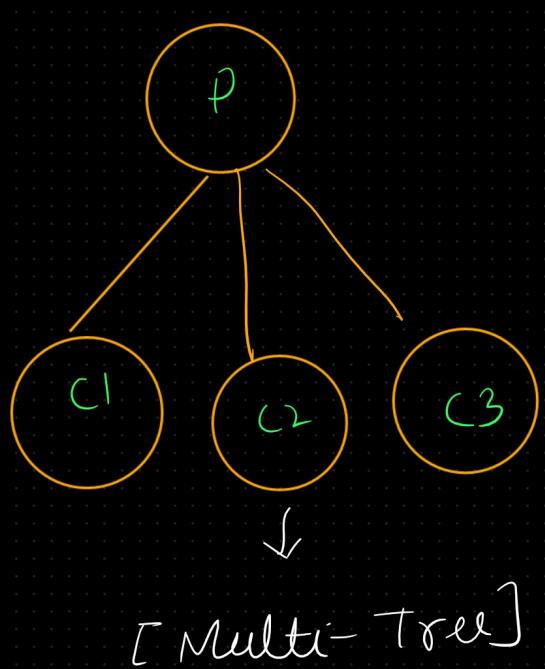
## Types of Decision Trees:

- ① Classification Trees  
(Decision Tree Classifier)

Eg → Predicting whether an email is spam or not

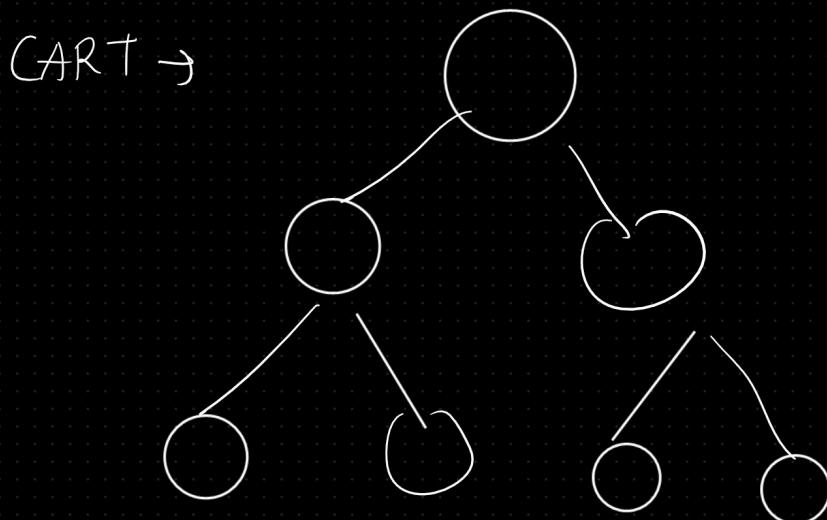
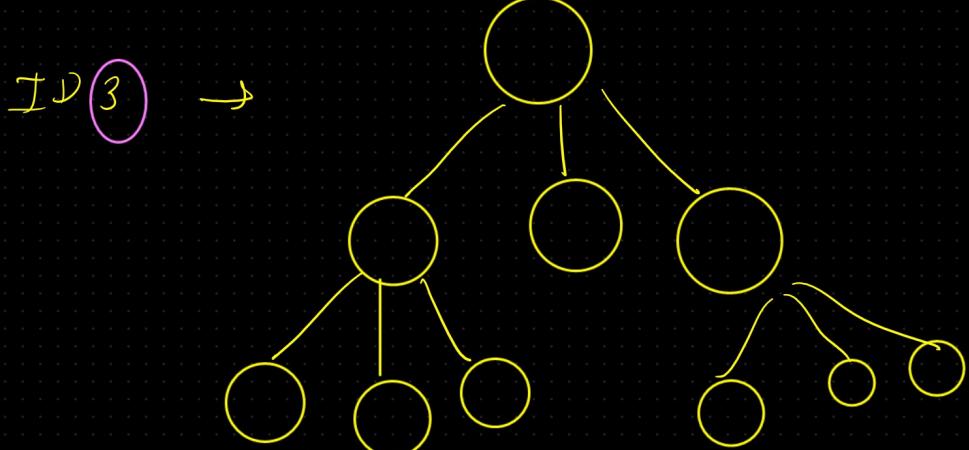
- ② Regression Trees [Decision Tree Regressor]

Eg → Predicting the price of the house based on its features.



