

## Unit 4-Tree-Based Methods

Tree-based methods are a class of machine learning algorithms used for both classification and regression tasks. These methods create decision trees to make predictions based on input features. Decision trees are hierarchical structures where each node represents a feature, each branch represents a decision based on that feature, and each leaf node represents the final prediction.

There are several tree-based methods, and some of the most popular ones include:

1. **Decision Trees:** Decision trees partition the feature space into subsets by asking binary questions based on feature values. Each internal node corresponds to a feature, and each leaf node represents a class label (in classification) or a continuous value (in regression). The goal is to create the simplest decision tree that accurately classifies or predicts the target variable.
2. **Random Forest:** Random Forest is an ensemble method that combines multiple decision trees to improve prediction accuracy and reduce overfitting. It works by training several decision trees on different subsets of the training data (randomly sampled with replacement) and averaging their predictions for improved robustness.
3. **Gradient Boosting:** Gradient Boosting is another ensemble method that builds multiple decision trees sequentially. Each tree is trained to correct the errors of the previous tree, gradually reducing the prediction errors and improving overall accuracy. Popular implementations include XGBoost, LightGBM, and AdaBoost.
4. **Bagging:** Bagging (Bootstrap Aggregating) is an ensemble technique that trains multiple decision trees independently on different bootstrapped samples of the training data. The final prediction is obtained by averaging the predictions of each individual tree, reducing variance and improving stability.

### Advantages of Tree-Based Methods:

They can handle both numerical and categorical data without requiring extensive data preprocessing.

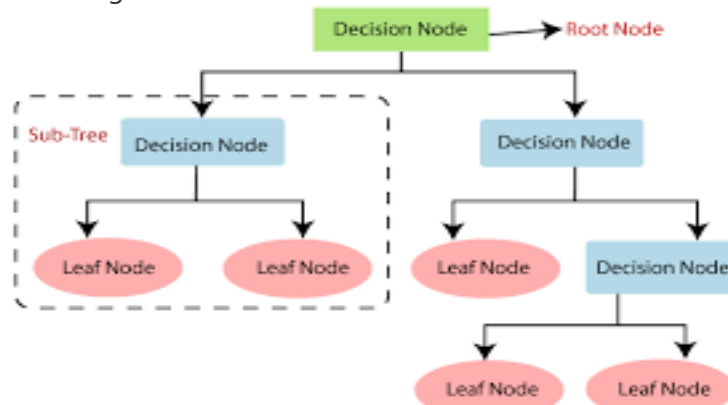
- **Interpretability:** Decision trees are relatively easy to interpret and visualize, making them useful for understanding the decision-making process.
- **Non-linear relationships:** Trees can capture non-linear relationships between features and the target variable.
- **Robustness:** Tree-based methods are robust to outliers and missing values.
- **Feature importance:** They provide information about feature importance, helping to identify influential features in the prediction process.

### Disadvantages of Tree-Based Methods:

- **Overfitting:** Decision trees can easily overfit the training data, especially if they are deep and complex.
- **Instability:** Small changes in the data can lead to significantly different tree structures.
- **Bias:** Decision trees tend to be biased towards features with more levels (in classification) or higher variance (in regression).
- **Lack of smoothness:** Prediction surfaces in tree-based models can be piecewise constant, which may not capture smooth relationships in the data.

### The Basics of Decision Trees

Decision trees are a fundamental concept in machine learning and data mining, widely used for both classification and regression tasks. They are simple yet powerful models that can capture complex relationships between input features and the target variable. Here are the basics of decision trees:



1. **Tree Structure:** A decision tree is a hierarchical structure that resembles a flowchart. It consists of nodes, branches, and leaves. Each node represents a decision based on a specific feature, each branch represents an outcome of that decision, and each leaf node represents the final prediction or outcome.

