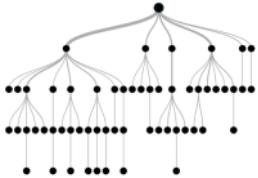


1 Introduction

- A **Decision Tree** is a supervised learning algorithm.

Decision Tree → classification
Decision Tree → regression

- It splits data into branches based on **feature values** to make predictions, resembling a flowchart or tree structure.



Types of Decision Tree:

- ① ID3
- ② CART

Key idea:

- At each step, it chooses the **feature and threshold** that best separates the data into pure (homogeneous) subsets.

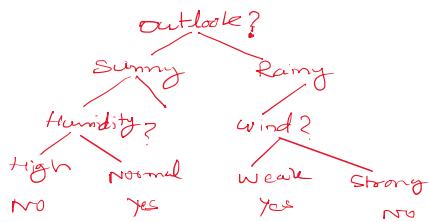
2 Why Decision Trees?

- Easy to understand and interpret
- No need for feature scaling
- Works with numerical and categorical data
- Can handle non-linear relationships
- ⚠ But:
- Can overfit easily
- Sensitive to small changes in data
- Greedy nature — may not find the global optimum

3 Anatomy of a Decision Tree

- **Root Node:** The first node where splitting starts (entire dataset).
- **Internal Nodes:** Nodes where the data is split further.
- **Leaf Nodes (Terminal Nodes):** Represent class labels (for classification) or predicted values (for regression).

Example (Classification):



4 How Does a Decision Tree Work?

It recursively partitions data using the **best feature** based on a chosen splitting criterion.

For Classification:

Common metrics:

- ① Entropy and Gini Impurity → pure split
- ② Information Gain → Feature selection

For Regression:

Metric: Variance Reduction (mean squared error reduction)

1. Entropy

Entropy measures the impurity or randomness:

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

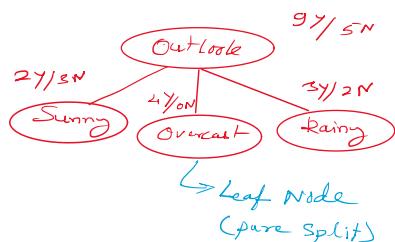
Example : Play Tennis Examples Dataset

Outlook	Temperature	PlayTennis
Sunny	Hot	No
Sunny	Hot	No
Overcast	Hot	Yes
Rainy	Mild	Yes
Rainy	Cool	Yes
Rainy	Cool	No
Overcast	Cool	Yes
Sunny	Mild	No
Sunny	Cool	Yes
Rainy	Mild	Yes
Sunny	Mild	Yes
Overcast	Mild	Yes
Overcast	Hot	Yes
Rainy	Mild	No

Frequency Table

Outlook		PlayTennis	
Sunny	2	3	No
Overcast	4	0	Yes
Rainy	3	2	No
	9	5	

Temperature	Y	N
hot	2	2
mild	4	2
cool	3	1
	9	5



$$\begin{aligned}
 \text{Entropy } H(S) &= - \sum_{i=1}^c p_i \log_2(p_i) \\
 &= - (P_S \log_2 P_S + P_N \log_2 P_N) \\
 &= - [\frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5}] \\
 &= 0.159 + 0.132 \\
 &= 0.291
 \end{aligned}$$

Entropy of node Hot ($\frac{2}{5}$)

$$\begin{aligned}
 H(Hot) &= - \left[\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right] \\
 &= - [0.5 \times -1 + 0.5 \times -1] \\
 &= - [0.5 - 0.5] \\
 &= 0
 \end{aligned}$$

Entropy of node Overcast ($\frac{4}{5}$)

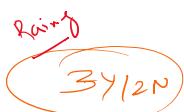
$$\begin{aligned}
 H(On) &= - \frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} \\
 &= 0
 \end{aligned}$$

Entropy of node Rainy ($\frac{3}{5}$)

$$\begin{aligned}
 H(Rainy) &= - \frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \\
 &= 0.291
 \end{aligned}$$