## Decision Tree Classifier

- ① ID3 → [Iterative Dichotomiser 3]
- ② CART → [Classification And Regression Trees]
- ✓ ③ Random Forest
- ✓ ④ Gradient Boosted Trees (GBT)
- ✓ ⑤ XGBoost (Extreme Gradient Boosting)



$Age = 14 \text{ years}$

if ( $age \leq 15$ )

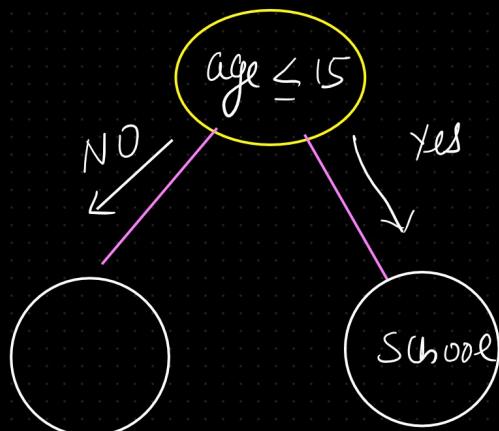
print ("School")

elif ( $age > 15$  and  $age \leq 21$ )

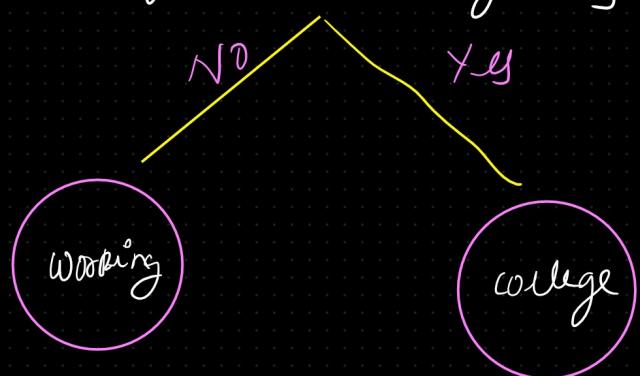
print ("College")

else:

print ("Working")