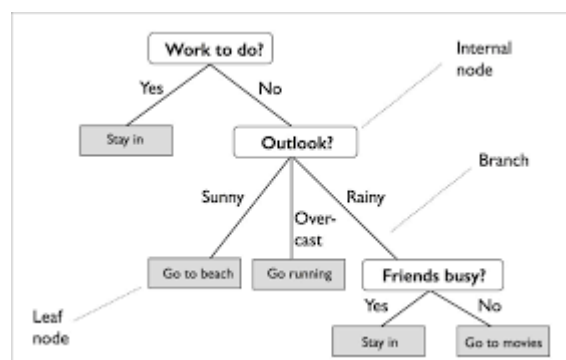
## 2. Nodes:

- **Root Node:** The topmost node of the tree, representing the first decision to be made. It corresponds to the most important feature or the feature that best splits the data.
  - **Internal Nodes:** Nodes that have child nodes. Each internal node represents a decision based on a specific feature and splits the data into subsets based on the feature's values.
  - **Leaf Nodes:** Terminal nodes with no child nodes. They provide the final prediction or classification label for the input data.
3. **Branches:** The branches represent the possible outcomes of a decision made at each internal node. For instance, if a feature's value is less than a certain threshold, the left branch might be taken; otherwise, the right branch might be chosen.
  4. **Splitting Criteria:** When building a decision tree, the algorithm looks for the best feature and its corresponding value to split the data at each internal node. The goal is to minimize impurity or maximize information gain, depending on the task (classification or regression).
  5. **Classification:**
    - In classification tasks, decision trees predict categorical or discrete class labels for the input data.
    - The tree is constructed by repeatedly partitioning the data into subsets based on the values of features, aiming to create homogeneous subsets with respect to the target class labels at the leaf nodes.
  6. **Regression:**
    - In regression tasks, decision trees predict continuous values for the input data.
    - The tree is constructed in a similar way as for classification, but the predicted value at each leaf node is the average (or another measure) of the target values in that subset.
  7. **Predictions:** To make predictions using a decision tree, you start at the root node and traverse the tree based on the feature values of the input data. At each internal node, you follow the appropriate branch until you reach a leaf node, where the final prediction is obtained.
  8. **Pruning:** Decision trees are prone to overfitting, especially if they are deep and complex. Pruning is a technique used to avoid overfitting by removing nodes and branches that do not contribute significantly to the model's performance. Pruned trees are simpler and more generalizable.
  9. **Interpretability:** Decision trees are highly interpretable models. The paths from the root to the leaves can be easily understood, making them useful for gaining insights into the decision-making process.
  10. **Ensemble Methods:** To improve prediction accuracy and reduce overfitting, ensemble methods like Random Forest and Gradient Boosting use multiple decision trees in combination.

---

## Regression Trees

Regression trees are a specific type of decision tree used for predicting continuous numerical values rather than categorical class labels. They are a variant of decision trees designed to handle regression tasks, where the goal is to model the relationship between input features and a continuous target variable.



**The process** of building a regression tree is similar to that of a standard decision tree, but instead of predicting class labels at the leaf nodes, regression trees predict continuous values based on the average (or other measure) of the target variable within each leaf node's subset.

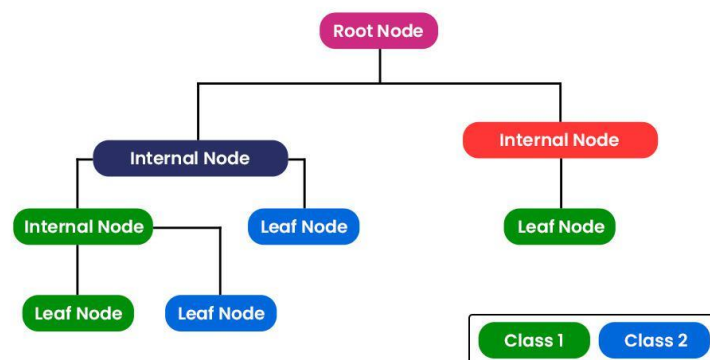
Here are the key points and steps involved in constructing a regression tree:

- 1) **Tree Construction:** The tree construction process begins with the entire dataset represented by the root node. The algorithm searches for the best feature and its corresponding value to split the data into two subsets (left and right child nodes). The splitting criterion typically aims to minimize the variance of the target variable within each subset.
- 2) **Splitting Criterion:** In regression trees, commonly used splitting criteria include:
  - Mean Squared Error (MSE): Minimizing the MSE at each split, which corresponds to the average squared difference between the predicted values and the actual target values.
  - Mean Absolute Error (MAE): Minimizing the MAE at each split, which corresponds to the average absolute difference between the predicted values and the actual target values.
- 3) **Recursive Splitting:** The process of finding the best feature and value to split the data is repeated for each child node. This recursive splitting continues until a stopping criterion is met, such as a predefined maximum tree depth or a minimum number of samples required in a leaf node.
- 4) **Leaf Nodes and Predictions:** Once the recursive splitting is complete, the tree will have several leaf nodes, each representing a subset of the data. The predicted value for each leaf node is determined by aggregating the target values within that subset, commonly using the mean or median value.
- 5) **Prediction:** To make predictions using a regression tree, you start at the root node and traverse the tree based on the feature values of the input data. At each internal node, you follow the appropriate branch based on the feature's value until you reach a leaf node. The predicted continuous value is then obtained from the leaf node.
- 6) **Pruning:** Similar to decision trees for classification, regression trees are also prone to overfitting. Pruning techniques are used to avoid overfitting by removing nodes and branches that do not contribute significantly to the model's accuracy.

---

## Classification Trees

Classification trees are a type of decision tree used for solving classification tasks. They are machine learning models that partition the feature space into distinct regions or leaves, each associated with a specific class label. Classification trees are particularly useful when dealing with categorical target variables and are widely used in various applications, including pattern recognition, medical diagnosis, and customer segmentation.



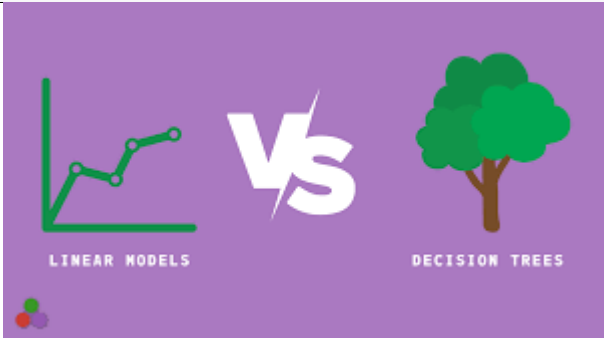
Here are the key points and steps involved in constructing a classification tree:

1. **Tree Construction:** The tree construction process starts with the entire dataset represented by the root node. The algorithm searches for the best feature and its corresponding value to split the data into two or more subsets (left and right child nodes) based on class labels. The goal is to maximize the purity of each subset, ensuring that samples within each leaf node belong predominantly to one class.
2. **Splitting Criterion:** In classification trees, commonly used splitting criteria to measure impurity and guide the tree's construction include:
3. **Gini impurity:** It measures the probability of incorrectly classifying a randomly chosen sample in a node if it were randomly labeled according to the class distribution in that node.

- 4. **Entropy:** It measures the uncertainty or disorder in a node's class distribution. Lower entropy corresponds to more pure nodes with more samples belonging to the same class.
- 5. **Classification Error:** It simply measures the misclassification rate in a node.
- 6. **Recursive Splitting:** The process of finding the best feature and value to split the data is repeated for each child node. This recursive splitting continues until a stopping criterion is met, such as a predefined maximum tree depth, a minimum number of samples required in a leaf node, or when all leaf nodes contain samples from the same class.
- 7. **Leaf Nodes and Predictions:** Once the recursive splitting is complete, the tree will have several leaf nodes, each representing a subset of the data belonging predominantly to a specific class. The majority class in each leaf node becomes the predicted class for any new input falling into that region.
- 8. **Prediction:** To make predictions using a classification tree, you start at the root node and traverse the tree based on the feature values of the input data. At each internal node, you follow the appropriate branch based on the feature's value until you reach a leaf node. The predicted class is then determined by the majority class in the corresponding leaf node.
- 9. **Pruning:** Similar to regression trees, classification trees are also susceptible to overfitting. Pruning techniques are used to avoid overfitting by removing nodes and branches that do not contribute significantly to the model's accuracy. Pruned trees are simpler and more generalizable.