② Gini Impurity

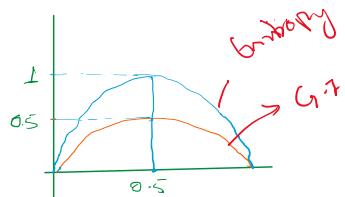
$$G.I. = 1 - \sum_{i=1}^c (p_i)^2$$



$$\begin{aligned}
 G.I. &= 1 - [(p_S)^2 + (p_N)^2] \\
 &= 1 - \left[\left(\frac{2}{5}\right)^2 + \left(\frac{3}{5}\right)^2 \right]
 \end{aligned}$$

$$= 1 - [0.36 + 0.16]$$

$$= 1 - 0.52 = 0.48$$



34/BN

$$H(S) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$= 1$$

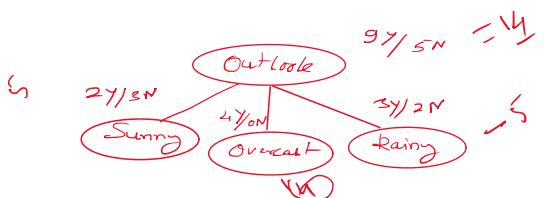
$$\text{G.I.} = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2$$

$$= 1 - \frac{1}{4} - \frac{1}{4} = 0.5$$

(b) Information Gain

$$\text{Information Gain (IG)} = H(S)_{\text{parent}} - \sum_{k \in \text{children}} \frac{n_k}{n} \times H(S)_k$$

Goal: Maximize Information Gain.



$$H(\text{parent}) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$= -(0.643)(-0.632) - (0.357)(-1.485)$$

$$= 0.4045 + 0.5301$$

$$= 0.94$$

$$H(\text{Sunny}) = 0.291$$

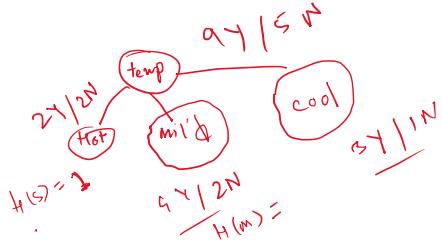
$$H(\text{outlook}) = 0$$

$$H(\text{Rainy}) = 0.291$$

$$\text{Information Gain} = 0.94 - \left[\frac{5}{14} \times 0.291 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.291 \right]$$

$$= 0.94 - 0.2078$$

$$= 0.732$$



$$H(\text{hot}) = 1$$

$$H(\text{mild}) = 0.918$$

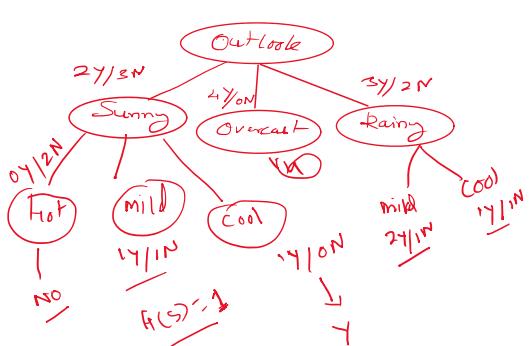
$$H(\text{cool}) = 0.8$$

$$H(\text{Temp}) = 0.94$$

$$\begin{aligned} \text{Info Gain} &= 0.94 - \left[\frac{4}{14} \times 1 + \frac{6}{14} \times 0.918 + \frac{4}{14} \times 0.8 \right] \\ &= 0.94 - [0.28 + 0.42 \times 0.918 + 0.28 \times 0.8] \\ &= \underline{0.832} \end{aligned}$$

$$I.G.(\text{Outlook}) = 0.752 \quad \checkmark$$

$$I.G.(\text{Temp}) = 0.832$$



	Outlook	Temperature	PlayTennis
1	Sunny	Hot	No
2	Sunny	Hot	No
3	Overcast	Hot	Yes
4	Rainy	Mild	Yes
5	Rainy	Cool	Yes
6	Rainy	Cool	No
7	Overcast	Cool	Yes
8	Sunny	Mild	No
9	Sunny	Cool	Yes
10	Rainy	Mild	Yes
11	Sunny	Mild	Yes
12	Overcast	Mild	Yes
13	Overcast	Hot	Yes
14	Rainy	Mild	No