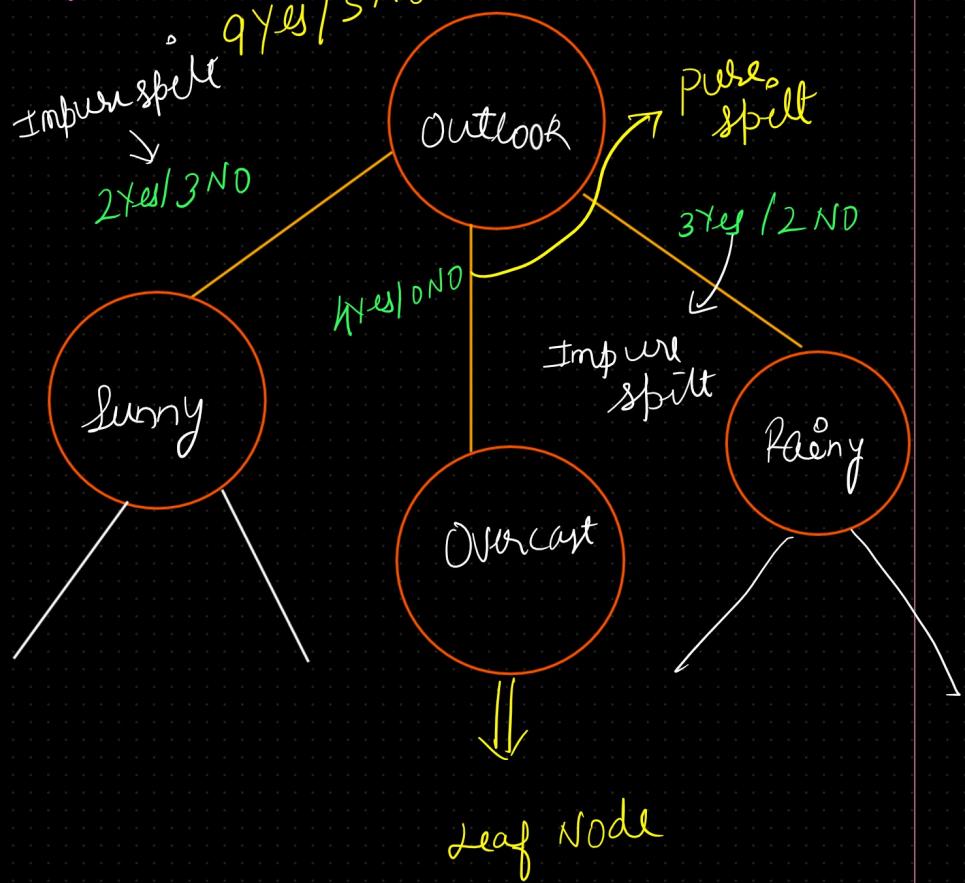


[ $age > 15$  and  $age \leq 21$ ]

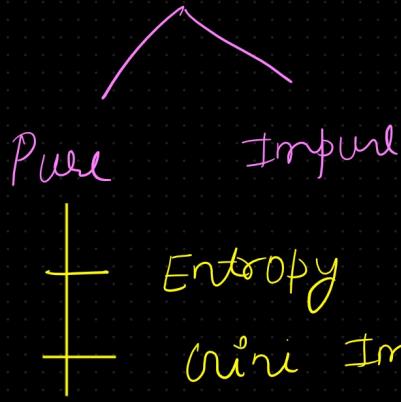


Dataset: Predict Play Tennis or Not

Outlook	Temperature	Humidity	Wind	Play Tennis
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rainy	Mild	High	Weak	Yes
Rainy	Cool	Normal	Weak	Yes
Rainy	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rainy	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes
Rainy	Mild	High	Strong	No



## ① Purity Check



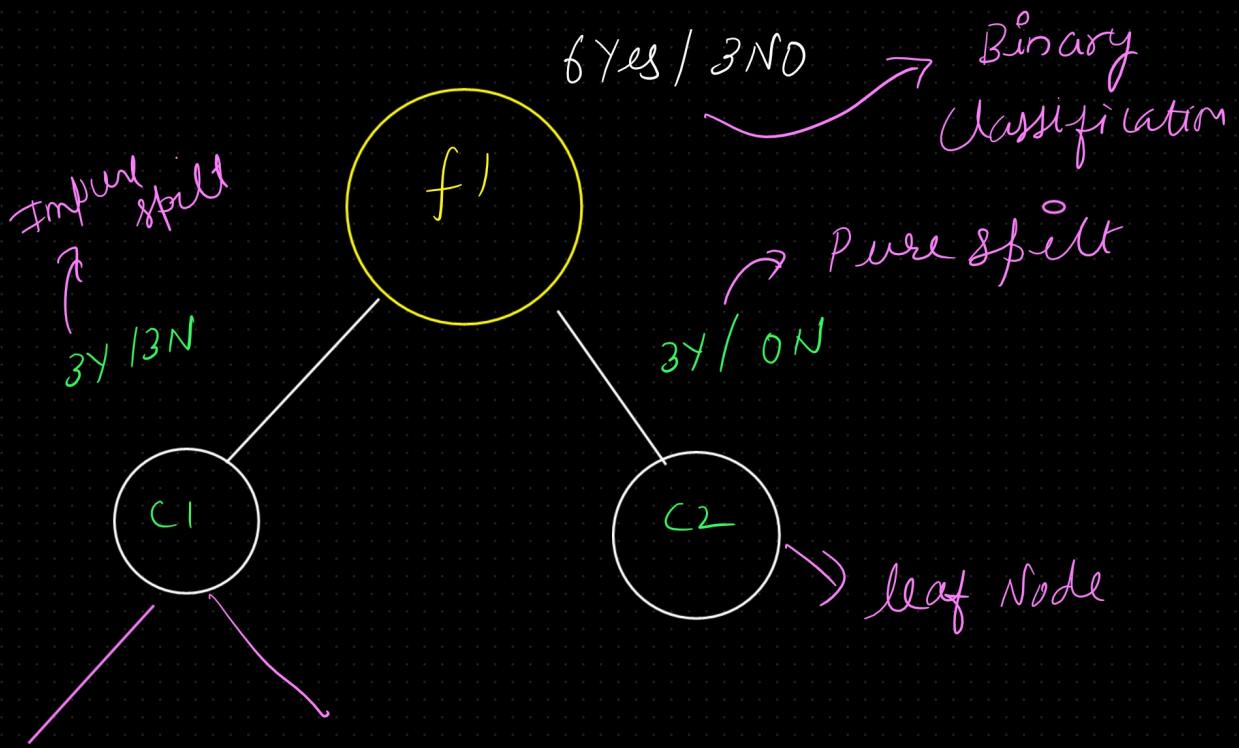
Decision Trees split nodes based on the concept of purity. A node is considered pure if all its data belongs to the same class (for classification) or has similar values (for regression).