Trees Versus Linear Models

Aspect	Linear Models	Trees (Decision Trees)
Linearity	Assume linear relationship between features and target.	Can handle both linear and non-linear relationships.
Interpretability	More interpretable due to coefficient insights.	Offer interpretable decision rules through tree structure.
Robustness to Outliers	Sensitive to outliers affecting coefficients.	Less sensitive to outliers as they split data in branches.
Handling Missing Data	Require data imputation for missing values.	Can handle missing data in features directly.
Feature Scaling	Often require feature scaling for consistent impact.	Not sensitive to feature scaling due to comparisons.
Handling Categorical Data	Often need one-hot encoding for categorical data.	Can handle categorical features directly.
Overfitting	Prone to overfitting complex data or underfitting.	Prone to overfitting, but ensembles help mitigate.
Complexity and Computation	Less computationally expensive.	Can be computationally expensive, especially for deep trees.
Use Cases	Well-suited for linear relationships and interpretability.	Suitable for non-linear patterns and complex feature interactions.



## **Advantages and Disadvantages of Trees**

Trees, especially decision trees, have various advantages and disadvantages that are important to consider when choosing them as a machine learning model for a particular problem. Let's explore the advantages and disadvantages of trees:

### **Advantages of Trees:**

**Interpretability:** Decision trees are highly interpretable. The decision-making process can be visualized as a flowchart with clear decision rules, making it easy to understand how the model arrives at its predictions.

**Handling Non-linearity:** Trees can capture non-linear relationships between input features and the target variable, making them suitable for problems with complex interactions and non-linear patterns.

**No Feature Scaling Required:** Trees are not sensitive to the scale of the features, so there's no need to perform feature scaling (e.g., mean normalization or standardization) before training the model.

**Handling Categorical Features:** Decision trees can handle categorical features directly without the need for one-hot encoding or other transformations.

**Robustness to Outliers:** Trees are less sensitive to outliers in the data since they recursively split the data, and outliers typically get isolated in separate branches.

**Feature Importance:** Trees provide a measure of feature importance, indicating which features are more influential in the decision-making process.

**Mixed Data Types:** Trees can handle both numerical and categorical data, making them versatile for a wide range of datasets.

**Ensemble Methods:** Decision trees can be combined into ensemble methods like Random Forest and Gradient Boosting, which often yield improved prediction accuracy and robustness.

### **Disadvantages of Trees:**

**Overfitting:** Decision trees are prone to overfitting, especially when the tree becomes too deep and complex. Overfitting occurs when the model fits the noise in the training data rather than the underlying patterns.

**Instability:** Small changes in the data can lead to significantly different tree structures, making the model unstable.

**Bias Towards Features with More Levels:** Decision trees tend to be biased towards features with more levels (e.g., categorical features with many categories).

**Lack of Smoothness:** Prediction surfaces in tree-based models can be piecewise constant, which may not capture smooth relationships in the data.

**Limited Expressiveness for Simple Relationships:** Trees may struggle to represent simple relationships effectively, and linear models might be more appropriate in such cases.

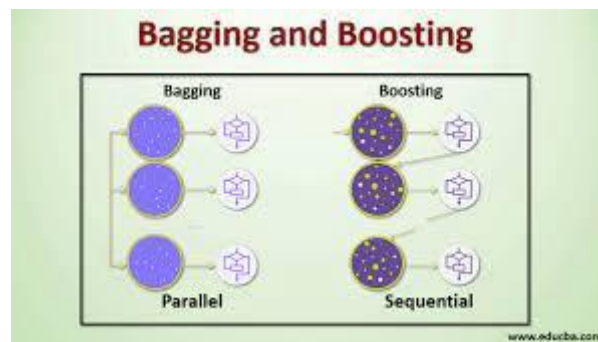
**Extrapolation:** Trees are not suitable for extrapolating beyond the range of data seen during training, and predictions can be unreliable outside this range.

**High Variance:** Individual decision trees can have high variance, meaning they can produce significantly different results on different subsets of the data.

---

## **Bagging**

Bagging, short for Bootstrap Aggregating, is an ensemble learning technique used to improve the accuracy and robustness of machine learning models. It involves creating multiple copies of the same model, each trained on a different random subset of the training data. The predictions from these individual models are then combined or aggregated to obtain the final prediction.



