

# Decision Trees

## Introduction

A decision tree is a supervised learning algorithm that uses a tree-like structure to make decisions. It splits data into smaller subsets based on the most significant attributes, creating a tree of decisions that can be used for both classification and regression tasks.

## How Decision Trees Work

### 1. Tree Structure

- Root Node: Starting point, contains entire dataset
- Internal Nodes: Decision points based on features
- Branches: Possible outcomes of each decision
- Leaf Nodes: Final predictions/outcomes

### 2. Splitting Criteria

- For Classification:
  - Gini Impurity:  $1 - \sum(p_i)^2$
  - Entropy:  $-\sum(p_i \times \log_2(p_i))$
  - Information Gain: Parent Entropy - Weighted Child Entropy
- For Regression:
  - Mean Squared Error (MSE)
  - Mean Absolute Error (MAE)
  - Root Mean Squared Error (RMSE)

## Advantages and Disadvantages

### Advantages

1. Easy to understand and interpret
2. Requires minimal data preprocessing
3. Can handle both numerical and categorical data
4. Can handle multi-output problems
5. Validates model using statistical tests

### Disadvantages

1. Can create overly complex trees (overfitting)
2. Can be unstable (small variations in data might result in different trees)
3. Biased toward dominant classes
4. May create biased trees if classes are imbalanced

## **Best Practices**

### **1. Preventing Overfitting**

- Use max\_depth to limit tree growth
- Set minimum samples for splits
- Implement pruning techniques
- Use cross-validation

### **2. Handling Missing Values**

- Create a new category for missing values
- Use surrogate splits
- Impute missing values

### **3. Feature Engineering**

- Bin continuous variables
- Create interaction features
- Handle categorical variables appropriately

### **4. Model Evaluation**

- Use cross-validation
- Check feature importance
- Analyze confusion matrix
- Consider multiple metrics