## Entropy vs Gini Impurity

→ Entropy : Measures the Impurity or uncertainty in a dataset

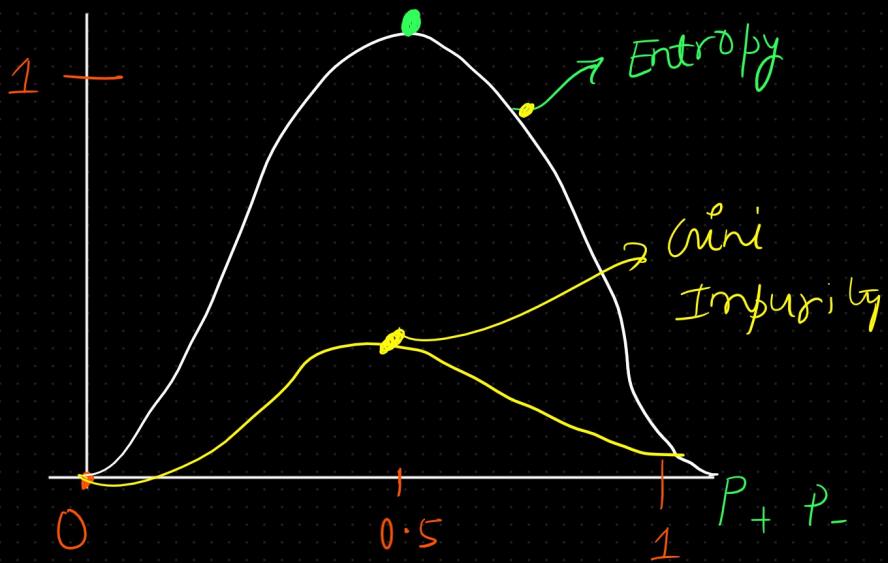
$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$P_+$  = Prob. of +ve category  
 $P_-$  = Prob. of -ve category



$$\begin{aligned}
 H(C_1) &= -P_+ \log_2 P_+ - P_- \log_2 P_- \\
 &= -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \left(\frac{3}{6}\right) \log_2 \left(\frac{3}{6}\right) \\
 &= -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \left(\frac{1}{2}\right) \log_2 \left(\frac{1}{2}\right) \\
 &= 1 \Rightarrow \text{Impure split}
 \end{aligned}$$

$$\begin{aligned}
 H(C_2) &= -\left(\frac{3}{3}\right) \log_2 \left(\frac{3}{3}\right) - 0/3 \log_2 (0/3) \\
 &= -\log_2 = 0 \Rightarrow \text{Pure split}
 \end{aligned}$$



$$\text{Mutual Information} \rightarrow -P_{C_1} \log P_{C_1} - P_{C_2} \log P_{C_2} - P_C \log P_C$$