Here's a step-by-step explanation of how bagging works:

### 1) Bootstrapping:

- The process begins by creating multiple random subsets (with replacement) from the original training data.
- Each subset contains a random sample of the same size as the original data.
- Some examples may appear multiple times in a subset, while others may not be included at all.

### 2) Model Training:

- A single machine learning model (e.g., a decision tree, neural network, or regression model) is trained independently on each subset of the data.
- Since each subset contains different samples, the resulting models will have some variations.

### 3) Prediction Aggregation:

- During the prediction phase, each of the individual models makes predictions on the test data.
- For classification tasks, the final prediction is obtained by combining the predictions using majority voting (the most common predicted class).
- For regression tasks, the final prediction is obtained by averaging the predictions from all individual models.

The key idea behind bagging is that by training multiple instances of the model on different subsets of the data, the variance of the prediction errors is reduced. This reduction in variance leads to a more robust and accurate prediction, especially when dealing with noisy or complex datasets.

### Advantages of Bagging:

1. **Improved Prediction Accuracy:** Bagging can significantly improve the predictive performance of machine learning models, especially when dealing with high-variance models prone to overfitting.
2. **Robustness:** Bagging reduces the variance of the model's predictions, making it more resistant to noise and outliers in the data.
3. **Versatility:** Bagging can be applied to various types of models, including decision trees, random forests, neural networks, and more.
4. **Simplicity:** Bagging is a relatively simple and easy-to-implement ensemble technique.

### Disadvantages of Bagging:

1. **Lack of Interpretability:** The ensemble of models may not be as interpretable as a single model, especially when using complex models like neural networks.
2. **Computational Cost:** Training multiple models can be computationally expensive, especially for large datasets and complex models.
3. **Memory Overhead:** Storing multiple models in memory can increase memory usage, particularly when working with large ensembles.

## Random Forests

Random Forests is an ensemble learning method that builds upon the bagging technique and extends it to decision trees. It is a popular and powerful algorithm for both classification and regression tasks. Random Forests combine the predictions of multiple decision trees to make more accurate and robust predictions. The name "Random Forests" comes from the idea of creating an ensemble of decision trees, each trained on a random subset of the data and using random subsets of features at each split.

Here's how Random Forests work:

### 1. Bootstrapping:



- Random Forests employ bootstrapping to create multiple random subsets (with replacement) from the original training data. Each subset contains a random sample of the same size as the original data.
  - Each subset is used to train a separate decision tree.
2. **Feature Randomness:**
    - During the training of each individual decision tree, a random subset of features is selected at each split to find the best split. This introduces feature randomness.
    - The number of features considered at each split is a hyperparameter, typically set to the square root of the total number of features or a fixed fraction of them.
  3. **Model Training:**
    - Multiple decision trees are trained independently on their respective bootstrap samples and using the random subsets of features.
    - Each tree can grow to its maximum depth or until a predefined stopping criterion is met.
  4. **Prediction Aggregation:**
    - During the prediction phase, each of the individual decision trees makes predictions on the test data.
    - For classification tasks, the final prediction is obtained by combining the predictions using majority voting (the most common predicted class).
    - For regression tasks, the final prediction is obtained by averaging the predictions from all individual trees.

The key advantages of Random Forests lie in their ability to handle complex data, reduce overfitting, and provide an estimate of feature importance. By averaging predictions from multiple trees, Random Forests are less prone to overfitting compared to a single decision tree and are more robust in handling noisy or high-dimensional datasets.

#### **Advantages of Random Forests:**

1. **Improved Prediction Accuracy:** Random Forests often provide better prediction accuracy than individual decision trees, especially when dealing with complex and non-linear data.
2. **Robustness:** Random Forests are more robust against overfitting due to the averaging of multiple trees and the use of feature randomness.
3. **Feature Importance:** Random Forests provide a measure of feature importance, which helps in identifying influential features in the prediction process.
4. **Versatility:** Random Forests can handle both classification and regression tasks and can be used with various types of data.
5. **Parallelization:** The training and prediction processes of individual trees in Random Forests can be easily parallelized, making them scalable for large datasets.

