# Decision Trees

## Decision Trees

## A. Understanding Decision Trees

### Definition and Structure of a Decision Tree

A decision tree is a powerful and intuitive machine learning model that mimics human decision-making. It is a flowchart-like structure where each internal node represents a "test" or "decision" on an attribute (feature), each branch represents the outcome of that decision, and each leaf node represents a class label (for classification) or a continuous value (for regression).

- **Root Node**: This is the topmost node in a decision tree. It represents the entire dataset and is the starting point for splitting the data based on the most significant feature.
- **Internal Nodes**: These nodes represent the features of the dataset. Each internal node performs a test on one of the input features and splits the data accordingly. The test is usually a condition, such as whether a feature's value is above or below a certain threshold.
- **Leaf Nodes**: The terminal nodes of the tree, where the decision-making process ends. Each leaf node represents a class label (in classification tasks) or a predicted value (in regression tasks).

### How Decision Trees Make Predictions

- **Classification Tasks**: In classification, the decision tree divides the dataset into subsets based on the value of input features. Starting from the root node, the model checks the value of a particular feature and moves down the tree following the branch corresponding to that value. This process continues until it reaches a leaf node, which holds the final prediction—the class label for the input data.
  For example, in a decision tree that predicts whether a customer will purchase a product, the root node might split customers based on their income level. Internal nodes might further split based on age or previous purchasing history, and the leaf nodes would indicate "Yes" or "No" for the purchase prediction.
- **Regression Tasks**: In regression, the tree works similarly but instead of predicting a class label, each leaf node represents a continuous value. The model splits the data based on feature values, and the final prediction is the average value of the target variable within that leaf.

### The Concept of Splitting (Features, Thresholds)

- **Feature Selection**: At each internal node, the decision tree algorithm chooses the feature that best splits the data into distinct classes or values. The selection is based on certain criteria that measure the "purity" of the split, such as Information Gain or Gini Index.
- **Thresholds**: For continuous features, the model determines an optimal threshold value that divides the data into two parts, maximizing the homogeneity of the resulting subsets. For categorical features, the model considers different possible groupings of categories.
- **Recursive Splitting**: This process of splitting the data is repeated recursively at each node, creating branches of the tree until a stopping condition is met. The aim is to continue splitting the data to create the most homogeneous subsets possible within the tree's constraints.

# B. Decision Tree Construction

**Algorithm for Building a Decision Tree (e.g., ID3, CART)**

Building a decision tree involves recursively selecting the best feature to split the data at each node. The most common algorithms used are:

- **ID3 (Iterative Dichotomiser 3)**: This algorithm builds a decision tree by employing a top-down, greedy approach. It selects the feature that maximizes Information Gain (a measure of the effectiveness of a feature in classifying data) at each step. ID3 is primarily used for categorical data.
- **CART (Classification and Regression Trees)**: CART is a versatile algorithm used for both classification and regression tasks. It builds binary trees, where each internal node has exactly two branches. The algorithm splits nodes using the Gini Index for classification tasks and Mean Squared Error (MSE) for regression tasks. CART can handle both categorical and continuous data.

**Criteria for Splitting Nodes**

The effectiveness of a split is determined by various criteria, which measure how well the split separates the data:

- **Gini Index**: Used by the CART algorithm, the Gini Index measures the impurity of a dataset. The lower the Gini Index, the purer the subset is after the split. The formula is:

$$Gini = 1 - \sum_{i=1}^{n} p_i^2$$

  where $p_i$ is the probability of class $i$ in the subset.
- **Information Gain**: Employed by the ID3 algorithm, Information Gain measures the reduction in entropy (disorder or uncertainty) after a dataset is split on a feature. The feature with the highest

Information Gain is chosen for splitting. The formula is:

$$IG(T, X) = Entropy(T) - \sum_{v \in Values(X)} \frac{|T_v|}{|T|} Entropy(T_v)$$

where $T$ is the dataset, $X$ is the feature, and $T_v$ is the subset of $T$ for each value $v$ of $X$.

- **Chi-Square**: This criterion is used to measure the statistical significance of the splits. It compares the observed distribution of data points with the expected distribution if the data were evenly split. Higher chi-square values indicate more significant splits.

### Stopping Criteria

To prevent a decision tree from growing too complex, stopping criteria are applied. These criteria determine when the tree should stop splitting:

- **Maximum Depth**: Limits how deep the tree can go. The depth of a tree is the number of nodes from the root to the farthest leaf. Limiting depth helps prevent overfitting by stopping the tree from learning overly specific patterns in the data.
- **Minimum Samples per Leaf**: Specifies the minimum number of samples that must be present in a leaf node. If a split results in a leaf with fewer samples than this threshold, the split is discarded.
- **Minimum Samples per Split**: Sets the minimum number of samples required to split an internal node. This prevents the tree from making splits that result in very small subsets, which could lead to overfitting.

## C. Pruning Techniques

### What is Overfitting in Decision Trees?

Overfitting occurs when a decision tree becomes too complex, capturing noise and specific details in the training data that do not generalize well to unseen data. An overfitted model has high accuracy on the training data but performs poorly on test data due to its excessive sensitivity to the training dataset's nuances.

### Pruning Methods to Prevent Overfitting

Pruning is a technique used to reduce the size of a decision tree by removing sections of the tree that provide little power in predicting target variables. There are two main types of pruning:

- **Pre-Pruning (Early Stopping)**: This approach involves stopping the growth of the tree before it reaches full depth. It applies constraints like maximum depth, minimum samples per split, or

minimum samples per leaf to control the growth. By limiting the tree's complexity, pre-pruning reduces the risk of overfitting.

- **Post-Pruning**: In post-pruning, the tree is allowed to grow to its full depth without any early stopping. Afterward, the tree is trimmed by removing branches that have little importance. The tree is evaluated, typically using cross-validation, and branches are pruned if they do not improve the model's performance on validation data. Methods include cost-complexity pruning, which removes nodes if the increase in misclassification error is minimal, and reduced error pruning, which prunes nodes if their removal does not increase the error rate on the validation set.

# D. Advantages and Disadvantages of Decision Trees

**Strengths of Decision Trees**

- **Interpretability**: Decision trees are easy to understand and interpret, even for non-experts. They provide a clear visual representation of the decision-making process, where each path from the root to a leaf represents a specific decision rule.
- **Handling of Categorical and Numerical Data**: Decision trees can handle both categorical and numerical features without the need for pre-processing such as normalization or dummy encoding.
- **No Assumption of Linear Relationships**: Unlike linear models, decision trees do not assume a linear relationship between the input features and the target variable. This allows them to model complex, non-linear relationships effectively.
- **Automatic Feature Selection and Interaction**: Decision trees naturally perform feature selection by choosing the most significant features for splits. They also automatically handle feature interactions without needing explicit specification.
- **Robustness to Outliers**: Decision trees are less sensitive to outliers compared to linear models. Outliers do not heavily influence the splitting process as they might in models that assume linearity.

**Limitations of Decision Trees**

- **Prone to Overfitting**: Without proper tuning or pruning, decision trees can easily overfit the training data, capturing noise and leading to poor generalization on new data.
- **Sensitivity to Small Variations in Data**: Small changes in the training data can result in a completely different tree structure, making decision trees less stable compared to some other models.
- **Bias in Splits with Many Levels**: Decision trees can favor features with many levels (like a categorical feature with many categories), as they can create more potential splits, sometimes leading to biased or less meaningful splits.

- **Limited Expressiveness**: Simple decision trees might not capture very complex patterns in data. Although they handle non-linearity well, their expressiveness is limited by their depth and the quality of the splits.