## Gini Impurity

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

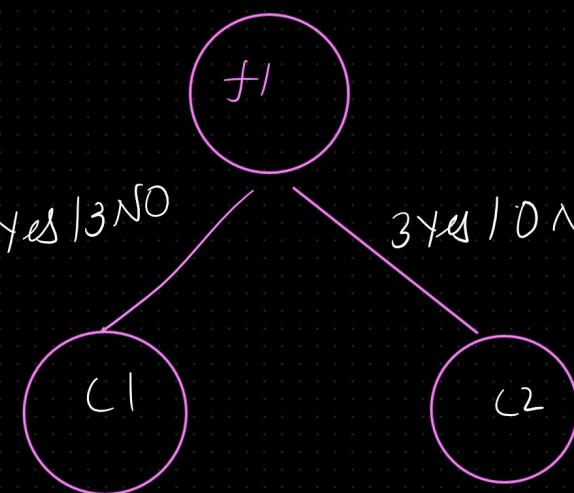
$$G.I = 1 - \sum_{i=1}^n (P_i)^2$$

$$CI = 1 - [(P_+)^2 + (P_-)^2]$$

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

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

$\therefore \frac{1}{2} = 0.5 \Rightarrow \text{Impure split}$



$$G.I(c_2) = 1 - [1^2 + 0]$$

$$= 1 - 1$$

$$= 0$$

$\Downarrow$

pure split

# Information gain

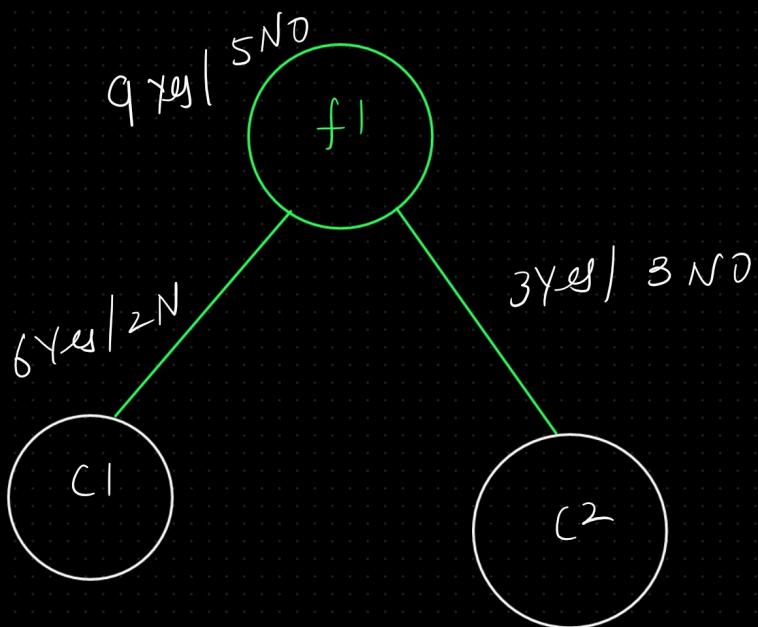
Information gain measures the effectiveness of a particular attribute in classifying the training data. It is calculated by comparing the entropy (or Gini impurity) before and after splitting the dataset based on that attribute.

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

where  $H(S)$  = Entropy

$f_1 \rightarrow$  feature

$f_1, f_2, f_3$   
↓ ↓ ↓



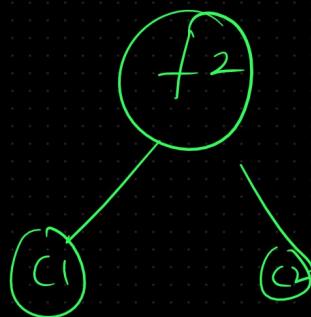
$$\begin{aligned}
 H(S) &= -P_+ \log_2 P_+ - P_- \log_2 P_- \\
 &= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \approx 0.94
 \end{aligned}$$

$$H(C_1) = 0.81$$

$$H(C_2) = 1$$

$$\text{Gain}(S, f_1) = 0.94 - \left[ \frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$\boxed{\text{Gain}(S, f_1) = 0.049}$$



$$\text{Gain}(S, f_2) = 0.051$$

Information gain is more when we split using  $f_2$

Entropy  $\rightarrow$  Dataset  $\rightarrow$  small

Min Impurity  $\rightarrow$  Dataset  $\rightarrow$  large

When to Use Decision Trees:

a. Pros:

Easy to understand and interpret.

Can handle both numerical and categorical data.

Nonlinear relationships between features and target variable are well captured.

b. Cons:

→ Prone to overfitting, especially with deep trees.

Can be sensitive to small variations in the data.

Biased towards features with more levels.

#### Mathematical Intuition:

- Splitting Decision: At each node, Decision Trees select the feature that provides the best split, maximizing information gain or minimizing impurity.
- Entropy/Gini Calculation: These metrics quantify the uncertainty or impurity in the dataset. Lower values indicate higher purity.
- Information Gain Calculation: Measures the reduction in entropy (or decrease in impurity) after splitting the dataset based on a particular attribute.
- Recursive Process: The decision-making process continues recursively until a stopping criterion is met, such as reaching a maximum depth or purity threshold.