#### **Disadvantages of Random Forests:**

**Lack of Interpretability:** The ensemble of decision trees in Random Forests may not be as interpretable as a single decision tree.

**Computational Complexity:** Training multiple decision trees can be computationally expensive, especially for large ensembles and complex datasets.

---

#### **Boosting and Bayesian Additive Regression Trees**

Boosting and Bayesian Additive Regression Trees (BART) are both advanced machine learning techniques used for regression tasks. While they both aim to improve predictive accuracy, they have different approaches and characteristics. Let's explore each method:

##### **Boosting:**

Boosting is an ensemble learning technique that combines multiple weak learners (typically decision trees) to create a strong predictive model. It works by iteratively building weak models and giving more emphasis to the misclassified or mispredicted data points in subsequent iterations. The key idea behind boosting is to focus on the mistakes made by the previous models, and in doing so, gradually improve the overall model's performance.

##### **1. Model Training:**

- Boosting starts by training an initial weak learner on the training data.

- Subsequent weak learners are then trained, giving more weight to the misclassified data points from the previous models. This means that the subsequent models focus on the mistakes made by the previous ones.

## 2. Weight Updates:

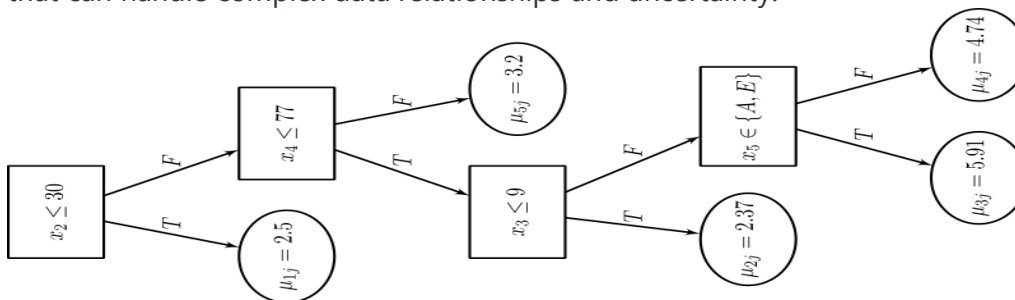
- After each iteration, the weights of the misclassified data points are updated, and this updated dataset is used to train the next weak learner.
- The process is repeated for a predefined number of iterations or until a stopping criterion is met.

## 3. Model Combination:

- The final prediction is obtained by combining the predictions from all the weak learners. The combination is typically weighted based on the performance of each weak learner.

## Bayesian Additive Regression Trees (BART):

BART is a Bayesian nonparametric modeling technique that uses a sum of regression trees to model the relationship between the input features and the target variable. BART is a flexible and powerful approach that can handle complex data relationships and uncertainty.



## 1. Model Construction:

- BART constructs the final model as a sum of regression trees, where each tree captures a particular subset of the data relationships.
- The model is built using a Bayesian framework, where prior distributions and posterior distributions are used to estimate model parameters.

## 2. MCMC Sampling:

- BART uses Markov Chain Monte Carlo (MCMC) sampling techniques to explore the posterior distribution of model parameters and make predictions.
- MCMC methods allow BART to capture uncertainty and provide credible intervals for predictions.

## 3. Model Combination:

- BART combines the predictions from multiple regression trees to create the final prediction for a given input.

## Advantages and Differences:

- **Boosting Advantages:** Boosting can achieve high predictive accuracy and is less prone to overfitting compared to individual decision trees. It is a flexible and powerful ensemble technique that works well with various weak learners.
- **BART Advantages:** BART is a Bayesian method, allowing it to capture uncertainty and provide credible intervals for predictions. It is particularly useful when the relationship between features and the target is highly non-linear and complex.
- **Differences:** Boosting is an ensemble method that focuses on iteratively correcting the mistakes of the previous models, while BART is a Bayesian nonparametric model that estimates uncertainty and provides credible intervals. Boosting typically uses decision trees as weak learners, while BART is based on additive regression trees.